

School Districts and Student Achievement

Matthew M. Chingos
Fellow, Brown Center on Education Policy
Brookings Institution

Grover J. Whitehurst
Senior Fellow and Director, Brown Center on Education Policy
Brookings Institution

Michael R. Gallaher
Former Research Analyst, Brown Center on Education Policy
Brookings Institution

Education Finance and Policy, Forthcoming

Final pre-publication draft: January 2014

Abstract

School districts are a focus of education reform efforts in the U.S., but there is very little existing research about how important they are to student achievement. We fill this gap in the literature using 10 years of student-level, statewide data on fourth- and fifth-grade students in Florida and North Carolina. A variance decomposition analysis based on hierarchical linear models indicates that districts account for only a small share (1 to 2 percent) of the total variation in student achievement. But the differences between lower and higher performing districts are large enough to be of practical and policy significance, with a one standard deviation difference in district effectiveness corresponding to about 0.11 standard deviations in student achievement (about 9 weeks of schooling). District performance is generally stable over time, but there are examples of districts that have shown significant increases or decreases in performance.

Introduction

School districts are at the center of public attention and public policy on education reform. Many of the most popular and aggressively promoted school reform efforts are focused at the district level. Performance-based teacher evaluation is a notable case in point. As a condition of competing for funding under the Obama administration's \$4.3 billion Race to the Top program, states promised to establish policies requiring school districts to put in place teacher evaluation systems that would heavily weight student achievement gains on state tests. Districts were expected to tie decisions on tenure, promotion, and salary for individual teachers to the resulting evaluation scores. States around the country are now in the process of requiring districts to implement such teacher evaluation systems, often with short time frames and much of the decisions on design and implementation left to each school district (Klein 2013). Presumably individual differences among school districts, certainly including the quality and skills of their management teams, will influence the results.

Many other reform initiatives are focused at the district level in the sense that they are intended to disrupt the school district's monopoly in delivering publicly funded K-12 education services. These include charter schools, vouchers, on-line education, and school portfolio management models. In some sense, these disruptive reforms proceed on the premise that school districts are the irremediable problem rather than the lever for reform, and that the focus on strong leadership from the top misconstrues where the action is on student achievement, which is individual schools, teachers, curriculum, and parental choice of where to educate their children.

The presumed importance of districts is also implied by the media attention given to prominent school superintendents such as Michelle Rhee in Washington, D.C. and Joel Klein in New York City. Rhee, for example, was on the cover of Time Magazine in 2008 with the lead,

“Michelle Rhee ... head of the D.C. public schools ... could transform public education.” The pay scale for superintendents also indicates their perceived importance. In New York State, for example, 63 district leaders each received over \$300,000 in salary and benefits for the 2011-12 school year, with the superintendent at the top of the list receiving a salary and benefits package of \$541,000 (New York State Education Department 2011).

Private philanthropy has invested heavily in district-level reforms on the premise that districts are a powerful fulcrum for change. One of those philanthropies, the Eli and Edythe Broad Foundation, has led the way in other initiatives that are predicated on the importance of school districts and their leadership. Their annual Broad Prize for Education bestows \$1 million on the school district that has shown the best performance and improvement. The Broad Superintendents Academy was founded in 2002 with the goal of finding leaders from both inside and outside education, training them, and having them fill superintendent positions in a third of the 75 largest school districts in the nation. The foundation has not reached that goal, but it has been remarkably successful in placing its graduates in high-level positions: as of 2011, 21 of the nation’s 75 largest districts had superintendents or other highly placed central-office executives who have undergone Broad training (Samuels 2011).

When we turn from perceptions and intuitions about the importance of school districts to empirical evidence on their impact, we move from a rich to a sparse landscape. Little is known about the impact of school districts on student achievement. And almost all of the existing research on district effectiveness is riddled with methodological shortcomings. The most important problem, which plagues most of this literature, is selection on the dependent variable, in which atypically high-performing districts are examined in order to identify the factors that enabled those districts to be successful. But without examining a broad group of districts, it is

impossible to determine whether those factors are correlated with student outcomes, much less whether any of the observed relationships are causal.

Two recent reviews of this research ignore this fundamental problem, which appears to plague most of the studies that are reviewed (Rorrer, Skrla, and Scheurich 2008; Leithwood 2010). However, Trujillo's (2013) review of 50 studies quantifies the scale of the problem. Of these studies, 40 percent selected districts for study based on anecdotal evidence about their effectiveness, the reputation of the districts, or convenience, and a majority (56 percent) chose districts based on unusually high effectiveness. Fully 63 percent of the studies only examined districts judged to be high performing, and less than 15 percent included districts that spanned a cross-section of performance. Roughly 90 percent of the studies examined used a case study approach, and quantitative techniques were rarely used.

The small body of existing quantitative research rests on a highly questionable set of methods and assumptions. Among the handful of studies that have addressed the importance of school districts, most have focused on district leadership. The most recent review of the quantitative research literature on district leadership, produced by the Mid-Continent Regional Educational Lab (McRel), comes to the conclusion that "district-level leadership matters." This conclusion is based on the authors' finding of "a statistically significant relationship (a positive correlation of 0.24) between district leadership and student achievement" in a meta-analysis of "studies conducted since 1970 that used rigorous, quantitative methods to study the influence of school district leaders on student achievement" (Waters and Marzano 2006).

This review is a notable example of the misuse of meta-analysis to draw causal conclusions. In particular, nearly all of the 14 studies the authors use for the meta-analysis of the impact of district leaders employed survey methodologies in which samples of superintendents

answered questions about their management practices and philosophies. A typical study in the meta-analysis is an unpublished doctoral dissertation in which the means for districts on a state assessment of student academic achievement were regressed on a linear combination of answers provided by superintendents to survey questions. The finding, for example, that a combination of answers to questions about leadership style by district superintendents is associated with differences in student achievement scores is taken as evidence that the leadership style of district leaders is causally related to student outcomes.

By and large the primary studies and the meta-analysis conflate correlation and causation, in particular in that they fail to consider that much of the variation in district performance that is attributed to variables at the district level such as leadership style could be due to differences among districts in the characteristics of the students and families that are served or in the characteristics of the teachers the district employs. This almost surely leads to false conclusions, including overstating the extent to which student achievement varies across districts by combining variation in actual district effectiveness with variation in the characteristics of the enrolled students. In other words, even if one ignores the problems in causal interpretation that arise from models of the effects of district organization that do not include appropriate controls for student ability, there is still the problem that what the models are accounting for may be such a small portion of the overall variance in student achievement as to be educationally unimportant.

This paper begins to fill this significant gap in the literature by exploring how student outcomes vary across districts using two statewide, student-level longitudinal databases. We do not aim to obtain rigorous causal estimates of the effect of individual school districts on student achievement. Rather, we explore the associations between school districts and student

achievement using the kinds of databases that are now readily available but did not exist when many of the existing studies were conducted. We believe that policy decisions made at the state or federal levels having to do with the resources and effort to be placed on district-level reform strategies can be informed by carefully examining the extent to which current variation in student achievement is associated with school districts as organizational units versus schools and teachers.

We apply variance decomposition techniques based on hierarchical linear modeling to administrative data from the states of Florida and North Carolina in order to measure how much student achievement varies across observationally similar districts, and put this in the context of variation at the school, classroom, and student levels. This study differs from other recent research on district performance (see, e.g., Bowers 2010) in that it uses statewide student-level data—instead of data aggregated at the school or district level—thereby enabling us to more accurately estimate the share of variance in student achievement associated with district-level factors.

We find that district-level variation in Florida and North Carolina accounts for a relatively small fraction of the variation in student achievement, on an order of magnitude of less than 2 percent of the total variation. But even though district effects are only a small piece of the total variation in student achievement, there are still differences among the academic achievement of demographically similar students in higher and lower performing districts in North Carolina and Florida that are large enough to be of practical and policy significance. A one standard deviation increase in the estimated district effect is associated with an increase in student achievement of 0.10-0.14 standard deviations in math and 0.07-0.11 standard deviations in reading. There are also districts that have displayed noteworthy patterns of performance in

terms of student achievement over the last decade, including districts with consistently high or low performance and districts that saw significant growth or declines.

Data

For our analysis we constructed two student-level datasets for Florida and North Carolina and matched students with teachers, schools, and districts for each dataset. Our extract from the Florida Department of Education's K-20 Education Data Warehouse (EDW) contains observations for every student who took state assessments in math and reading from 1998-99 to 2009-10. In addition to student test scores from the Florida Comprehensive Assessment Test (FCAT), the EDW contains information on student demographics, attendance, and program participation—such as the gifted and talented, free and reduced lunch, and English language learner programs. As of 2000-01, students in grade 3-10 took the FCAT in both reading and writing. Therefore, we limit our analysis to the ten years of data between 2000-01 and 2009-10 for both Florida and North Carolina. However, most of our analyses are based only on the most recent year of data (2009-10).

In Florida, we used the EDW's course records and matched students with one teacher for math and one teacher for reading, with most students matching to the same teacher for both subjects. Students were only matched to teachers if the student spent 40 percent or more of their total academic time with that teacher. Additionally, because the focus of our analysis was on student test-score performance, only students that took the FCAT in at least one subject were included in the analysis. Students with duplicate test scores for the same subject in the same semester were excluded from the analysis, though these students made up a small portion of the initial dataset (less than one half of one percent). Of the students who took the FCAT during our

period of study, 90 percent were matched to a math teacher and 92 percent were matched to a reading teacher (in both cases, with whom they spent at least 40 percent of their time).

Similarly, our extract from the North Carolina Education Research Data Center (NCERDC) contains observations for every student who took End of Grade (EOG) assessments in math and reading through 2009-10. Like Florida, student data from North Carolina include test scores on math and reading assessments and student demographics. Unlike Florida, however, we matched students to teachers based on the identity of the EOG assessment proctor for that year, which according to NCERDC accurately identify 95 percent of all classroom teacher assignments (Hoxby and Weingarth 2005).

For both Florida and North Carolina, we standardize test scores, separately by state, subject, grade, and year, to have a mean of zero and standard deviation of one. We pool results for the two grades we examine (fourth and fifth) in order to increase the precision of the results (as compared to estimating results separately for each grade), especially given the fact that some rural districts enroll a relatively small number of students.

Even though districts are the focus of our analysis, the decision of which teachers to associate with students is important because the portion of the variance in outcomes attributable to teachers affects the portion that could be attributable to districts. Consequently, we limit our analyses to students in grades 4-5, primarily because students in these grades usually have a single classroom teacher whereas older students have multiple subject-specific teachers (and we exclude third-grade students in order to be able to control for prior-year test scores in some models).¹

¹ For the later grades, we could use teachers of a given subject in the analysis of test scores in that subject, but that would ignore any effects of teachers in other subjects (e.g., the effect of the English teacher on math scores, which might partly reflect students' ability to read word problems).

Methods

We first measure the variation in student achievement associated with the district, school, and teacher using variance decomposition techniques based on hierarchical linear models (Raudenbush and Bryk 2002). These models allow us to measure how much student achievement varies at different levels, namely students within classrooms, classrooms within schools, schools within districts, and districts within Florida and North Carolina.

Specifically, we use Stata's `xtmixed` command to estimate a four-level hierarchical linear model of test scores, with students (level 1) nested within classrooms (level 2) nested within schools (level 3) nested within districts (level 4). Following the notation of Raudenbush et al. (2011), the models are:

$$\text{Level-1 (students): } Y_{ijkl} = \pi_{0jkl} + \sum_{p=1}^P \pi_{pjkl} \alpha_{pijkl} + e_{ijkl},$$

where Y_{ijkl} is the test score of student i in classroom j in school k in district l , π_{0jkl} is a constant, π_{pjkl} are level-1 coefficients, α_{pijkl} is level-1 predictor p for student i in classroom j in school k in district l , and e_{ijkl} is the level-1 random error.

$$\text{Level-2 (classrooms): } \pi_{pjkl} = \beta_{p0kl} + \sum_{q=1}^{Q_p} \beta_{pqkl} X_{qjkl} + r_{pjkl},$$

where π_{pjkl} are the coefficients from the level-1 model, β_{p0kl} is a constant, β_{pqkl} are level-2 coefficients, X_{qjkl} are level-2 predictors, and r_{pjkl} are level-2 random effects.

$$\text{Level-3 (schools): } \beta_{pqjkl} = \gamma_{pq0} + \sum_{s=1}^{S_{pq}} \gamma_{pqsl} W_{skl} + u_{pqkl},$$

where β_{pqjkl} are the coefficients from the level-2 model, γ_{pq0} is a constant, γ_{pqsl} are level-3 coefficients, W_{skl} are level-3 predictors, and u_{pqkl} are level-3 random effects.

$$\text{Level-4 (districts): } \gamma_{pqsl} = \delta_{pq0} + \sum_{g=1}^{G_{pq}} \delta_{pqsg} Z_{gl} + v_{pqsl},$$

where γ_{pqsl} are the coefficients from the level-3 model, δ_{pqso} is a constant, δ_{pqsg} are level-4 coefficients, Z_{gl} are level-4 predictors, and v_{pqsl} are level-4 random effects.

The simplest implementation of these equations involves no predictors (other than the constants) and allows us to estimate the proportion of the variance in reading and math achievement test scores at fourth and fifth grades that is associated with differences between students vs. classrooms vs. schools vs. districts. In this simple case, we calculate the variance at each level as the variance of the random effects, and divide this by the total variance (sum of variances at each level) to obtain the proportion of the variance at that level. We implement the analysis separately by state, subject (math and reading), and school year.

Previous studies of leadership effects at the school level have debated the proper conceptual models and methodological approaches for addressing causal research questions on this subject. For example, Hallinger and Heck (2011) discuss how school leadership might be modeled as part of a reciprocal process rather than as only a driver of student learning. It stands to reason that this argument could be applied to the district level as well, but we use the simpler HLM approach because the aim of this study is to partition the variance in student achievement associated with each of the four levels in the model, not to estimate the causal impact of district leadership or any other specific factor. But the important issue of how to model reciprocal interactions in school districts certainly represents fertile ground for future research. For example, the variance decomposition estimates reported below set the stage for future work of topics such as how the variance shares associated with each level might change when contexts are more or less favorable to reciprocal interactions.

Our use of simple variance decomposition techniques does not imply that we view districts as only having direct (non-mediated) impacts on student achievement. To the contrary,

district effects likely work through effects at the school and classroom levels. For example, a more effective district may have policies that lead to the recruitment and retention of more effective principals and teachers. We would understate the importance of districts if we were to directly control for outcome or process variables at the school and classroom levels. However, our HLM analysis does not have this disadvantage because it will apportion to the district level variation in student outcomes that results from classroom- or school-level policies to the extent that such variation is common across all classrooms or schools. District-level policies may also influence the variation in student achievement accounted for at the school and classroom levels. Our HLM analysis essentially averages across all districts in calculating the share of the variance associated with each level.

We also estimate versions of the HLM models that include control variables at the student level to account for variation in student characteristics that is correlated with the teacher, school, and district effects. The controls include age, race/ethnicity, cognitive disability status, free and reduced lunch program status, limited English proficiency status, and, for Florida only, whether the parent and student are native English speakers and whether the student was born in the U.S. These models also include aggregate characteristics of classrooms, schools, and districts to account for the correlation between the student-level covariates and the random effects.²

Finally, we estimate models that also control for students' test scores from the prior year. We largely estimate this model because it is the most appropriate model for current-year teacher effects, which we compare to school and district effects. It is not our preferred model for district

² We selected which aggregate characteristics to include using a method developed by Mundlak (1978). We ran models at each level (district, school, and classroom) to determine which aggregate controls were significantly correlated with random effects at each nested level. In order to correct for this correlation, we included the following aggregate characteristics in the final models: district-level free and reduced lunch program status and race; school-level age, race, cognitive disability status, free and reduced lunch program status, whether the parent/student are native English speakers, and whether the student was born in the U.S.; and classroom-level age, ethnicity/race, cognitive disability status, free and reduced lunch program status, whether the parent/student are native English speakers, and whether the student was born in the U.S.

and school effects because it is likely that districts that have an impact on student achievement do so in all prior grades as well as the current grade. If we measure district effectiveness only by the gains generated for students in grades four and five, we would not capture variation in performance in earlier grades. We would also net out any district-level variables such as teacher quality that are associated with prior-year test scores. Consequently, although prior-year scores can serve as a proxy for other unmeasured student characteristics, in this case controlling for them will likely cause us to understate the influence of districts on student achievement.

Because our preferred analysis of district effects does not condition on prior-year scores, it potentially captures all prior years of district influence—not the single year typical of most “value-added” type models. In other words, the analysis estimates the extent to which districts serving similar student populations produce better or worse outcomes in fourth and fifth grade, which will reflect impacts both in those grades and persistent impacts from prior grades.

It is important to emphasize that variance decomposition, no matter how cleverly applied, does not lead to point estimates of causal effects. It is the underlying research design, such as random assignment, that permits causal inference. We frequently use terminology that suggests causal effects because alternate phrasing would be convoluted. However, our methods are observational and do not allow us to make rigorous causal conclusions. When we write, for example, about “differences in student achievement attributable to school districts” we might more accurately write about “associations among student test scores and the districts in which students are educated that remain after accounting for variation in student achievement within districts that is associated with teachers and schools, and with the inclusion of statistical controls for demographic characteristics of students.”

HLM is related to but not identical to approaches used by economists to deal with estimation of effects in hierarchical data, specifically fixed effects. We prefer HLM for our present purposes because it is descriptive whereas the econometric models are focused on the causal effect of one level in a multi-level design, e.g., what is the causal effect of teachers on student achievement having fixed the effects of schools? However, as a test of the robustness of the HLM results we also implement fixed effects regressions that calculate average achievement by district adjusted for student characteristics (but not taking into account the nested structure of the data). Fixed effect estimates, like the results of variance decomposition techniques, are not rigorous causal estimates.

We estimate district fixed effects models using the following specification:

$$Y_{ijkl} = \beta_0 + \alpha X_{ijkl} + v_l + \epsilon_{ijkl} ,$$

where Y_{ijkl} is the test score of student i in classroom j in school k in district l , β_0 is a constant, X_{ijkl} is a vector of student characteristics (identical to those used in the variance decomposition analysis) with coefficient vector α , v_l is a vector of district fixed effects, and ϵ_{ijkl} is a standard zero-mean error term. The coefficients on the district fixed effects are our district effect estimates for the state, test, subject, and year included in the estimation. Below we show that the district fixed effect estimates are highly correlated with the random effect estimates.

Results

The results of the variance decomposition analysis are presented in Tables 1a and 1b for fourth- and fifth-grade students in 2009-10, the most recent year in both the Florida and North Carolina datasets.³ The first column of Table 1a shows, for a model of math achievement with no control variables, the variance of the random effects at each level, as well as the corresponding share of the total variance in Florida. For example, the results for the district level indicate that the district random effects have an estimated variance of 0.015, which is 1.3 percent of the total variance in math achievement. The share of variance explained increases at the lower levels, to 9 percent at the school level, 31 percent at the teacher level, and the balance of 58 percent at the student level or unexplained (i.e. the residual variance, which includes differences across students within classrooms as well as measurement error).

The second column of Table 1a shows results that include controls for student-level demographic covariates, which explain about one-third of the variance in student outcomes.⁴ The shares of the variance explained at the district, school, and teacher levels drop to 1, 2, and 12 percent, respectively. This is not surprising given the well-documented correlation between students' demographic characteristics and their districts, schools, and teachers. The relative drop is smaller at the district level than the school level, probably because there is greater within-district sorting of families across schools than across-district sorting given that Florida school districts are coterminous with counties (i.e. geographically large).

³ The full output of the HLM models, including standard errors of the variance estimates, coefficients and standard errors for control variables, and fit statistics, are available from the authors upon request.

⁴ We calculate the variance explained by the controls as the difference between the total variance in a null model (i.e. a model with no controls) and the variance explained at the district, school, teacher, and student levels. The null model is estimated separately for the sample of students that can be included in each non-null model.

In column three we add controls for prior-year scores, which more than doubles the share of variance explained by the controls to 71 percent.⁵ It is unsurprising that prior-year scores are the strongest predictor of current-year scores, and that the share of variance explained by districts, schools, and teachers falls once student's prior achievement is taken into account. The share of variance at the district level falls to one-tenth of one percent (i.e. one thousandth of the total variance). But as we discuss above, these estimates likely understate the importance of districts and schools (but not teachers, who usually only instruct a student for a single year).

The analysis of reading scores, which is reported in columns four through six, yields a similar pattern of results. In our preferred model (column 5), the share of variance explained at the district, school, and teacher levels is smaller for reading than for math, which is consistent with prior research showing that formal education has a larger impact on math achievement than on reading achievement because the latter is more influenced by activities in the home. We also obtain similar results from each of the other school years going back to 2000-01, an interesting finding given the various reforms implemented in Florida during this decade and the substantial improvement in overall student achievement that occurred (Chingos 2012).

Table 1b shows corresponding results for the 115 districts in North Carolina. As in Florida, teachers account for more of the variation in student achievement than schools and districts and the three institutional components account for more variance in math than in reading. But districts explain more than twice the share of variance in North Carolina as they do in Florida: 1.9 and 1.3 percent in math and reading, respectively, compared to 0.9 and 0.5 percent in Florida (using our preferred model that controls only for student demographics). This may be due to the fact that North Carolina districts are smaller than Florida districts, on average.

⁵ The sample of students changes from columns (2) to (3) due to missing data on prior-year scores. When we run the model in column (2) on the sample of students included in column (3), we obtain qualitatively similar results (not shown).

Consequently, superintendents of smaller districts may more easily be able to change education policies and practices than their counterparts in larger districts.⁶ There may also be more idiosyncratic variability in smaller districts, such as the departure of a highly effective principal of a school that accounts for a significant share of enrollment in the district.

A policymaker may be more interested in how important each of the three “institutional” levels is relative to the total variance across just these three levels (ignoring the variation explained by the control variables and the variation across students within classrooms). Table 2 shows the variance shares rescaled in this way, and indicates that the relative importance of the three levels is much less sensitive to model specification than implied by Tables 1a and 1b, especially in Florida. Teachers in Florida consistently account for 75-86 percent of the institutional variance. Once demographics are accounted for, schools explain 9-13 percent of the variance and districts explain the remaining 4-6 percent. In North Carolina, our preferred estimates apportion 20-23 percent of the variance to the district level, 22-27 percent to the school level, and 53-55 percent to the teacher level.

Districts explain a relatively small share of the total variation in student achievement, but are there differences among districts in their contribution to student achievement that are large enough to be relevant for policy? Table 3 converts the variances reported in our preferred models in Tables 1a and 1b to standard deviations (recall that the standard deviation is equal to the square root of the variance). These results indicate that a one standard deviation move in the distribution of district effects is associated with an increase in student achievement of 0.07-0.14 standard deviations. Such a difference corresponds to about 7 weeks of schooling in Florida (about one-fifth of a school year) and 10-11 weeks in North Carolina (one-quarter to one-third of

⁶ Some districts in North Carolina contain only one elementary school. We obtain qualitatively similar results when we exclude these schools from the analysis.

a school year).⁷ It is worth noting that this additional learning is measured in the tested grades, but potentially reflects cumulative effects from earlier grades as well.

The finding that districts explain a small share of variance but are potentially important for student achievement may seem counterintuitive at first, but is consistent with the finding that teachers also explain a relatively paltry share of test scores (less than 4 percent in models that control for prior-year scores) but are hugely important for student achievement. For example, a teacher-level variance of 0.039 in math in Florida, about 3.7 percent of the total variance, corresponds to a standard deviation of 0.20 standard deviations (more than 40 percent of a year of learning). Of course this variation corresponds to differences in teacher quality within schools, which averages out to be roughly the same for any given student over time, as compared to the district and school effects which are relative to other districts and schools.

We also compute the best linear unbiased predictions (BLUPs) of the district-level random effects for each state, subject, and year from the hierarchical linear model. The BLUP calculation procedure shrinks the noisier estimates—those based on a level that explains less variance or those based on smaller clusters—toward the mean. The BLUPs and their 95 percent confidence intervals for math scores in 2009-10 are shown in Figures 1 and 2.

There are a number of districts in both states that perform at levels that are above or below the average for districts in the state by a statistically significant margin (using a 5 percent significant level). In Florida, 15 percent of districts are statistically significantly above average and 12 percent are statistically below average. The comparable numbers in North Carolina are 16 percent above and 16 percent below. This means there are districts that are over- or under-

⁷ Standard deviations are converted to weeks of schooling using the average learning gains for fourth and fifth grade reported by Hill et al. (2008), assuming a 180-day (36-week) school year. Consequently, the weeks of schooling correspond to these two grades, a useful fact to bear in mind given that our preferred estimates reflect cumulative learning from prior years as well.

performing on student achievement relative to what might be expected of them given the characteristics of their students.

The difference in performance between districts is quite large at the extremes: 0.40 student-level standard deviations separate the highest and lowest performing district in the Florida math data whereas the difference is 0.61 standard deviations for North Carolina. Using the same conversion from Table 3, these ranges correspond to about 80 percent of a school year in Florida and more than 120 percent of a school year in North Carolina. The random effect estimates for reading (not shown) follow a similar pattern, as we would expect given the high correlation between math and reading performance aggregated to the district level (correlation coefficients of 0.80 in 2009-10 in both states).

These results are robust to replacing the HLM random effects specification with a fixed effects specification.⁸ The standard deviations of the district fixed effects estimates for 2009-10 are 0.13 and 0.10 in Florida for math and reading, respectively, and 0.15 and 0.12 in North Carolina.⁹ These standard deviations are modestly larger than those from the random effects model because the fixed effects model does not shrink noisier estimates toward the mean (as the random effects model does). The random and fixed effects produce similar rank orderings of districts, with correlation coefficients in the 0.70-0.88 range.

School district performance, as measured by fourth- and fifth-grade reading and math scores, is quite stable over time. Table 4 shows the year-to-year correlations of the estimated random effects (using our preferred model with demographic controls only) for the ten-year

⁸ The HLM and fixed effects specifications use the same student-level control variables, but the HLM specification also includes additional control variables aggregated to the classroom, school, and district levels, as described above (note that it is not possible to include district-level controls in a single-year fixed effects specification).

⁹ These standard deviations are weighted by the number of students contributing to each district's fixed effect estimate. The unweighted standard deviations are 0.19 and 0.12 in math and reading, respectively, in Florida, and 0.17 and 0.14 in North Carolina.

period from 2001-01 through 2009-10. The correlation coefficients are in the 0.85-0.90 range in both Florida and North Carolina. In other words, districts with high or low student achievement (conditional on demographics) in a given year tend to also have high or low achievement the following year.

There are some notable exceptions to this pattern in our data. In Whitehurst, Chingos, and Gallaher (2013), we show examples of districts that, instead of having consistently high or low performance, saw steep declines or impressive gains. For example, one North Carolina district that was at or above the 90th percentile of math performance every year from 2000-01 to 2004-05 was below the 10th percentile in every year from 2007-08 through 2009-10. One way to formalize this analysis is to ask how many districts in Florida and North Carolina experienced statistically significant changes (at the 10 percent level) in performance over the decade covered by our data. In other words, how many had 90 percent confidence intervals in 2000-01 that do not overlap with their 2009-10 intervals?

Out of the 67 Florida districts, five experienced significant gains in math over this period and four experienced significant declines. In reading, six Florida districts experienced significant gains and there were seven significant declines. Among the 115 North Carolina districts, eight experienced gains and seven experienced declines in math, whereas seven gained and four declined in reading. These results indicate that significant changes were uncommon, generally occurring in fewer than 10 percent of districts, but not unheard of.

Discussion

It is unsurprising that student achievement varies more at the levels close to students—teachers and schools—than at the district level. But our results suggest that there are important differences in student achievement across school districts after taking into account differences in student characteristics. Moving from approximately the 30th to the 70th percentile corresponds to 0.07-0.14 standard deviations in fourth- and fifth-grade student test scores, or 20-33 percent of a year of learning.

A few caveats to these results are worth noting. First, it could be the case that each and every Florida district has a sizeable and similar impact on student achievement through a set of very similar practices. If, for example, every district in Florida added about 6 months of academic growth to the students it served, our analysis would not pick this up, depending as it does on examining variation in outcomes across districts. Such a “main effect” for districts could conceivably be wrought by relieving those closer to instructional interactions, such as building principals, from the time they would have to spend on non-instructionally relevant tasks such as providing for student transportation and meals if there were no district administration above them. We note, however, that this model of a school district effect in which the major function of the district is to provide efficient business services, which nearly all do, is very different from the model in which districts compete for great leaders to drive education reform and enhance student achievement.

Second, our model only considers student performance on state math and reading exams in fourth and fifth grades. Districts may well have effects on performance in other grades, other subjects, and on skills not captured by standardized exams. Our analysis will not capture these effects to the extent that they are unrelated to math and reading test scores in elementary school.

And as discussed earlier, our observational methodology generates exploratory findings regarding the importance of districts relative to other levels but does not produce rigorous estimates of causal effects.

Finally, it should be noted that the results from Florida and North Carolina cannot necessarily be extrapolated to states with numerous smaller districts. Smaller districts may be more likely to vary in their practices given the greater ease with which a superintendent can change the course of a small district containing a handful of schools, as compared to her counterpart in a large district. This theory is supported by a comparison of the results for Florida and North Carolina. North Carolina has more districts than Florida despite being a less populous state, and districts account for more of the variation in student achievement in North Carolina than in Florida. But data from these two states do not allow us to measure whether districts are more important in states with more numerous districts, such as Texas which has over 1,000 districts.

These limitations aside, these results represent the first effort to measure the importance of school districts using student-level databases and modern statistical techniques. As we discuss in Whitehurst, Chingos, and Gallaher (2013), the decade of data from Florida and North Carolina is rife with examples of consistently high and low performing districts, districts that have seen precipitous declines in achievement, and districts that have made transformative progress. Whether certain characteristics of districts, such as policies and leadership, help explain these different patterns is a ripe subject for future research.

However, it is important that the next generation research on district effectiveness not repeat the mistakes of the first. Previous studies have purported to identify the features of more effective districts, including the importance of standards and curriculum, organizational

structures, instructional leadership, monitoring and evaluation, and professional development (Trujillo 2013). But as we discuss above, the existing literature largely consists of case studies of districts judged to be effective, ignoring the fact that it is difficult to extrapolate beyond a single district or handful of districts, and it is impossible to judge even the correlation between district practices and effectiveness without examining data on districts across the performance spectrum. It is our hope that future research on this subject will test the theories suggested by the largely anecdotal existing evidence using more credible methodologies.

The first stage in such research would be to examine whether there are associations between observed district characteristics, such as the superintendent, and measures of district performance, such as student achievement adjusted for demographics. In addition to taking student characteristics into account, it is important that this research examine a wide range of districts with varying performance rather than a handful of unusually high-performing ones. The second stage would be to evaluate promising district-level reforms using randomized controlled trials. For example, Slavin et al. (2013) evaluate a data-driven reform strategy using random assignment across 59 districts. A key challenge facing such work is the need to recruit a large enough number of districts to have sufficient statistical power. But if successful district-level policies can be rigorously identified, then they can be put in place to the benefit of many more students than are typically enrolled in any single school or classroom.

Acknowledgments

We thank the Florida Department of Education and the North Carolina Education Research Data Center for sharing the data used in this paper. Helpful comments were provided by participants at the APPAM 2012 fall research conference and two anonymous referees. An anonymous foundation provided financial support for this project.

References

- Bowers, Alex J. 2010. Toward addressing the issues of site selection in district effectiveness research: A two-level hierarchical linear growth model. *Educational Administration Quarterly* 46(3): 395-425.
- Chingos, Matthew M. 2012. The impact of a universal class-size reduction policy: Evidence from Florida's statewide mandate. *Economics of Education Review* 31(5): 543-562.
- Hallinger, Philip, and Ronald H. Heck. 2011. Conceptual and methodological issues in studying school leadership effects as a reciprocal process. *School Effectiveness and School Improvement: An International Journal of Research, Policy, and Practice* 22(2): 149-173.
- Hill, Carolyn J., Howard S. Bloom, Alison R. Black, and Mark W. Lipsey. 2008. Empirical benchmarks for interpreting effect sizes in research. *Child Development Perspectives* 2(3): 172-177.
- Hoxby, Caroline M., and Gretchen Weingarth. 2005. Taking race out of the equation: School reassignment and the structure of peer effects. Harvard University, working paper.
- Klein, Alyson. 2013. GAO: Race to Top states have mixed record on teacher evaluation. *Education Week Politics K-12* blog, September 18, 2013. Available http://blogs.edweek.org/edweek/campaign-k-12/2013/09/gao_race_to_the_top_states_hav.html.
- Leithwood, Kenneth. 2010. Characteristics of school districts that are exceptionally effective in closing the achievement gap. *Leadership and Policy in Schools* 9(3): 245-291.
- Mundlak, Yair. 1978. On the pooling of time series and cross section data. *Econometrica* 46(1): 69-85.
- New York State Education Department, Administrative Compensation Information for 2011-2012. Available http://www.p12.nysed.gov/mgtserv/admincomp/docs/2011-12_AdminSalDisc_5_11_11_Post_r.xls. Accessed 13 July 2011.
- Raudenbush, Stephen W., and Anthony S. Bryk. 2002. *Hierarchical linear models: Applications and data analysis methods*. 2nd edition. Thousand Oaks, CA: SAGE Publications.
- Raudenbush, Stephen W., Anthony S. Bryk, Yuk F. Cheong, Richard T. Congdon, and Mathilda du Toit. 2011. *HLM 7: Hierarchical linear and nonlinear modeling*. Lincolnwood, IL: Scientific Software International.
- Rorrer, Andrea K., Linda Skrla, and James J. Scheurich. 2008. Districts as institutional actors in educational reform. *Educational Administration Quarterly* 44(3): 307-357.
- Samuels, Christina A. 2011. Critics target growing army of Broad leaders. *Education Week*, June 8, 2011.

Slavin, Robert E., Alan Cheung, GwenCarol Holmes, Nancy A. Madden, and Anne Chamberlain. 2013. Effects of a data-driven district reform model on state assessment outcomes. *American Educational Research Journal* 50(2): 371-396.

Trujillo, Tina. 2013. The reincarnation of the effective schools research: Rethinking the literature on district effectiveness. *Journal of Educational Administration* 51(4): 426-452.

Waters, J. Timothy, and Robert J. Marzano. 2006. School district leadership that works: The effect of superintendent leadership on student achievement. Denver, CO: Mid-continent Research for Education and Learning.

Whitehurst, Grover J., Matthew M. Chingos, and Michael R. Gallaher. 2013. Do school districts matter? Washington, DC: Brown Center on Education Policy, Brookings Institution.

Figure 1. Random-Effect Estimates and 95 Percent Confidence Intervals, Math Achievement, Student-Level Standard Deviation Units, Florida, 2009-10

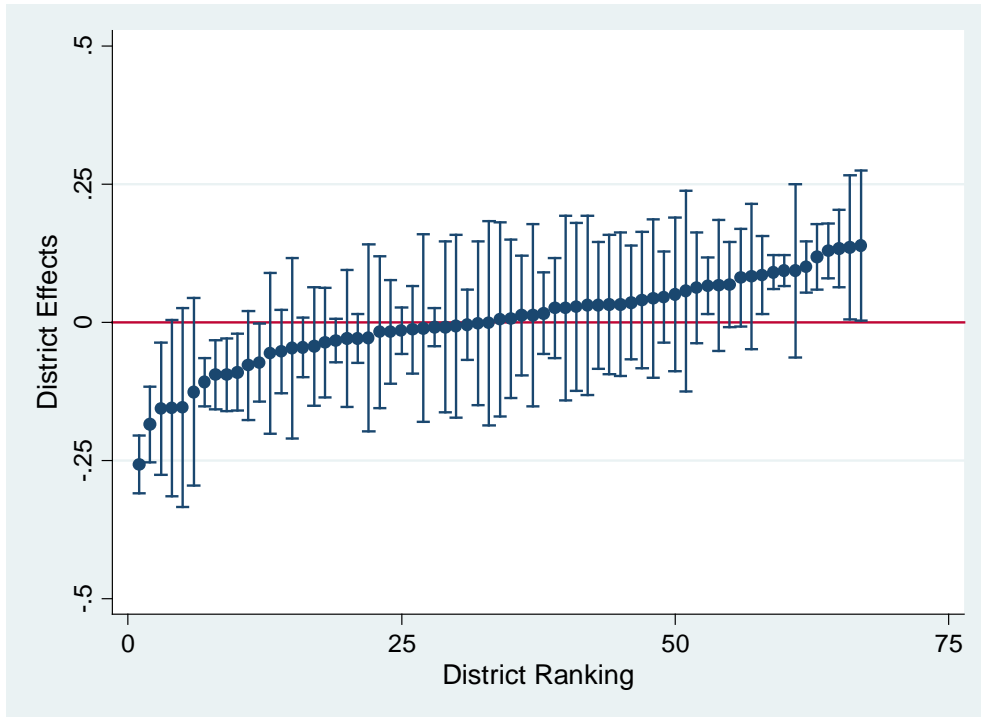


Figure 2. Random-Effect Estimates and 95 Percent Confidence Intervals, Math Achievement, Student-Level Standard Deviation Units, North Carolina, 2009-10

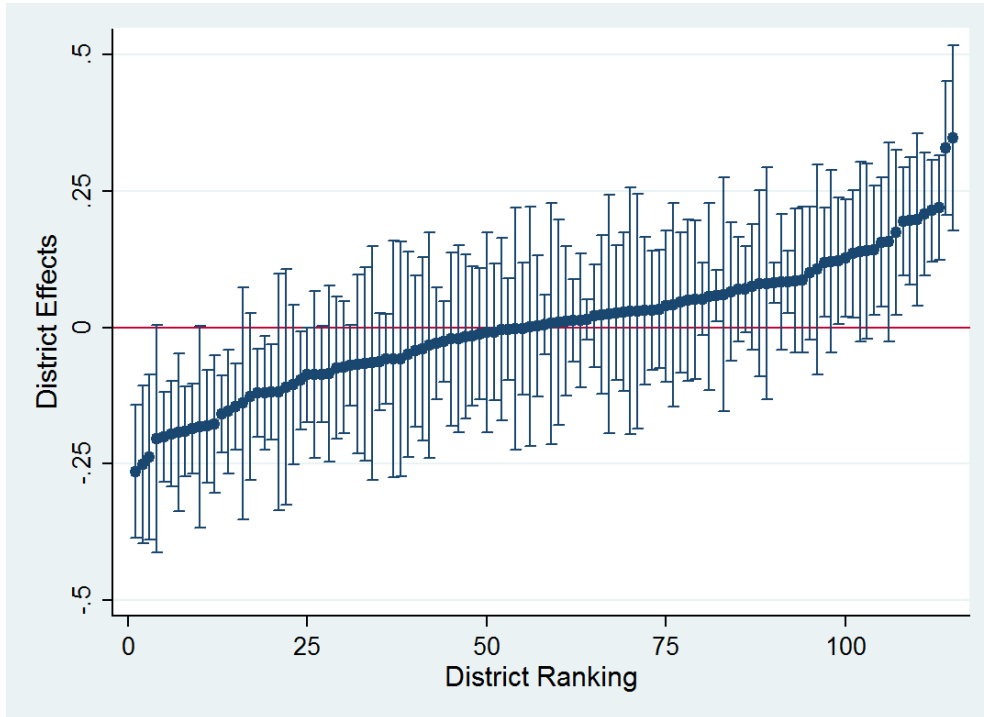


Table 1a. Variance Decomposition of 4th- and 5th-Grade Student Achievement, Florida, 2009-10

| | Math | | | Reading | | |
|-----------------------------------|----------------|----------------|----------------|----------------|----------------|----------------|
| | (1) | (2) | (3) | (4) | (5) | (6) |
| District level | 0.015 1.3% | 0.010 0.9% | 0.001 0.1% | 0.017 1.6% | 0.005 0.5% | 0.001 0.1% |
| School level | 0.099 9.0% | 0.021 1.9% | 0.007 0.7% | 0.089 8.2% | 0.011 1.0% | 0.003 0.3% |
| Teacher level | 0.347 31.4% | 0.136 12.3% | 0.039 3.7% | 0.328 30.1% | 0.104 9.5% | 0.018 1.7% |
| Student level (residual variance) | 0.646 58.3% | 0.596 53.9% | 0.255 24.1% | 0.657 60.2% | 0.605 55.4% | 0.306 29.3% |
| Controls | | 0.34 31.0% | 0.76 71.4% | | 0.37 33.6% | 0.72 68.6% |
| Total variance, null model | 1.11 | 1.11 | 1.06 | 1.09 | 1.09 | 1.04 |
| Demographic controls? | No | Yes | Yes | No | Yes | Yes |
| Prior-year scores? | No | No | Yes | No | No | Yes |
| Observations | 304,168 | 304,168 | 286,677 | 339,783 | 339,783 | 320,554 |

Notes: Estimates of variance of random effects are calculated using hierarchical linear models. Share of variance is calculated as the variance divided by the total variance from the "null model" which includes no control variables. Demographic controls include student-level variables for grade, age, race/ethnicity, cognitive disability status, free and reduced lunch program status, limited English proficiency status, whether the parent/student are native English speakers, and whether the student was born in the U.S.; district-level free and reduced lunch program (FLRP) status and race; school-level age, race, cognitive disability, free and reduced lunch program status, whether the parent/student are native English speakers, and whether the student was born in the U.S; and classroom-level age, race, cognitive disability, FLRP status, English language parent/student, and U.S. born. Prior-year scores include test scores in both math and reading, dummy variables for the student's grade in the prior year, and interactions between these dummies and the prior-year scores.

Table 1b. Variance Decomposition of 4th- and 5th-Grade Student Achievement, North Carolina, 2009-10

| | Math | | | Reading | | |
|-----------------------------------|----------------|----------------|----------------|----------------|----------------|----------------|
| | (1) | (2) | (3) | (4) | (5) | (6) |
| District level | 0.047 4.6% | 0.019 1.9% | 0.003 0.3% | 0.041 4.1% | 0.013 1.3% | 0.002 0.2% |
| School level | 0.098 9.7% | 0.026 2.6% | 0.011 1.1% | 0.093 9.3% | 0.012 1.2% | 0.003 0.3% |
| Teacher level | 0.082 8.2% | 0.052 5.2% | 0.032 3.3% | 0.059 5.9% | 0.030 3.0% | 0.011 1.1% |
| Student level (residual variance) | 0.781 77.5% | 0.635 63.2% | 0.236 24.1% | 0.815 80.8% | 0.652 64.9% | 0.278 28.1% |
| Controls | | 0.27 27.1% | 0.70 71.2% | | 0.30 29.6% | 0.69 70.3% |
| Total variance, null model | 1.01 | 1.00 | 0.98 | 1.01 | 1.00 | 0.99 |
| Demographic controls? | No | Yes | Yes | No | Yes | Yes |
| Prior-year scores? | No | No | Yes | No | No | Yes |
| Observations | 208,485 | 205,572 | 191,479 | 206,649 | 203,780 | 190,955 |

Notes: Estimates of variance of random effects are calculated using hierarchical linear models. Share of variance is calculated as the variance divided by the total variance from the "null model" which includes no control variables. Demographic controls include student-level variables for grade, age, race/ethnicity, cognitive disability status, free and reduced lunch program status, and limited English proficiency status; district-level free and reduced lunch program (FLRP) status and race; school-level age, race, cognitive disability, and FLRP status; and classroom-level age, race, cognitive disability, and FLRP status. Prior-year scores include test scores in both math and reading, dummy variables for the student's grade in the prior year, and interactions between these dummies and the prior-year scores.

Table 2. Institutional Variance Shares, 4th- and 5th-Grade Student Achievement, 2009-10

| Florida | | | | | | |
|-----------------------|---------|---------|---------|---------|---------|---------|
| | Math | | | Reading | | |
| | (1) | (2) | (3) | (4) | (5) | (6) |
| District level | 3.2% | 6.0% | 3.0% | 3.9% | 4.4% | 2.7% |
| School level | 21.5% | 12.7% | 15.2% | 20.6% | 9.4% | 13.3% |
| Teacher level | 75.3% | 81.3% | 81.8% | 75.5% | 86.2% | 84.0% |
| Demographic controls? | No | Yes | Yes | No | Yes | Yes |
| Prior-year scores? | No | No | Yes | No | No | Yes |
| Observations | 304,168 | 304,168 | 286,677 | 339,783 | 339,783 | 320,554 |
| North Carolina | | | | | | |
| | Math | | | Reading | | |
| | (1) | (2) | (3) | (4) | (5) | (6) |
| District level | 20.6% | 19.6% | 6.6% | 21.2% | 23.2% | 10.8% |
| School level | 43.1% | 27.0% | 23.3% | 48.3% | 21.7% | 19.0% |
| Teacher level | 36.4% | 53.4% | 70.1% | 30.5% | 55.1% | 70.2% |
| Demographic controls? | No | Yes | Yes | No | Yes | Yes |
| Prior-year scores? | No | No | Yes | No | No | Yes |
| Observations | 208,485 | 205,572 | 191,479 | 206,649 | 203,780 | 190,955 |

Notes: See notes to Tables 1a and 1b. Percentages are calculated as the variance of the random effects at each level divided by the total variance at the three institutional levels (district, school, and teacher).

Table 3. Distribution of District-Level Random Effect Estimates, 4th- and 5th-Grade Student Achievement, 2009-10

| | Florida | | North Carolina | |
|---|---------|---------|----------------|---------|
| | Math | Reading | Math | Reading |
| Standard deviation, in test scores | 0.10 | 0.07 | 0.14 | 0.11 |
| Standard deviation, in weeks of schooling | 7.4 | 7.3 | 10.3 | 11.3 |

Notes: Standard deviations, calculated in student-level test-score standard deviations, are calculated as square root of variances reported in columns 2 and 5 of Tables 1a and 1b. These statistics are converted to weeks of schooling by dividing by 0.485 in math and 0.36 in reading (the average learning gain between for 4th- and 5th-grade students reported by Hill et al. [2008]) and multiplying by 36 (the number of weeks in a typical 180-day school year).

Table 4. Year-to-Year Correlations of Random Effect Estimates, 2000-01 through 2009-10

| | Florida | | North Carolina | |
|------------|---------|---------|----------------|---------|
| | Math | Reading | Math | Reading |
| Unweighted | 0.85 | 0.84 | 0.86 | 0.87 |
| Weighted | 0.88 | 0.87 | 0.88 | 0.89 |

Notes: Coefficients indicate linear correlation between best linear unbiased predictions of district-level random effects (for 4th- and 5th-grade student achievement) in adjacent years. Weighted estimates are weighted by number of students contributing to each random effect estimate.