NBER WORKING PAPER SERIES

NON-COGNITIVE ABILITY, TEST SCORES, AND TEACHER QUALITY:
EVIDENCE FROM 9TH GRADE TEACHERS IN NORTH CAROLINA

C. Kirabo Jackson

Working Paper 18624
http://www.nber.org/papers/w18624

Non-Cognitive Ability, Test Scores, and Teacher Quality: Evidence from 9th Grade Teachers
in North Carolina
C. Kirabo Jackson
NBER Working Paper No. 18624
December 2012, Revised January 2013
JEL No. H0,I2,J0

## ABSTRACT

This paper presents a model where students have cognitive and non-cognitive ability, and a teacher's
effect on long-run outcomes is a combination of her effects on both ability-types. Conditional on cognitive
scores, an underlying non-cognitive factor associated with student absences, suspensions, grades, and
grade progression, is strongly correlated with long-run educational attainment, arrests, and earnings
in survey data. In administrative data, teachers have meaningful causal effects on both cognitive-scores
and this non-cognitive factor. Calculations indicate that teacher effects based on test scores alone fail
to identify many excellent teachers, and may greatly understate the importance of teachers for adult
outcomes.

C. Kirabo Jackson
Northwestern University
School of Education and Social Policy
2040 Sheridan Road
Evanston, IL 60208
and NBER
kirabo-jackson@northwestern.edu

# Non-Cognitive Ability, Test Scores, and Teacher Quality: Evidence from 9[1] Grade Teachers in North Carolina[1]

*C. Kirabo Jackson*,
Northwestern University and NBER

This paper presents a model where students have cognitive and non-cognitive ability, and a teacher's effect on long-run outcomes is a combination of her effects on both ability-types. Conditional on cognitive scores, an underlying non-cognitive factor associated with student absences, suspensions, grades, and grade progression, is strongly correlated with long-run educational attainment, arrests, and earnings in survey data. In administrative data, teachers have meaningful causal effects on both cognitive-scores and this non-cognitive factor. Calculations indicate that teacher effects based on test scores alone fail to identify many excellent teachers, and may greatly understate the importance of teachers for adult outcomes. (JEL I21, J00)

> *"The preoccupation with cognition and academic "smarts" as measured by test scores to the exclusion of social adaptability and motivation causes a serious bias in the evaluation of many human capital interventions"* (Heckman, 1999)

There is a growing consensus that non-cognitive skills not captured by standardized tests, such as adaptability, self-restraint, self-efficacy, and motivation are important determinants of adult outcomes (Lindqvist & Vestman, 2011; Heckman & Rubinstein, 2001; Borghans, Weel, & Weinberg, 2008; Waddell, 2006). At the same time, many interventions that have no effect on test-scores have meaningful effects on long-run outcomes such as education, earnings, and crime (Booker, Sass, Gill, & Zimmer, 2011; Deming, 2009; Deming, 2011). Also Heckman, Pinto, & Savelyev (forthcoming) find that changes in personality traits explain the positive effect of the Perry Preschool Program on adult outcomes, and Fredriksson, Ockert, & Oosterbeek, (forthcoming) find that class size effects on earnings are predicted by effects on *both* test scores and non-cognitive ability. This suggests that schooling inputs produce both cognitive skills (measured by standardized tests) and non-cognitive skills (reflected in socio-behavioral outcomes), both of which determine adult outcomes. Accordingly, evaluating interventions based on test scores may capture only one dimension of the skills required for adult success, and "*a more comprehensive evaluation of interventions would account for their effects on producing the noncognitive traits that are also valued in the market*" (Heckman & Rubinstein, 2001).

There is broad consensus among policy makers, educators, parents, and researchers that teachers are one of the most important components of the schooling environment. Studies show that having a teacher at the 85[th] percentile of the quality distribution (as measured by effects on student test scores) versus one at the 15[th] percentile is associated with between 8 and 20 percentile points higher scores in math and reading (Aaronson, Barrow, & Sander, 2007; Kane & Staiger, 2008; Chetty, Friedman, & Rockoff, 2011; Rivkin, Hanushek, & Kain, 2005). The focus on test scores is largely because they are typically the best available measure. However, research on non-cognitive skills provides strong reason to suspect that effects on test scores may fail to capture a teachers overall effect. Several districts publicly release estimates of teachers' average effects on student test scores (value-added) and use them in hiring and firing decisions. Accordingly, it is important that these measures reflect a teacher's effect on long-run outcomes and not *only* her effect on cognitive ability. To shed light on this issue, this paper tests whether teachers have causal effects on *both* cognitive ability (measured by test scores) and non-cognitive ability (measured by absences, suspensions, grades, and grade progression). It also investigates whether teachers who improve test scores also improve non-test score outcomes, and estimates the extent to which test score measures understate the overall importance of teachers. While there is a growing literature on the importance of non-cognitive skills, and a burgeoning literature on the effect of teachers on test scores, this paper presents the first comprehensive analysis of teacher effects on *both* cognitive and objective non-cognitive outcomes.

Opponents of using test-scores to infer teacher quality raise two concerns. The first concern is that improving test scores does not necessarily imply better long-run outcomes because teachers might teach to the test, test score improvements might reflect transitory student effort, and those skills measured by test-scores may not be those that are associated with improved long-term outcomes. However, Chetty, Friedman, & Rockoff (2011) assuage this concern by demonstrating that teachers who improve test scores also improve outcomes into adulthood. The second concern is that student ability is multidimensional while test-scores may measure only one dimension of ability. If teachers improve skills not captured by test-scores then (a) many excellent teachers who improve long-run outcomes may not raise test scores, (b) the ability to raise test scores may not be the best predictor of effects on long-run outcomes, and (c) a regime that emphasizes test scores might induce teachers to divert effort away from skills not captured by test scores to increase test score outcomes – potentially decreasing teacher quality

overall (Holmstrom & Milgrom, 1991). This paper speaks to this second critique by assessing whether teachers affect non-cognitive skills not captured by test scores.

In existing work, Alexander, Entwisle, & Thompson (1987), Ehrenberg, Goldhaber, & Brewer (1995) and Downey & Shana (2004) find that students receive better teacher evaluations of behavior when students and teachers are more demographically similar, and Jennings & DiPrete (2010) finds that certain kindergarten classrooms are associated with meaningful differences in teacher evaluations of student behavioral skills. However, these studies may reflect differences in teacher perception rather than actual student behavior.[2] Accordingly, the extent to which teachers affect non-cognitive outcomes remains unclear.[3]

This paper is organized into three sections. The first section presents a latent factor model following Heckman, Stixrud, & Urzua (2006) in which both student ability and teacher ability have a cognitive and a non-cognitive dimension. It shows that a teacher's effect on long-run outcomes can be expressed as a combination of her effect on predicted cognitive ability (a factor that is highly correlated with test scores) and her effect on predicted non-cognitive ability (a factor that is a weighted average of several non-test score outcomes). The second section uses longitudinal survey data to form predicted non-cognitive ability, and demonstrates the extent to which predicted non-cognitive ability in $8^{th}$ grade predicts adult outcomes conditional on test scores. The third section uses administrative data to estimate $9^{th}$ grade Algebra and English teacher effects on test-scores and predicted cognitive and non-cognitive ability. Using these estimates within the context of the model, this paper investigates how well test-score measures alone identify teachers that have large predicted effects on long-run adult outcomes.

There are two distinct empirical challenges to credibly identifying teacher effects in high-school: The first challenge is that differences in outcomes across teachers may be due to selection of students to teachers; The second challenge is that *even with random assignment of students to teachers*, in a high-school setting, students in different tracks may be exposed to both different teachers and different "track treatments" *outside the classroom* that influence classroom performance and confound teacher effects (e.g. students in the gifted track take Algebra with Mr. Smith and also take a study skills class that has a direct effect on their Algebra test

---

[2] Moreover, these studies are based on single cohorts of nationally representable samples that have too few student observations per school and teacher for credible identification of individual teacher effects.

[3] In related work Koedel (2008) estimates high school teacher effects on graduation. However, he does not differentiate between effects that are due to improved cognitive skills versus non-cognitive skills.

performance).[4] To address both these challenges, this paper follows Jackson (2012) and estimates models that condition on a student's school-track (the unique combination of school, courses taken, and level of courses taken) so that all comparisons are made among students *at the same school and in the same academic track* — precluding any bias due to student selection to tracks or treatments that vary across tracks. In such models, variation comes from comparing the outcomes of students in the same track and school but who are exposed to different teachers either due to (a) changes in the teachers for a particular course and track over time, or (b) schools having multiple teachers for the same course and track in the same year. Because personnel changes within schools over time may be correlated with other changes within schools, I estimate models that also include school-by-year fixed effects. The remaining concern is that comparisons among students within the same track may be susceptible to selection bias. While most plausible stories of student selection involve selection to tracks rather than to teachers within tracks, this paper presents several empirical tests that suggest little to no selection bias.

Using the National Educational Longitudinal Study 1988 (NELS-88), a standard deviation increase in estimated non-cognitive ability in 8th grade (a weighted average of attendance, suspensions, grades, and grade progression) is associated with larger improvements in arrests, college-going, and wages than a standard deviation increase in test scores or estimated cognitive ability. Non-cognitive ability is particularly important at the lower end of the earnings distribution. Using administrative data, 9th grade Algebra and English teachers have meaningful effects on test-scores, the same non-test score outcomes, and estimated cognitive and non-cognitive ability. Calculations suggest that teacher effects on college-going and wages may be as much as five times larger than that predicted based on test scores alone. As such, more than half of teachers who would improve long-run outcomes may not be identified using test scores alone.

This paper presents the first evidence that teachers have meaningful effects on non-cognitive outcomes that are strongly associated with adult outcomes *and are not captured by test scores*. This has important implications regarding measuring teacher quality, and suggests that one might worry that test-based accountability may induce teachers to divert effort away from improving students' non-cognitive skills in order to improve test scores. Also, the finding that teachers have effects on ability unmeasured by test scores offers a potential explanation for the

---

[4] This is not a student selection problem, but an Omitted Variables Bias problem that may exist when <u>teachers</u> are not randomly assigned to tracks, such that other classes, other inputs, and peer quality are not balanced across <u>teachers</u>. See Jackson (2012) for an illustration of this kind of teacher sorting and the resulting bias.

empirical regularity that interventions that have test score effects that "fade out" over time can have lasting effects on adult outcomes (Chetty et. al. 2011; Cascio & Staiger, 2012).

The remainder of this paper is organized as follows: Section II presents the theoretical framework. Section III presents evidence that both test score and non-test score outcomes in 8th grade predict important adult outcomes. Sections IV and V present the data and empirical strategy, respectively. Sections VI and VII present the results and Section VIII concludes.

## II    Theoretical Framework

*A Model of Multidimensional Teacher and Student Ability*

Following Heckman, Stixrud, & Urzua (2006), this section presents a latent two-factor model where all student outcomes are a linear combination of students' cognitive and non-cognitive ability. Teachers vary in their ability to increase student cognitive and non-cognitive ability. Unlike a model with uni-dimensional student ability where a teacher who improves one outcome (such as test scores) necessarily improves *all* student outcomes (such as crime and adult wages), a model with multidimensional student ability has predictions that differ in important ways.

**Student ability:** Student ability is two-dimensional. One dimension is cognitive skills (e.g. content knowledge and computation speed), and the other dimension is non-cognitive or socio-behavioral skills (e.g. agreeableness, motivation, and self-control). Each student $i$ has a two-dimensional ability vector $v_i = (v_{c,i}, v_{n,i})$, where the subscript $c$ denotes the cognitive dimension and the subscript $n$ denotes the non-cognitive dimension.

**Teacher ability:** Each teacher $j$ has a two-dimensional ability vector $\omega_j = (\omega_{c,j}, \omega_{n,j})$ such that $E[\omega] = (0,0)$, which describes how much teacher $j$ changes each dimension (cognitive or non-cognitive) of student ability. The *total* ability of student $i$ with teacher $j$ is thus $\alpha_{ij} = v_i + \omega_j$.

**Outcomes:** There are multiple outcomes $y_z$ observed for each student $i$. Each outcome $z$ is a particular linear function of the ability vector such that $y_{zij} = \alpha_{ij}'\beta_z = (v_i + \omega_j)'\beta_z + \varepsilon_{zi}$ where $\beta_z = (\beta_{c,z}, \beta_{n,z})$ and $E[\varepsilon_{zi} | \varepsilon_{z'i}, \alpha, \omega] = 0$. Vector $\beta_z$ captures the fact that while some outcomes may depend on cognitive ability (such as test scores) others may depend on non-cognitive skills (such as attendance). The variable $\varepsilon_{zi}$ captures the fact that outcomes are measured with random student-level error. Following Heckman, Stixrud, & Urzua (2006), the ability types are

5

uncorrelated so $\text{cov}(\alpha_{c,ij}, \alpha_{n,ij}) = 0$. In the factor model representation, the two factors are the *total* ability of student $i$ with teacher $j$ in cognitive and non-cognitive ability, and the parameter vector $\beta_z$ is the factor loadings for student outcome $z$. The two-factor model is illustrated in a path diagram in Figure 1.

**Teacher Effects:** The difference in student outcomes between teacher $j$ with $\omega_j = (\omega_{c,j}, \omega_{n,j})$ and an average teacher with $\omega = (0,0)$ is a measure of $j$'s effect (*relative to the average teacher*). Teacher $j$'s effect for outcome z is therefore $\theta_{zj} = \omega_j ' \beta_z$. Teacher $j$'s effect for long-run outcome $y_*$ is $\theta_{*j} = \omega_j ' \beta_* = \beta_{c,*}\omega_{c,j} + \beta_{n,*}\omega_{n,j}$ where $\beta_{c,*}\beta_{n,*} \neq 0$. The long-run outcome is not observed and policy-makers wish to predict teacher effects for long-run outcome $y_*$.

**Proposition 1:** *A teacher's effect on the long-run outcome may be correlated with her effect on multiple short-run outcomes, even if her effects on these short-run outcomes are not correlated with each other.*

Consider the case with two outcomes $y_1$ *and* $y_2$. Each outcome reflects only one dimension of ability so that $\theta_{1j} = \beta_{c,1}\omega_{c,j}$ and $\theta_{2j} = \beta_{n,2}\omega_{n,j}$ where $\beta_{c,1}\beta_{n,2} \neq 0$. The two dimensions of teacher ability are uncorrelated so $\text{cov}(\omega_{c,j}, \omega_{n,j}) = 0$. In this scenario, the covariance between teacher effects across all three outcomes are given by [1] through [3] below.

$$\text{cov}(\theta_1, \theta_2) = \text{cov}(\beta_{c,1}\omega_{c,j}, \beta_{n,2}\omega_{n,j}) = \beta_{c,1}\beta_{n,2}\,\text{cov}(\omega_{c,j}, \omega_{n,j}) = 0 \qquad [1]$$

$$\text{cov}(\theta_1, \theta_*) = \text{cov}(\beta_{c,1}\omega_{c,j}, \beta_{c,*}\omega_{c,j}) + \text{cov}(\beta_{c,1}\omega_{c,j}, \beta_{n,*}\omega_{n,j}) = \beta_{c,1}\beta_{c,*}\,\text{var}(\omega_{c,j}) \neq 0 \qquad [2]$$

$$\text{cov}(\theta_2, \theta_*) = \text{cov}(\beta_{n,2}\omega_{n,j}, \beta_{c,*}\omega_{c,j}) + \text{cov}(\beta_{n,2}\omega_{n,j}, \beta_{n,*}\omega_{n,j}) = \beta_{n,2}\beta_{n,*}\,\text{var}(\omega_{n,j}) \neq 0 \qquad [3]$$

**Proposition 2:** *With multiple short-run outcomes that reflect a mix of both ability types, one can uncover estimates of a students' two-dimensional ability vector.*

Because the outcomes are linear combinations of the underlying ability types, any two outcomes are correlated *iff* they share the same mix of ability types. As such, outcomes that are based mostly on cognitive skills will be highly correlated and these will be weakly correlated with outcomes that are based largely on non-cognitive skills. For simplicity, we can standardize the variance of each ability-type so that $\text{var}(\alpha_{c,ij}) = \text{var}(\alpha_{n,ij}) = 1$. From above, given that

$y_z = \beta_{c,z}\alpha_{c,ji} + \beta_{n,z}\alpha_{n,ji} + \varepsilon_{zi}$, and $\text{cov}(\alpha_{c,ij}, \alpha_{n,ij}) = 0$, it follows that all covariance across

outcomes are due to commonality in the underlying ability types across these outcomes so that

$\text{cov}(y_1, y_2) = \text{cov}(\beta_{c,1}\alpha_{c,j} + \beta_{n,1}\alpha_{n,j} + \varepsilon_{1i}, \beta_{c,2}\alpha_{c,j} + \beta_{n,2}\alpha_{n,j} + \varepsilon_{2i}) = \beta_{c,1}\beta_{c,2} + \beta_{n,1}\beta_{n,2}$. As such, the

covariance of the outcomes is given by [4], where E is a diagonal variance of the error terms.

$$
\text{var}\begin{bmatrix} y_1 \\ \\ \\ y_z \end{bmatrix} = \Sigma = \begin{bmatrix} \beta_{c,1}\beta_{c,1} + \beta_{n,1}\beta_{n,1} & \cdots & \cdots & \beta_{c,1}\beta_{c,z} + \beta_{n,1}\beta_{n,z} \\ \cdots & \cdots & \cdots & \cdots \\ \cdots & \cdots & \cdots & \cdots \\ \beta_{c,1}\beta_{c,z} + \beta_{n,1}\beta_{n,z} & \cdots & \cdots & \beta_{c,z}\beta_{c,z} + \beta_{n,z}\beta_{n,z} \end{bmatrix} + E = \beta\beta' + E \qquad [4]
$$

It also follows that $\text{cov}(y_z, \alpha_c) = \text{cov}(\beta_{c,z}\alpha_{c,j} + \beta_{n,z}\alpha_{n,j} + \varepsilon_{zi}, \alpha_{c,j}) = \beta_{c,z}$ and similarly

$\text{cov}(y_z, \alpha_n) = \text{cov}(\beta_{c,z}\alpha_{c,j} + \beta_{n,z}\alpha_{n,j} + \varepsilon_{zi}, \alpha_{n,j}) = \beta_{n,z}$. Accordingly, $\begin{bmatrix} Y \\ \alpha \end{bmatrix}$ is distributed with mean

$\begin{bmatrix} \mu \\ 0 \end{bmatrix}$ and variance $\begin{bmatrix} \beta\beta' + E & \beta \\ \beta' & I \end{bmatrix}$. From Johnson and Wichern (2007), the conditional

expectation of the underlying factor vector α is given by [5] below.[5]

$$
E(\alpha \mid Y) = \beta'(\beta\beta' + E)^{-1}(Y - \mu) \qquad [5]
$$

As such, the best linear unbiased predictor of each underlying factor is a linear combination of

each of the outcomes y such that $\hat{\alpha}_1 = m_1 Y$ and $\hat{\alpha}_2 = m_2 Y$, where $m_1$ and $m_2$ are defined from [5].

Simply put, one can predict student non-cognitive ability with a weighted average of their short-

run non-test-score outcomes and one can predict student cognitive ability with a weighted

average of their short-run test-score outcomes.

**Proposition 3:** *One can express a teacher's overall effect on the long-run outcome as a linear*

*combination of her effect on the predicted cognitive and non-cognitive factors.*

From above, $y_{*ij} = \beta_{c,*}\alpha_{c,ij} + \beta_{n,*}\alpha_{n,ij} + \varepsilon_{*i}$. Because $E(\hat{\alpha}_{c,ij}) = \alpha_{c,ij}$ and $E(\hat{\alpha}_{n,ij}) = \alpha_{n,ij}$, the

conditional expectation of long-run outcome based on all the available short-run outcomes is

---

[5] The general statement is that if $X = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$ is distributed with mean $\mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ and variance $\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}$ and

$|\Sigma_{22}| > 0$. Then the conditional expectation of $X_1$ given $X_2=x_2$ is $E(X_1 \mid x_2) = \mu_1 + \Sigma_{12}\Sigma_{22}^{-1}(x_2 - \mu_2)$. See Johnson and Wichern (2007) for a formal proof. A condition for identification is that the number of outcomes be equal to or greater than the number of factors. As such, the two factors can be only be identified with more than one outcome.

$\hat{y}_{*ij} = E[y_{*ij} | Y] = \beta_{c,*}\hat{\alpha}_{c,ij} + \beta_{n,*}\hat{\alpha}_{n,ij}$. From the law of iterated expectations, a teacher's predicted effect on the long-run outcome is her effect on the predicted long-run outcome, which is a weighted average of her effects on the predicted cognitive and non-cognitive factors. Specifically, $\hat{\theta}_{*j} = \beta_{c,*}\theta_{\hat{\alpha}_c j} + \beta_{n,*}\theta_{\hat{\alpha}_n j}$ where $\theta_{\hat{\alpha}_c j}$ and $\theta_{\hat{\alpha}_n j}$ are teacher $j$'s effects on predicted cognitive and predicted non-cognitive ability, respectively.

*Discussion of the model:* The model highlights that if test scores largely measure cognitive skills, improvements in test scores will be a good measure of how much a teacher improves cognitive skills but may not reflect the extent to which a teacher improves non-cognitive skills. The model further indicates that non-test score outcomes that are sensitive to non-cognitive skills might allow one to predict students' non-cognitive ability and thus identify teacher effects on non-cognitive skills. Because long-run outcomes reflect a combination of cognitive and non-cognitive skills, if we know the relationship between cognitive skills, non-cognitive skills, and long-run outcomes *and we have a mix of test score and non-test score outcomes*, one can identify teachers that are excellent at improving long-run student outcomes even if their effect on long-run outcomes are not directly observed and such teachers have small effects on test scores. I aim to use the implications of the model to (a) form predictions of student cognitive and non-cognitive ability, (b) estimate teacher effects on both cognitive and non-cognitive ability, (c) determine the extent to which teachers who improve test scores also improve non-cognitive outcomes, and (d) form estimates of a teacher's effect on long-run outcomes to determine the extent to which teacher effects on test scores may understate the overall effect of teachers on long-run outcomes.

**Section III     *Test Scores, Non-Cognitive Ability, and Long-Run Outcomes.***

This section uses data from the National Educational Longitudinal Survey of 1988 (NELS-88) to empirically test the underlying assumptions and implications of the model. Specifically, evidence is presented that (a) there are two underlying factors that explain most of the covariance between test score and non-test score outcomes in 8[th] grade that can be identified as cognitive and non-cognitive ability; (b) test scores are a good proxy for cognitive ability; (c) *both* higher cognitive and non-cognitive ability are independently associated with better adult outcomes; (d) as the model indicates, teacher effects on long-run outcomes should reflect the combination of her effects on the cognitive factor and the non-cognitive factor.

The NELS-88 is a nationally representative sample of respondents who were eighth-graders in 1988. These data contain information on short-run outcomes (absences, suspensions, GPA, grade repetition, and math and reading test scores) and long-run outcomes (being arrested or having a close friend who has been arrested, attending a post-secondary institution, and income in 1999 when most respondents were 25 years old). The left side of the top panel of Table 1 presents the correlations between the short-run outcomes in these data.

The first notable pattern is that test score outcomes are relatively strongly correlated with each other (math scores and reading scores have correlation≈0.8), are moderately correlated with grade point average (correlation≈0.38 for math and correlation≈0.34 for reading), and are weakly correlated with other non-test-score outcomes (the average absolute value of the individual correlations is 0.098). This is suggestive of an underlying cognitive factor that is highly predictive of standardized test scores, is moderately related to grades, and is largely unrelated to socio-behavioral outcomes. The other pattern is that the non-test score outcomes are more strongly correlated with each other (average correlation of 0.16) than with the test score outcomes (the average absolute value of the individual correlations is 0.098), and are moderately correlated with GPA (correlation≈0.35). This suggests that there is a non-cognitive factor that explains correlations between the non-test score outcomes, is an important determinant of grades, and is unrelated to test scores. The fact that GPA is correlated with both sets of variables is consistent with research (e.g. Howley, Kusimo, & Parrott, 2000; Brookhart, 1993) finding that most teachers base their grading on some combination of student product (exam scores, final reports etc.), student process (effort, class behavior, punctuality, etc.) and student progress — so that grades reflect a combination of cognitive and non-cognitive skill.

To assess the degree to which the variables can be classified as measuring cognitive and non-cognitive skills, I estimate the factor analytical model outlined in section II. From the model we know that $E(\alpha \mid Y) = \beta'(\beta\beta' + E)^{-1}(Y - \mu)$ and the covariance of the short-run outcomes is $\beta\beta' + E$. With a sufficiently large sample, the sample covariance matrix $S$ is a consistent estimate of $\beta\beta' + E$. If the error vector is normally distributed, one can estimate the factor

loadings $\beta$ by maximum likelihood.[6] As suggested in Johnson and Wichern (2007), I use $\hat{\alpha} = \hat{\beta}' \hat{S}^{-1}(Y - \bar{y})$ as my estimate of the underlying factor vector. If the model is reasonable, there should be two factors that are identifiable as measuring cognitive and non-cognitive skills.

Linear regressions reveal that the first estimated factor explains about 90 percent of the variability in math and reading scores, 18 percent of the variation in GPA, and less than 5 percent of the variability of all the other outcomes (right side of Table 1) — consistent with this factor measuring cognitive skills and test scores being a good proxy for cognitive skills. The second factor explains 45 percent of the variability in suspensions, 28 percent of the variation in absences, 47 percent of the variation in GPA, 20 percent of the variation in grade repetition, and less then five percent of the variation in math and reading test scores. This second factor explains much variability in socio-behavioral outcomes and is weakly correlated with the standardized test scores. It is reasonable to call this factor a measure of non-cognitive skill.

While I am agnostic about what exact set of skills are captured by this factor, studies suggest that each of these non-test score outcomes is associated with the same personality traits. Psychologists typically classify non-cognitive traits in terms of five dimensions; Neuroticism, Extraversion, Openness to Experience, Agreeableness, and Conscientiousness. Low levels of agreeableness and high neuroticism are associated with higher school absences, higher externalizing behaviors, greater juvenile delinquency and lower educational attainment (Lounsbury, Steel, Loveland, & Gibson, 2004; Barbaranelli, Caprara, Rabasca, & Pastorelli, 2003; John, Caspi, Robins, Moffit, & Stouthamer-Loeber, 1994; Carneiro, Crawford, & Goodman, 2007). High conscientiousness is associated with fewer absences, fewer externalizing behaviors, higher grades, more on-time grade progression, and higher educational attainment (Duckworth, Peterson, Matthews, & Kelly, 2007). This suggests that the factor explaining covariation between absences, grade progression, suspensions, and grades reflects a skill-set that is associated with high conscientiousness, high agreeableness, and low neuroticism. Consistent with this factor being largely unrelated to test scores, cognitive ability is associated with high openness, but largely unrelated to the other traits (Furnham, Monsen, & Ahmetoglu, 2009).

---

[6] The likelihood of the data is $L(\mu, \Sigma) = (2\pi)^{\frac{-np}{2}} |\Sigma|^{\frac{-n}{2}} e^{-\left(\frac{1}{2}\right) tr\left[\Sigma^{-1}\left(\sum_{j=1}^{n}(x_j - \bar{x})(x_j - \bar{x})' + n(\bar{x} - \mu)(\bar{x} - \mu)'\right)\right]}$, where $p$ is the number of short-run outcomes. Because the covariance matrix $\Sigma = \beta\beta' + E$ is a function of the factor loading vector $\beta$, the maximum likelihood estimate $\hat{\beta}$ of $\beta$ is the one that maximizes this likelihood (Johnson and Wichern 2007).

In Table 2, I show that both cognitive and non-cognitive ability independently predict long-run outcomes. I regress various adult outcomes on these two ability measures. Test scores and the two factors are standardized to be mean zero unit variance. To remove correlation between these ability measures and adult outcomes due to differences in socioeconomic status or demographic differences, *all models include controls for household income, whether English is the primary language at home, whether the student lives with their mother or father or both, the mothers and father highest level of education, family size, student race and student gender.* Columns 1 through 4 show that while test scores and cognitive ability have little relation to whether a student is arrested or has had a close friend who was arrested, a one standard deviation increase in the non-cognitive factor is associated with a 5.6 percentage point reduction (a 25 percent reduction relative to the sample mean). Column 1 presents results for the individual non-test score outcomes, and columns 2 and 3 show results using math and English test scores as a proxy for cognitive ability. Columns 5 through 8 predict college attendance based on a linear probability model. Column 8 indicates that a one standard deviation increase in cognitive ability is associated with a 6.3 percentage point increase in college going, while a one standard deviation increase in non-cognitive ability is associated with a 10.6 percentage point increase in college going (the baseline level is 0.79). Columns 6 and 7 indicate that using English or math test scores as proxies for cognitive ability yield similar results.

Columns 9 through 12 show marginal effects on the natural log of earnings based on unconditional quantile regression (Firpo, Fortin, & Lemieux, 2009) evaluated at the median of the earnings distribution.[7] The coefficients in column 12 indicate that a one standard deviation increase in cognitive ability is associated with 4.17 percent higher earnings, while a one standard deviation increase in non-cognitive ability is associated with an even larger 6.14 percent higher earnings. In recent findings, Lindqvist & Vestman (2011) and Heckman, Stixrud, & Urzua (2006) find that non-cognitive ability is particularly important at the lower end of the earnings distribution. To ensure that my measure of non-cognitive ability is valid, it is important to show similar relationships in these data. To test for this, I estimate quantile regressions to obtain the marginal effect on log wages at different points in the earnings distribution and present the estimated coefficients in Figure 2. The results are strikingly similar to Lindqvist & Vestman,

---

[7] I look at income conditional on any income data. Because roughly 90 percent of observations have income data the quantile regression results are similar for any reasonable treatment of missing data.

(2011). Specifically, cognitive ability is similarly important over all parts of the earnings distribution, while non-cognitive ability is very important at the lower end of the earnings distribution and becomes less important at higher ends of the distribution. This is consistent with the effects on arrests and echoes findings in the literature— thereby suggesting that this factor is a reasonable measure of non-cognitive ability.

The results indicate that student outcomes are indeed a reflection of both cognitive and non-cognitive ability, and in many cases the effect of non-cognitive ability is larger than that of cognitive ability. While the results indicate that students with higher test scores do tend to have higher non-cognitive ability *on average*, there is much variation in non-cognitive ability that is unmeasured by test scores and has important effects on long-run adult outcomes. Accordingly, teachers who improve non-cognitive ability may have important effects on adult outcomes even if they have no effects on test scores. In light of these relationships, the remainder of the paper aims to estimate the extent to which 9[th] grade teachers improve test scores, non-test score outcomes, and predicted cognitive and non-cognitive ability. I then aim to use these estimates to uncover teacher effects on long-run outcomes through *both* dimensions of ability, and determine the extent to which test score measures of quality may understate the importance of teachers.

## IV    Data:

To estimate the effect of teachers on student outcomes, this paper uses data on all public middle- and high-school students in North Carolina from 2005 to 2010 from the North Carolina Education Research Data Center (NCERDC). The data include demographics, transcript data on all courses taken, middle-school test scores, end of course scores for Algebra I and English I and codes allowing one to link students' end of course test-score data to individual teachers who administered the test.[8] I limit the analysis to students who took either the Algebra I or English I course (the two courses for which standardized tests have been consistently administered over time). Over 90 percent of all 9th graders take at least one of these courses so that the resulting sample is representative of 9th graders as a whole. To avoid endogeneity bias that would result from teachers having an effect on repeating 9th grade, the master data is based on only the first observation for when a student is in 9th grade. Summary statistics are presented in Table 3.

---

[8] Because the teacher identifier listed is not always the student's teacher, I use an algorithm to ensure high quality matching of students to teachers. I detail this in Appendix note 1.

These data cover 348,547 9th grade students in 619 secondary schools in classes with 4296 English I teachers, and 3527 Algebra I teachers. While roughly half of the students are male, about 58 percent are white, 29 percent are black, 7.5 percent are Hispanic, 2 percent are Asian, and the remaining one percent is Native American, mixed race, or other. About 7.5 percent of students have the highest parental education level (i.e. the highest level of education of the student's two parents) below high-school, 40 percent with a high school degree, about 15 percent with a junior college or trade school degree, 20 percent with a four year college degree or greater, and 6.4 percent with an advanced degree (about 10 percent of students are missing data on parental education). The test score variables have been standardized to be mean zero with unit variance for each cohort and test. Incoming 7th and 8th grade test scores in the final 9th grade sample are approximately 8 percent of a standard deviation higher than that of the average in 7th or 8th grade. This is because the sample of 9th grade students is less likely to have repeated a grade and to have dropped out of the schooling system.

Looking to the outcomes, the average number of absent days is 3 and the average of the log of days absent is 0.585. About 85 percent of students were in 10th grade the following year, and about 5.6 percent of all $9^{th}$ graders had a school suspension. The lower panel of Table 1 shows the correlations of these outcomes in the NCERDC data. While not identical, these variables have a similar correlational structure to the NELS-88 data such that students with better test scores do not necessarily have better non-test score outcomes. Specifically, test scores are strongly correlated with each other (corr = 0.616) and with GPA (corr ≈ 0.55), have a moderate correlation with being in $10^{th}$ grade on time (corr ≈ 0.32), and weak correlations with absences, and being suspended. To obtain measures of cognitive and non-cognitive skills that improve adult outcomes, I recreated the factors in the NCERDC data using the factor weights derived from the NELS-88 data. As in the NELS-88 data, cognitive ability explains much more variability in test scores than the non-cognitive factor, and the non-cognitive factor explains much more variation in GPA, on time grade progression, absences, and suspensions than the cognitive factor.[9] While there are some differences, the correlations between the non-cognitive

---

[9] The correlation between the non-cognitive factor using the correlational structure in the NCERDC data and the non-cognitive factor as predicted using the correlational structure in the NELS data is 0.98. Even though the results are similar, because I aim to use the same factor that is demonstrated to improve adult outcomes in the NELS-88 data, I create the factor using the weights derived from the correlational structure in the NELS-88 data.

factor and the outcomes are similar to those in the NELS-88 data such that there is much variation in non-cognitive ability that is unrelated to test scores. [10]

## *Measuring Tracks*

Even though schools may not have explicit labels for tracks, most do practice de-facto tracking by placing students of differing levels of perceived ability into distinct groups of courses (Sadker & Zittleman, 2006; Lucas & Berends, 2002). As highlighted in Jackson (2012) and Harris & Anderson (2012), it is not just the course that matters but also the levels at which students take a course. Indeed, even among students taking Algebra I or English I courses, there are three different levels of instruction (advanced, regular, and basic). As such, I exploit the richness of the data and take as my measure of a school-track, the unique combination of the 10 largest academic courses, the level of Algebra I taken, and the level of English I taken in a particular school.[11] As such, *only students who take the same academic courses, and take the same level of English I, and the same level of Algebra I, all at the same school are in the same school-track.*[12] Defining tracks flexibly at the school-by-course-group-by-course-level level allows for different schools to have different selection models and treatments for each track. Because many students pursue the same course of study, only 3.7 percent of all student observations are in singleton tracks, most students are in school-tracks with more than 50 students, and the average student is in a school-track with 117 other students.

## V    Empirical Strategy

This section outlines the strategy to estimate teacher effects on test score outcomes, non-test-score outcomes, and predicted cognitive and non-cognitive ability. The empirical approach is to model student outcomes as a function of lagged student achievement and student covariates, with the additional inclusion of controls for student selection to tracks and any treatments that are specific to tracks that might affect student outcomes directly. This removes the influence of

---

[10] One difference between the two datasets is that the non-cognitive factor explains more of the variation in test scores in the NCERDC data than in the NELS-88 data. This is likely due to the fact that the NELS test data were zero stakes, while the test data in the NCERDC has some stakes and may reflect a greater amount of non-cognitive ability (such as motivation). This is evidenced by test scores being more strongly correlated with on time grade progression and GPA in the NCERDC data than in the NELS.
[11] While there are hundreds of courses that students can take (including special topics and reading groups), there 10 academic courses that constitute two-thirds of all courses taken. They are listed in Appendix Table A1.
[12] Students taking the same courses at different schools are in different school-tracks. Students at the same school taking a different number of courses or at least one different course are in different school-tracks. Students at the same school taking the same courses but taking Algebra or English at different levels are in different school-tracks.

track-level treatments and selection to tracks on estimated teacher effects by comparing student outcomes across teachers within groups of students *in the same track at the same school*. Specifically, I model the outcomes $Y_{icjgys}$ of student $i$ in class $c$ with teacher $j$ in school-track $g$, at school $s$, in year $y$ with [6] below (note: most teachers are observed in multiple classes).

$$Y_{icjgys} = A_{iy\text{-}1}\delta + X_i\beta + I_{ji}.\theta_j + I_{gi}\,\theta_g + I_{sy}\,\theta_{sy} + \phi_c + \varepsilon_{icjgys} \qquad [6]$$

$A_{iy\text{-}1}$ is a matrix of incoming achievement of student $i$ (8$^{th}$ grade and 7$^{th}$ grade math and reading scores), $X_i$ is a matrix of student-level covariates (parental education, ethnicity, and gender), $I_{ij}$ is an indicator variable equal to 1 if student $i$ has teacher $j$ and equal to 0 otherwise so that $\theta_j$ is a time-invariant fixed effect for teacher $j$, $I_{gi}$ is an indicator variable equal to 1 if student $i$ is in school-track $g$ and $0$ otherwise so that $\theta_g$ is a time-invariant fixed effect for school-track $g$,[13] $I_{sy}$ is an indicator variable denoting whether the student is in school $s$ in year $y$ so that $\theta_{sy}$ is a school-by-year fixed effect, $\phi_c$ is a random classroom-level shock, and $\varepsilon_{ijgy}$ is mean zero random error term. In these models, by conditioning on school-tracks, one can obtain consistent estimates of the teacher effects $\theta_j$ as long as there is no selection to teachers *within* a school-track. In these models, the teacher effects are teacher-level means of the outcome after adjusting for incoming student characteristics, school-by-year level shocks and school-by-track effects. For test score outcomes, this model is just a standard value-added model with covariate adjustments.

Because the main models include school-by-track effects, all inference is made within school-tracks so that identification of teacher effects comes from two sources of variation; (1) comparisons of teachers at the same school teaching students in the same track *at different points in time*, and (2) comparisons of teachers at the same school teaching students in the same track *at the same time*. To illustrate these sources of variation, consider the simple case illustrated in Table 4. There are two tracks A and B in a single school. There are two math teachers at the school at all times, but the identities of the teachers change from year to year due to staffing changes. The first source of variation is due to changes in the identities of teachers over time due to staffing changes within schools. For example, between 2000 and 2005 teacher 2 is replaced by teacher 3. Because, teachers 2 and 3 both teach in track B (in different years) one can estimate the effect of teacher 2 relative to teacher 3 by comparing the outcomes of students in track B with teacher 2 in 2000 with those of students in tracks B with teacher 3 in 2005. To account for

---

[13] Note: In expectation, the coefficient on the school-track indicator variable reflects a combination of *both* the unobserved treatment specific and selection to school-track g.

differences in outcomes between 2000 and 2005 that might confound comparisons within tracks over time (such as school-wide changes that may coincide with the hiring of new teachers), one can use the change in outcomes between 2000 and 2005 for teacher 1 (who is in the school in both years) as a basis for comparison. In a regression setting this is accomplished with the inclusion of school-by-year fixed effects (Jackson & Bruegmann, 2009). This source of variation is valid as long as students do not select across cohorts (e.g. skip a grade) or schools in response to changes in Algebra I and English I teachers. Tests in section VI provide little evidence of such selection. The second source of variation comes from having multiple teachers for the same course in the same track at the same time. In the example, because both teachers 1 and 2 taught students in track B in 2000, one can estimate the effect of teacher 1 relative to teacher 2 by comparing the outcomes of teachers 1 and 2 among students in track B in 2000. This source of variation is robust to student selection to school-tracks and is valid as long as students do not select to teachers *within* school-tracks. Tests in section VI show that the findings are not driven by student selection within school-tracks.

To illustrate how much variation there is within school-tracks during the same year versus how much variation there is within school-tracks across years (cohorts), I computed the number of teachers in each non singleton school-track-year-cell for both Algebra I and English I (Appendix Table A2). About 63 and 51 percent of all school-track-year cells include one teacher in English and Algebra, respectively, so that for more than half the data the variation is based on comparing single teachers across cohorts within the same school-track. It follows that more than one-third of the variation is within school-track-year cells. Section V.4 shows that results using variation within school-track-cohort cells are similar to those obtained using only variation within school-tracks but across cohorts.

### *Illustrating the importance of accounting for tracks*

Because including school-track effects is nonstandard, it is important to illustrate the importance of conditioning on tracks. I do this by showing how conditioning on school-tracks affects the relationship between teacher experience and test scores. To do this, I regress English I and Algebra I scores on 8[th] and 7[th] grade math and reading scores and indicator variables for each year of teacher experience. I estimate models with school fixed effects, and then with track-by-school fixed effects. I plot the estimated coefficients on the years of experience indicator variables in Figure 3. Across all models, students perform better with Algebra teachers who have

16

more years of experience (top panel), but models with track-by-school effects are about 25 percent smaller than those without. The results for English are even more dramatic. Models with track-by-school effects are 66 percent smaller than those without. This suggests that either (a) more experienced teachers select into tracks with unobservably higher achievement students, or (b) tracks with more experienced teachers provide other supports (outside the teachers own classroom) that directly improve outcomes— indicating that failing to account for school-tracks may lead to much bias. For both subjects, about 60 percent of the variation in teacher experience occurs within school-tracks. Accordingly, a lack of variability within tracks *cannot* explain the reduced effects— underscoring the importance of accounting for school-tracks.

### *Estimating the Variance of Teacher Effects*

It is known that the variance of the estimated teacher effects $\hat{\theta}$ from [6] will overstate the variance of true teacher quality because of sampling variation and classroom-level shocks. As such, I follow Kane & Staiger (2008), Jackson (forthcoming) and Chetty, Friedman, & Rockoff (2011) and use the covariance between mean classroom-level residuals for the same teacher as my measure of the variance of teacher effects. This is done in two steps:

*Step 1:* Estimate equation [7] below.

$$Y_{icjgys} = A_{iy-1}\delta + X_i\beta + I_{gi}\theta_g + I_{sy}\theta_{sy} + \phi_c + \theta_j + \varepsilon_{icjgys} \qquad [7]$$

There are no teacher (or classroom) indicator variables so the total error term is $\varepsilon^* = \phi_c + \theta_j + \varepsilon_{igjy}$ (i.e. a teacher effect, a classroom effect, and the error term). I then compute mean residuals from [7] for each classroom $\bar{e}_c^* \equiv \theta_j + \phi_c + \hat{\varepsilon}_c$ where $\hat{\varepsilon}_c$ is the classroom-level mean error term.

*Step 2:* Link every classroom-level mean residual and pair it with a random different classroom level mean residual for the same teacher and compute the covariance of these mean residuals. That is, compute $cov(\bar{e}_c^*, \bar{e}^*_{c'} \mid J = j)$. To ensure that the estimate is not driven by any particular random pairing of classrooms within teachers, replicate this calculation 100 times and take the median of the estimated covariance as the parameter estimate. If the classroom errors $\phi_c$ are uncorrelated with each other (recall that the model includes school-by-year fixed effects) and are uncorrelated with teacher quality $\theta_j$, the covariance of mean residuals *within teachers* but *across classrooms* is a consistent measure of the true variance of persistent teacher quality (Kane &

Staiger, 2008). That is, $cov(\bar{e}_{c}^{*}, \bar{e}*_{c'} | J = j) = cov(\theta_{j}, \theta_{j}) = var(\theta_{j}) \longrightarrow \sigma_{\theta_{j}}^{2}$.[14]

As discussed in Jackson (2012) because the variance of the true teacher effects can be estimated by a sample covariance, one can compute confidence intervals for the sample covariance. I use the empirical distribution of 100 randomly computed "placebo" covariances (i.e. sample covariance across classrooms for *different* teachers) to form an estimate of the standard deviation of the sampling covariance across classrooms for the same teacher. I use this "bootstrapped" standard deviation of the covariance for normal-distribution-based hypothesis testing. Because most studies report the standard deviation of teacher effects, I report the square root of the sample covariance and the square root of the confidence bounds.

## VI    Main Results

Table 5 presents the estimated covariance across classrooms for the same teachers under three different models. Because standard deviations are positive by definition, when the sample covariance is negative (*note that none of the negative covariance estimates is statistically significantly different from zero at the 5 percent level*), I report the standard deviation to be zero. I present naïve models that include school fixed effects and attempt to account for tracking with peer characteristics (mean peer $8^{th}$ and $7^{th}$ grade math and reading scores in addition to mean peer demographics). I then present models that account for tracking with school-track fixed effects. And finally I estimate the preferred models that include both school-track fixed effects and school-by-year effects to account for both bias due to tracking and any school-wide shocks that might be confounded with teacher effects.

In models that include school effects and peer characteristics, the estimated standard deviation of Algebra teacher effects on Algebra test scores (left panel) is $0.12\sigma$ (in student outcome units). This is similar to estimates found in other studies of high school teachers that do not account for tracking explicitly. In this naïve model, Algebra teachers have statistically significant effects on all outcomes except for English scores, and the standard deviations of Algebra teacher effects on cognitive and non-cognitive ability are $0.085\sigma$ and $0.146\sigma$, respectively. For English teachers (right panel), the models indicate that the standard deviation of English teacher effects on English scores is $0.049\sigma$. In this naïve model, English teachers have

---

[14] Note: $cov(\theta_{j}, \phi_{c'}) = cov(\theta_{j}, \bar{e}_{jgyc'}) = cov(\phi_{c}, \theta_{j}) = cov(\phi_{c}, \phi_{c'}) = cov(\phi_{c}, \bar{e}_{jgyc'}) = cov(\bar{e}_{jgyc}, \theta_{j}) = cov(\bar{e}_{jgyc}, \phi_{c'}) = cov(\bar{e}_{jgyc}, \bar{e}_{jgyc'}) = 0$

significant effects on all non-test score outcomes, and the standard deviations of English teacher effects on cognitive and non-cognitive ability are 0.052σ and 0.106σ, respectively.

The middle panel presents results that include track-by-school fixed effects instead of peer characteristics to account for tracking. In such models, the estimated effects fall by about 20 percent relative to only including peer level characteristics. Adding additional controls for school-by-year effects reduces the effects by a further 30 to 50 percent— suggesting that there may be bias associated with omitting both school-track effects and not accounting for school-by year shocks. In the preferred model (lowest panel), the standard deviation of the Algebra teacher effect on Algebra test scores is 0.066σ; an estimate about half the size as that obtained in the naïve model that does not account for tracking or school specific shocks. The estimated teacher effects are statistically significantly different from zero for some non-cognitive outcomes such that the standard deviation of teacher effects on GPA is 0.045σ and that on enrolling in 10[th] grade is 0.025. Algebra teachers do appear to influence both students' cognitive and non-cognitive ability such that the standard deviation of the effects on cognitive and non-cognitive ability are 0.046σ and 0.069σ, respectively. To assuage concerns that the effects on non-cognitive ability are due to the correlation with test scores, I took the estimated non-cognitive ability and removed any correlation with test scores so that any effects on this measure will not be detected by changes in test scores. The correlation between this orthonogalized non-cognitive measure and the raw measure is 0.87. The lower panel shows that the standard deviation of the effects on non-cognitive ability (that is not captured by test scores) is 0.064σ.

Looking to English I teachers, in the preferred model (right lower panel), the standard deviation of English teacher effects on English test scores is 0.03σ (about half the magnitude of most estimates in the literature). While there is no effect of English teachers on Algebra scores, the estimated teacher effects are statistically significantly different from zero for all of the non-test score outcomes such that the standard deviation of teacher effects on the likelihood of being suspended is 0.014, the effect on log of absences is 0.037, that on GPA is 0.027, and that on enrolling in 10[th] grade is 0.024. Summarizing these effects, the standard deviation of the English teacher effects on the cognitive and non-cognitive factors are 0.023σ (*p*-value of 0.083) and 0.054σ (*p*-value of 0.007), respectively. Similar to Algebra, the effects on non-cognitive ability are largely undetected by test scores.

To put the non-test score estimates into perspective, having an Algebra or English teacher

at the 85[th] percentile of GPA quality versus one at the 15[th] percentile would be associated with 0.09 and 0.054 higher GPA, respectively. For both subjects, a teacher at the 85 percentile of on-time grade progression quality versus one at the 15[th] percentile would be associated with being 5 percentage points ($0.14\sigma$) more likely to enroll in 10[th] grade on time. Given that not enrolling in 10[th] grade is a strong predictor of dropout, this suggests significant teacher effects on dropout—consistent with Koedel (2008). Having an English teacher at the 85[th] percentile of absences and suspensions quality versus one at the 15[th] percentile would be associated with being 2.8 percentage points ($0.12\sigma$) less likely to be suspended, and having 7.4 percent fewer days absent. Teachers in both subjects have larger effects on non-cognitive skills than on cognitive skills. Overall, having an Algebra teacher at the 85[th] percentile of improving non-cognitive ability versus one at the 15[th] percentile would be associated with $0.14\sigma$ higher non-cognitive ability, while the same calculation for English teachers is $0.11\sigma$. If non-cognitive ability is as important, or possibly more important, in determining long-run outcomes as cognitive ability, and test score effects are a better measure of a teacher's effect on cognitive ability than her effects on non-cognitive ability, then test-score-based measures of teacher quality may drastically understate the importance of teachers for long-run outcomes. I investigate this in section VII.

### *Tests for Bias due to Selection*

Several papers indicate that conditioning on a single lag of achievement removes most selection bias (Koedel & Betts, 2011; Kinsler, 2012; Kane & Staiger, 2008) and others indicate that the inclusion of two lags of achievement effectively removes bias due to dynamic student sorting to teachers (Rothstein, 2010). However, it is worth testing these assumptions in these data. I present a test for selection on observables and a test for selection on unobservables.

To test for selection on observables, I show that conditioning on lagged test scores eliminates selection to teacher quality with regards to other observable student characteristics. Following Chetty, Friedman, & Rockoff (2011), I predict each outcome (based on 7th grade math and reading scores, parental education, gender, and ethnicity) and regress predicted outcomes on school-track fixed effects, year fixed effects, the estimated effect of the student's teacher (estimated out of sample)[15], and 8th grade test scores. If students with better

---

[15] To remove any endogeneity, for each observation year, I estimate teacher effects using all *other* years of data. For example for observations in 2005, the estimated teacher effects are based on teacher performance in 2006, 2007, 2008, 2009, and 2010. For estimates in 2008, estimates are based on 2005, 2006, 2007, 2009, and 2010. I follow

characteristics select to classrooms based on teacher effectiveness, there would be a systematic relationship between estimated teacher quality and predicted outcomes. The results (top row of lower panel of Table 6) indicate that where there are multiple teachers for a given cohort of students in the same school-track, teachers with higher estimated effects do not receive students in the cohort with better or worse predicted outcomes on average— suggesting no selection.

To test for selection on unobservables within school-track-years, I follow Chetty, Friedman, & Rockoff (2011) and exploit the statistical fact that the effects of any selection among students *within* the same school-track and cohort will be eliminated by aggregating the treatment to the school-track-year level and relying only on cohort-level variation across years within school-tracks. That is, if the estimated teacher effects merely capture student selection to teachers, then the arrival of a teacher with a large positive estimated effect (who increases the average estimated teacher effect for a cohort) should have no effect on student outcomes for that cohort. Conversely, if the estimated effects are real, differences in average estimated teacher quality across cohorts (driven by changes in teaching personnel within schools over time) should be associated with similar differences across cohorts in average cohort level outcomes as the same difference in estimated teacher quality across individual students (due to there being multiple teachers in the same school-track at the same time) within the same cohort.

To test this, I estimate equations [8] and [9] where $\hat{\theta}_j$ is the estimated (out of sample) effect of teacher $j$, $\overline{\hat{\theta}}_{sgy}$ is the mean estimated teacher effect in school-track $g$ in year $y$, $\theta_g$ is a school-track effect, $\theta_{sy}$ is a school-year effect, and $\theta_{sgy}$ is a school-track-year effect.

$$Y_{ijcy} = A_{iy-1}\delta + \psi_1\hat{\theta}_j + X_{iy}\beta + \theta_{sgy} + \varepsilon_{ijcy} \tag{8}$$

$$Y_{ijcy} = A_{iy-1}\delta + \psi_2\overline{\hat{\theta}}_{sgy} + X_{iy}\beta + I_{gi}\theta_g + \theta_{sy} + \varepsilon_{ijcy} \tag{9}$$

Equations [8] and [9] both estimate the effect of having teachers with higher estimated effects, but each use a distinct source of variation. In [8], teacher quality is defined at the student level and the model includes a track-school-year fixed effect, so that it only makes comparisons among students with different teachers in the same school-track and year (removing all variation due to personnel changes over time). In contrast, by defining teacher quality at the school-track-cohort level in [9], one no longer compares students within the same school-track-year (where selection is likely), and only compares *cohorts* of students in the same school-track over time

(where selection is unlikely because variation in this aggregate measure is due *only* to changes in the identities of teachers in the school-track over time). To control for school-level changes that could affect the cohort-level results, all models include school-by-year fixed effects. Relating the predictions to the equations directly, if there is no sorting $\psi_1$ should be similar to $\psi_2$, and if the effects are due to sorting $\psi_2$ will be equal to 0. Note that because the teacher effects are estimated with noise, the coefficients $\psi_1$ and $\psi_2$ will be less than 1. The results are presented in the top panel of Table 6. Despite no relationship between estimated teacher quality and *predicted* outcomes, there are economically and statistically significant effects of estimated teacher quality on *actual* outcomes for both subjects. For both subjects, marginal effects obtained using variation within school-track-cohorts are similar to those obtained using variation across cohorts within school-tracks. For all outcomes, mean school-track-cohort-level teacher quality has a statistically significant effect on outcomes so that the null hypothesis that estimated teacher effects are driven by selection within school-track-cohorts can be rejected at the 5 percent level.

## VII    *Relationship between Teacher Effects across Outcomes*

Having established that teachers have real causal effects on test scores, non-test score outcomes, and predicted cognitive and non-cognitive ability, this section documents the extent to which test score based value-added understates the importance of teachers on long-run outcomes. To gain a sense of whether teachers who improve test scores also improve other outcomes, I regress the estimated teacher effects for all the outcomes on the Algebra test score, English test score, cognitive ability and non-cognitive ability effects and report the $R^2$ in Table 7. The reported $R^2$ measures the fraction of teacher effects on each outcome that can be explained by (or is associated with) effects on test scores, cognitive ability, or non-cognitive ability. The top panel presents effects for Algebra teachers. For all outcomes, Algebra teachers with higher test score value-added are associated with better non-test score outcomes, but the relationships are weak. Algebra teacher effects on Algebra test scores explain 1.15 percent of the variance in estimated teacher effects on suspensions,  2.11 percent of the estimated effect on absences, 10.57 percent of the effect on GPA, and 5.03 percent of the effect on on-time 10th grade enrollment (top panel top row). This indicates that while teachers who raise test score may also be associated with better non-test-score outcomes, most of effects on non-test score outcomes are unrelated to effects on test scores. The results for cognitive ability are consistent with this. Algebra teacher

effects on cognitive ability explain 75 percent of the variation in Algebra test score value added, but only 3.74 percent of the variance in estimated teacher effects on suspensions, 6.3 percent of the estimated effect on absences, 28.3 percent of the effect on GPA, and 11.3 percent of the effect on on-time 10th grade enrollment (top panel second row). In contrast to effects on test scores and cognitive ability, effects on non-cognitive ability explain much of the estimated effects on the non-test score outcomes. Specifically, Algebra teacher effects on non-cognitive ability explain 18.4 percent of the estimated teacher effects on suspensions, 47.6 percent of the estimated effect on absences, 80 percent of the effect on GPA, and 40 percent of the effect on on-time 10th grade enrollment (top panel third row). However, teacher effects on non-cognitive ability explain only 10 percent of the variance in estimated teacher effects on Algebra test scores.

Results for English teachers follow a similar pattern. English teacher effects on English test scores explain little of the estimated effects on non-test score outcomes. Specifically, teacher effects on English test scores explain only 1.16, 2.47, 5.16, and 5.8 percent of the variance in teacher effects on suspensions, absences, GPA, and on-time 10th grade enrollment, respectively (lower panel top row). Unlike with Algebra test score value added, English teacher effects on cognitive ability explain only 25 percent of the variation in English test score value-added. This would imply that *some* teachers who raise English test scores are not also improving GPA and on time grade enrollment (which both form part of the cognitive measure). This could be an artifact of the English test such that English test scores are noisy measures of cognitive ability. Consistent with this view, several studies find that scores on high school math tests are stronger predictors of adult success than scores on high school English tests (findings echoed in Table 2). English teacher effects on cognitive ability explain only 2.62, 4.2, 13.9, and 6.98 percent of the variance in teacher effects on suspensions, absences, GPA, and on-time 10th grade enrollment, respectively. However, English teacher effects on non-cognitive ability explain only 6.8 percent of the variance in estimated teacher effects on English test scores, but explain 18.24 percent of the variance in estimated teacher effects on suspensions, 46.58 percent of the estimated effect on absences, 79.68 percent of the effect on GPA, and 48.67 percent of the effect on on-time 10th grade enrollment (bottom panel third row).

In sum, for both subjects, teacher effects on cognitive scores measure an important set of skills and teacher effects on non-cognitive scores measure a *different* and equally important set of skills. For both subjects, a teacher's effect on test scores is a better measure of her effect on

cognitive ability than her effect on non-cognitive ability. To show this visually, Figure 4 presents a scatterplot of the estimated effects on cognitive ability and non-cognitive ability against her effect on test scores. It is clear that (a) a teacher's effect on test scores captures much of the effect on cognitive ability (particularly for math), (b) a teacher's effect on test scores captures some of the effect on cognitive ability, and (c) teacher effects on test scores in both subjects leave much variability in effects on non-cognitive ability unexplained.

Because variability in outcomes associated with individual teachers that is unexplained by test scores is not just noise, but is systematically associated with their ability to improve typically unmeasured non-cognitive skills, classifying teachers based on their test score value-added will likely lead to large shares of excellent teachers being deemed poor and *vice versa*. This also implies that there could be considerable gains associated with measuring teacher effects on both test score and non-test score outcomes over using test score measures alone. Another implication is that if teachers must expend less effort improving non-cognitive ability in order to improve cognitive ability, regimes that increase the external rewards for test scores (such as paying teachers for test score performance or test-based accountability) may undermine the creation of students' non-cognitive skills (Holmstrom & Milgrom, 1991). In light of the large estimated benefits to higher non-cognitive skills (particularly for students at the lower end of the earnings distribution) in Table 2, this may be cause for concern. Finally, the results suggest that the effect of teachers on student long-run outcomes may be much larger than that suggested by comparing the adult outcomes of students with high versus low value-added teachers.

### *How much do Test Score Measures Understate the Importance of Teachers?*

The results thus far have shown that (a) variation in non-cognitive outcomes may be more determinative of adult outcomes than variation in cognitive ability, (b) teachers have effects on non-cognitive ability that are larger than their effects on cognitive ability, and (c) teacher effects on test scores explain a modest fraction of their effects on non-cognitive skills. Taken together, this implies that test score based measures of teacher quality will understate the true importance of teachers for long-run outcomes and many of the most effective teachers will not be identified based on test score based measures. To see this point, consider this thought exercise. If cognitive ability and non-cognitive ability were uncorrelated, were equally important in determining wages, test scores were a perfect proxy for cognitive ability, and teachers had the exact same effects sizes on cognitive and non-cognitive ability, then teacher effects on student wages as

measured by cognitive ability would reflect roughly 50 percent of a teachers' overall effect. If now test score effects capture only 75 percent of the effect on cognitive ability (as is the case for Algebra) then test score effects will reflect roughly 50*0.75=37.5 percent of a teachers' overall effect. If teacher effects are larger on non-cognitive ability than for cognitive outcomes (as is the case for both subjects) test scores effects would explain even less of the overall effect. If the marginal effect of non-cognitive ability is larger than that for cognitive ability (as is the case for arrests, college going, and wages) then the fraction of the overall variability in teacher effects on the long-run outcomes captured by test score value-added will be smaller still. However, if teacher effects on cognitive ability and non-cognitive ability were positively (negatively) correlated, then the fraction of the overall variability in teacher effects on the long-run outcomes captured by test score value-added will be larger (smaller). This discussion highlights that the extent to which test score value-added explains variability in teachers' overall effects depends on (a) the relationship between the importance of cognitive- and non-cognitive ability in determining adult outcomes, (b) the extent to which improvements in test scores measure improvements in cognitive and non-cognitive ability, (c) the relative magnitude of teacher effects on cognitive and non-cognitive ability, and (d) the relationship between effects on both dimensions of ability. This section presents estimates that take all these parameters into account.

As discussed in Section II, one can express a teacher's overall predicted effect on long-run outcomes as a linear combination of her effect on predicted cognitive ability and her effect on predicted non-cognitive ability. Accordingly, we can predict a teacher's effect on college-going or wages based on her effect on predicted cognitive and non-cognitive ability as long as we know the marginal effect of cognitive and non-cognitive ability on wages. While there is no way to know for certain, the results in Table 2 present some guidance. Specifically, column 8 of Table 2 indicates that the marginal effect of increasing cognitive ability by 1sd is to increase college going by 6.32 percentage points, and the marginal effect of increasing non-cognitive ability by 1sd is to increase college-going by 10.6 percentage points. One could compute a teacher's predicted effect on college-going based on both her contribution to cognitive skills and non-cognitive skills as $\hat{\theta}_j^{college} = (0.0632)\hat{\theta}_j^{cognitive} + (0.106)\hat{\theta}_j^{non-cognitive}$. The $R^2$ of a regression of $\hat{\theta}_j^{college}$ on $\hat{\theta}_j^{scores}$ will provide a measure of how much of the predicted effect on college-going can be explained by a teacher's effect on test scores. Similarly, using estimates from column 12, one

can compute a teacher's predicted effect on log wages based on both her contribution to cognitive and non-cognitive skills as $\hat{\theta}_j^{wages} = (0.0417)\hat{\theta}_j^{cognitive} + (0.0614)\hat{\theta}_j^{non-cognitive}$. The $R^2$ of a regression of $\hat{\theta}_j^{wages}$ on $\hat{\theta}_j^{scores}$ will provide a measure of how much of the predicted effect on log wages can be explained by a teacher's effect on test scores.

The problem with simply using the estimated relationships from the NELS-88 data is that the relationships are based on correlations and the variables are measured with error, so one cannot be certain that these represent causal relationships. As such, I construct teacher's predicted effects on college-going and wages under different assumptions about the relative importance of non-cognitive ability. That is, when the relative non-cognitive weight for college going is $r$ then, $\hat{\theta}_j^{college} = (0.0632) \cdot (\hat{\theta}_j^{cognitive} + r \cdot \hat{\theta}_j^{non-cognitive})$ and when the relative non-cognitive weight for log wages is $r$ then $\hat{\theta}_j^{\log wages} = (0.0417) \cdot (\hat{\theta}_j^{cognitive} + r \cdot \hat{\theta}_j^{non-cognitive})$. The relative weight for college-going from the empirical estimates is 0.106/0.0632=1.67 and that for log wages is 0.0614/0.0417=1.47. This is consistent with existing studies finding that non-cognitive abilities may be at least as important as cognitive skills in predicting adult outcomes (e.g. Bowles, Gintis, & Osborne, 2001; Jencks, 1979; Heckman, Stixrud, & Urzua, 2006).

To show how much of the variability in a teacher's predicted effect on college going and wages is explained or detected by her effect on test scores, Figure 5 plots the $R^2$ of a regression of $\hat{\theta}_j^{college}$ or $\hat{\theta}_j^{wages}$ on $\hat{\theta}_j^{scores}$ for different weights $r$. Because both the $R^2$ and the relative weight are invariant to the unit of measurement, the pattern in Figure 5 will hold for any long-run outcome. For Algebra teachers, when the relative weight is 0, the estimated teacher effect on test scores predicts 75 percent all of the predicted teacher effect on wages. This is because teacher effects on Algebra tests explain only 75 percent of the variability in teacher effects on cognitive ability. As one increases the contribution of non-cognitive ability (which test scores only weakly explain) to 0.5 (such that non-cognitive skills are half as important as cognitive skills in determining the outcome), Algebra test score effects explain about 55 percent of the full effect. If the weight on non-cognitive ability is 1 (i.e. non-cognitive skills are as important as cognitive skills), Algebra test score effects explain about 40 percent of the full effect. For the empirical weight for wages of 1.47, Algebra test score effects explain only 33 percent of the full effect, and at the empirical weight for college-going of 1.67, Algebra test score effects explain only 29

percent of the teacher's full effect. This drop stems from the fact that as the weight on non-cognitive ability increases there is more variability in the predicted effect that is weakly correlated with test score effects so that the variability explained by test score effects approaches 0.1 (the fraction of variability in non-cognitive ability predicted by Algebra test scores).

The results are even more dramatic for English teachers. Because teacher effects on English tests explain only 25 percent of the variability in teacher effects on cognitive ability, even when the relative non-cognitive weight is zero, the estimated teacher effect on test scores predicts only 25 percent all of the predicted teacher effect on college-going or wages. If the weight on non-cognitive ability is 1, the fraction of the overall effect predicted by English test score effects falls slightly to 20 percent. For the empirical weight for wages of 1.47, English test score effects explain only 17 percent of the full effect, and at the empirical weight for college-going of 1.67, English test score effects explain only about 17 percent of the teacher's full effect. To be clear, *this implies that English teachers are important, and may be much more so than their effects on test scores might suggest*.

These results suggest that under reasonable relative weights for the importance of non-cognitive ability, a teacher's effect on college-going and earnings as measured based on her effect on test scores alone may only reflect between 20 and 40 percent of her overall effect. In the context of recent findings that having a teacher at the 85[th] percentile of the test score value-added distribution versus a median teacher is associated with students being 0.5 percentage points more likely to attend college and earning $186 more per year (Chetty et. al. 2011), calculations suggest that having a teacher at the 85[th] percentile of the quality distribution for college-going and wages versus a median teacher may be associated with being as much as 2.5 percentage points more likely to attend college and earning as much as $930 more per year. While the plausible range of values is wide, the calculations illustrate that estimated teacher effects based on test scores may drastically understate the overall importance of teachers. This is consistent with Fredriksson, Ockert, & Oosterbeek (forthcoming) who find that the imputed wage effect of class-size based on test scores are substantially smaller than the direct estimates. This under-estimation of the importance of teachers may be particularly large for outcomes such as being arrested that are largely functions of non-cognitive rather than cognitive skills.

**VIII** *Conclusions*

This paper presents a two-factor model such that all student outcomes are a function of student cognitive and non-cognitive ability. The model shows that one can use a variety of short-run outcomes to construct a measure of student cognitive and non-cognitive ability and use this to estimate a teacher's predicted effect on long-run outcomes. Using longitudinal survey data, evidence indicates that non-cognitive ability (an index of non-test score socio-behavioral outcomes in 8[th] grade) is associated with sizable improvements in adult outcomes, and that non-cognitive ability is at least as important a determinant of adult outcomes as cognitive ability.

Using administrative data with students linked to individual teachers, 9[th] grade English and Algebra teachers have economically meaningful effects on test scores, absences, suspensions, on time 10[th] grade enrollment, and grades. The results indicate that teachers have larger effects on non-cognitive ability than they do on cognitive ability. A variety of additional tests suggest that these effects can be interpreted causally. Teacher effects on test scores and teacher effects on non-cognitive ability are weakly correlated such that many teachers in the top of test score value-added distribution will also be among the bottom of teachers at improving non-cognitive skills. This means that a large share of teachers thought to be highly effective based on test score performance will be no better than the average teacher at improving college-going or wages. Under reasonable assumptions about the importance of non-cognitive skills for long-run outcomes, calculations indicate that test score based measures may understate the importance of teachers on college-going and wages by between 50 and 80 percent. In fact, results indicate that test score based measures may understate the importance of teachers to an even greater extent for outcomes that depend heavily on non-cognitive skills such as being arrested.

This study highlights that a failure to account for the effect of educational interventions on non-cognitive skills can lead to biased estimates of the effect of such interventions on important long-run outcomes. The two-sample factor analytic framework put forth in this paper can be used in other settings to estimate the effects of educational interventions on both cognitive and non-cognitive ability. Finally, the results illustrate the large potential gains that may be associated with making policy decisions about educational interventions based on estimated effects on both cognitive and non-cognitive outcomes rather than just test scores alone.

# Bibliography

1. Aaronson, D., Barrow, L., & Sander, W. (2007). Teachers and Student Achievement in the Chicago Public High Schools. *Journal of Labor Economics, 25*, 95-135.
2. Alexander, K. L., Entwisle, D. R., & Thompson, M. S. (1987). School Performance, Status Relations, and the Structure of Sentiment: Bringing the Teacher Back in. *American Sociological Review, 52*, 665-82.
3. Barbaranelli, C., Caprara, G. V., Rabasca, A., & Pastorelli, C. (2003). A questionnaire for measuring the BigFive in late childhood. *Personality and Individual Differences, 34*(4), 645-664.
4. Booker, K., Sass, T. R., Gill, B., & Zimmer, R. (2011). The Effect of Charter High Schools on Educational Attainment. *Journal of Labor Economics, 29*(2), 377-415.
5. Borghans, L., Weel, B. t., & Weinberg, B. A. (2008). Interpersonal Styles and Labor Market Outcomes. *Journal of Human Resources, 43*(4), 815-58.
6. Bowles, S., Gintis, H., & Osborne, M. (2001). The Determinants of Earnings: A Behavioral Approach. *Behavioral Approach, 39*(4), 1137-76.
7. Brookhart, S. M. (1993). Teachers' Grading Practices: Meaning and Values. *Journal of Educational Measurement, 30*(2), 123-142.
8. Carneiro, P., Crawford, C., & Goodman, A. (2007). The Impact of Early Cognitive and Non-Cognitive Skills on Later Outcomes. *CEE Discussion Papers 0092*.
9. Cascio, E., & Staiger, D. (2012). Knowledge, Tests, and Fadeout in Educational Interventions. *NBER working Paper Number 18038*.
10. 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. *Quarterly Journal of Economics, 126*(4), 1593-1660.
11. Chetty, R., Friedman, J., & Rockoff, J. (2011). The Long-Term Impacts of Teachers: Teacher Value-Added and Student Outcomes in Adulthood. *unpublished manuscript*.
12. Deming. (2009). Early Childhood Intervention and Life-Cycle Skill Development: Evidence from Head Start. *American Economic Journal: Applied Economics, 1*(3), 111-134.
13. Deming, D. (2011). Better Schools, Less Crime? *The Quarterly Journal of Economics, 126*(4), 2063-2115.
14. Downey, D., & Shana., P. (2004). When Race Matters: Teachers' Evaluations of Students' Classroom Behavior. *Sociology of Education, 77*, 267-82.
15. Duckworth, A. L., Peterson, C., Matthews, M. D., & Kelly, D. R. (2007). Grit: Perseverance and Passion for Long-Term Goals. *Journal of Personality and Social Psychology, 92*(6), 1087–1101.
16. Ehrenberg, R. G., Goldhaber, D. D., & Brewer, D. J. (1995). Do teachers' race, gender, and ethnicity matter? : evidence from NELS88. *Industrial and Labor Relations Review, 48*, 547-561.
17. Firpo, S., Fortin, N. M., & Lemieux, T. (2009). Unconditional Quantile Regressions. *Econometrica, 77*(3), 953-973.
18. Fredriksson, P., Ockert, B., & Oosterbeek, H. (forthcoming). Long-Term Effects of Class Size. *Quartlerly Journal of Economics*.
19. Furnham, A., Monsen, J., & Ahmetoglu, G. (2009). Typical intellectual engagement, Big Five personality traits, approaches to learning and cognitive ability predictors of academic performance. *British Journal of Educational Phycology, 79*, 769-782.
20. Harris, D., & Anderson, A. (2012). Bias of Public Sector Worker Performance Monitoring: Theory and Empirical Evidence From Middle School Teachers. *Association of Public Policy Analysis & Management*. Baltimore.
21. Heckman, J. (1999). Policies to Foster Human Capital. *NBER Working Paper 7288*.
22. Heckman, J. J., & Rubinstein, Y. (2001). The Importance of Noncognitive Skills: Lessons from the GED Testing Program. *American Economic Review, 91*(2), 145-49.
23. 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), 411-82.
24. Heckman, J., Pinto, R., & Savelyev, P. (forthcoming). Understanding the Mechanisms Through Which an Influential Early Childhood Program Boosted Adult Outcomes. *American Economic Review*.
25. Holmstrom, B., & Milgrom, P. (1991). Multitask Principal-Agent Analysis: Incentive Contracts, Asset Ownership and Job Design. *Journal of Law, Economics and Organization, 7*(Special Issue), 24-52.
26. Howley, A., Kusimo, P. S., & Parrott, L. (2000). Grading and the ethos of effort. *Learning Environments Research, 3*, 229-246.

27. Jackson, C. K. (2012). Teacher Quality at the High-School Level: The Importance of Accounting for Tracks. *NBER Working Paper 17722*.
28. Jackson, C. K. (forthcoming). Match quality, worker productivity, and worker mobility: Direct evidence from teachers. *Review of Economics and Statistics*.
29. Jackson, C. K., & Bruegmann, E. (2009). Teaching Students and Teaching Each Other: The Importance of Peer Learning for Teachers. *American Economic Journal: Applied Economics, 1*(4), 85-108.
30. Jencks, C. (1979). *Jencks, Christopher.* New York: Basic Books.
31. Jennings, J. L., & DiPrete, T. A. (2010). Teacher Effects on Social and Behavioral Skills in Early Elementary School. *Sociology of Education, 83*(2), 135-159.
32. John, O., Caspi, A., Robins, R., Moffit, T., & Stouthamer-Loeber, M. (1994). The "Little Five": exploring the nomological network of the Five-Factor Model of personality in adolescent boys. *Child Development, 65*, 160–178.
33. Johnson, R. A., & Wichern, D. W. (2007). *Applied Multivariate Statistical Analysis.* Pearson.
34. Kane, T., & Staiger, D. (2008). Estimating Teacher Impacts on Student Achievement: An Experimental Evaluation. *NBER working paper 14607*.
35. Kane, T., & Staiger, D. (2008). Estimating Teacher Impacts on Student Achievement: An Experimental Evaluation. *NBER Working Paper # 14607*.
36. Kinsler, J. (2012). Assessing Rothstein's critique of teacher value-addedmodels. *Quantitative Economics, 3*, 333-362.
37. Koedel, C. (2008). An Empirical Analysis of Teacher Spillover Effects in Secondary School. *Department of Economics, University of Missouri Working Paper 0808*.
38. Koedel, C. (2008). Teacher Quality and Dropout Outcomes in a Large, Urban School District. *Journal of Urban Economics, 64*(3), 560-572.
39. Koedel, C., & Betts, J. (2011). Does Student Sorting Invalidate Value-Added Models of Teacher Effectiveness? *Education Finance and Policy, 6*(1), 18-42.
40. 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.
41. Lounsbury, J. W., Steel, R. P., Loveland, J. M., & Gibson, L. W. (2004). An Investigation of Personality Traits in Relation to Adolescent School Absenteeism. *Journal of Youth and Adolescence, 33*(5), 457–466.
42. Lucas, S. R., & Berends, M. (2002). Sociodemographic Diversity, Correlated Achievement, and De Facto Tracking. *Sociology of Education, 75*(4), 328-348.
43. Mansfield, R. (2012). Teacher Quality and Student Inequality. *Working Paper Cornell University*.
44. Rivkin, S. G., Hanushek, E. A., & Kain, J. F. (2005). Teachers, Schools, and Academic Achievement. *Econometrica, 73*(2), 417-458.
45. Rothstein, J. (2010). Teacher Quality in Educational Production: Tracking, Decay, and Student Achievement. *Quarterly Journal of Economics*.
46. Sadker, D. M., & Zittleman, K. (2006). *Teachers, Schools and Society: A Brief Introduction to Education.* McGraw-Hill.
47. Waddell, G. (2006). Labor-Market Consequences of Poor Attitude and Low Self-Esteem in Youth. *Economic Inquiry, 44*(1), 69-97.

# Tables

**Table 1:** *Correlations among Short-run Outcomes in the NELS-88 and NCERDC Data*

| | NELS-88 Data | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Raw correlations between outcomes | | | | | | Percentage of Variance Explained by Factors | |
| | Ln of # Days Absent | Suspended | Grade Point Average | In 9th grade on time | Math Score in 8th Grade | Reading Score in 8th Grade | Cognitive Factor | Non-cognitive Factor |
| Ln of # Days Absent | 1 | | | | | | 0.0045 | 0.2751 |
| Suspended | 0.148 | 1 | | | | | 0.0207 | 0.4531 |
| Grade Point Average | -0.171 | -0.243 | 1 | | | | 0.1799 | 0.4678 |
| In 9th grade on time | -0.051 | -0.127 | 0.127 | 1 | | | 0.0003 | 0.2026 |
| Math Score in 8th Grade | -0.075 | -0.128 | 0.383 | 0.023 | 1 | | 0.9087 | 0.0176 |
| Reading Score in 8th Grade | -0.049 | -0.126 | 0.342 | 0.024 | 0.803 | 1 | 0.8897 | 0.0045 |

| | NCERDC Data | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Raw correlations between outcomes | | | | | | Percentage of Variance Explained by Factors | |
| | Ln of # Days Absent | Suspended | Grade Point Average | In 10th grade on time | Algebra Score in 9th Grade | English Score in 9th Grade | Cognitive Factor | Non-cognitive Factor |
| Ln of # Days Absent | 1 | | | | | | 0.0348 | 0.3976 |
| Suspended | 0.252 | 1 | | | | | 0.0370 | 0.1558 |
| Grade Point Average | -0.232 | -0.192 | 1 | | | | 0.5466 | 0.8013 |
| In 10th grade on time | -0.167 | -0.16 | 0.482 | 1 | | | 0.1337 | 0.3146 |
| Algebra Score in 9th Grade | -0.098 | -0.13 | 0.592 | 0.321 | 1 | | 0.8358 | 0.2830 |
| English Score in 9th Grade | -0.082 | -0.13 | 0.539 | 0.323 | 0.616 | 1 | 0.7576 | 0.2311 |

Note: The cognitive and non-cognitive factors were uncovered using factor analysis and are linear combinations of all the short-run outcomes. The results were then standardized. Note that the factors in the NCERDC use the weights derived from the NELS-88 data. However, the factors using weights derived from the NCERDC have correlations greater than 0.95 with those derived using weights from the NELS-88.

**Table 2**: *Relationship between 8th Grade Outcomes and Adult Outcomes in the NELS-88 Survey.*

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  | Arrested, or Close Friend Arrested | | | | Attend College | | | | Log of Income in 1999 | | | |
| Ln of # Days Absent | 0.0145*** | | | | -0.0182*** | | | | -0.00965 | | | |
|  | [0.00528] | | | | [0.00496] | | | | [0.00956] | | | |
| Grade Point Average | -0.0451*** | | | | 0.117*** | | | | 0.0924*** | | | |
|  | [0.00617] | | | | [0.00592] | | | | [0.0111] | | | |
| Repeat 8th Grade | 0.0014 | | | | -0.177*** | | | | 0.042 | | | |
|  | [0.0393] | | | | [0.0407] | | | | [0.0640] | | | |
| Suspended | 0.131*** | | | | -0.162*** | | | | -0.0648*** | | | |
|  | [0.0138] | | | | [0.0135] | | | | [0.0221] | | | |
| Math Score in 8th Grade | 0.0109 | | 0.00194 | | 0.0325*** | | 0.0621** | | 0.0508*** | | 0.0472** | |
|  | [0.00709] | | [0.00424] | | [0.00605] | | [0.00410] | | [0.0128] | | [0.00781] | |
| Reading Score in 8th | 0.00347 | -0.00156 | | | 0.00299 | 0.0582** | | | -0.0315** | 0.0299** | | |
|  | [0.00695] | [0.00420] | | | [0.00602] | [0.00412] | | | [0.0124] | [0.00763] | | |
| Cognitive Factor | | | | -0.000146 | | | | 0.0632** | | | | 0.0417** |
|  | | | | [0.00411] | | | | [0.00408] | | | | [0.00755] |
| Non-cognitive Factor | | -0.0559*** | -0.0561*** | -0.0560*** | | 0.110*** | 0.107*** | 0.106*** | | 0.0643** | 0.0610** | 0.0614** |
|  | | [0.00440] | [0.00443] | [0.00443] | | [0.00425] | [0.00426] | [0.00426] | | [0.00736] | [0.00737] | [0.00737] |
| | | | | | | | | | | | | |
| Demographics | YES | YES | YES | | YES | YES | YES | | YES | YES | YES | YES |
| Observations | 10,792 | 10,792 | 10,792 | | 10,792 | 10,792 | 10,792 | | 9,956 | 9,956 | 9,956 | 9,956 |

Robust standard errors in brackets. + significant at 10%; * significant at 5%; ** significant at 1%.

**Notes:** All models include controls for household income while in 8th grade, whether English is the students primary language at home, whether the student lives with their mother, whether the student lives with their father, the mothers highest level of education, the fathers highest level of education, family size, student race and student gender. Regressions of log income are median regression (unconditional quantiles) rather than OLS to account for any bias due to very low levels of earnings. Also, note that all the independent variables were collected in 1988 when respondent were in 8th grade while the outcomes were collected in 1999 when respondents were approximately 25 years old. Models that include *both* 8th grade math and reading scores suffer from co-linearity.

**Table 3:** *Summary Statistics of Student Data*

| Variable | Mean | SD | SD within school-tracks | SD within schools |
|---|---|---|---|---|
| Math z-score 8th grade | 0.091 | (0.944) | (0.600) | (0.878) |
| Reading z-score 8th grade | 0.073 | (0.941) | (0.678) | (0.891) |
| Male | 0.510 | (0.50) | (0.482) | (0.498) |
| Black | 0.288 | (0.453) | (0.375) | (0.399) |
| Hispanic | 0.075 | (0.263) | (0.245) | (0.256) |
| White | 0.579 | (0.494) | (0.404) | (0.432) |
| Asian | 0.020 | (0.141) | (0.133) | (0.138) |
| Parental education: Some High-school | 0.075 | (0.263) | (0.25) | (0.259) |
| Parental education: High-school Grad | 0.400 | (0.49) | (0.454) | (0.474) |
| Parental education: Trade School Grad | 0.018 | (0.132) | (0.129) | (0.132) |
| Parental education: Community College Grad | 0.133 | (0.339) | (0.327) | (0.335) |
| Parental education: Four-year College Grad | 0.205 | (0.404) | (0.376) | (0.394) |
| Parental education: Graduate School Grad | 0.064 | (0.245) | (0.225) | (0.237) |
| Number of Honors classes | 0.880 | (1.323) | (0.575) | (1.163) |
| Algebra I z-Score (9th grade) | 0.063 | (0.976) | (0.775) | (0.889) |
| English I z-Score (9th grade) | 0.033 | (0.957) | (0.670) | (0.906) |
| Ln Absences | 0.586 | (1.149) | (0.927) | (0.984) |
| Suspended | 0.056 | (0.23) | (0.214) | (0.225) |
| GPA | 2.763 | (0.87) | (0.604) | (0.801) |
| In 10th grade | 0.856 | (0.351) | (0.305) | (0.339) |
| Observations | | | 348547 | |

**Notes:** These summary statistics are based on student who took the English I exam. Incoming math scores and reading scores are standardized to be mean zero unit variance. About 10 percent of students do not have parental education data so that the missing category is "missing parental education".

**Table 4:** *Illustration of the Variation at a Hypothetical School*

| | Year | Track A<br>Alg I (regular)<br>Eng I (regular)<br>Natural Sciences<br>US History | Track B<br>Alg I (regular)<br>Eng I (regular)<br>Biology<br>World History<br>Geometry |
|---|---|---|---|
| Math Teacher 1 | 2000 | X | X |
| Math Teacher 2 | 2000 | | X |
| | | | |
| Math Teacher 1 | 2005 | X | X |
| Math Teacher 2* | 2005 | - | - |
| Math Teacher 3 | 2005 | | X |

**Table 5:** *Estimated Covariance across Classrooms for the Same Teacher under Different Models*

| | | Algebra Teachers | | | | English Teachers | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | SD | Prob Cov≤0 | 95% CI Upper bound | 95% CI Lower bound | SD | Prob Cov≤0 | 95% CI Upper bound | 95% CI Lower bound |
| **School Effects with Peer Characteristics** | Algebra Score 9th Grade | **0.12** | **0.000** | **0.128** | **0.112** | **0.068** | **0.000** | **0.079** | **0.054** |
| | English Score 9th Grade | 0.029 | 0.056 | 0.043 | 0.000 | **0.049** | **0.000** | **0.059** | **0.037** |
| | Suspended | **0.023** | **0.001** | **0.029** | **0.015** | 0.032 | 0.000 | 0.036 | 0.028 |
| | Log of # Absences | **0.138** | **0.000** | **0.154** | **0.119** | **0.132** | **0.000** | **0.143** | **0.12** |
| | GPA | **0.095** | **0.000** | **0.103** | **0.086** | **0.081** | **0.000** | **0.089** | **0.072** |
| | On time enrollment in 10th grade | **0.058** | **0.000** | **0.064** | **0.052** | **0.047** | **0.000** | **0.053** | **0.041** |
| | Cognitive Factor | **0.085** | **0.000** | **0.093** | **0.075** | **0.052** | **0.000** | **0.065** | **0.038** |
| | Non-cognitive factor | **0.146** | **0.000** | **0.157** | **0.134** | **0.106** | **0.000** | **0.116** | **0.094** |
| **Track-by-School and Year Effects** | Algebra Score 9th Grade | **0.1** | **0.000** | **0.107** | **0.093** | **0.039** | **0.005** | **0.051** | **0.019** |
| | English Score 9th Grade | 0.013 | 0.842 | 0.034 | 0.000 | **0.043** | **0.000** | **0.05** | **0.035** |
| | Suspended | 0.009 | 0.645 | 0.022 | 0.000 | **0.021** | **0.000** | **0.025** | **0.017** |
| | Log of # Absences | **0.093** | **0.000** | **0.109** | **0.073** | **0.095** | **0.000** | **0.108** | **0.08** |
| | GPA | **0.065** | **0.000** | **0.075** | **0.053** | **0.049** | **0.000** | **0.059** | **0.037** |
| | On time enrollment in 10th grade | **0.035** | **0.000** | **0.042** | **0.027** | **0.031** | **0.000** | **0.037** | **0.025** |
| | Cognitive Factor | **0.071** | **0.000** | **0.794** | **0.0594** | 0.0327 | 0.024 | 0.0455 | 0.009 |
| | Non-cognitive factor | **0.103** | **0.000** | **0.1161** | **0.088** | **0.085** | **0.000** | **0.0944** | **0.074** |
| **Track-by-School and School-by-Year Effects** | Algebra Score 9th Grade | **0.066** | **0.000** | **0.074** | **0.056** | 0.000 | 0.968 | 0.008 | 0.000 |
| | English Score 9th Grade | 0.000 | 0.917 | 0.009 | 0.000 | **0.03** | **0.003** | **0.04** | **0.014** |
| | Suspended | 0.000 | 0.887 | 0.008 | 0.000 | **0.014** | **0.003** | **0.019** | **0.007** |
| | Log of # Absences | 0.000 | 0.798 | 0.03 | 0.000 | **0.037** | **0.024** | **0.054** | **0.009** |
| | GPA | **0.045** | **0.000** | **0.057** | **0.028** | **0.027** | **0.007** | **0.039** | **0.008** |
| | On time enrollment in 10th grade | **0.025** | **0.002** | **0.033** | **0.012** | **0.024** | **0.000** | **0.031** | **0.015** |
| | Cognitive Factor | **0.046** | **0.000** | **0.055** | **0.0359** | 0.023 | 0.083 | 0.0347 | 0.000 |
| | Non-cognitive factor | **0.069** | **0.000** | **0.056** | **0.0811** | **0.054** | **0.007** | **0.068** | **0.033** |
| | Non-cognitive factor (orthogonal) | **0.064** | **0.0055** | **0.083** | **0.0391** | **0.051** | **0.001** | **0.062** | **0.038** |

Notes: The estimated covariances are computed by taking the classroom level residuals from equation 7 and computing the covariance of mean residuals across classrooms for the same teacher. Specifically, I pair each classroom with a randomly chosen other classroom for the same teacher and estimate the covariance. I replicate this 50 times and report the median estimated covariance as my sample covariance. To construct the standard deviation of this estimated covariance, I pair each classroom with a randomly chosen other classroom for a different teacher and estimate the covariance. The standard deviation of 50 replications of these "placebo" covariances is my bootstrap estimate of the standard deviation of the estimated covariance. These two estimates can then be used to form confidence intervals for the covariance which can be used to compute estimates and confidence intervals for the standard deviation of the teacher effects (by taking the square root of the sample covariance and the estimated upper and lower bounds). When the estimated covariance is negative, I report a value of zero for the standard deviation. Note that none of the negative covariances is statistically significant at the five percent level.

**Table 6:** *Effect of Out of Sample Estimated Teacher Effects and School-Track-Year-Level mean Teacher Effects on Outcomes and Predicted Outcomes*

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| | Algebra Teachers | | | | English Teachers | | | |
| | Algebra | Cognitive | Non-cognitive | Non-cognitive (orth)[a] | English | Cognitive | Non-cognitive | Non-cognitive (orth)[a] |
| Estimated Effect (*variation within cohorts*) | **0.274**\*\* [0.0349] | **0.148**\*\* [0.0257] | **0.0684**\*\* [0.0238] | **0.0693**\* [0.0292] | **0.174**\*\* [0.0240] | **0.0983**\*\* [0.0269] | **0.0569**\*\* [0.0217] | **0.0466**\* [0.0221] |
| Mean Estimated Effect (*variation across cohorts*) | **0.260**\*\* [0.0611] | **0.225**\*\* [0.0656] | **0.140**\* [0.0552] | **0.0975**\* [0.0536] | **0.243**\*\* [0.0440] | **0.0739**\* [0.0373] | **0.173**\*\* [0.0379] | **0.176**\*\* [0.0426] |
| | Predicted Algebra | Predicted Cognitive | Predicted Non-cognitive | Predicted Non-cognitive (orth) | Predicted English | Predicted Cognitive | Predicted Non-cognitive | Predicted Non-cognitive (orth) |
| Estimated Effect (*all exogenous variation*) [b] | 0.00473 [0.00600] | 0.00382 [0.00780] | -0.00376 [0.00355] | -0.00604+ [0.00361] | 0.00956 [0.00685] | 0.00672 [0.00824] | 0.00133 [0.00621] | 0.00369 [0.00299] |
| Mean Estimated Effect (*variation across cohorts*) | 0.0491 [0.0582] | 0.0564 [0.0779] | 0.0145 [0.0267] | -0.00128 [0.0126] | -0.0271 [0.0522] | -0.0284 [0.0651] | -0.03 [0.0202] | -0.00113 [0.00983] |
| Observations | 137,600 | 136,127 | 139,129 | 137,573 | 137,600 | 136,127 | 139,129 | 137,573 |

Standard errors in brackets. \*\* p<0.01, \* p<0.05, + p<0.1

All models include school-year effects and all models include school-track fixed effects. The independent variable in the within cohort models is the estimated effect of a student's teacher (from all other years of data) for that outcome. The independent variable in the across cohort models is the mean estimated effect (from all other years of data) of all students in the same school-track and the same cohort as the students for that outcome.

a. The outcome "Non-cognitive (orth)" is the non-cognitive factor that has been purged of any correlation with test scores.

b. Note that these models are conditional on 8th grade math and reading scores. Accordingly the predicted outcome reflects the effects of gender, ethnicity, 7th grade test scores, and parental education.

**Table 7**:    *Proportion of the Variability in Estimated Effects Explained by Estimated Effects on Test Scores and Effects on the Non-cognitive Factor*
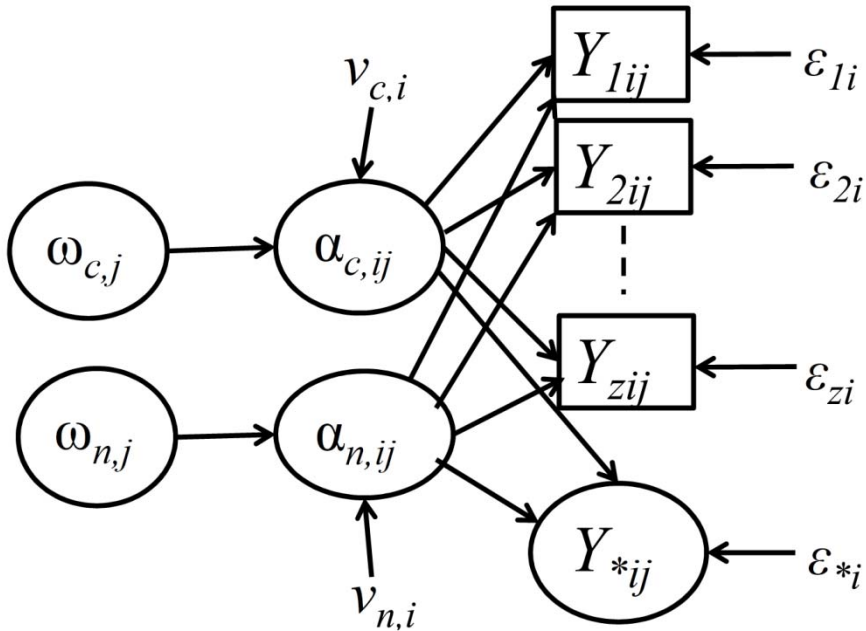
| | Algebra Score | English Score | Suspended | Log of # Absences | GPA | On time enrollment in 10th grade | Cognitive Factor | Non-cognitive factor |
|---|---|---|---|---|---|---|---|---|
| Algebra FX | **1** | - | 0.0115 | 0.0211 | **0.1057** | 0.0503 | **0.75** | **0.1** |
| Cognitive FX | **0.75** | - | 0.0374 | 0.063 | **0.2833** | **0.113** | 1 | **0.269** |
| Non-cog FX | 0.1 | - | **0.1839** | **0.4766** | **0.8007** | **0.407** | **0.269** | 1 |

| | Algebra Score | English Score | Suspended | Log of # Absences | GPA | On time enrollment in 10th grade | Cognitive Factor | Non-cognitive factor |
|---|---|---|---|---|---|---|---|---|
| English FX | - | 1 | 0.0116 | 0.0247 | 0.0516 | 0.0585 | **0.2523** | 0.0683 |
| Cognitive FX | - | | 0.0262 | 0.042 | **0.1392** | 0.0698 | 1 | **0.1458** |
| Non-cog FX | - | 0.0186 | **0.1824** | **0.4658** | **0.7968** | **0.4867** | **0.1458** | 1 |

This table presents the estimated R-squared from separate regressions of a teacher's effect on each outcome on her effect on test scores and her effect on the non-cognitive factor. Estimates greater than 10 percent are in bold.
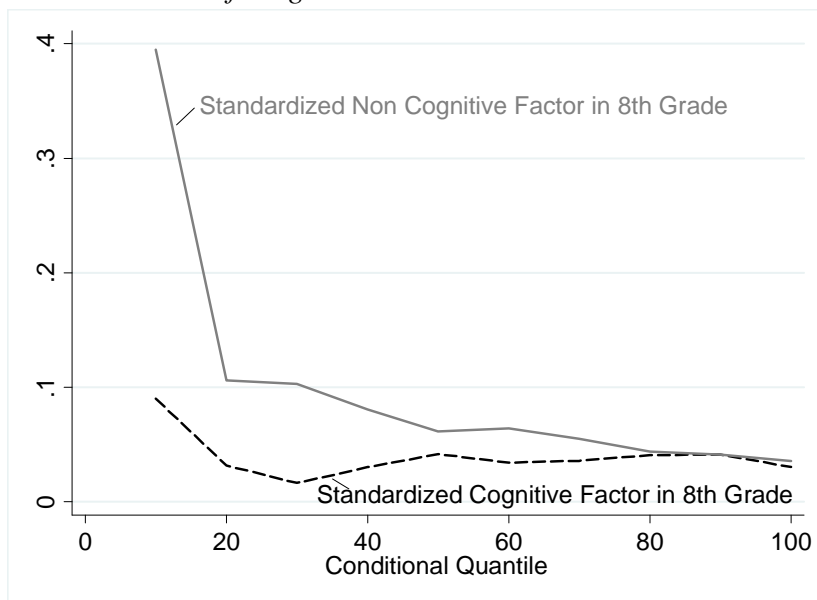
# Figures

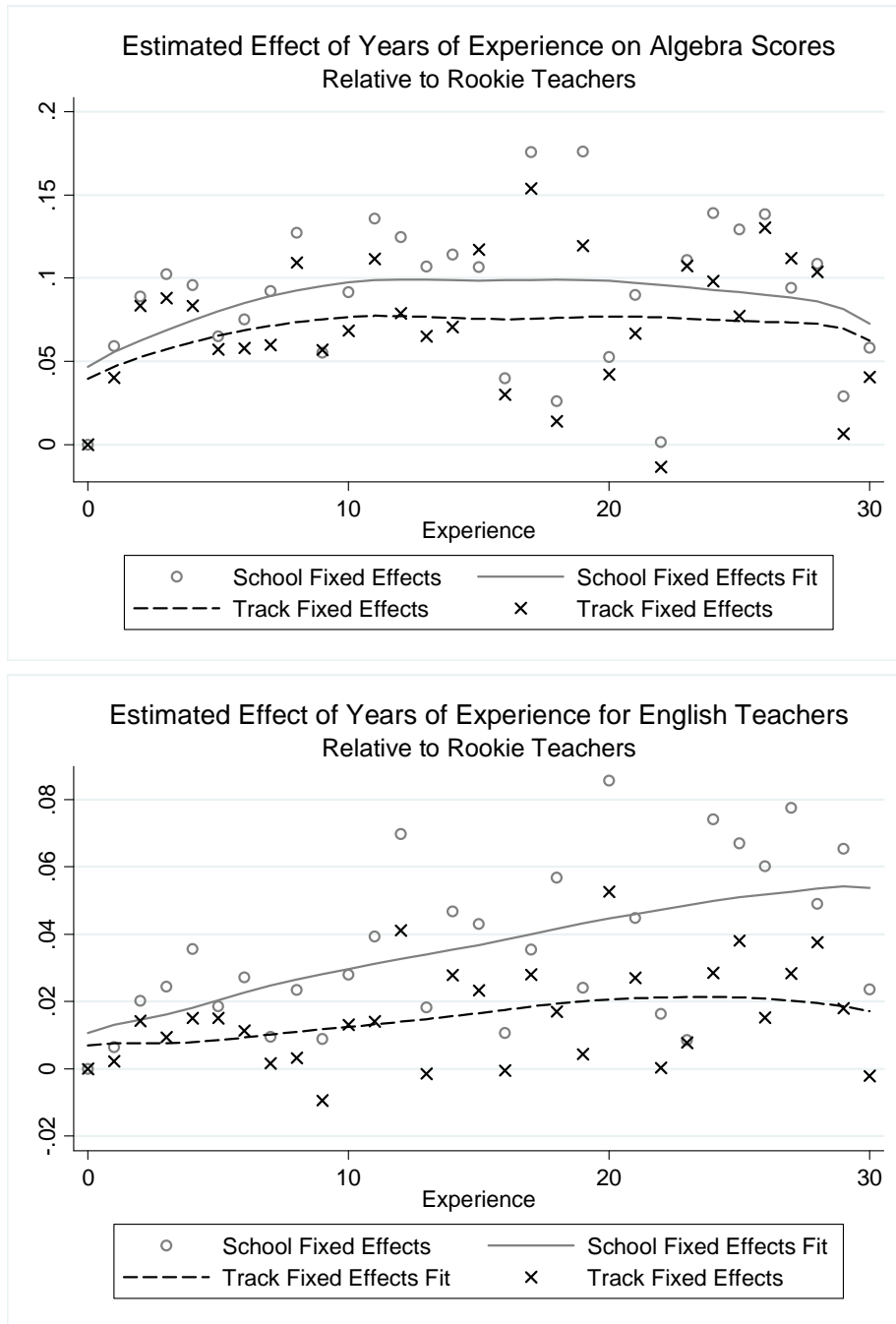**Figure 1:** *Path Diagram of the Two-Factor Model*



**Note:** An arrow from a to b indicates that variable b is a linear function of variable a. Square boxes denote observed variables; while ovals denote unobserved or latent variables.

**Figure 2:** *The Marginal Effect of Cognitive and Non-cognitive Skills at Different Points of the Distribution of Wages*



**Note:** This figure presents estimates from models that control for demographics. The effects are even more pronounced when they are not included.

**Figure 3:** *The Marginal Effect of Teacher Experience under Different Models*



These figures present the estimated coefficients on indicator variables denoting each year of experience on English and Math test scores. The figures show the estimated point estimates for each year of experience and a lowess fit of the point estimates for models with school fixed effects, and track-by-school fixed effects.

**Figure 4:** *Relationship between Teacher Effects on Test Scores and Non-cognitive Factor*
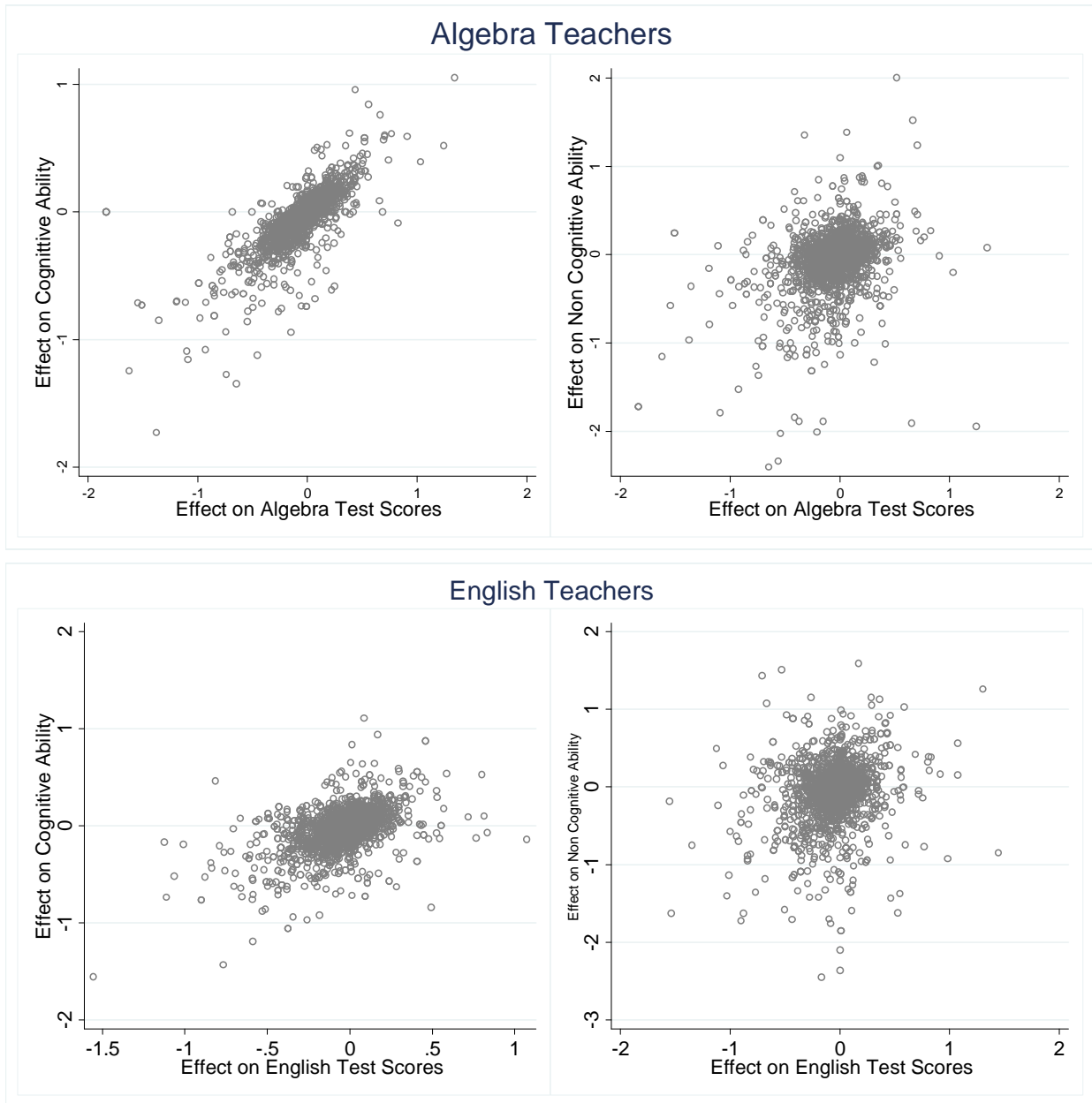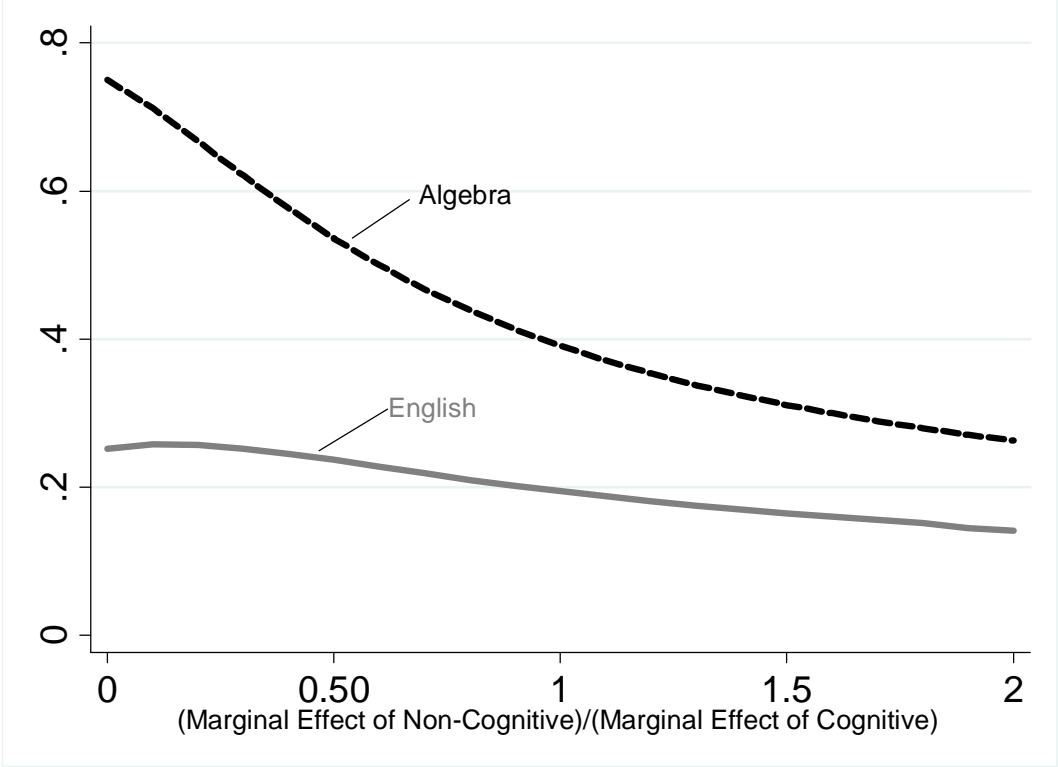
**Figure 5:** *The Proportion of the Teacher Effect on Predicted Adult Outcomes Explained by Test Score Value Added for Different Relative Weights on Non-cognitive Ability.*

# Appendix

**Table A1:**   *Most common academic courses*

| Academic course rank | Course Name | % of 9th graders taking | % of all courses taken |
|:---:|:---:|:---:|:---:|
| 1 | English I* | 90 | 0.11 |
| 2 | World History | 84 | 0.11 |
| 3 | Earth Science | 63 | 0.09 |
| 4 | Algebra I* | 51 | 0.06 |
| 5 | Geometry | 20 | 0.03 |
| 6 | Art I | 16 | 0.03 |
| 7 | Biology I | 15 | 0.02 |
| 8 | Intro to Algebra | 14 | 0.02 |
| 9 | Basic Earth Science | 13 | 0.01 |
| 10 | Spanish I | 13 | 0.02 |

**Table A2:**   *Distribution of Number of Teachers in Each School-Track Year Cell*

| Number of Teachers in Track-Year-School Cell | English | Algebra |
|:---:|:---:|:---:|
| | Percent | |
| 1 | 63.37 | 51.07 |
| 2 | 18.89 | 26.53 |
| 3 | 9.12 | 11 |
| 4 | 5.6 | 6.38 |
| 5 | 3.03 | 3.25 |
| 6 | 0 | 1.77 |

Note: This is after removing singleton tracks.

**Appendix Note 1:**     *Matching Teachers to Students*


   The teacher ID in the testing file corresponds to the teacher who administered the exam, who is not always the teacher that taught the class (although in many cases it will be). To obtain high quality student-teacher links, I link classrooms in the End of Course (EOC) testing data with classrooms in the Student Activity Report (SAR) files (in which teacher links are correct). The NCERDC data contains The End of Course (EOC) files with test score level observations for a certain subject in a certain year. Each observation contains various student characteristics, including, ethnicity, gender, and grade level. It also contains the class period, course type, subject code, test date, school code, and a teacher ID code. Following Mansfield (2012) I group students into classrooms based in the unique combination of class period, course type, subject code, test date, school code, and the teacher ID code. I then compute classroom level totals for the student characteristics (class size, grade level totals, and race-by-gender cell totals). The Student Activity Report (SAR) files contain classroom level observations for each year. Each observation contains a teacher ID code (the actual teacher), school code, subject code, academic level, and section number. It also contains the class size, the number of students in each grade level in the classroom, and the number of students in each race-gender cell.
   To match students to the teacher who taught them, unique classrooms of students in the EOC data are matched to the appropriate classroom in the SAR data. To ensure the highest quality matches, I use the following algorithm:

(1) Students in schools with only one Algebra I or English I teacher are automatically linked to the teacher ID from the SAR files. These are perfectly matched. Matched classes are set aside.
(2) Classes that match exactly on all classroom characteristics and the teacher ID are deemed matches. These are deemed perfectly matched. Matched classes are set aside.
(3) Compute a score for each potential match (the sum of the squared difference between each observed classroom characteristics for classrooms in the same school in the same year in the same subject, and infinity otherwise) in the SAR file and the EOC data. Find the best match in the SAR file for each EOC classroom. If the best match also matches in the teacher ID, a match is made. These are deemed imperfectly matched. Matched classes are set aside.
(4) Find the best match (based on the score) in the SAR file for each EOC classroom. If the SAR classroom is also the best match in the EOC classroom for the SAR class, a match is made. These are deemed imperfectly matched. Matched classes are set aside.
(5) Repeat step 4 until no more high quality matches can be made.


This procedure leads to a matching of approximately 60 percent of classrooms. All results are similar when using cases when the matching is exact so that error due to the fuzzy matching algorithm does not generate any of the empirical findings.

**Appendix Note 2:** *Estimating Efficient Teacher Fixed Effects*

      I follow the procedure outlined in Kane and Staiger (2008) to compute efficient teacher fixed effects. This approach accounts for the fact that (1) teachers with larger classes will tend to have more precise estimates and (2) there are classroom level disturbances so that teachers with multiple classrooms will have more precise estimates. As before, I compute mean residuals from [7] for each classroom $\bar{e}_{cj}^{*} \equiv \theta_{j} + \phi_{c} + \hat{\varepsilon}_{c}$. Since the classroom error is randomly distributed, I use the covariance between the mean residuals of classrooms for the same teacher $\text{cov}(\bar{e}_{cj}^{*}, \bar{e}_{c'j}^{*}) = \hat{\sigma}_{\theta_{j}}^{2}$ as an estimate of the variance of true teacher quality. I use the variance of the classroom demeaned residuals as an estimate of $\hat{\sigma}_{\varepsilon}^{2}$. Because the variance of the residuals is equal to the sum of the variances of the true teacher effects, the classroom effects, and the student errors, I compute the variance of the classroom errors $\sigma_{c}^{2}$ by subtracting $\sigma_{\varepsilon}^{2}$ and $\hat{\sigma}_{\theta_{j}}^{2}$ from the total variance of the residuals. For each teacher I compute [A1], a weighted average of their mean classroom residuals, where classrooms with more students are more heavily weighted in proportion to their reliability.

$$\hat{\theta}_{j} = \sum_{t=1}^{T_{j}} z_{jt} \cdot \frac{(1/(\sigma_{c}^{2} + (\sigma_{\varepsilon}^{2}/N_{jt})))}{\sum_{t=1}^{T_{j}}(1/(\sigma_{c}^{2} + (\sigma_{\varepsilon}^{2}/N_{jt})))} \qquad \text{[A1]}$$

Where $N_{jt}$ is the number of students in classroom $c$, and $T_{j}$ is the total number of classrooms for teacher $j$. This is a more efficient estimate of the teacher fixed effect that the simple teacher average.