

education sector reports

Growth Models and Accountability: A Recipe for Remaking ESEA

By Kevin Carey and Robert Manwaring



EDUCATIONSECTOR

www.educationsector.org

ACKNOWLEDGEMENTS

This report was funded by the Stuart Foundation. Education Sector thanks the foundation for its support. The views expressed in the paper are those of the authors alone.

ABOUT THE AUTHORS

KEVIN CAREY is the policy director at Education Sector. He can be reached at kcarey@educationsector.org.

ROBERT MANWARING is a fiscal and policy consultant. He can be reached at robert.manwaring@msn.com.

ABOUT EDUCATION SECTOR

Education Sector is an independent think tank that challenges conventional thinking in education policy. We are a nonprofit, nonpartisan organization committed to achieving measurable impact in education, both by improving existing reform initiatives and by developing new, innovative solutions to our nation's most pressing education problems.

© Copyright 2011 Education Sector

Education Sector encourages the free use, reproduction, and distribution of our ideas, perspectives, and analyses. Our Creative Commons licensing allows for the noncommercial use of all Education Sector authored or commissioned materials. We require attribution for all use. For more information and instructions on the commercial use of our materials, please visit our website, www.educationsector.org.

1201 Connecticut Ave., N.W., Suite 850, Washington, D.C. 20036
202.552.2840 • www.educationsector.org

There are nearly 100,000 public schools in the United States, but President Barack Obama praised just one of them in his 2011 State of the Union address. It was Bruce Randolph School in Denver, Colorado. The Colorado Department of Education had identified Bruce Randolph as the worst-performing middle school in the state just four years before. But, after firing most of the teachers, expanding to grades six-12, and being liberated from district and teachers union regulations on spending and hiring, Bruce Randolph made rapid progress. Student test scores grew rapidly, and in May 2010, 97 percent of seniors graduated. Nearly nine out of 10 went on to college. “That’s what good schools can do,” said the President to Congress and the nation, “and we want good schools all across the country.”

To achieve that vision, the Obama administration has proposed major changes to the federal Elementary and Secondary Education Act (ESEA) created in 1965 and last reauthorized by Congress in 2001 as the No Child Left Behind Act (NCLB). Rock-bottom performers like Bruce Randolph should be aggressively reconstituted, according to the administration, and judged by how much academic progress, or achievement growth, individual students make each year. Such “growth model” systems of evaluating school performance stand in contrast to the NCLB system of judging schools, which is based strictly on the percentage of students who pass standardized tests, regardless of how well or poorly those students had performed in previous years. According to the Colorado Department of Education, the rate of achievement growth among middle and high school students at Bruce Randolph has consistently outpaced most other students statewide.

But growth model systems also bring complications. While the state found that individual students at Bruce Randolph had improved more than their peers, the state’s data also indicated that *overall* achievement at Bruce Randolph was not good. Forty-three percent of its students scored “proficient” in reading in 2010, near the state average. But only 16 percent were

proficient in writing, and only 13 percent hit the mark in math. The state also acknowledged that although achievement growth at Bruce Randolph was above average in every subject, those growth rates were inadequate to put students on pace to catch up and learn what they needed to know before graduating. Nearly every student in Bruce Randolph’s first class of freshmen earned a diploma and went to college, a remarkable achievement. But it’s likely that many of them arrived on campus with serious learning deficits that will hamper their ability to stay in college and earn a degree.

Bruce Randolph epitomizes the challenge of incorporating information about student growth into educational accountability—a challenge that every state and school district in America will face if ESEA is revised as the administration proposes. Measuring growth is a delicate balancing act. Policymakers need to be fair and constructive with educators working in immensely difficult school environments. But public officials must also hold fast to the end goal of helping students thrive in a world that makes ever-higher demands on workers and citizens. As the political will and technical capacity to hold schools accountable for student academic progress converge, growth models appear to be an idea whose time has come.

LOOKING BACK: GROWTH AND ACCOUNTABILITY

The modern standards- and testing-based school accountability movement began in the late 1980s and accelerated in 1994 when President Clinton and a bipartisan group of legislators in Congress reauthorized ESEA. That version was called the Improving America's Schools Act (IASA). For the first time, the federal government required states to create common academic standards for all students and hold schools accountable for student scores on standardized tests. It wasn't easy work. In 1998, the National Education Goals Panel (a nonprofit group originally created by President George H.W. Bush and a bipartisan collection of reform-minded governors) recognized the limitations of relying solely on bottom-line measures of academic proficiency and spoke to the promise of measuring annual growth:

“A key issue faced by states in establishing systems of accountability is how to take into account the strong correlation of test scores with the socio-economic status (SES) of the students. Perceived unfairness in the system of rankings and rewards can seriously erode the trust necessary for effective incentives. If actual scores were primarily utilized to rank schools and give rewards, the schools in higher SES school districts would currently dominate the top rankings. However, year-to-year gains in scores can provide a potential advantage to schools with lower SES students since gains can be greater for lower scoring students.”¹

Educational accountability, in other words, isn't just a matter of identifying which schools have the most failing students. It also requires some response to that information that will help fewer students fail. It's unfair to blame educators for test scores that are substantially caused by external SES factors. And while the Goals Panel didn't say so explicitly, it's also unfair to blame educators for the failures and shortcomings of other educators who previously taught their students. Unfair accountability systems are unlikely to spur improvement.

To date, responsibility for wrestling with this dilemma has fallen primarily to the states. IASA mandated standards, tests, and accountability, but it also gave states a great deal of flexibility in deciding how to

implement such a system. Some took to the project with more enthusiasm than others. Then-Tennessee Gov. Lamar Alexander had been an early standards proponent in the 1980s before becoming U.S. Secretary of Education in 1991. In North Carolina, four-term Gov. James Hunt pushed his state toward standards-based reform. And most prominently, standards and tests were enthusiastically backed in Texas by then-Gov. George W. Bush.

These early adopter states made two decisions that were crucial to the development of growth models. First, they tested students annually, allowing for the calculation of year-to-year growth in student achievement. Second, they created sophisticated statewide repositories of student data, allowing them to calculate annual learning growth in an accurate, consistent manner for every school. These large data systems also allowed states to estimate learning growth for students who moved among different schools, something beyond the capacity of local districts.

In the early 1990s, William Sanders, an agricultural statistics professor at the University of Tennessee, used the state's recently created annual test data to gauge the effectiveness of individual teachers by comparing an estimate of how their students' test scores were expected to grow, based on the students' previous performance history, to how much their students' test scores actually grew. These so-called “value-added” estimates slowly spread across the country as more states created annual tests and data systems. (They are now at the center of a raging controversies in Los Angeles, New York City, and elsewhere, as education reformers and teachers unions debate the use of standardized test-score data in determining teacher tenure, firing, and compensation policies.² The use of such estimates for individual schools has been less controversial.) Researchers employed by the Dallas Independent School District were among the first to create measures similar to the Tennessee value-added model, with the backing of a local school board member named Sandy Kress. When Gov. Bush became president in 2001, he brought Kress to Washington, D.C., as his chief education adviser.

Kress dived into the 2001 reauthorization of ESEA and was enthusiastic about value-added data and the potential of measuring growth. But he knew that most states were far behind Texas and Tennessee in

developing the annual tests and data systems on which growth models depend. “It became clear that it was not viable at the time because it was so ahead of common usage,” Kress said recently.³

Growth models had a political problem as well. “The civil rights community had concerns about it,” Kress said, “and wanted to make sure that all students were held to the same expectations.” Advocates for the rights of traditionally underserved children were concerned that schools would be judged as high-performing (and therefore not be held accountable for helping low-performing students) as long as academically deficient low-income and minority students made a year’s worth of growth—even if they never actually caught up and achieved proficiency in math and reading. Growth models, they feared, could institutionalize what President Bush memorably described as “the soft bigotry of low expectations.”

The final version of the law, No Child Left Behind, held schools almost exclusively accountable for absolute levels of student performance—the percentage who passed state standardized tests. In a small concession to growth, low-performing schools could escape potential sanctions if the percentage of students who failed the test in a given grade declined enough relative to the percentage of students who had failed the test in the same grade in the previous year. This so-called “cohort” growth measure—this year’s fourth-graders compared to last year’s fourth-graders, for example—was distinct from, and arguably inferior to, growth models that tracked the progress of *the same* students from year to year. Individual classes of students vary in aptitude and myriad other factors, making valid comparisons among them statistically tricky. But most states didn’t have the testing and data infrastructure to calculate anything else.

NCLB passed Congress with broad bipartisan support, and President Bush signed it into law in 2002. But it wasn’t long before good feelings about the law began to evaporate, and the lack of a true growth model played a significant role. Educators felt it was inherently unfair to label a school that had made great strides with low-performing students as “failing” just because the students had not yet made it all the way to a “proficient” level of achievement. Support for NCLB among parents and influential policymakers began to decline, and major interest

groups such as the National Education Association, the nation’s largest teachers union, called for it to be revised or repealed.

Worried that its signature domestic policy initiative was faltering, the Bush administration moved to incorporate more growth measures into state accountability systems. In 2005, U.S. Secretary of Education Margaret Spellings announced that states would be allowed to apply for permission to incorporate growth models into their accountability systems. There was a catch, however. States couldn’t use just *any* growth model. The proposed models would be evaluated by a group of education experts to ensure that they met certain strict criteria. The most important was a concept called “growth to proficiency.”

Growth models, they feared, could institutionalize what President Bush memorably described as “the soft bigotry of low expectations.”

The NCLB accountability model was based on tests tied to academic standards—“criterion-referenced” tests, in education-speak. In such a system, the government decides that students need to know some things—how to factor polynomials, that World War I ended in 1918—and administers a test of such knowledge and skills. The passing score, or “proficiency” level, indicates whether students had learned enough. This was a change from the common practice in states of using so-called “norm-referenced” tests, which indicated where students stood *relative to one another*. The widely used Stanford 10 test, for example, yields scores in percentiles. An 80th percentile score means the tester knows more than four out of five other students. It doesn’t indicate whether they know the year World War I drew to a close.

Supporters of criterion-referenced tests were leery of the relativity inherent to norm-referenced scores. Certain things had to be learned, they believed, irrespective of what other students know. And growth

models were just another kind of relativity. Instead of showing where students stood relative to other students, like the Stanford 10, growth models showed where students stood relative to themselves at an earlier time. This left open the question of how much growth was sufficient to label a student—and thus, his or her school—a failure or a success. This question of how to *interpret* growth measures, as opposed to merely calculate them—to decide how much growth is *enough* growth—would come to dominate the growth model debate.

Secretary Spellings decided that the accountability system had to remain anchored to a criterion-referenced proficiency measure. Therefore, states were only allowed to interpret growth as enough growth if they could show that underperforming students were on track to become proficient within a relatively short time period—three or four years. Critics of NCLB asserted that many schools were being unfairly labeled as failures despite achieving phenomenal growth. The growth model pilot projects would put that assertion to the test.

LEARNING FROM THE PILOTS: HOW MUCH GROWTH IS ENOUGH?

Since 2005, 15 states have been approved to implement a growth model pilot. They have adopted four distinct models, each with virtues and drawbacks.

The simplest and most common strategy is the “Trajectory” model employed by Alaska, Arizona, Arkansas, Florida, Missouri, and North Carolina. Using the U.S. Department of Education’s growth model pilot restrictions as a guide, these states examine the growth in test scores for individual students and calculate the achievement level each student would reach in the future if his growth continued at the same pace that occurred in the most recent year. If this linear trajectory leads to proficiency within the three- or four-year window, the student is deemed to have made enough growth that year.

Table 1. Four Types of Growth Models Under the Federal Pilot Program

Growth Model	States Using Model	How It Works
Trajectory	Alaska, Arizona, Arkansas, Florida, Missouri, and North Carolina	First a state determines the gap between a student’s current achievement level and proficient. Then a student must close a portion of that gap each year over a three- or four-year period. The simplest trajectory model is a linear trajectory. In Florida, for example, a student makes enough growth (“adequate yearly growth” or AYG) if the student closes one third of the gap each year. Some states require the gap to be closed over four years.
Transition Tables	Delaware, Iowa, Michigan, and Minnesota	States have several achievement categories below the proficiency level. In Iowa, for example, a student can score weak, low marginal, or high marginal. A student is determined to have made AYG if he or she moves up at least one category (e.g., from weak to low marginal or from low marginal to high marginal).
Student Growth Percentiles	Colorado (never implemented federal pilot) and Pennsylvania	A student’s year-to-year growth is compared to other students with similar test scores in past years. The amount of growth that a student made is converted into a percentile (from 0 to 100). The state then figures out whether students in the past at similar growth percentiles were able to make it to the state’s proficiency target within the next three years. So students whose growth percentiles are high enough are deemed on the track to proficient and have passed AYG.
Projection	Ohio, Tennessee, and Texas	Through a complex statistical analysis, the state develops a “projection” or prediction for each student based on how students with similar achievement patterns have done in the past. If the model predicts that a student with similar achievement in the past reached the state’s proficiency level within a three-year period, then the student is deemed to be on track to proficiency and makes AYG.

The second strategy, used by Delaware, Iowa, Michigan, and Minnesota, employs “Transition Tables” that identify certain thresholds of achievement below the “proficient” level. In Iowa, for example, non-proficient students can score, in ascending order, as “weak,” “low marginal,” or “high marginal.” If a student crosses one of these thresholds—moving from “weak” to “low marginal,” for example—he has made enough growth. Delaware has a more complicated system. There are four achievement levels below proficiency: 1A, 1B, 2A, and 2B. Schools get a certain number of points depending on how many thresholds each student crosses in a year: 150 points for moving from 1A to 1B and 225 points for moving from 1A to 2A, for example, and 300 points for proficiency. Students in Delaware schools must achieve a certain average point value for their schools to make “adequate yearly progress,” or AYP, under NCLB.

The third strategy, proposed by Colorado and Pennsylvania, was the most relativistic of the four. The “Student Growth Percentiles” model starts with a norm-referenced measure of growth, converting student growth measures to percentiles. The state then identifies the growth percentiles that, in the past, were high enough such that students were likely to become proficient within three or four years. Students who meet or exceed that growth percentile are deemed to have made enough growth.

The fourth and most sophisticated growth model, “Projection,” was used by three states that had made major investments in testing and data systems over the last two decades: Ohio, Tennessee, and Texas. Taking advantage of their sophisticated student data systems, these states were able to create models that use multiple years of past achievement data—not just for the individual students in question but for whole cohorts of similar students—to make a more accurate prediction of how individual students were likely to score in the future. Some projections, for example, use “hierarchical linear modeling,” an advanced technique that accounts for statistical effects occurring at multiple levels of aggregation (e.g., classrooms, schools, and districts) in predicting future student achievement.

The growth model pilots were implemented over the course of several years. Test scores were tallied, growth rates estimated, and new school achievement

levels calculated. In 2010, the U.S. Department of Education published a report designed to answer the question of how much growth models had changed NCLB.⁴

The answer: not much.

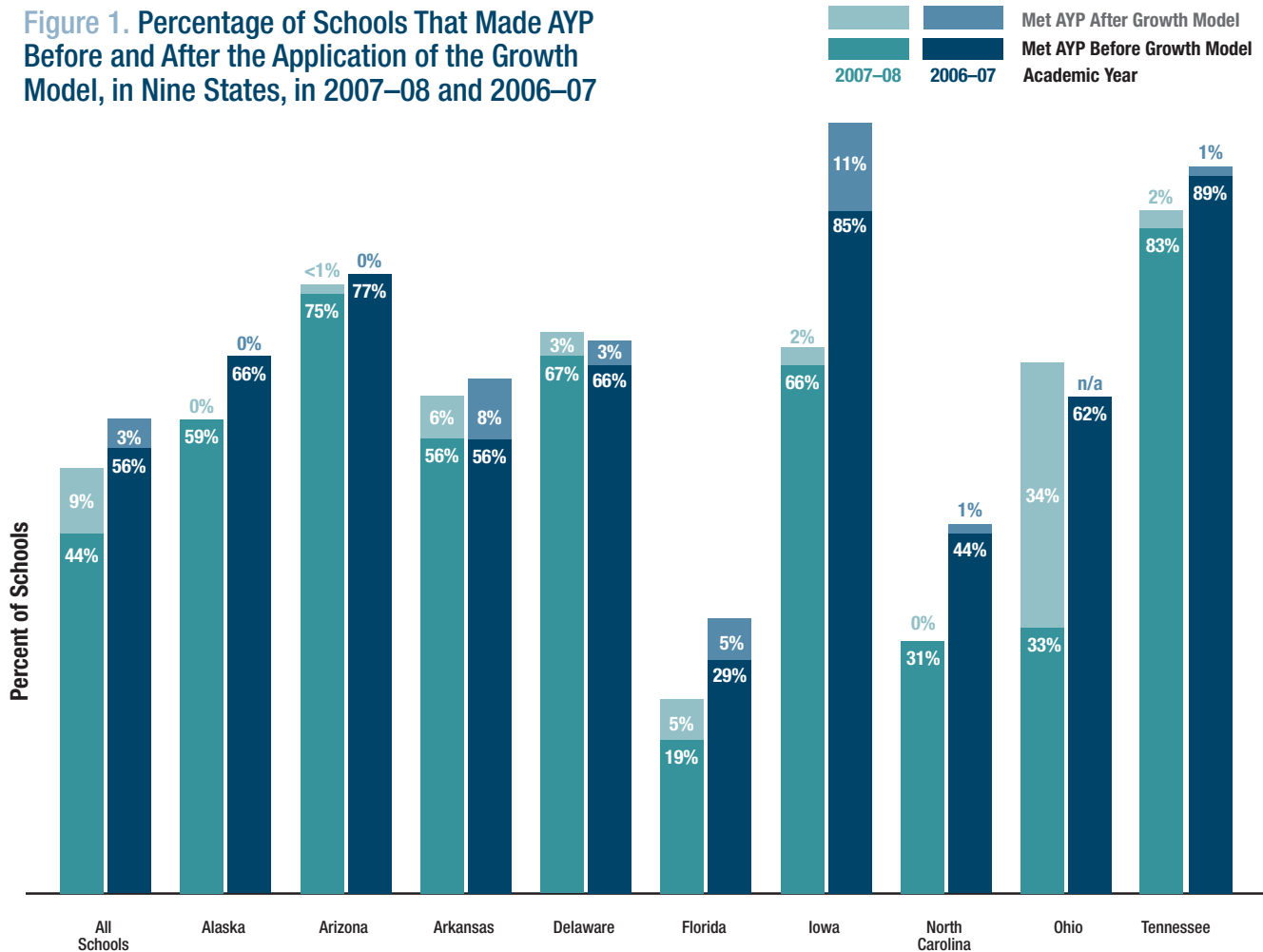
Figure 1 shows the percentage of schools making AYP in the nine states that had implemented growth model pilots in the 2007-08 school year. (Six more states were approved in subsequent years.) For each state, the darker bars show the percentage of schools making AYP under the original provisions of NCLB. The lighter bars show the percentage of *additional* schools making AYP due to their salutary levels of growth. On average, 56 percent of schools made AYP under the old model in the 2006-07 school year. The growth model pilots increased that amount by only 3 percentage points. The difference was larger, but still modest, in 2007-08: 44 percent under the old system, 53 percent after adding growth.

There were a number of reasons that the growth model pilots had little effect on AYP. Tennessee had the most sophisticated growth model in the country. But it also had unusually lax academic standards—the criterion against which student proficiency and school performance were judged.⁵ When the results of the growth model pilots were tallied, 89 percent of Tennessee schools were *already* making the grade under the traditional NCLB system. That left few schools—19, to be exact—to benefit from the growth model pilot. The percentage of schools making AYP in Tennessee rose by only a single percentage point in 2006-07 and two points in 2007-08.

Other states, like neighboring North Carolina, had much tougher standards than Tennessee. Only 44 percent of schools made the grade in the Tar Heel State in 2006-07. Yet the growth model pilot increased that amount to just 45 percent. The reason? First, it turned out that a lot of schools that were bad at helping students reach the proficiency bar were also bad at helping students grow. They were just bad all around. Only 8 percent of students in North Carolina were found to be below proficiency but on track to get there. Thirty-seven percent, by contrast, were neither proficient nor on-track.⁶

Second, the three- to four-year time window mandated by the U.S. Department of Education

Figure 1. Percentage of Schools That Made AYP Before and After the Application of the Growth Model, in Nine States, in 2007–08 and 2006–07



Source: U.S. Department of Education, Office of Planning, Evaluation and Policy Development, Policy and Program Studies Service, *Final Report on the Evaluation of the Growth Model Pilot Project* (Washington, D.C., U.S. Department of Education, 2011), Exhibit ES.1.

meant that only a lot of growth was enough growth. A student who was three years behind—not uncommon in the middle and high school grades—would have to make double the normal progress every year. Such progress is unusual, particularly for multiple years in a row.

The more sophisticated projection models presented additional complications. Forecasting the future means analyzing the past. If the last 20 low-pressure systems in January yielded an average snowfall of 8 inches, meteorologists act accordingly when forecasting the 21st. States like Tennessee and Texas looked at vast amounts of old student achievement data to predict the most likely future path of students who demonstrated certain patterns of achievement. This brought a welcome dose of realism to the exercise of deciding how much growth was enough growth.

If 95 percent of past students with similarly low test scores ultimately failed to learn what they needed to know, it's a safe bet that future such students will probably fail too. Merely hoping otherwise is not a plan.

Yet, there's a conundrum at the heart of such forecasting. Projection models are based on test-score data from an education system in which, despite significant progress in some grades and subjects over time, low-income, minority, and other at-risk students continue to fall short in large numbers. Fixing this national crisis was the point of passing NCLB in the first place. The premise of the law, therefore, was that the problem can be fixed, that humans and human institutions are not immutable forces like the weather but fundamentally changeable things.

If 95 percent of past students with similarly low test scores ultimately failed to learn what they needed to know, it's a safe bet that future such students will probably fail too. Merely hoping otherwise is not a plan.

But the statistical analysis used in the forecasting models made no such allowances. It treated years of low achievement like a nasty storm front. Even when previously failing students in states like Tennessee made unexpected upward progress, the models tended to treat those numbers like a statistical blip, an outlier likely to regress to the mean. Schools in the present weren't given credit for new progress that schools in the past had been unable to maintain, and so the number of Tennessee schools making AYP under the growth model pilot barely budged.

Texas, another projection model state, encountered the flip side of this problem. Just as winters in Boston tend to yield nor'easters, meteorologists in San Antonio tend to be on safe ground when predicting dry sunny heat in July. For each student, the Texas projection model combined scores on a given test with scores from the same student on other standardized tests (e.g., reading, writing, and math), along with scores of other students at the same school, to predict whether students who failed the test in one year would pass the test in the next year. If the model predicted such success in the future, the student was deemed to have passed in the present, even if he failed.

Texas' state-specific, non-NCLB accountability system was also unusually generous in the way it interpreted the results. Students who passed a test in Texas but were statistically predicted to fail in the future weren't counted as failures, unlike states that used projection models to discount both positive and negative deviations from past trends. Texas also left open the possibility of kicking the proficiency can all the way to the end of the road. If an elementary

school is given credit because a non-proficient third-grader is projected to be proficient in the sixth grade, and he or she doesn't actually make it, there is no retroactive penalty. The same is true for sixth to ninth grade, and so on.

As a result, hundreds of districts and thousands of schools across Texas improved their ratings under the Texas school accountability system.⁷ (Like a number of other states, Texas maintains two distinct K-12 accountability systems, one mandated by NCLB and another specific to Texas.) During contentious 2010 legislative hearings about the growth model, Texas state legislator Scott Hochberg noted that a student could be deemed as "passing" the state's writing test even if he got *every single question on the test wrong*, as long as his scores in reading and math were high enough.⁸ The State Department of Education responded by showing that, statistically speaking, schools that had been given the statistical benefit of the doubt deserved it nearly all of the time. Most of those predicted to succeed in the future, did.⁹

Ohio, meanwhile, joined the growth model pilot in 2007-08 after not participating in 2006-07. This had the effect of more than doubling the number of Ohio schools making AYP from the number who would have under normal NCLB rules (33 percent to 77 percent), a result that was far different from the other eight participants studied and was substantially responsible for the increase in the percentage of schools affected by the growth model pilot between 2006-07 and 2007-8. This was not because Ohio had an unusually large number of low-proficiency, fast-growing schools. Instead, like Texas, Ohio chose to interpret growth results in an unusually lenient way. Ohio added the equivalent of two standard deviations of performance to each student's score to determine whether students were on track to reach proficiency and based school ratings on these "augmented predictions." Such artificial augmentations of actual students' scores have been used by other states to manipulate the interpretation of regular, proficiency-based NCLB ratings.¹⁰

In sum, the growth model pilot system implemented in 2005 provided numerous examples of how growth models could fail to meaningfully change NCLB-style accountability systems. More than anything, they highlighted how the public policy questions around growth models are less an issue of measurement than

interpretation of measurement. The question, it turns out, is not how much growth a student has made, or is likely to make, or even how much growth is “enough.” The real question is how should growth of any kind be interpreted in a way that will plausibly lead to more growth? A good place to begin answering that question is Colorado, home of Bruce Randolph School and the one state that successfully applied for a growth model pilot only to change its mind.

COMBINING GROWTH AND PROFICIENCY: THE COLORADO MODEL

Colorado had planned to use a “Student Growth Percentiles” model, which combines elements of the “Projection” and “Trajectory” models by determining whether a student’s relative level of growth matched historical patterns of students who successfully grew toward proficiency within a certain amount of time. But when education officials there saw the results come in from other states, they realized that it wasn’t worth the effort—the new system would likely identify almost exactly the same schools as the old system. So Colorado officials scrapped their growth model pilot plan and focused on creating a state-specific accountability system that puts a premium on communicating information to the public and making meaningful distinctions between different kinds of schools.

Figure 2 shows 2010 performance results for the 182 public elementary, middle, and high schools in Denver. Each circle is a school. The circles are proportional to school size: The more students, the larger the circle. The vertical axis on Figure 2 shows the percentage of students who scored “proficient” on the state standardized test, the standard NCLB metric. The horizontal axis shows the “median student growth percentile.” That means that the Colorado Department of Education calculated how much growth each student made in math since the previous year. They compared that growth to other students with similar academic performance histories, yielding a percentile for each student. The horizontal axis on Figure 2 is the median such percentile for all students in a school.

One of the advantages of the Colorado system is that it provides more information than simple indicators

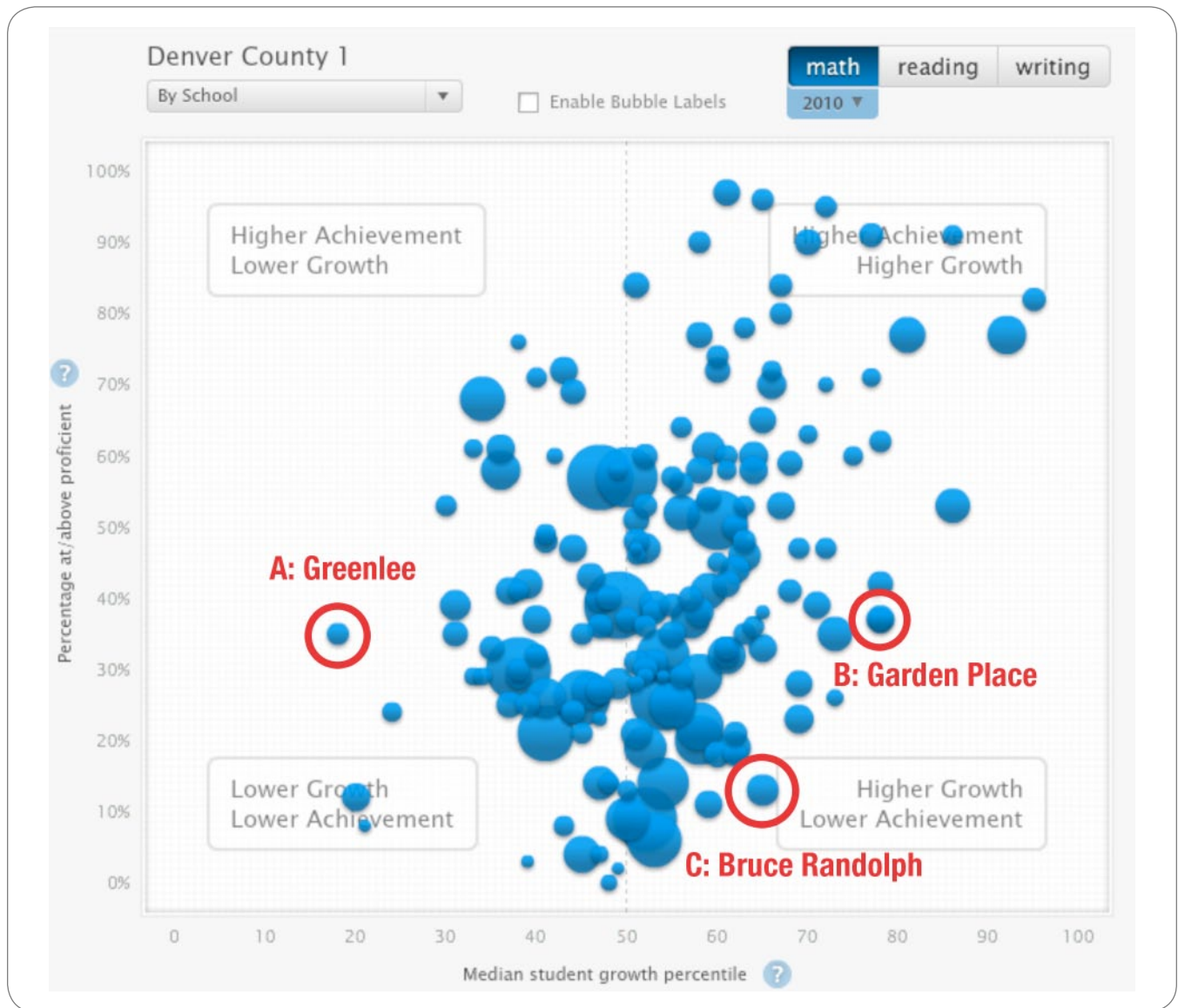
of whether a school has reached a certain threshold of performance. Under NCLB, the percentage of students who need to be proficient in order to make AYP rises steadily until it reaches 100 percent in 2014. In other words, the threshold of “enough” proficiency changes regularly over time. Under the growth model pilots, states also focused on a threshold: whether growth was enough to reach proficiency within three or four years.

Thresholds have the advantage of being decisive, but they also carry the disadvantage of discarding useful information. It matters whether a school is far above or below a given proficiency level, as opposed to near the margin. It also matters whether student growth is above, below, or equal to the growth achieved by similar students.

Figure 2 illustrates the magnitude of these distinctions. There is a visible trend within the 182 schools, sloping upward from the lower left-hand quadrant (low proficiency, low growth) to the upper right-hand quadrant (high proficiency, high growth). This pattern is common when growth and proficiency are plotted together and suggests that proficiency and growth are not independent of one another. Schools that achieve proficiency also tend to achieve growth; schools that fail to achieve proficiency also tend to fail to achieve growth. This correlation is one of the reasons that states utilizing the NCLB growth model pilots failed to identify large numbers of additional schools as good enough.

But there is still much meaningful variation to be found. At Greenlee K-8 Elementary School, the circle labeled “A” in the lower-left quadrant, 35 percent of students scored proficient in mathematics, well below the state threshold of 70 percent. What’s worse, only 18 percent of students displayed growth higher than was typical among similar students. Greenlee displayed a similar low / low combination on tests of reading and writing. Such low-proficiency, low-growth schools are prime candidates for the kind of aggressive “turnaround” interventions currently being championed by the Obama administration. And indeed, Greenlee is one of the bottom 5 percent of all schools, as identified by the U.S. Department of Education’s School Improvement Grant (SIG) program. Greenlee will receive SIG funding to implement a comprehensive turnaround strategy.¹¹

Figure 2. Denver School Performance—2010



Source: https://edx.cde.state.co.us/growth_model/public/index.htm#/year-2010, accessed May 3, 2011.

Garden Place Elementary School, the “B” circle in the lower-right quadrant, has a 37 percent proficiency level in math, almost the same as Greenlee. But 78 percent of students displayed growth higher than is typical for similar students. According to the state, those students are on pace to catch up and become proficient over time. Not coincidentally, Garden Place is not a bottom 5 percent SIG school and is not in the process of being aggressively reconstituted. This illustrates the value of adding growth information to proficiency information when considering school performance.

Interpreting growth information can be more complicated, however, for other kinds of schools. At the school labeled “C” in the lower-right quadrant, 65 percent of students had growth above the median for similar students. Only one other high school in the city, the Denver School of Science and Technology, had better growth scores.

But according to the same Colorado Department of Education data, students at the “C” school aren’t growing fast enough. Only if the typical student there was in the 99th percentile of growth would the growth

rate be enough for that student to catch up and reach proficiency before finishing his or her education. Only 13 percent are proficient in math, and high school ends for them in a few short years. It would take astronomical levels of growth for students to end up where they need to be.

The “C” school is Bruce Randolph, the school singled out by President Obama in his State of the Union. The leader of the free world thinks that Bruce Randolph is a model for the nation to follow. The Colorado Department of Education says its growth is unusually good, but still not good enough.

How should we understand schools like Bruce Randolph? And having understood them, what should we do? These are the questions that Congress must answer as it reauthorizes ESEA.

REMAKING ESEA: HOW TO ACCOMMODATE GROWTH

There are a number of specific challenges to confront and opportunities to take in remaking ESEA to accommodate and promote growth models.

Better Testing

Growth models, like all test-based measurement systems, are only as good as the test on which they rely. Many of the standardized tests used in K-12 education are inadequate, and their flaws can be magnified by growth calculations. A key problem is the scope of what tests try to assess. Unsurprisingly, tests designed for eighth-graders focus on eighth-grade standards and eighth-grade skills. But a significant number of eighth-graders aren’t learning at that level, or even close to it. For students who are far below grade level, grade-level tests can provide little or no information about *how far* below. An eighth-grader reading at the second-grade level could get exactly the same test result as a classmate reading at the fourth-grade level: the worst possible score. The second-grade reader needs to grow faster than the fourth-grade reader to catch up, but growth models using tests of eighth-grade standards might not recognize the distinction. The same problem exists on the other end of the achievement spectrum, among students who score far above the norm.

The key is to broaden the range of achievement that tests can detect. Two assessment consortia—the Partnership for the Assessment of Readiness for College and Careers (PARCC) and the SMARTER Balanced Assessment Consortium (SBAC)—are currently working to design tests aligned to the Common Core State Standards developed by a consortia of governors and nonprofit organizations.

The key is to broaden the range of achievement that tests can detect.

Both assessment consortia are working to accommodate the needs of growth models. The SBAC, for example, is tackling this challenge by using so-called “adaptive testing.” When a group of students sits down for traditional tests, like the SAT, they all take the same test with the same questions. Adaptive tests, by contrast, are administered via computers that change the questions students are given based on their answers to previous questions. Students who get questions right are given progressively more difficult problems to solve, while students who get questions wrong are given easier questions in turn. In this way, adaptive tests extend the scope of knowledge and skills assessed beyond a single grade range, providing more useful information for growth models.

The consortia are also developing interim tests that will be given during the middle of the school year. The interim SBAC test will use the same scale as tests given at the end of the year, allowing states to potentially give schools more fine-grained, actionable estimates of ongoing student growth during the year.¹² In addition, the consortia are considering aligning their testing scales *across* grades. It’s possible to estimate growth between non-aligned tests through statistical correlation (e.g., if students who get a 600 on the SAT math exam in their junior year are most likely to score a “3” on the Advanced Placement (AP) calculus exam in their senior year, a 600-scoring student who gets a “5” on the AP exam may be inferred to have made an unusual amount of growth in math.) But aligned tests may allow for more accurate estimates

of growth over time.¹³ If a score of 700 denotes grade-level proficiency on the seventh-grade math test, for example, a score of 700 on the eighth-grade math test would indicate that the student is one year behind.

More Grades

Growth models also need to incorporate information about students beyond what is mandated by NCLB—testing in grades three through eight and once in high school. Growth calculations require at least two points of time to compare. That means that under the current testing regime, growth can't be calculated until grade four, because there is no grade two test to use as a baseline. For elementary schools that go up to grades four through six, this could create perverse incentives to neglect grades K-three, resulting in low achievement, and a concentration of resources in grades four through six, where growth would be measured. Because standardized tests are less accurate and less developmentally appropriate in the early grades, this problem can't be solved by simply extending the testing window all the way down to toddlers. Instead, ESEA should require states to use multiple measures to evaluate elementary school quality, such as inspections by trained observers and the observation-based Classroom Assessment Scoring System (CLASS) developed at the University of Virginia.

ESEA should require states to use multiple measures to evaluate elementary school quality....

The same problem exists for older students. NCLB mandates only one accountability test in high school, which is typically given in grade 10 and doesn't address advanced secondary subjects like chemistry, calculus, history, economics, and other courses needed to properly prepare students for college and careers. States designing accountability systems should be required to administer a 12th-grade test in reading and math, as well as include results from standardized "end-of-subject" tests that

states are increasingly requiring students to pass in order to graduate from high school. States should also incorporate information about what happens to students after they finish high school. The best way to know if a student has been adequately prepared to succeed in college is to see if he or she actually succeeds in college. States like Florida that have linked their K-12 and higher education data systems can extend their growth model projections across the administrative and conceptual chasm that often separates high school and college.¹⁴ Data about college enrollment, first-year retention, college grades, and student placement in remedial courses can be used to assess whether student growth in high school is enough growth.

Tougher Standards

The Common Core standards were designed to identify what students need to know and be able to do in order to succeed in college and careers. While the new tests that will assess student mastery of the common standards are still being developed by the two consortia, it is widely expected that they will be more rigorous and difficult than what is typical among states today. In 2009, for example, 79 percent of Alabama fourth-graders scored as "proficient" in mathematics, based on Alabama standards and the Alabama test.¹⁵ In the same year, only 25 percent of Alabama fourth-graders scored proficient in math on the U.S. Department of Education-administered National Assessment of Educational Progress. Like most states, Alabama's adoption of the Common Core standards and related tests will result in fewer students making the grade.

ACT, publisher of the widely used college admission test of the same name, recently conducted an analysis of how high school students might fare on a test based on the Common Core standards.¹⁶ By matching ACT test questions to similar elements of the Common Core and examining hundreds of thousands of actual ACT test results, researchers found that only about one-third to one-half of 11th-graders were college- and career-ready in reading, writing, and math, as defined by the Common Core. Passing rates for minority students were substantially worse.

Raising state standards to meet the level of rigor established by the Common Core will increase the

...researchers found that only about one-third to one-half of 11th-graders were college- and career-ready in reading, writing, and math, as defined by the Common Core.

challenge of balancing growth and proficiency. Fewer students will score as proficient and the growth trajectory of underperforming students toward proficiency will be even steeper. Projection models will deem more students unlikely to succeed. Schools like Bruce Randolph will have to do even better in order to achieve adequate growth.

Different Models for Different Things

The growth model pilot program launched by the U.S. Department of Education in 2005 has provided valuable information about growth models. The experiences of the pilot states show the consequences of different approaches to identifying schools as making enough growth. More importantly, they demonstrate that “How much growth is enough growth?” is a necessary but insufficient question to ask. Growth model information is only useful if interpreted along with other perspectives on student success. And the way that information should be used depends on what it is to be used for.

One of those uses is public information. The major innovation of the Colorado growth model is not the method of estimating growth or judging whether growth is enough. Colorado stands out for the ease with which policymakers, principals, school board members, parents, and other stakeholders can access the information. The charts, available on the Internet, allow people to examine growth for subgroups of students—low-income, minority, English language learners, students with disabilities—to see how traditionally disadvantaged students are performing relative to their peers. The colorful arrays of circles show exactly where each school stands compared to all others. Colorado developed its system with

open-source software so other states could quickly and inexpensively present their growth data in a similar way. As of early 2011, 14 states had formed a consortium to do exactly that, at a per-state cost of as little as \$250,000.

Growth model information is also used to make specific policy choices: Should a school be identified as failing? Should it be given more money? Should it be forced to reorganize or reform? Again, the type of decision dictates the type of model. As researchers like Harvard University’s Andrew Ho have demonstrated, the “Projection” and “Trajectory” models can yield radically different results when used for accountability purposes.¹⁷ Both are valuable, but only when matched to the right perspective and the right use.

“Projection” models make sense from the 40,000-foot perspective. In the aggregate, past is often prologue. This is particularly true in public education, which has proved to be remarkably immune to external shocks, both positive and negative, over the years. State policymakers deciding how to distribute funding among school districts or where to concentrate intensive reform efforts should take projections based on long-term statistical trends very seriously. If the projections strongly suggest that students in a distressed urban district are collectively not on track to reach proficiency, that information should be treated with deadly seriousness and acted on accordingly.

As the perspective narrows and descends, however, projections have less value. There is an aspiration at the center of public education. It’s a bet on human potential, an idea that institutions are improvable, and a faith that the best means of helping young people learn and grow into their fullest selves have yet to be discovered. It is crucial that these ideas aren’t crushed by the weight of aggregate statistics. The only reason to hold schools accountable for student learning is to improve student learning. That won’t happen if accountability systems presuppose that such improvement can never occur.

Growth information about individual teachers and students should be considered with a particular sensitivity to the fact that all risks are not equal. Imagine a parent receiving her child’s standardized writing test results in the mail. The scores are horrific,

There is an aspiration at the center of public education. It's a bet on human potential, an idea that institutions are improvable, and a faith that the best means of helping young people learn and grow into their fullest selves have yet to be discovered.

the worst possible. But imagine there's a note attached, from a statistician at the State Department of Education. "Our records indicate that other students like your daughter did much better on the writing test," the hypothetical note says. "And your daughter did much better on tests of things other than writing. As a result, our statistical model predicts that your daughter will do better in writing next year. So don't worry, nothing to see here."

Would a reasonable parent cancel the upcoming parent-teacher conference and stop checking written homework assignments? Probably not. Parents would be sensitive to the risk that their child is an exception to the statistical trend and in danger of an academic catastrophe. The risk of failing to intervene on behalf of that child is substantial when weighed against the risk of providing extra assistance to a child who is actually fine.

Similarly, teachers who achieve unusually good results in a given year should be recognized and rewarded for their success. Any instance of deviation from the educational norm might be a statistical anomaly—but it might not. And the aspirational education idea depends on believing that such successes can be learned from, replicated, and spread out to the world at large. Teachers who experience an unusual drop, by contrast, should be given the benefit of the doubt, with multiple consecutive years of failure or success given steadily more weight, and other factors like expert observation and peer review taken into account. Teachers' rights, reputations, and livelihoods are important, and the risk of damaging them unnecessarily should be minimized, even as

schools also weigh the countervailing risk of assigning vulnerable students to poor-performing teachers.

At the school level, Colorado's two-dimensional combination of proficiency and growth strikes a reasonable balance. Greenlee and Garden Place elementary schools aren't the same and shouldn't be treated as such. The only responsible action on behalf of young children trapped in a low-proficiency, low-growth school like Greenlee is immediate, radical change. Students at Garden Place also need more help than most. But schools that can achieve unusual growth with disadvantaged students are hard to come by. The best strategy in such schools is often to invest in doing more of what is making them successful, not doing something else with different people.

Making such distinctions isn't always easy. Most elementary schools in Denver are in between Greenlee and Garden Place. They are simply ordinary when it comes to growth. No single mathematical formula can adequately capture all the distinctions among schools. Add the odds of graduating from high school and enrolling in college to the model and Colorado's neat two-dimensional array becomes a three-dimensional cube. Add the odds of *succeeding* in college and the data inhabit four dimensions, beyond visualization. Layer on achievement differences among different student groups, and the complexity level shoots up like a rocket. It will take a strong dose of wise human judgment among state and federal policymakers to synthesize this information and decide how and where to intervene.

The next class of students at Bruce Randolph School might not graduate in such high numbers. The last class may have trouble in college, victim of educational failures that occurred long before their high school teachers ever knew their names. But in 2010, the school did something extraordinary. That is worth understanding.

Notes

1. David Grissmer and Ann Flanagan, *Exploring Rapid Achievement Gains in North Carolina and Texas*, (Washington, D.C.: National Education Goals Panel, 1998).
2. "Grading the Teachers: Value Added Analysis," Los Angeles Times, <http://www.latimes.com/news/local/teachers-investigation/>

3. Sandy Kress, in discussion with Robert Manwaring, 2010.
4. U.S. Department of Education, Office of Planning, Evaluation and Policy Development, Policy and Program Studies Service, *Final Report on the Evaluation of the Growth Model Pilot Project* (Washington, D.C., U.S. Department of Education, 2011). The department released a preliminary version of the report in 2010.
5. Charles Barone, *Are We There Yet? What Policymakers Can Learn from Tennessee's Growth Model*, (Washington, D.C.: Education Sector, March 2009).
6. *Final Report on the Evaluation of the Growth Model Pilot Project*, Exhibit 29.
7. Jeffrey Weiss, "The Scores Are in, and the Texas Projection Measure Results Are Mixed," *Dallas Morning News*, August 3, 2010.
8. Hochberg also noted that the projections weren't technically based on "growth" at all, in the sense of a linear sequence of scores over time, but rather just a set of test scores associated with a certain student and group of similar students with no attention given (statistically speaking) to the sequence in which those scores occurred.
9. Consistent with these findings, the *Final Report on the Evaluation of the Growth Model Pilot Project* found that "[The] projection model has the highest correct classification rates for future proficiency: over 80 percent. These rates are 5 to 20 percentage points higher than trajectory and transition matrix models, depending on the grade level and proximity to the growth model time limit."
10. See for example Kevin Carey, *The Pangloss Index: How States Game the No Child Left Behind Act*, (Washington, D.C.: Education Sector, November 2007).
11. Padmini Jambulapati, *A Portrait of School Improvement Grantees* (Washington, D.C.: Education Sector, April 2011).
12. The PARCC approach is somewhat different; they are pursuing a "through-course" design where the material currently assessed at the end of the year in a single high-stakes annual test will be broken up into pieces and tested at four points throughout the year, with the results then aggregated into a final score. Since the first three PARCC tests will use "performance tasks" and the fourth will be based on a different "selected response" assessment method, it may be more difficult to calculate growth levels at each point along the way.
13. There is not a consensus in the assessment community as to whether the considerable complication and expense of creating so-called "vertical" alignment of test scales across grades is worth the additional growth estimate accuracy such alignment can potentially create. The PARCC consortia has not formerly committed to creating such alignment.
14. Chad Aldeman, *College- and Career-Ready: Using Outcomes Data to Hold High Schools Accountable for Student Success* (Washington, D.C.: Education Sector, January 2010).
15. <http://educationnext.org/files/ProficiencyData.pdf>
16. *A First Look at the Common Core and College and Career Readiness* (Iowa City, IA: ACT, 2010).
17. Andrew Ho, *Supporting Growth Interpretations Using Through-Course Assessments*, (Austin, TX: Center for K-12 Assessment and Performance Management at ETS, March 2011), http://www.k12center.org/rsc/pdf/TCSA_Symposium_Final_Paper_Ho.pdf