

Validation methods for aggregate-level test scale linking: A case study mapping school district test score distributions to a common scale

AUTHORS

Sean F. Reardon

Stanford University

Demetra Kalogrides

Stanford University

Andrew D. Ho

Harvard University

ABSTRACT

Linking score scales across different tests is considered speculative and fraught, even at the aggregate level (Feuer et al., 1999; Thissen, 2007). We introduce and illustrate validation methods for aggregate linkages, using the challenge of linking U.S. school district average test scores across states as a motivating example. We show that aggregate linkages can be validated both directly and indirectly under certain conditions, such as when the scores for at least some target units (districts) are available on a common test (e.g., the National Assessment of Educational Progress). We introduce precision-adjusted random effects models to estimate linking error, for populations and for subpopulations, for averages and for progress over time. In this case, we conclude that the linking method is accurate enough to be used in analyses of national variation in district achievement, but that the small amount of linking error in the methods renders fine-grained distinctions among districts in different states invalid. We discuss how this approach may be applicable when the essential counterfactual question—"what would means/variance/progress for the aggregate units be, had students taken the other test?"—can be answered directly for at least some of the units.

Acknowledgements: The research described here was supported by grants from the Institute of Education Sciences (R305D110018), the Spencer Foundation, and the William T. Grant Foundation. Some of the data used in this paper were provided by the National Center for Education Statistics (NCES). The paper would not have been possible without the assistance of Ross Santy, Michael Hawes, and Marilyn Seastrom, who facilitated access to the EDFacts data. Additionally, we are grateful to Yeow Meng Thum at NWEA, who provided the NWEA data used in some analyses. This paper benefitted substantially from ongoing collaboration with Erin Fahle, Ken Shores, and Ben Shear. The opinions expressed here are our own and do not represent views of NCES, NWEA, the Institute of Education Sciences, the Spencer Foundation, the William T. Grant Foundation, or the U.S. Department of Education. Direct correspondence and comments to Sean F. Reardon, sean.reardon@stanford.edu, 520 CERAS Building #526, Stanford University, Stanford, CA 94305.

VERSION

July 2018

Suggested citation: Reardon, S.F., Kalogrides, D., & Ho, A. (2018). Linking U.S. School District Test Score Distributions to a Common Scale (CEPA Working Paper No.16-09). Retrieved from Stanford Center for Education Policy Analysis: <http://cepa.stanford.edu/wp16-09>

Validation methods for aggregate-level test scale linking:

A case study mapping school district test score distributions to a common scale

Sean F. Reardon
Demetra Kalogrides
Stanford University

Andrew D. Ho
Harvard Graduate School of Education

July, 2018

The research described here was supported by grants from the Institute of Education Sciences (R305D110018), the Spencer Foundation, and the William T. Grant Foundation. Some of the data used in this paper were provided by the National Center for Education Statistics (NCES). The paper would not have been possible without the assistance of Ross Santy, Michael Hawes, and Marilyn Seastrom, who facilitated access to the *EDFacts* data. Additionally, we are grateful to Yeow Meng Thum at NWEA, who provided the NWEA data used in some analyses. This paper benefitted substantially from ongoing collaboration with Erin Fahle, Ken Shores, and Ben Shear. The opinions expressed here are our own and do not represent views of NCES, NWEA, the Institute of Education Sciences, the Spencer Foundation, the William T. Grant Foundation, or the U.S. Department of Education. Direct correspondence and comments to Sean F. Reardon, sean.reardon@stanford.edu, 520 CERAS Building #526, Stanford University, Stanford, CA 94305.

Validation methods for aggregate-level test scale linking:

A case study mapping school district test score distributions to a common scale

Abstract

Linking score scales across different tests is considered speculative and fraught, even at the aggregate level (Feuer et al., 1999; Thissen, 2007). We introduce and illustrate validation methods for aggregate linkages, using the challenge of linking U.S. school district average test scores across states as a motivating example. We show that aggregate linkages can be validated both directly and indirectly under certain conditions, such as when the scores for at least some target units (districts) are available on a common test (e.g., the National Assessment of Educational Progress). We introduce precision-adjusted random effects models to estimate linking error, for populations and for subpopulations, for averages and for progress over time. In this case, we conclude that the linking method is accurate enough to be used in analyses of national variation in district achievement, but that the small amount of linking error in the methods renders fine-grained distinctions among districts in different states invalid. We discuss how this approach may be applicable when the essential counterfactual question—“what would means/variance/progress for the aggregate units be, had students taken the other test?”—can be answered directly for at least some of the units.

Keywords: *linking, scaling, multilevel modeling, achievement testing, NAEP*

Introduction

As educational testing programs proliferate, non-overlapping populations and incomparable scales can limit the scope of research about the correlates and causes of educational achievement. Linking is the psychometric solution to this problem. Common persons, common populations, or common items across tests form the basis for estimated linking functions (Kolen & Brennan, 2014). These functions can enable mappings of scores from various tests to a common scale, enabling large-scale research about educational achievement. However, the bases for these linkages—common persons, populations, or items—are not always available at a large scale. When they are available, methods for evaluating the linkage for the purpose of large-scale research, rather than student-level uses like diagnosis and selection, are still in development (Thissen, 2007).

Dorans and Holland (2000) outline five requirements for equating: (1) equal constructs, (2) equal reliability, (3) examinee indifference between tests, and (4) a symmetrical linking function that is (5) invariant across populations. These requirements are only realistically met within testing programs, not across them. For linkages that do not meet the stringent conditions of equating, the appropriateness of the linkage becomes dependent on the interpretations and uses of the linked scores.

We present a case of aggregate-level linking whose purpose is to support educational research. First, we show how a common assessment at one level of aggregation (the state, in our example) can serve as the basis for a common-population linkage. Second, we demonstrate how the assessment can directly validate the linkage on which it is based, if the assessment also reports scores at a lower level of aggregation (the school district, here). Third, we show how to validate inferences about progress over time in addition to inferences about relative achievement. Fourth, we show how additional assessments that are common across a subset of the lower-level units can provide indirect validation of the linking. Although none of the methods we present is new on its own, the logic and methods in this validation approach is likely to be useful in other aggregate linking scenarios.

A case comparing U.S. school district achievement across states

To understand how a “patchwork” administration of tests can support aggregate linking, we present the case of linking U.S. school district average scores to a common scale. U.S. school districts differ dramatically in their socioeconomic and demographic characteristics (Reardon, Yun, & Eitle, 1999; Stroub & Richards, 2013), and districts have considerable influence over instructional and organizational practices that may affect academic achievement (Whitehurst, Chingos, & Gallaher, 2013). Nonetheless, we have relatively little rigorous large-scale research describing national patterns of variation in achievement across districts, let alone an understanding of the factors that cause this variation. Such analyses generally require district-level test score distributions that are comparable across states. No such nation-wide, district-level achievement dataset currently exists, because school districts do not administer a common set of assessments to all districts across states.

Existing assessments enable some comparisons of academic performance across states or school districts, but none provides comprehensive comparisons across grades, years, and all school districts. At the highest level, the National Assessment of Educational Progress (NAEP) provides comparable state-level scores in odd years, in reading and mathematics, in grades 4 and 8. NAEP also provides district-level scores, but only for a small number of 27 large urban districts under the Trial Urban District Assessment (TUDA) initiative: 18 districts in 2009, and 21 districts in 2011 and 2013. Within individual states, we can compare district achievement using state math and reading/English Language Arts (ELA) tests federally mandated by the No Child Left Behind (NCLB) act, administered annually in grades 3-8. Comparing academic achievement across state lines requires either that districts administer a common test, or that the scores on the state tests can be linked to a common scale. However, state accountability tests generally differ across states. Each state develops and administers its own tests; these tests may assess somewhat different content domains; scores are reported on different, state-determined scales; and proficiency thresholds are set at different levels of achievement. Moreover, the content, scoring, and

definition of proficiency may vary within any given state over time and across grades.

As a result, direct comparisons of average scores or percentages of proficient students across states (or in many cases within states, across grades and years) are unwarranted and misleading. Average scores may differ because scales differ and because performance differs. Proficiency rates may differ because proficiency thresholds differ (Bandeira de Mello, Blankenship, & McLaughlin; Braun & Qian, 2007) and because performance differs. The ongoing rollout of common assessments developed by multistate assessment consortia (such as the Partnership for Assessment of Readiness for College and Careers, PARCC, and the Smarter Balanced Assessment Consortium, SBAC) is certainly increasing comparability across states, but only to the extent that states use these assessments. Customization of content standards by states may also discourage the reporting of results on a common scale across states (Gewertz, 2015; U.S. Department of Education, 2009). Given the incomplete, divided, and declining state participation in these consortia, comprehensive, directly comparable district-level test score data in the U.S. remain unavailable.

In some cases, districts also administer voluntarily-chosen assessments, often for lower-stakes purposes. When two districts adopt the same such assessments, we can compare test scores on these assessments among districts. One of the most widely used assessments, the Measures of Academic Progress (MAP) test from Northwest Evaluation Association (NWEA), is voluntarily administered in several thousand school districts, over 20% of all districts in the country. However, the districts using MAP are neither a representative nor comprehensive sample of districts.

In this paper, we present a validation strategy for comparisons of district-level test scores across states, years, and grades. We rely on a combination of a) population-level state test score data from NAEP and state tests; b) linear transformations that link state test scores to observed and interpolated NAEP scales; and c) a set of validation checks to assess the accuracy of the resulting linked estimates. Together, this represents a suite of approaches for constructing and evaluating linked estimates of test score

distributions.

We use data from the *EDFacts* Initiative (U.S. Department of Education, 2015), NAEP, and NWEA. We obtain population-level state testing data from *EDFacts*; these data include counts of students in ordered proficiency categories for each district-grade-year-subject combination. We fit heteroskedastic ordered probit (HETOP) models to these district proficiency counts, resulting in estimated district means and variances on a state standardized (zero mean and unit variance) scale (Reardon, Shear, Castellano, & Ho, 2016). We then apply linear linking methods that adjust for test reliability (reviewed by Kolen and Brennan, 2014) to place each district's estimated score distribution parameters on a common national scale. Our linking methods are similar to those that Hanushek and Woessman (2012) used to compare country-level performance internationally. At the district level (Greene & McGee, 2011) and school level (Greene & Mills, 2014), the *Global Report Card* (GRC) maps scores onto a national scale using proficiency rates, using a somewhat different approach than ours.¹ What we add to these standard linear linking methods are direct and indirect validation methods that take advantage of patchwork reporting of test scores at the target levels of aggregation.

Although some have argued that using NAEP as a basis for linking state accountability tests is both infeasible and inappropriate for high-stakes student-level reporting (Feuer, Holland, Green, Bertenthal, & Hemphill, 1999), our goal here is different. We do not attempt to estimate student-level scores, and we do not intend the results to be used for high-stakes accountability. Rather, our goal is to estimate transformations that render aggregate test score distributions roughly comparable across districts in different states, so that the resulting district-level distributions can be used in aggregate-level

¹ Our data and methods are more comprehensive than those used in the GRC (Greene & McGee, 2011; Greene & Mills, 2014; <http://globalreportcard.org/>). First, we provide grade-specific estimates (by year), allowing for estimates of measures of progress. Second, instead of the statistical model we describe below (Reardon, Shear, Castellano, & Ho, 2016), which leverages information from three cut scores in each grade, the GRC uses only one cut score and aggregates across grades. This assumes that stringency is the same across grades and that district variances are equal. Third, our methods allow us to provide standard errors for our estimates. Fourth, we provide both direct and indirect validation checks for our linkages.

research. We grant that NAEP and state tests may differ in many respects, including content, testing dates, motivation, accommodations for language, accommodations for disabilities, and test-specific preparation. While accepting these sources of possible linking error, we focus on the counterfactual question that linking asks: How well do our linked district scores from state tests recover the NAEP scores that these districts would have received, had their students taken NAEP? In this way, we treat the issue of feasibility empirically, by using validation checks to assess the extent to which our methods yield unbiased estimates of aggregate parameters of interest.

Data

We use state accountability test score data and state NAEP data to link scores, and we use NAEP TUDA data and NWEA MAP data to evaluate the linkage. Under the *EDFacts* Initiative (U.S. Department of Education, 2015), states report frequencies of students scoring in each of several ordered proficiency categories for each tested school, grade, and subject (mathematics and reading/ELA). The numbers of ordered proficiency categories vary by state, from 2 to 5, most commonly 4. We use *EDFacts* data from 2009 to 2013, in grades 3-8, provided to us by the National Center for Education Statistics under a restricted data use license. These data are not suppressed and have no minimum cell size. We also use reliability estimates collected from state technical manuals and reports for these same years and grades, imputing when they are not reported.²

Average NAEP scores and their standard deviations are reported for states and participating TUDA districts in odd years, in grades 4 and 8, in reading and mathematics. In each state and TUDA district, these scores are based on an administration of the NAEP assessments to representative samples

² From 2009-2012, around 70% of 2,400 state(50)-grade(6)-subject(2)-year(4) reliability coefficients were available. Missing reliabilities were imputed as predicted values from a linear regression of reliability on state, grade, subject, and year. Reliabilities from 2013, which were not yet available when these data were gathered, were assumed to be the same as corresponding reliabilities from 2012. As Reardon and Ho (2015) show, reliabilities are almost always within a few hundredths of 0.90, so imputation errors are not likely to be consequential.

of students in the relevant grades and years. We use years 2009, 2011, and 2013 as a basis for linking; we use additional odd years from 2003 through 2007 as part of some validation analyses. The NAEP state and district means and standard deviations, as well as their standard errors, are available from the NAEP Data Explorer (U.S. Department of Education, n.d.). To account for NAEP initiatives to expand and standardize inclusion of English learners and students with disabilities over this time period, we rely on the Expanded Population Estimates (EPE) of means and standard deviations provided by the National Center of Education Statistics (see Braun, Zhang, & Vezzu, 2008; McLaughlin, 2005; National Institute of Statistical Sciences, 2009).³

Finally, we use data from the NWEA MAP test that overlap with the years, grades, and subjects available in the *EDFacts* data: 2009-2013, grades 3-8, in reading/ELA and mathematics. Student-level MAP test score data (scale scores) were provided to us through a restricted-use data sharing agreement with NWEA. Several thousand school districts chose to administer the MAP assessment in some or all years and grades that overlap with our *EDFacts* data. Participation in the NWEA MAP is generally binary in districts administering the MAP; that is, in participating districts, either very few students or essentially all students are assessed. We exclude cases in any district's grade, subject, and year, where the ratio of assessed students to enrolled students is lower than 0.9 or greater than 1.1. This eliminates districts with scattered classroom-level implementation as well as very small districts with accounting anomalies. Excluded districts comprise roughly 10% of the districts using the NWEA MAP tests. After these exclusions, we estimate district-grade-subject-year means and standard deviations from student-level data reported on the continuous MAP scale.

³ Note that the correlation of EPE and regular NAEP estimates are near unity; as a result, our central substantive conclusions are unchanged if we use the regular NAEP estimates in the linking.

Linking Methods

The first step in linking the state test scores to the NAEP scale is to estimate district-level score means and standard deviations from the coarsened proficiency count data available in *EDFacts*. These distributions are standardized relative to the state-wide student distribution of scores on the state assessment. We do this in each state, separately for each grade, year, and subject, using the methods described in detail by Reardon, Shear, Castellano, and Ho (2016).

In brief, Reardon, Shear, Castellano, and Ho (2016) demonstrate that a heteroskedastic probit (HETOP) model can be used to estimate group (district) test score means and standard deviations from coarsened data. The HETOP model assumes that there is some monotonic transformation of a state's test scale in which each district's test score distribution is normal (with a district-specific mean and standard deviation) and that the observed ordered proficiency counts in each district are the result of coarsening the districts' normal score distributions using a common set of proficiency thresholds. Given these assumptions, the HETOP model provides estimates of each districts' score mean and standard deviation. The resulting estimates are generally unbiased and are only slightly less precise than estimates obtained from (uncoarsened) student-level scale score data in typical state and national educational testing contexts. We refer readers to their paper (Reardon, Shear, Castellano, & Ho, 2016) for technical specifics. Because most states do not report district-level means and standard deviations, the ability to estimate these distributional parameters from coarsened proficiency category data is essential, given that such categorical data are much more readily available (e.g., *EDFacts*). Of course, if individual scale score data or district-level means and standard deviations were readily available, this step would be unnecessary.

Fitting the HETOP model to *EDFacts* data yields estimates of each district's mean test score, where the means are expressed relative to the state's student-level population mean of 0 and standard deviation of 1, within each grade, year, and subject. We denote these estimated district means and standard deviations as $\hat{\mu}_{dygb}^{\text{state}}$ and $\hat{\sigma}_{dygb}^{\text{state}}$, respectively, for district d , year y , grade g , and subject b . The

HETOP model estimation procedure also provides standard errors of these estimates, denoted $se(\hat{\mu}_{dygb}^{\text{state}})$ and $se(\hat{\sigma}_{dygb}^{\text{state}})$, respectively (Reardon, Shear, Castellano, & Ho, 2016).

The second step of the linking process is to estimate a linear transformation linking each state/year/grade/subject scale (standardized to a student-level mean of 0 and standard deviation of 1—the scale of $\hat{\mu}_{dygb}^{\text{state}}$) to its corresponding NAEP distribution. Recall that we have estimates of NAEP means and standard deviations at the state (denoted s) level, denoted $\hat{\mu}_{sygb}^{\text{naep}}$ and $\hat{\sigma}_{sygb}^{\text{naep}}$, respectively, as well as their standard errors. To obtain estimates of these parameters in grades (3, 5, 6, and 7) and years (2010 and 2012) in which NAEP was not administered, we interpolate and extrapolate linearly. First, within each NAEP-tested year, 2009, 2011, and 2013, we interpolate between grades 4 and 8 to grades 5, 6, and 7 and extrapolate to grade 3. Next, for all grades 3-8, we interpolate between the NAEP-tested years to estimate parameters in 2010 and 2012. We illustrate this below for means, and we apply the same approach to standard deviations. Note that this is equivalent to interpolating between years first and then interpolating and extrapolating to grades.

$$\begin{aligned} \hat{\mu}_{sygb}^{\text{naep}} &= \hat{\mu}_{sy4b}^{\text{naep}} + \frac{g-4}{4} (\hat{\mu}_{sy8b}^{\text{naep}} - \hat{\mu}_{sy4b}^{\text{naep}}), \quad \text{for } g \in \{3, 5, 6, 7\}; y \in \{2009, 2011, 2013\}; \text{ and } \forall s, b \\ \hat{\mu}_{sygb}^{\text{naep}} &= \frac{1}{2} (\hat{\mu}_{s[y-1]gb}^{\text{naep}} + \hat{\mu}_{s[y+1]gb}^{\text{naep}}), \quad \text{for } g \in \{3, 4, 5, 6, 7, 8\}; y \in \{2010, 2012\}; \text{ and } \forall s, b \end{aligned} \quad (1)$$

We evaluate the validity of linking to interpolated NAEP grades and years explicitly later in this paper.

Because the estimated district test score moments $\hat{\mu}_{dygb}^{\text{state}}$ and $\hat{\sigma}_{dygb}^{\text{state}}$ are already expressed on a state scale with mean 0 and unit variance, the estimated mapping of the standardized test scale in state s , year y , grade g , and subject b to the NAEP scale is given by Equation (2) below, where $\hat{\rho}_{sygb}^{\text{state}}$ is the estimated reliability of the state test. Given $\hat{\mu}_{dygb}^{\text{state}}$, this mapping yields an estimate of the of the district average performance on the NAEP scale; denoted $\hat{\mu}_{dygb}^{\text{naep}}$. Given this mapping, the estimated standard deviation, on the NAEP scale, of scores in district d , year y , grade g , and subject b is given by Equation (3).

$$\hat{\mu}_{dygb}^{naep} = \hat{\mu}_{sygb}^{naep} + \frac{\hat{\rho}_{dygb}^{state}}{\sqrt{\hat{\rho}_{sygb}^{state}}} * \hat{\sigma}_{sygb}^{naep} \quad (2)$$

$$\hat{\sigma}_{dygb}^{naep} = \left[\frac{(\hat{\sigma}_{dygb}^{state})^2 + \hat{\rho}_{sygb}^{state} - 1}{\hat{\rho}_{sygb}^{state}} \right]^{1/2} \cdot \hat{\sigma}_{sygb}^{naep} \quad (3)$$

The intuition behind Equation (2) is straightforward: districts that belong to states with relatively high NAEP averages, $\hat{\mu}_{sygb}^{naep}$, should be placed higher on the NAEP scale. Within states, districts that are high or low relative to their state (positive and negative on the standardized state scale) should be relatively high or low on the NAEP scale in proportion to that state's NAEP standard deviation, $\hat{\sigma}_{sygb}^{naep}$.

The reliability term, $\hat{\rho}_{sygb}^{state}$, in Equations (2) and (3) is necessary to account for measurement error in state accountability test scores. Recall that district means and standard deviations on the state scale, $\hat{\mu}_{dygb}^{state}$ and $\hat{\sigma}_{dygb}^{state}$, are expressed in terms of standard deviation units of the state student-level observed score distribution. The standardized means are attenuated toward zero due to measurement error. They must be disattenuated before being mapped to the NAEP scale, given that the NAEP scale accounts for measurement error due to item sampling. We disattenuate the means by dividing them by the square root of the state test score reliability estimate, $\hat{\rho}_{sygb}^{state}$. The district standard deviations on the state scale, $\hat{\sigma}_{dygb}^{state}$, are biased toward 1 due to measurement error; we adjust them before linking them to the NAEP scale, as shown in Equation (3).

From Equations (2) and (3), we can derive the (squared) standard errors of the linked means and standard deviations for non-interpolated grades and years. For simplicity in these derivations, we assume $\hat{\mu}_{sygb}^{naep}$ and $\hat{\sigma}_{sygb}^{naep}$ are independent random variables,⁴ which yields:

⁴ This is not strictly true, since $\hat{\mu}_{sygb}^{naep}$ and $\hat{\sigma}_{sygb}^{naep}$ are estimated from the same sample. However, the NAEP samples are large within each state-year-grade-subject, so the covariance of the estimated means and standard deviations will generally be small relative to other sources of sampling variance in Equation (4).

$$\begin{aligned}
\text{var}(\hat{\mu}_{dygb}^{\widehat{\text{naep}}}) &= \text{var}(\hat{\mu}_{sygb}^{\text{naep}}) + \frac{\text{var}(\hat{\sigma}_{sygb}^{\text{naep}})\text{var}(\hat{\mu}_{dygb}^{\text{state}})}{\hat{\rho}_{sygb}^{\text{state}}} \\
&+ \frac{(\hat{\sigma}_{sygb}^{\text{naep}})^2 \text{var}(\hat{\mu}_{dygb}^{\text{state}})}{\hat{\rho}_{sygb}^{\text{state}}} + \frac{(\hat{\mu}_{dygb}^{\text{state}})^2 \text{var}(\hat{\sigma}_{sygb}^{\text{naep}})}{\hat{\rho}_{sygb}^{\text{state}}},
\end{aligned} \tag{4}$$

for $g \in \{4,8\}$ and $y \in \{2009,2011,2013\}$;

$$\begin{aligned}
\text{var}(\hat{\sigma}_{dygb}^{\widehat{\text{naep}}}) &= \frac{(\hat{\sigma}_{dygb}^{\text{state}})^2 \cdot \text{var}(\hat{\sigma}_{dygb}^{\text{state}})}{\hat{\rho}_{sygb}^{\text{state}} [(\hat{\sigma}_{dygb}^{\text{state}})^2 + \hat{\rho}_{sygb}^{\text{state}} - 1]} \left[\text{var}(\hat{\sigma}_{sygb}^{\text{naep}}) + (\hat{\sigma}_{sygb}^{\text{naep}})^2 \right] \\
&+ \text{var}(\hat{\sigma}_{sygb}^{\text{naep}}) \left(\frac{(\hat{\sigma}_{dygb}^{\text{state}})^2 + \hat{\rho}_{sygb}^{\text{state}} - 1}{\hat{\rho}_{sygb}^{\text{state}}} \right),
\end{aligned} \tag{5}$$

for $g \in \{4,8\}$ and $y \in \{2009,2011,2013\}$.

For interpolated grades and years, the sampling variances differ, because interpolated and extrapolated values are essentially weighted averages. For example, it follows from Equation (1), and assuming uncorrelated terms, that for grades 3, 5, 6, and 7 in odd years the sampling variances of interpolated and extrapolated means are:

$$\text{var}(\hat{\mu}_{dygb}^{\widehat{\text{naep}}}) = \left(\frac{8-g}{4}\right)^2 \text{var}(\hat{\mu}_{dy4b}^{\widehat{\text{naep}}}) + \left(\frac{g-4}{4}\right)^2 \text{var}(\hat{\mu}_{dy8b}^{\widehat{\text{naep}}}), \text{ for } g \in \{3, 5, 6, 7\} \tag{6}$$

Validation Checks and Results

The linking method we use here, on its own, is based on the untested assumption that districts' distributions of scores on the state accountability tests have the same relationship to one another (i.e., the same relative means and standard deviations) as they would if the NAEP assessment were administered in lieu of the state test. Implicit in this assumption is that differences in the content, format, and testing conditions of the state and NAEP tests do not differ in ways that substantially affect aggregate relative distributions. This is, on its face, a strong assumption.

Rather than assert that this assumption is valid, we empirically assess it, using the patchwork reporting and administration of district results by NAEP and NWEA. We do this in several ways. First, for the districts participating in the NAEP TUDA assessments over these years, we compare $\hat{\mu}_{dygb}^{naep}$ —the estimated district mean based on our linking method—to $\hat{\mu}_{dygb}^{naep}$ —the mean of NAEP TUDA scores from the district. This provides a direct validation of the linking method, since the TUDA scores are in the metric that the linking method attempts to recover but are not themselves used in any way in the linking process. We repeat this linkage for demographic subgroups to assess the population invariance of the link.

Second, we assess the correlation of our linked district estimates with district mean scores on the NWEA MAP tests. This provides the correlation across a larger sample of districts. However, the NWEA MAP test has a different score scale, so it does not provide direct comparability with the NAEP scale that is the target of our linking.

Third, for the 20 relevant TUDA districts, we assess whether within-district differences in linked scores across grades and cohorts correspond to those differences observed in the NAEP data. That is, we assess whether the linking provides accurate measures of changes in scores across grades and cohorts of students, in addition to providing accurate means in a given year.

Fourth, we conduct a set of validation exercises designed to assess the validity of the interpolation of the NAEP scores in non-NAEP years and grades. For all of these analyses, we present evidence regarding the district means; corresponding results for the standard deviations are in the appendices.

Validation Check 1: Recovery of TUDA means

The NAEP TUDA data provide estimated means and standard deviations on the actual “naep”

scale, $\hat{\mu}_{dygb}^{naep}$ and $\hat{\sigma}_{dygb}^{naep}$ for 17 large urban districts in 2009 and 20 in 2011 and 2013.⁵ For these particular large districts, we can compare the NAEP means and standard deviations to their linked means and standard deviations. For each district, we obtain discrepancies $\widehat{\hat{\mu}}_{dygb}^{naep} - \hat{\mu}_{dygb}^{naep}$ and $\widehat{\hat{\sigma}}_{dygb}^{naep} - \hat{\sigma}_{dygb}^{naep}$. If there were no sampling or measurement error in these estimates, we would report the average of these discrepancies as the bias, and would report the square root of the average squared discrepancies as the Root Mean Squared Error (RMSE). We could also report the observed correlation between the two as a measure of how well the linked estimates align linearly with their reported TUDA values. However, because of imprecision in both the NAEP TUDA and linked estimates, the RMSE will be inflated and the correlation will be attenuated as measures of recovery. Instead, we report measurement-error corrected RMSEs and correlations that account for imprecision in both the linked and TUDA parameter estimates. To estimate the measurement-error corrected bias, RMSE, and correlation in a given year, grade, and subject, we fit the model below, using the sample of districts for which we have both estimates $\widehat{\hat{\mu}}_{dygb}^{naep}$ and $\hat{\mu}_{dygb}^{naep}$ (or $\widehat{\hat{\sigma}}_{dygb}^{naep}$ and $\hat{\sigma}_{dygb}^{naep}$ as the case may be; the model is the same for the means or standard deviations):

$$\begin{aligned}\hat{\mu}_{idygb} &= \alpha_{0dygb}(LINKED_i) + \alpha_{1dygb}(TUDA_i) + e_{idygb} \\ \alpha_{0dygb} &= \beta_{00} + u_{0dygb} \\ \alpha_{1dygb} &= \beta_{10} + u_{1dygb} \\ e_{idygb} &\sim N(0, \omega_{idygb}^2); \mathbf{u}_{dygb} \sim MVN(\mathbf{0}, \boldsymbol{\tau}^2),\end{aligned}\tag{7}$$

where i indexes source (linked or NAEP TUDA test), ω_{idygb}^2 is the estimated sampling variance of $\hat{\mu}_{idygb}$,

⁵ In 2009, the 17 districts are Atlanta, Austin, Baltimore, Boston, Charlotte, Chicago, Cleveland, Detroit, Fresno, Houston, Jefferson County, Los Angeles, Miami, Milwaukee, New York City, Philadelphia, and San Diego. Albuquerque, Dallas, and Hillsborough County joined in 2011 and 2013. Washington, DC is not included for validation, as it has no associated state for linking. California districts (and Texas districts in 2013) did not have a common Grade 8 state mathematics assessment, so the California and Texas districts lack a linked district mean for Grade 8 mathematics.

and $\boldsymbol{\tau}^2 = \begin{bmatrix} \tau_{00}^2 & \tau_{01}^2 \\ \tau_{01}^2 & \tau_{11}^2 \end{bmatrix}$ is the variance-covariance matrix of the linked and TUDA parameter values which must be estimated. Given the model estimates, we estimate the bias, $\hat{B} = \hat{\beta}_{00} - \hat{\beta}_{10}$, and $\widehat{RMSE} = [\hat{B}^2 + \mathbf{b}\hat{\boldsymbol{\tau}}^2\mathbf{b}']^{1/2}$ where $\mathbf{b} = [\mathbf{1} \quad -\mathbf{1}]$ is a 1×2 design matrix. Finally, we estimate the correlation of α_{0dygb} and α_{1dygb} as $\hat{r} = \frac{\hat{\tau}_{01}^2}{\hat{\tau}_{00}\hat{\tau}_{11}}$.

Table 1 reports the results of these analyses in each subject, grade, and year in which we have TUDA estimates (see Table A1 for the corresponding table for standard deviations). Although we do not show the uncorrected estimates here, we note that the measurement error corrections have a negligible impact on bias and reduce the (inflated) RMSE by around 8% on average. On average, the linked estimates overestimate actual NAEP TUDA means by roughly 1.8 points on the NAEP scale, or around 0.05 of a standard deviation unit, assuming a NAEP scale standard deviation of 35 (NAEP standard deviations vary from roughly 30 to 40 across subjects, years, and grades). The bias is slightly greater in earlier years and in mathematics.

Table 1: Recovery of NAEP TUDA means following state-level linkage of state test score distributions to the NAEP scale, measurement error adjusted.

Subject	Grade	Year	n	Recovery		
				RMSE	Bias	Correlation
Reading	4	2009	17	3.95	2.12	0.96
		2011	20	3.69	1.25	0.96
		2013	20	2.62	0.20	0.98
	8	2009	17	2.92	1.12	0.95
		2011	20	2.20	0.63	0.97
		2013	20	3.62	1.67	0.93
Math	4	2009	17	6.09	4.10	0.93
		2011	20	4.97	2.60	0.94
		2013	20	3.60	1.46	0.95
	8	2009	14	5.21	3.40	0.95
		2011	17	3.77	2.09	0.96
		2013	14	4.54	1.47	0.94
Average		2009	14-17	4.70	2.69	0.95
		2011	17-20	3.79	1.64	0.96
		2013	14-20	3.66	1.20	0.95
		Reading	17-20	3.23	1.17	0.96
		Math	14-20	4.77	2.52	0.95
		All	14-20	4.07	1.84	0.95
Subgroup Average		Male	14-20	4.14	1.84	0.97
		Female	14-20	3.95	1.70	0.98
		White	11-19	3.89	0.66	0.98
		Black	13-19	4.11	1.80	0.96
		Hispanic	12-20	4.07	2.08	0.94

Source: Authors calculations from ED*Facts* and NAEP TUDA Expanded Population Estimates data. RMSE and bias are measured in NAEP scale score points. Estimates are based on Equation 7 in text. Subgroup averages are computed from a model that pools across grades and years within subject (like Equation 9 in text); the subject averages are then pooled within subgroup.

This positive bias indicates that the average scores of students in the TUDA districts are systematically higher in the statewide distribution of scores on the state accountability tests than on the NAEP test. This leads to a higher-than-expected NAEP mapping. Table 1 also shows that the average estimated precision-adjusted correlation (disattenuated to account for the imprecision in the observed means) is 0.95 (note that the simple unadjusted correlation is 0.94; measurement error in the means is

relatively minor relative to the true variation in the means of the TUDA districts). Figure 1 shows scatterplots of the estimated linked means versus the observed TUDA means, separately for grades and subjects, with the identity lines displayed as a reference.

Figure 1 here

Note that under a linear linking such as Equation 2, our definition of bias implies that the weighted average bias, among all districts within each state, and across all states, is 0 by design. If we had all districts, the bias in Table 1 would be 0; it is not 0 because Table 1 summarizes the bias for only the subset of NAEP urban districts for which we have scores. The RMSE similarly describes the magnitude of error (the square root of average squared error) for these districts and may be larger or smaller than the RMSE for other districts in the state.

We review here four possible explanations for discrepancies between a district's average scores on the state accountability test and on the NAEP assessments. First, the population of students assessed in the two instances may differ. For example, a positive discrepancy may result if the target district excluded low scoring students from state tests but not from NAEP. If this differential exclusion were greater in the target district, on average, than in other districts in the state, the target district would appear higher in the state test score distribution than it would in the NAEP score distribution, leading to a positive discrepancy between the district's linked mean score and its NAEP mean scores. Likewise, a positive discrepancy would result if the NAEP assessments excluded high scoring students more in the TUDA assessment than in the statewide assessment, or if there were differential exclusion of high-scoring students in other districts on the state test relative to the target district and no differential exclusion on NAEP. In other words, the discrepancies might result from a target district's scores being biased upward on the state test or downward on the NAEP assessment relative to other districts in the state, and/or from other districts' scores being biased downward on the state test or upward on the NAEP assessment relative to the target district.

Second, the discrepancies may result from differential content in NAEP and state tests. If a district's position in the state distribution of skills/knowledge measured by the state test does not match its position in the statewide distribution of skills measured by the NAEP assessment, the linked scores will not match those on NAEP. The systematic positive discrepancies in Table 1 and Figure 1 may indicate that students in the TUDA districts have disproportionately higher true skills in the content areas measured by their state tests than the NAEP assessments relative to other districts in the states. In other words, if large districts are better than other districts in their states at teaching their students the specific content measured by state tests, relative to their effectiveness in teaching the skills measured by NAEP, we would see a pattern of positive discrepancies like that in Table 1 and Figure 1.

Third, relatedly, students in the districts with a positive discrepancy may have relatively high motivation for state tests over NAEP, compared to other districts. Fourth, the bias evident in Table 1 and Figure 1 may indicate differential score inflation or outright cheating. For example, some of the largest positive discrepancies among the 20 TUDA districts illustrated in Figure 1 are in Atlanta in 2009, where there was systematic cheating on the state test in 2009 (Wilson, Bowers, & Hyde, 2011). The discrepancies in the Atlanta estimates are substantially smaller (commensurate with other large districts) in 2011 and 2013, after the cheating had been discovered. In this way, we see that many possible sources of bias in the linking are sources of bias with district scores on the state test itself, rather than problems with the linking per se.

We also address the population invariance of the linking (e.g., Kolen & Brennan, 2014; Dorans & Holland, 2000) by reporting the average direction and magnitude (RMSE) of discrepancies, $\hat{\mu}_{dygb}^{naep} - \hat{\mu}_{dygb}^{naep}$, for selected gender and racial/ethnic subgroups in Table 1. The number of districts is lower in some grade-year cells due insufficient subgroup samples in some districts.⁶ The RMSEs are only slightly

⁶ Our model for subgroups pools across grades (4 and 8) and years (2009, 2011, and 2013) to compensate for smaller numbers of districts in some grade-year cells. We describe this model in Validation Check 3 below. On

larger for subgroups than the RMSE for all students, and bias is around the same magnitude, albeit smaller for White students than for the other subgroups. We conclude from these comparable values that the linking functions recover NAEP district means similarly, on average, across subgroups.

Validation Check 2: Association with NWEA MAP means

The NWEA MAP test is administered in thousands of school districts across the country. Because the MAP tests are scored on the same scale nationwide, district average MAP scores can serve as a second audit test against which we can compare the linked scores. As noted previously, in most tested districts, the number of student test scores is very close to the district's enrollment in the same subject, grade, and year. For these districts, we estimate means and standard deviations on the scale of the MAP test, which we designate "map". The scale differs from that of NAEP, so absolute discrepancies are not interpretable. However, strong correlations between linked district means and standard deviations and those on MAP represent convergent evidence that the linking is appropriate.

We calculate disattenuated correlations between the observed MAP means deviations and both the HETOP estimated means (prior to linking them to the NAEP scale) and the linked means from Equation 2. The improvement from the correlation of MAP means and HETOP estimates to the correlation of MAP means and NAEP-linked estimates is due solely to the move from the "state" to the "naep" scale, shifting all districts within each state according to NAEP performance.

Table 2 shows that correlations between the linked district means and MAP district means are 0.93 on average when adjusting for imprecision (see Table A2 for the corresponding table for standard deviations). This is larger than the average correlation of 0.87 between the MAP means and the

average across grades and years, the results are similar to a model that does not use pooling. We also calculate bias and RMSE for Asian student populations but do not report them due to small numbers of districts: 5-10 per cell. However, bias and RMSE of linked district estimates were higher for Asians, suggesting caution against conducting a separate linkage for Asian students.

(unlinked) HETOP estimates. Figure 2 shows a bubble plot of district MAP scores on linked scores for Grade 4 mathematics in 2009, as an illustration of the data underlying these correlations. Note that the points plotted in Figure 2 are means estimated with imprecision. The observed (attenuated) correlations are generally .03 to .10 points lower than their disattenuated counterparts.

Table 2: Precision-adjusted correlations of linked district means with NWEA MAP district means before and after state-level linkage of state test score distributions to the NAEP scale.

Subject	Grade	Year	N	Precision-Adjusted Correlations	
				...With HETOP Estimates	...With Linked Estimates
Reading	4	2009	1139	0.90	0.95
		2011	1472	0.87	0.93
		2013	1843	0.92	0.95
	8	2009	959	0.84	0.91
		2011	1273	0.87	0.91
		2013	1597	0.88	0.92
Math	4	2009	1128	0.86	0.93
		2011	1467	0.82	0.90
		2013	1841	0.87	0.93
	8	2009	970	0.83	0.93
		2011	1279	0.84	0.92
		2013	1545	0.87	0.95
Average		2009	4196	0.86	0.93
		2011	5491	0.85	0.91
		2013	6826	0.88	0.94
	All Years		16513	0.87	0.93

Source: Authors' calculations from ED*Facts* and NWEA MAP data. NWEA MAP = Northwest Evaluation Association Measures of Academic Progress. Sample includes districts with reported NWEA MAP scores for at least 90% of students.

Figure 2 here

Validation Check 3: Association of between-grade and -cohort trends

An additional assessment of the extent to which the linked state district means match the corresponding NAEP district means compares not just the means in a given grade and year, but the within-district differences in means across grades and years. If the discrepancies evident in Figure 1 are

consistent across years and grades within a district, then the linked state estimates will provide accurate measures of the within-district trends across years and grades, even when there is a small bias in the average means.

To assess the accuracy of the across-grade and -year differences in linked mean scores, we use data from the grades and years in which we have both linked means and corresponding means from NAEP. We do not use the NAEP data from interpolated years and grades in this model. We fit the same model for both means and standard deviations, and separately by subject. For each model, we fit precision-weighted random coefficients models of this form:

$$\begin{aligned}
\hat{\mu}_{idygb} &= \alpha_{0dygb}(LINKED_i) + \alpha_{1dygb}(TUDA_i) + e_{idygb} \\
\alpha_{0dygb} &= \beta_{00d} + \beta_{01d}(year_{dygb} - 2011) + \beta_{02d}(grade_{dygb} - 6) + u_{0dygb} \\
\alpha_{1dygb} &= \beta_{10d} + \beta_{11d}(year_{dygb} - 2011) + \beta_{12d}(grade_{dygb} - 6) + u_{1dygb} \\
\beta_{00d} &= \gamma_{00} + v_{00d} \\
\beta_{01d} &= \gamma_{01} + v_{01d} \\
\beta_{02d} &= \gamma_{02} + v_{02d} \\
\beta_{10d} &= \gamma_{10} + v_{10d} \\
\beta_{11d} &= \gamma_{11} + v_{11d} \\
\beta_{12d} &= \gamma_{12} + v_{12d} \\
e_{idygb} &\sim N(0, \omega_{idygb}^2); \mathbf{u}_{dygb} \sim MVN(\mathbf{0}, \boldsymbol{\sigma}^2); \mathbf{v}_d \sim MVN(\mathbf{0}, \boldsymbol{\tau}^2),
\end{aligned}
\tag{8}$$

where i indexes source (linked or NAEP TUDA test) and ω_{idygb}^2 is the sampling variance of $\hat{\mu}_{idygb}$ (which we set equal to the square of its estimated standard error). The vector $\boldsymbol{\Gamma} = \{\gamma_{00}, \dots, \gamma_{12}\}$ contains the average intercepts, year slopes, and grade slopes (in the second subscript, 0, 1, and 2, respectively) for the linked values and the target values (in the first subscript, 0 and 1, respectively). The differences between the corresponding elements of $\boldsymbol{\Gamma}$ indicate average bias (i.e., the difference between γ_{00} and γ_{10}

indicates the average deviation of the linked means and the NAEP TUDA means, net of district-specific grade and year trends). Unlike Table 1 above, where we estimated bias separately for each year and grade and descriptively averaged them, the bias here is estimated by pooling over all years and grades of TUDA data, with district random effects. If the linking were perfect, we would expect this to be 0.

The matrix of random parameters $\boldsymbol{\tau}^2$ includes, on the diagonal, the between-district variances of the average district means and their grade and year trends; the off-diagonal elements are their covariances. From $\boldsymbol{\tau}^2$ we can compute the correlation between the within-district differences in mean scores between grades and years. The correlation $\text{corr}(v_{01d}, v_{11d})$, for example, describes the correlation between the year-to-year trend in district NAEP scores and the trend in the linked scores. Likewise the correlation $\text{corr}(v_{02d}, v_{12d})$ describes the correlation between the grade 4-8 differences in district NAEP scores and the corresponding difference in the linked scores. Finally, the correlation $\text{corr}(v_{00d}, v_{10d})$ describes the correlation between the NAEP and linked intercepts in the model—that is, the correlation between linked and TUDA mean scores. This correlation differs from that shown in Table 1 above because the former estimates the correlation separately for each grade and year; the model in Equation 8 estimates the correlation from a model in which all years and grades are pooled.

Table 3 shows the results of fitting this model separately by subject to the district means (see Table A3 for the corresponding table for standard deviations). When comparing the linked estimates to the NAEP TUDA estimates, several patterns are evident. First, the estimated correlation of the TUDA and linked intercepts is 0.98 (for both math and reading) and the bias in the means (the difference in the estimated intercepts in Table 3) is small and not statistically significant. The linked reading means are, on average 1.1 points higher (s.e. of the difference is 3.0; n.s.) than the TUDA means; and the linked mathematics means are, on average, 2.4 points higher (s.e. of the difference is 3.3, n.s.) than the TUDA means. These are, not surprisingly, similar to the average bias estimated from each year and grade separately and shown in Table 1.

Table 3. Estimated comparison of linked and TUDA district means, pooled across grades and years, by subject.

	Reading		Math	
Linked ED<i>Facts</i> Parameters				
Intercept (γ_{00})	228.53	***	250.59	***
	(2.00)		(2.10)	
Year (γ_{01})	0.91	***	0.43	*
	(0.17)		(0.19)	
Grade (γ_{02})	10.81	***	9.58	***
	(0.27)		(0.29)	
TUDA Parameters				
Intercept (γ_{10})	227.41	***	248.14	***
	(2.12)		(2.49)	
Year (γ_{11})	1.03	***	0.91	***
	(0.10)		(0.11)	
Grade (γ_{12})	10.84	***	9.68	***
	(0.22)		(0.17)	
L2 Intercept Variance - Linked (σ_0^2)	2.51		2.66	
L2 Intercept Variance - TUDA (σ_1^2)	0.82		1.26	
Correlation - L2 Residuals	1.00		0.36	
L3 Intercept Variance - Linked (τ_1^2)	8.87		9.27	
L3 Intercept Variance - TUDA (τ_4^2)	9.79		11.10	
L3 Year Slope Variance - Linked (τ_2^2)	--		--	
L3 Year Slope Variance - TUDA (τ_5^2)	--		--	
L3 Grade Slope Variance - Linked (τ_3^2)	1.06		1.03	
L3 Grade Slope Variance - TUDA (τ_6^2)	0.93		0.61	
Correlation - L3 Intercepts	0.98		0.98	
Correlation - L3 Year Slopes	--		--	
Correlation - L3 Grade Slopes	0.85		0.98	
Reliability L3 Intercept - Linked	0.98		0.98	
Reliability L3 Year Slope - Linked	--		--	
Reliability L3 Grade Slope - Linked	0.76		0.73	
Reliability L3 Intercept - TUDA	1.00		1.00	
Reliability L3 Year Slope - TUDA	--		--	
Reliability L3 Grade Slope - TUDA	0.87		0.72	
N – Observations	228		204	
N - Districts	20		20	

Source: Authors' calculations from ED*Facts* and NAEP TUDA Expanded Population Estimates data. Estimates are based on Equation 9 in text. *Note:* the level 3 random errors on the year slope were not statistically significant, and so were dropped from the model. L2 = "Level 2"; L3 = "Level 3".

Second, the estimated average linked and TUDA grade slopes ($\hat{\gamma}_{02}$ and $\hat{\gamma}_{12}$, respectively) are

nearly identical to one another in both math and reading. The estimated bias in grade slopes (-0.04 in reading and -0.10 in math) is only 1% as large as the average grade slope. The implied RMSE from the model is 0.56 in reading and 0.46 in math, roughly 5% of the average grade slope.⁷ The estimated correlation of the TUDA and linked grade slopes is 0.85 for reading and 0.98 for math. Finally, the reliability of the grade slopes of the linked estimates is 0.76 in reading and 0.73 in math.⁸ Together these indicate that the linked estimates provide unbiased estimates of the within-district differences across grades, and that these estimates are precise enough to carry meaningful information about between-grade differences.

Third, there is little or no variation in the year trends in the TUDA districts; for both math and reading, the estimated variation of year trends is small and not statistically significant. As a result, neither the TUDA nor the linked estimates provide estimates of trends across years that are sufficiently reliable to be useful (in models not shown, we estimate the reliabilities of the TUDA year trends to be 0.28 and 0.53 and of the linked year trends to be 0.45 and 0.72 in reading and math, respectively). As a result, we dropped the random effects on the year trends and do not report in Table 3 estimates of the variance, correlation, or reliability of the year trends.

Validation Check 4: Recovery of estimates under interpolation within years

Using the interpolated state means and standard deviations in Equation (1) for the linking establishes an assumption that the linkage recovers district scores that would have been reported in years 2010 and 2012 and grades 3, 5, 6, and 7. Although we cannot assess recovery of linkages in interpolated grades with only grades 4 and 8, we can check recovery for an interpolated year, specifically, 2011, between 2009 and 2013. By pretending that we do not have 2011 NAEP state data, we can assess

⁷ We compute the RMSE of the grade slope from the model estimates as follows. Let $\hat{C} = \hat{\gamma}_{02} - \hat{\gamma}_{12}$ be the bias in the grade slopes; then the RMSE of the grade slope will be: $\widehat{RMSE} = [\hat{C}^2 + \mathbf{d}\hat{\tau}^2\mathbf{d}']^{1/2}$, where $\mathbf{d} = [0\ 0\ 1\ 0\ 0\ -1]$.

⁸ The reliability of the level-3 slopes and intercepts is computed as described in Raudenbush and Bryk (2002).

performance of our interpolation approach by comparing linked estimates to actual 2011 TUDA results.

For each of the TUDAs that participated in both 2009 and 2013, we interpolate, for example,

$$\hat{\mu}_{s2011gb}^{naep'} = \frac{1}{2} (\hat{\mu}_{s2009gb}^{naep} + \hat{\mu}_{s2013gb}^{naep})$$

$$\hat{\sigma}_{s2011gb}^{naep'} = \frac{1}{2} (\hat{\sigma}_{s2009gb}^{naep} + \hat{\sigma}_{s2013gb}^{naep})$$
(9)

Applying Equations 2-5, we obtain estimates, for example, $\hat{\mu}_{d2011gb}^{naep'}$, and we compare these to actual TUDA estimates from 2011. We estimate bias and RMSE for discrepancies $\hat{\mu}_{d2011gb}^{naep'} - \hat{\mu}_{d2011gb}^{naep}$ using the model from Validation Check 1. Table 4 shows results in the same format as Table 1 (see Table A4 for the corresponding table for standard deviations). We note that the average RMSE of 3.8 and bias of 1.4 in Table 4 are approximately the same as the average RMSE of 3.8 and bias of 1.6 shown for 2011 in Table 1. Note that the interpolations we use in our proposed linking are those between observed scores that are only two years apart, rather than four years apart as in the validation exercise here. The two-year interpolations should be more accurate than the four-year interpolation, which itself is accurate enough to show no degradation in our recovery of estimated means. We conclude that the between-year interpolation of state NAEP scores adds no appreciable error to the linked estimates for TUDA districts.

Table 4. Recovery of reported 2011 NAEP TUDA means following state-level linkage of state test score distributions to a NAEP scale interpolated between 2009 and 2013, measurement error adjusted.

Subject	Grade	Year	n	Recovery		
				RMSE	Bias	Correlation
Reading	4	2011	20	3.78	0.89	0.95
	8	2011	20	2.14	1.47	0.99
Math	4	2011	20	4.67	2.25	0.94
	8	2011	14	3.81	1.66	0.96
Average	Reading		20	3.07	1.18	0.97
	Math		14-20	4.26	1.95	0.95
	All		14-20	3.72	1.57	0.96

Note: Using NAEP Expanded Population Estimates. Adjusted correlations account for imprecision in linked and target estimates.

We next investigate the viability of interpolation by comparing correlations of linked district estimates with MAP scores at different degrees of interpolation. Some grade-year combinations need no interpolation, others are singly interpolated, and others are doubly interpolated. Table 5 shows that, on average, precision-adjusted correlations between linked NAEP means and MAP means are almost identical across different degrees of interpolation, around 0.93 (see Table A5 for the corresponding table for standard deviations). This lends additional evidence that interpolation adds negligible aggregate error to recovery.

Table 5: Precision-adjusted correlations between NWEA MAP district means and NAEP-linked estimates.

Subject	Grade	2009	2010	2011	2012	2013
Reading	3	0.95	0.94	0.93	0.93	0.94
	4	0.95	0.94	0.93	0.94	0.95
	5	0.94	0.94	0.93	0.93	0.94
	6	0.92	0.94	0.92	0.93	0.93
	7	0.92	0.93	0.92	0.92	0.92
	8	0.91	0.91	0.91	0.91	0.92
Math	3	0.91	0.89	0.91	0.91	0.91
	4	0.93	0.92	0.90	0.92	0.93
	5	0.91	0.90	0.92	0.91	0.93
	6	0.93	0.93	0.94	0.93	0.94
	7	0.94	0.95	0.95	0.95	0.95
	8	0.93	0.93	0.92	0.94	0.95
No interpolation		0.93				
Single interpolation		0.93				
Double interpolation		0.93				

Note. Linked using NAEP Expanded Population Estimates. NWEA MAP = Northwest Evaluation Association Measures of Academic Progress. Sample includes districts with reported NWEA MAP scores for at least 90% of students.

Discussion

We present validation methods for aggregate-level linking, motivated by the goal of constructing a U.S.-wide district-level dataset of test score means and standard deviations. We demonstrate that test score distributions on state standardized tests can be linked to a national NAEP-linked scale in a way that

yields district-level distributions that correspond well—but not perfectly—to the absolute performance of districts on NAEP and the relative performance of districts on MAP. The correlation of district-level mean scores on the NAEP-linked scale with scores on the NAEP TUDA and NWEA MAP assessments is generally high (averaging 0.95 and 0.93 across grades, years, and subjects). Nonetheless, we find some evidence that NAEP-linked estimates include some small, but systematically positive, bias in large urban districts (roughly +0.05 standard deviations, on average). This implies a corresponding small downward bias for other districts in the same states, on average.

Are these discrepancies a threat to the validity of the linked estimates of district means? The answer depends on how the estimates will be used. Given evidence of the imperfect correlation and small bias, the linked estimates should not be used to compare or rank school districts' performance when the estimated means are close and when the districts are in different states. As we noted, there are several possible sources of error in a cross-state comparison, including differences in content, motivation, sampling, and inflation. Our methods cannot identify the presence of any one type of error, but do allow us to quantify the total amount of error in cross-state comparisons. This error is small relative to the distribution of test scores and the variation in average district scores. Of course, relative comparisons within states do not depend on the linking procedure, so these are immune to bias and variance that arises from the linking methods.

On the basis of these results, we believe the linked estimates are accurate enough to be used to investigate broad patterns in the relationships between average test performance and local community or schooling conditions, both within and between states. The validation exercises suggest that the linked estimates can be used to examine variation among districts and across grades within districts. It is unclear whether the estimates provide unbiased estimates of within-grade trends over time, given that there is little or no variation in the NAEP district trends over time against which to benchmark the linked trend estimates. This is true more generally even of within-grade national NAEP trends, which are often

underpowered to detect true progress over shorter time spans of 2- to 4-years.

Validation methods must begin with an intended interpretation or use of scores (Kane, 2013). An operational interpretation of the linked aggregate estimates is the result of monotonic transformations of district score distributions on state tests. They are state score distributions with NAEP-based adjustments, with credit given for being in a state with relatively high NAEP performance and, for districts within the states, greater discrimination among districts when a state's NAEP standard deviation is high. Our contribution is to provide a strategy and methods for answering the essential counterfactual question: What would district results have been, had district scores on NAEP or MAP been available? When patchwork administrations of district tests are available, we can obtain direct and indirect answers to this question.

Because the testing conditions, purpose, motivation, and content of NAEP and state tests differ, we find that district results do differ across tests. But our validation checks suggest that these differences are generally small relative to the variation among districts. This is evident in the high correspondence of the linked and NAEP TUDA estimates and of the linked and NWEA MAP estimates. This suggests that our set of estimated NAEP-linked district test score results may be useful in empirical research describing and analyzing national variation in local academic performance. When data structures with patchwork administrations of tests are available in other U.S. and international testing contexts, our strategy and methods are a roadmap to not only link scores at the aggregate level but to validate interpretations and uses of linked scores.

References

- Bandeira de Mello, V., Blankenship, C., and McLaughlin, D.H. (2009). *Mapping state proficiency standards onto NAEP scales: 2005-2007 (NCES 2010-456)*. National Center for Education Statistics, Institute of Education Sciences, U.S. Department of Education. Washington, DC.
- Braun, H. and Qian, J. (2007). An Enhanced Method for Mapping State Standards onto the NAEP Scale. In: Dorans, N.J., Pommerich M., Holland P.W. (eds) *Linking and Aligning Scores and Scales*. Statistics for Social and Behavioral Sciences. Springs, New York, NY
- Braun, H., Zhang, J., and Vezzu, S. (2010). An investigation of bias in reports of the National Assessment of Educational Progress. *Educational Evaluation and Policy Analysis*, 32, 24-43.
- Dorans, N.J., and Holland, P.W. (2000) Population invariance and the equatability of tests: Basic theory and the linear case. *Journal of Educational Measurement*, 37, 281-306.
- Feuer, M. J., Holland, P. W., Green, B. F., Bertenthal, M.W., & Hemphill, F. C. (1999). *Uncommon measures: Equivalence and linkage among educational tests*. Washington, DC: National Academy Press.
- Gewertz, C. (2015) PARCC Restructures, Allows States to Customize Test. *Education Week*. Retrieved from http://blogs.edweek.org/edweek/curriculum/2015/11/parcc_allows_states_to_customize_test.html
- Greene, J. P., & McGee, J. M. (2011). George W. Bush Institute Global Report Card: Technical Appendix. Report prepared for the George W. Bush Institute. Retrieved from <http://www.globalreportcard.org/docs/AboutTheIndex/Global-Report-Card-Technical-Appendix-9-28-11.pdf>.
- Greene, J. P., & Mills, J. N. (2014). George W. Bush Institute Global Report Card 3.0: Technical Appendix. Report prepared for the George W. Bush Institute. Retrieved from http://www.bushcenter.org/stateofourcities/data/GWBI_GRC_TechnicalAppendix.pdf.

- Hanushek, E. A., & Woessmann, L. (2012). Do better schools lead to more growth? Cognitive skills, economic outcomes, and causation. *Journal of Economic Growth*, 17, 267-321.
- Kane, M.T. (2013) Validating the interpretations and uses of test scores. *Journal of Educational Measurement*, 50(1), 1-73.
- Kolen, M. and Brennan, R. (2014) *Test Equating, Scaling, and Linking*. Springer-Verlag, New York
- McLaughlin D. (2005). Properties of NAEP Full Population Estimates. Unpublished report, American Institutes for Research.
- http://www.schooldata.org/Portals/0/uploads/reports/NSA_T1.5_FPE_Report_090205.pdf
- National Institute of Statistical Sciences (2009). NISS/NESSI Task Force on Full Population Estimates for NAEP. Technical Report #172.
- http://www.niss.org/sites/default/files/technical_reports/tr172.pdf.
- Raudenbush, S. W., & Bryk, A. S. (2002). *Hierarchical Linear Models: Applications and Data Analysis Methods* (2nd ed.). Thousand Oaks, CA: Sage Publications.
- Reardon, S. F., Shear, B. R., Castellano, K. E., & Ho, A. D. (2016). *Using heteroskedastic ordered probit models to recover moments of continuous test score distributions from coarsened data*. Retrieved from <https://cepa.stanford.edu/sites/default/files/wp16-02-v201601.pdf>.
- Reardon, S. F., Yun, J. T., & Eitle, T. M. (1999). *The changing context of school segregation: Measurement and evidence of multi-racial metropolitan area school segregation, 1989-1995*. Paper presented at the annual meeting of the American Educational Research Association. Montreal, Canada.
- Stroub, K. J., & Richards, M. P. (2013). From resegregation to reintegration: Trends in the racial/ethnic segregation of metropolitan public schools, 1993–2009. *American Educational Research Journal*, 50, 497-531.
- Thissen, D. (2007). Linking assessments based on aggregate reporting: Background and issues. In N. J. Dorans, M. Pommerich, & P. W. Holland (Eds.), *Linking and aligning scores and scales* (pp. 287–

312). New York, NY: Springer.

Thissen, D. (2012). *Validity issues involved in cross-grade statements about NAEP results*. NAEP Validity Studies Panel. Washington, DC: National Center for Education Statistics.

U.S. Department of Education. (2015). *EDFacts Submission System User Guide V11.2 (SY 2014-2015)*. Washington, DC: EDFacts. Retrieved from <http://www.ed.gov/EDFacts>.

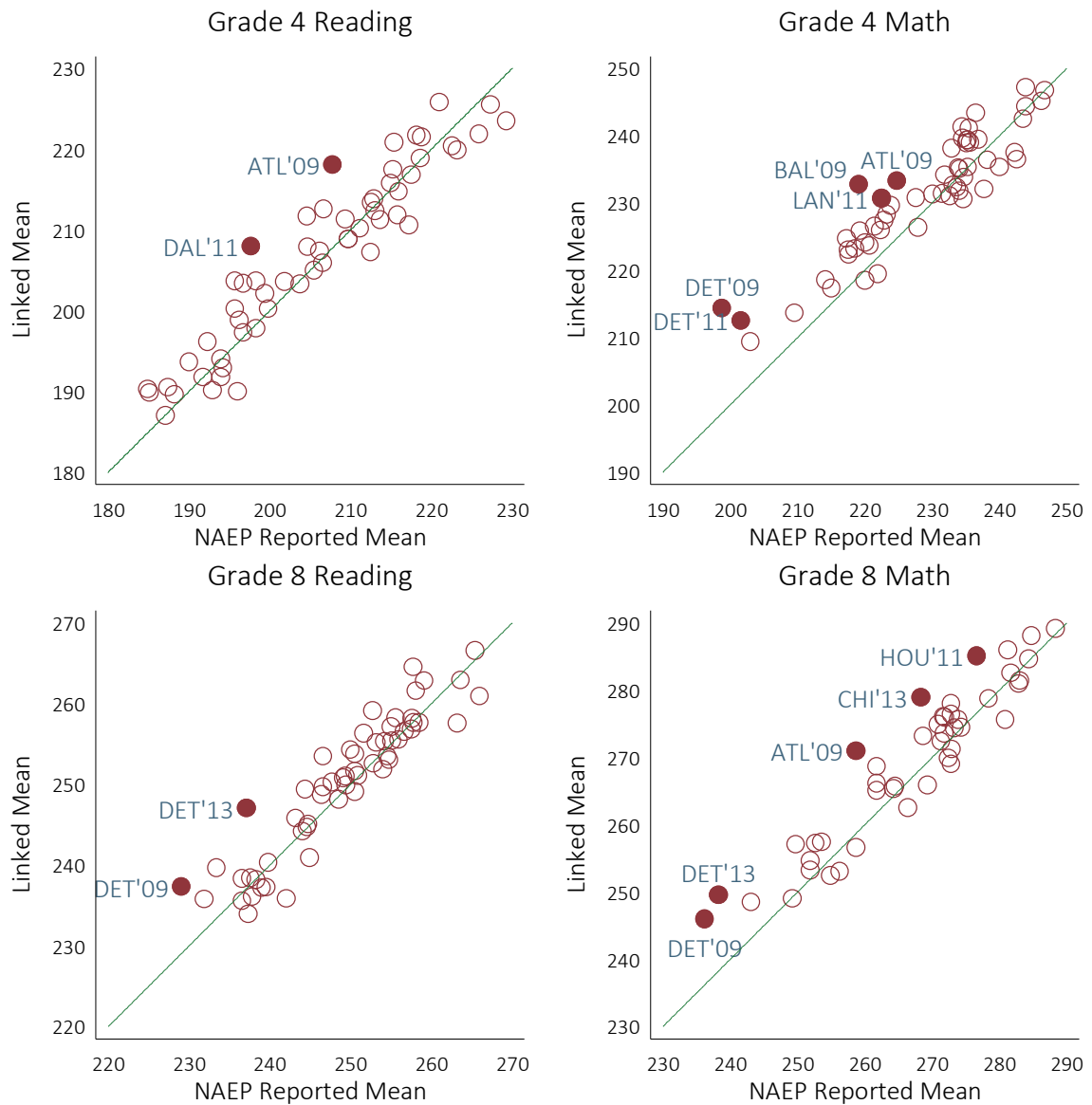
U.S. Department of Education (n.d.), *NAEP Data Explorer*, Washington, D.C.: National Center for Education Statistics, Institute of Education Sciences.

U.S. Department of Education (2009). *Race to the Top Executive Summary*. Washington, DC. Retrieved from <https://www2.ed.gov/programs/racetothetop/executive-summary.pdf>

Whitehurst, G. J., Chingos M. M., & Gallaher, M. R. (2013). *Do School Districts Matter?* Washington, DC: Brookings Institution.

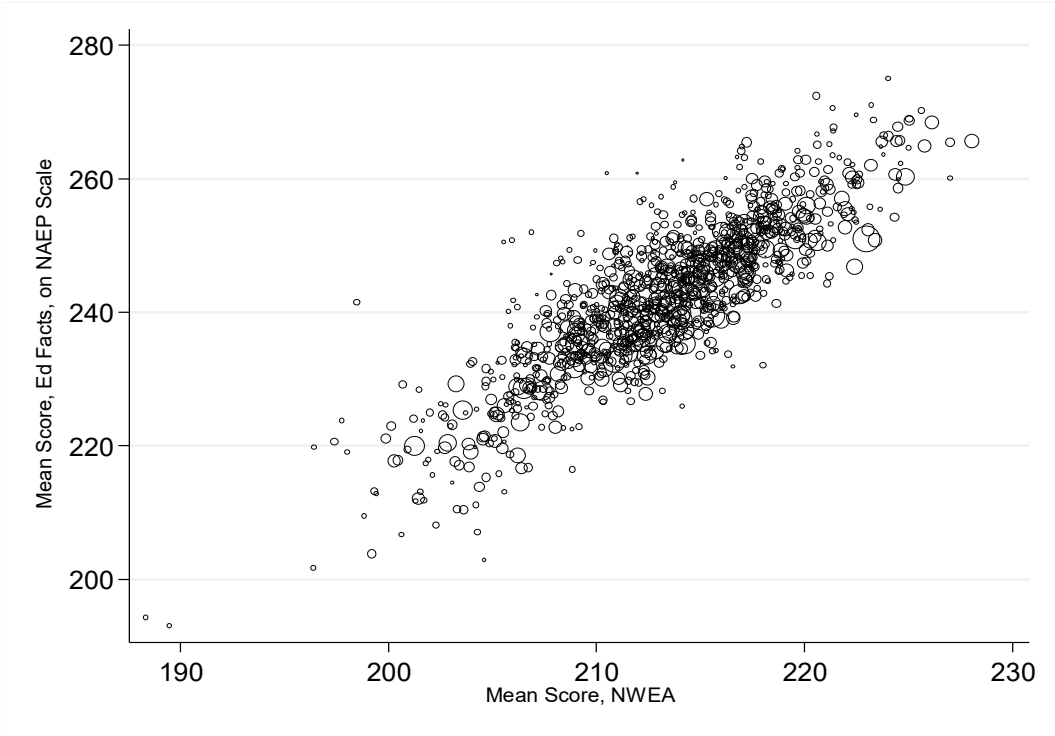
Wilson, R. E., Bowers, M. J., & Hyde, R. L. (2011). *Special Investigation into Test Tampering in Atlanta's School System*. Retrieved from <https://ia801000.us.archive.org/32/items/215252-special-investigation-into-test-tampering-in/215252-special-investigation-into-test-tampering-in.pdf>.

Figure 1. Comparison of reported means from NAEP TUDA and NAEP-linked state test score distributions, grades 4 and 8, Reading and Mathematics, in 2009, 2011, and 2013.



Note: DAL = Dallas; ATL = Atlanta; DET = Detroit; BAL = Baltimore; LAN = Los Angeles; CHI = Chicago; HOU = Houston. District-years with a greater than 8-point discrepancy are labeled.

Figure 2. Example of an association between linked means and NWEA MAP means, grade 4 math, 2009.



Note: Correlation of .87; precision-adjusted correlation of .93. Bubble size corresponds to district enrollment.

Appendix Tables

Table A1: Recovery of NAEP TUDA standard deviations following state-level linkage of state test score distributions to the NAEP scale, measurement error adjusted.

Subject	Grade	Year	n	Recovery		
				RMSE	Bias	Correlation
Reading	4	2009	17	0.89	-1.31	0.76
		2011	20	2.28	-0.14	0.46
		2013	20	1.07	0.34	0.96
	8	2009	17	1.84	-1.33	0.88
		2011	20	1.80	-1.00	0.41
		2013	20	1.01	-1.27	0.62
Math	4	2009	17	1.40	-0.09	0.72
		2011	20	1.71	0.83	0.68
		2013	20	1.64	0.99	0.71
	8	2009	14	2.15	-0.62	0.77
		2011	17	1.86	0.24	0.79
		2013	14	2.49	0.15	0.56
Average		2009	14-17	1.64	-0.84	0.78
		2011	17-20	1.93	-0.02	0.59
		2013	14-20	1.66	0.05	0.71
		Reading	17-20	1.57	-0.79	0.68
		Math	14-20	1.91	0.25	0.71
		All	14-20	1.75	-0.27	0.69
Subgroup Average	Male	14-20	1.60	-0.13	0.72	
	Female	14-20	1.76	0.05	0.65	
	White	11-19	4.43	2.99	0.42	
	Black	13-19	1.90	-0.24	0.40	
	Hispanic	12-20	2.51	-1.24	0.34	

Source: Authors calculations from ED*Facts* and NAEP TUDA Expanded Population Estimates data. Estimates are based on Equation 7 in text. Subgroup averages are computed from a model that pools across grades and years within subject (like Equation 9 in text); the subject averages are then pooled within subgroup.

Table A2: Precision-adjusted correlations of linked district standard deviations with NWEA MAP district standard deviations before and after state-level linkage of state test score distributions to the NAEP scale.

Subject	Grade	Year	n	Precision-Adjusted Correlations	
				...With HETOP Estimates	...With Linked Estimates
Reading	4	2009	1139	0.52	0.58
		2011	1472	0.55	0.61
		2013	1843	0.59	0.64
	8	2009	959	0.54	0.57
		2011	1273	0.51	0.58
		2013	1597	0.52	0.51
Math	4	2009	1128	0.65	0.75
		2011	1467	0.58	0.66
		2013	1841	0.63	0.69
	8	2009	970	0.65	0.70
		2011	1279	0.55	0.62
		2013	1545	0.65	0.70
Average		2009	4196	0.59	0.65
		2011	5491	0.55	0.62
		2013	6826	0.60	0.64
		All Years	16513	0.58	0.63

Source: Authors' calculations from ED Facts and NWEA MAP data. NWEA MAP = Northwest Evaluation Association Measures of Academic Progress. Sample includes districts with reported NWEA MAP scores for at least 90% of students.

Table A3. Estimated comparison of linked and TUDA district standard deviations, pooled across grades and years, by subject.

	Reading		Math	
Linked EDFacts Parameters				
Intercept (γ_{00})	36.21	***	33.17	***
	(0.42)		(0.35)	
Year (γ_{01})	0.18	+	0.39	***
	(0.10)		(0.11)	
Grade (γ_{02})	-0.68	***	1.66	***
	(0.13)		(0.18)	
TUDA Parameters				
Intercept (γ_{10})	36.85	***	32.75	***
	(0.39)		(0.36)	
Year (γ_{11})	0.07		0.17	*
	(0.10)		(0.08)	
Grade (γ_{12})	-0.38	**	1.81	***
	(0.13)		(0.14)	
L2 Intercept Variance - Linked (σ_0^2)	0.96		1.21	
L2 Intercept Variance - TUDA (σ_1^2)	0.74		0.43	
Correlation - L2 Residuals	1.00		1.00	
L3 Intercept Variance - Linked (τ_1^2)	1.73		1.30	
L3 Intercept Variance - TUDA (τ_4^2)	1.60		1.43	
L3 Year Slope Variance - Linked (τ_2^2)	--		--	
L3 Year Slope Variance - TUDA (τ_5^2)	--		--	
L3 Grade Slope Variance - Linked (τ_3^2)	0.44		0.67	
L3 Grade Slope Variance - TUDA (τ_6^2)	0.46		0.54	
Correlation - L3 Intercepts	0.60		0.70	
Correlation - L3 Year Slopes	--		--	
Correlation - L3 Grade Slopes	0.77		0.89	
Reliability L3 Intercept - Linked	0.86		0.76	
Reliability L3 Year Slope - Linked	--		--	
Reliability L3 Grade Slope - Linked	0.63		0.77	
Reliability L3 Intercept - TUDA	0.84		0.87	
Reliability L3 Year Slope - TUDA	--		--	
Reliability L3 Grade Slope - TUDA	0.63		0.79	
N - Observations	228		204	
N - Districts	20		20	

Source: Authors' calculations from EDFacts and NAEP TUDA Expanded Population Estimates data. Estimates are based on Equation 9 in text. Note: the level 3 random errors on the year slope were not statistically significant, and so were dropped from the model. L2 = "Level 2"; L3 = "Level 3".

Table A4. Recovery of reported 2011 NAEP TUDA standard deviations following state-level linkage of state test score distributions to a NAEP scale interpolated between 2009 and 2013, measurement error adjusted.

Subject	Grade	Year	n	Recovery		
				RMSE	Bias	Correlation
Reading	4	2011	20	2.12	0.39	0.53
	8	2011	20	2.53	-0.70	0.00
Math	4	2011	20	2.15	1.43	0.53
	8	2011	14	1.52	0.65	0.93
Average		Reading	20	2.34	-0.15	0.26
		Math	14-20	1.86	1.04	0.73
		All	14-20	2.11	0.44	0.50

Note: Using NAEP Expanded Population Estimates. Adjusted correlations account for imprecision in linked and target estimates.

Table A5: Precision-adjusted correlations between NWEA MAP district standard deviations and NAEP-linked estimates.

Subject	Grade	2009	2010	2011	2012	2013
Reading	3	0.53	0.57	0.63	0.62	0.61
	4	0.58	0.57	0.61	0.62	0.64
	5	0.61	0.54	0.58	0.64	0.64
	6	0.60	0.58	0.61	0.64	0.57
	7	0.60	0.57	0.56	0.55	0.53
	8	0.57	0.56	0.58	0.52	0.51
Math	3	0.71	0.70	0.69	0.67	0.65
	4	0.75	0.77	0.66	0.71	0.69
	5	0.73	0.69	0.73	0.72	0.74
	6	0.73	0.75	0.69	0.71	0.73
	7	0.75	0.65	0.70	0.71	0.70
	8	0.70	0.64	0.62	0.71	0.70
No interpolation		0.63				
Single interpolation		0.58				
Double interpolation		0.64				

Note. Linked using NAEP Expanded Population Estimates. NWEA MAP = Northwest Evaluation Association Measures of Academic Progress. Sample includes districts with reported NWEA MAP scores for at least 90% of students.