



Getting Teacher Evaluation Right: A Background Paper for Policy Makers

Linda Darling-Hammond
Stanford University

Audrey Amrein-Beardsley
Arizona State University

Edward H. Haertel
Stanford University

Jesse Rothstein
University of California, Berkeley

Research Briefing

Getting Teacher Evaluation Right: A Challenge for Policy Makers

***September 14, 2011
Dirksen Senate Office Building***

Convening Organizations

***American Educational Research Association
National Academy of Education***

This Capitol Hill research briefing was jointly held by the American Educational Research Association and the National Academy of Education as part of these organizations' commitment to the sound use of scientific research and data in education and education policy. The views in this background paper are those of the authors who presented at the briefing and not of the organizations planning and convening this event.

Executive Summary

Consensus that current teacher evaluation systems often do little to help teachers improve or to support personnel decision making has led to a range of new approaches to teacher evaluation. This brief looks at the available research about teacher evaluation strategies and their impacts on teaching and learning.

Prominent among these new approaches are value-added models (VAM) for examining changes in student test scores over time. These models control for prior scores and some student characteristics known to be related to achievement when looking at score gains. When linked to individual teachers, they are sometimes promoted as measuring teacher “effectiveness.”

Drawing this conclusion, however, assumes that student learning is measured well by a given test, is influenced by the teacher alone, and is independent of other aspects of the classroom context. Because these assumptions are problematic, researchers have documented problems with value-added models as measures of teachers’ effectiveness. These include the facts that:

1. Value-Added Models of Teacher Effectiveness Are Highly Unstable: Teachers’ ratings differ substantially from *class to class* and from *year to year*, as well as from one *test* to the next.

2. Teachers’ Value-Added Ratings Are Significantly Affected by Differences in the Students Who Are Assigned to Them: Even when models try to control for prior achievement and student demographic variables, teachers are advantaged or disadvantaged based on the students they teach. In particular, teachers with large numbers of new English learners and others with special needs have been found to show lower gains than the same teachers when they are teaching other students.

3. Value-Added Ratings Cannot Disentangle the Many Influences on Student Progress: Many other home, school, and student factors influence student learning gains, and these matter more than the individual teacher in explaining changes in scores.

Other tools have been found to be more stable. Some have been found both to predict teacher effectiveness and to help improve teachers’ practice. These include:

- Performance assessments for licensure and advanced certification that are based on professional teaching standards, such as National Board Certification and beginning teacher performance assessments in states like California and Connecticut.
- On-the-job evaluation tools that include structured observations, classroom artifacts, analysis of student learning, and frequent feedback based on professional standards.

In addition to the use of well-grounded instruments, research has found benefits of systems that recognize teacher collaboration, which supports greater student learning.

Finally, systems are found to be more effective when they ensure that evaluators are well-trained, evaluation and feedback are frequent, mentoring and coaching are available, and processes, such as Peer Assistance and Review systems, are in place to support due process and timely decision making by an appropriate body.

Getting Teacher Evaluation Right: A Background Paper for Policy Makers

There is a widespread consensus among practitioners, researchers, and policy makers that current teacher evaluation systems in most school districts do little to help teachers improve or to support personnel decision making. For this reason, new approaches to teacher evaluation are being developed and tested.

There is also a growing consensus that evidence of teachers' contributions to student learning should be a component of teacher evaluation systems, along with evidence about the quality of teachers' practice. *Value-added models* (VAMs) for examining gains in student test scores from one year to the next are promoted as tools to accomplish this goal. Policy makers can benefit from research about what these models can and cannot do, as well as from research about the effects of other approaches to teacher evaluation. This background paper addresses both of these important concerns.

Research on Value-Added Models of Teacher "Effectiveness"

Researchers have developed value-added models for examining gains in student achievement by using statistical methods that allow them to measure changes in student scores over time, while taking into account student characteristics and other factors often found to influence achievement. In large-scale studies, these methods have proved valuable for looking at a range of factors affecting achievement and measuring the effects of programs or interventions.¹

When applied to individual teacher evaluation, the use of value-added modeling (VAM) assumes that measured student achievement gain, linked to a specific teacher, reflect that teacher's "effectiveness." Drawing this conclusion, however, assumes that student learning is measured well by a given test, is influenced by the teacher alone, and is independent from the growth of classmates and other aspects of the classroom context.

But research reveals that a student's achievement and measured gains are influenced by much more than any individual teacher. Others factors include:

- school factors such as class sizes, curriculum materials, instructional time, availability of specialists and tutors, and resources for learning (books, computers, science labs, and more);
- home and community supports or challenges;
- individual student needs and abilities, health, and attendance;
- peer culture and achievement;
- prior teachers and schooling, as well as other current teachers;
- differential summer learning loss, which especially affects low-income children; and
- the specific tests used, which emphasize some kinds of learning and not others, and which rarely measure achievement that is well above or below grade level.

Most of these factors are not actually measured in value-added models, and the teacher's effort and skill, while important, constitute a relatively small part of this complex equation. As a

consequence, researchers have documented a number of problems with value-added models as accurate measures of teachers' effectiveness.

1. Value-Added Models of Teacher Effectiveness Are Highly Unstable

Researchers have found that teachers' effectiveness ratings differ substantially from *class to class* and from *year to year*, as well as from one statistical model to the next, as Table 1 shows.²

Table 1. Percent of Teachers Whose Effectiveness Rankings Change

	By 1 or More Deciles	By 2 or More Deciles	By 3 or More Deciles
Across models ^a	56–80%	12–33%	0–14%
Across courses ^b	85–100%	54–92%	39–54%
Across years ^b	74–93%	45–63%	19–41%

^aDepending on pair of models compared.

^bDepending on the model used.

Source: Newton, Darling-Hammond, Haertel, and Thomas (2010).

A study examining data from five separate school districts found, for example, that of teachers who scored in the bottom 20% of rankings in one year, only 20–30% had similar ratings the next year, while 25–45% of these teachers moved to the top part of the distribution, scoring well above average. (See Figure 1.) The same volatility occurred for those who scored at the top of the distribution in one year: A small minority (about 25%) stayed in the same rating band the following year, while most scores moved to other parts of the distribution.

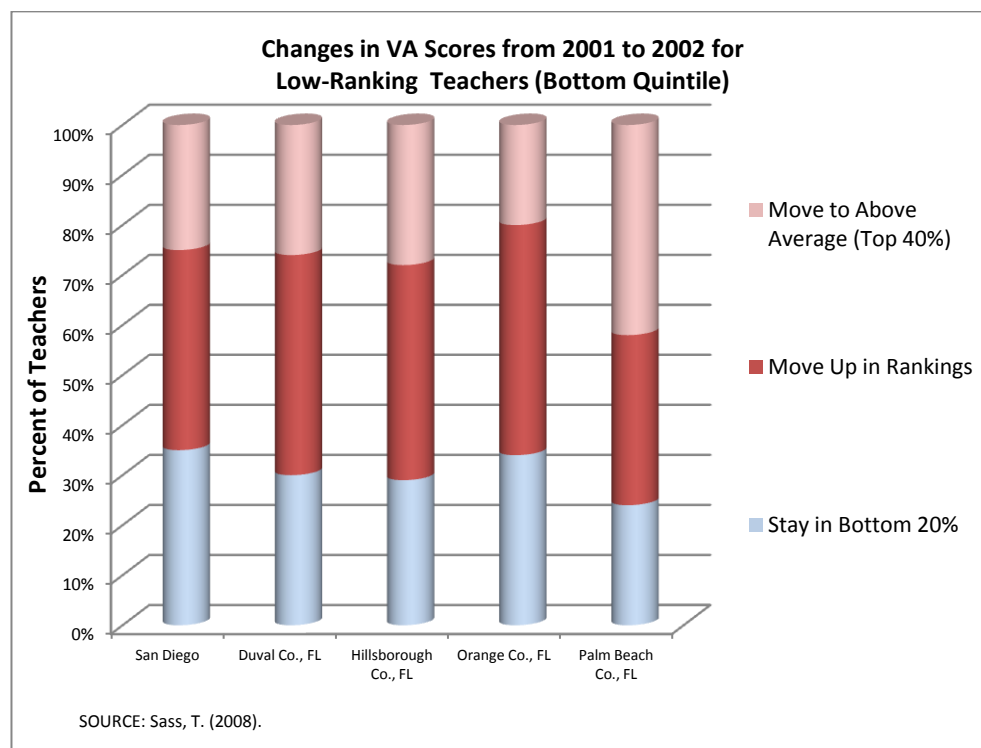


Figure 1. Changes in Value-Added Scores from 2001 to 2002 for Low-Ranking Teachers

Teachers' measured effectiveness varies significantly *when different statistical methods are used*.³ For example, when researchers used a different model to recalculate the value-added scores for teachers that were published in the *Los Angeles Times* in 2011, they found that 40–55% of them would get noticeably different scores using an alternative statistical model that accounted for student assignments in a different way.⁴

Teachers' value-added scores also differ significantly *when different tests are used*, even when these are within the same content area.⁵ For example:

- In a study using two tests measuring basic skills and higher order skills, 20–30% of teachers who ranked in the top quartile in terms of their impacts on state tests ranked in the bottom half of impacts on more conceptually demanding tests (and vice versa).⁶
- Teachers' estimated effectiveness is very different for "Procedures" and "Problem Solving" subscales of the same math test.⁷
- Teacher effects on high-stakes tests are not highly related to their effects on low stakes tests and dissipate more quickly.⁸

This raises concerns both about measurement error and, when teacher evaluation results are tied to student test scores, about the effects of emphasizing "teaching to the test" at the expense of other kinds of learning, especially given the narrowness of most tests currently used in the United States.

2. Teachers' Value-Added Ratings Are Significantly Affected by Differences in the Students Who Are Assigned to Them

VAMs require that students be assigned to teachers randomly. But students are not randomly assigned to teachers. Furthermore, statistical models cannot fully adjust for the fact that some teachers will have a disproportionate number of students who have greater challenges (students with poor attendance, who are homeless, who have severe problems at home, etc.) and those whose scores on traditional tests may not accurately reflect their learning (eg., those who have special education needs or who are new English-language learners). These factors can create both misestimates of teachers' effectiveness and disincentives for teachers to want to teach the students who have the greatest needs.

Even when the model includes controls for prior achievement and student demographic variables, teachers are advantaged or disadvantaged based on the students they teach. Several studies have shown this by conducting tests that look at a teacher's "effects" on their students in grade levels *before* or *after* the grade level in which he or she teaches them. Logically, for example, 5th grade teachers cannot influence their students' 3rd grade test scores. So a VAM that identifies teachers' true effects should show *no* effect of 5th grade teachers on their students' 3rd grade test scores two years earlier. But studies that have looked at this have shown large "effects"—a phenomenon suggesting that other factors associated with students have at least as much bearing on the value-added measure as the teachers who actually teach them in a given year.⁹

One study that found considerable instability in teachers' value-added scores from class to class and year to year examined changes in student characteristics associated with the changes in teacher ratings.¹⁰ After controlling for prior test scores of students *and* student characteristics, the study still found significant correlations between teachers' ratings and their students' race/ethnicity,

income, language background, and parent education. Figure 2 illustrates this finding for an experienced English teacher in the study whose rating went from the very lowest category in one year to the very highest category the next year (a jump from the 1st to the 10th decile). In the second year, this teacher had many fewer English learners, Hispanic students, and low-income students, and more students with well-educated parents, than in the first year.

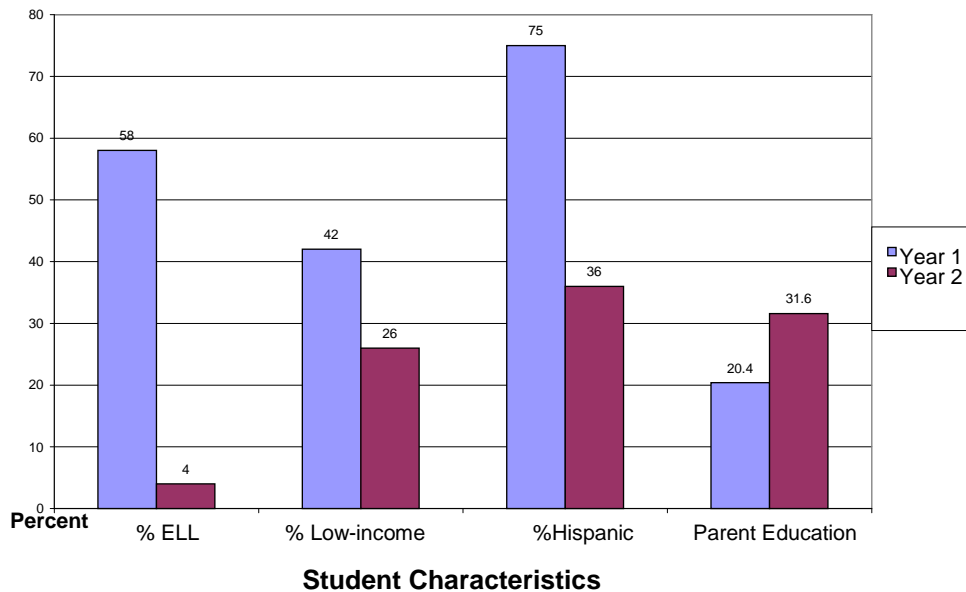


Figure 2. Student Characteristics in Years 1 and 2 for a Teacher Whose Ranking Changed From the 1st to the 10th Decile

This variability raises concerns that use of such ratings for evaluating teachers could create disincentives for teachers to serve high-need students. This could inadvertently reinforce current inequalities, as teachers with options would be well-advised to avoid classrooms or schools serving such students or to seek to prevent such students from being placed in their classes.

3. Value-Added Ratings Cannot Disentangle the Many Influences on Student Progress

It is impossible to fully separate out the influences of students' other teachers, as well as of school conditions, on their reported learning. No single teacher accounts for all of a student's learning. Prior teachers have lasting effects, for good or ill, on students' later learning, and current teachers also interact to produce students' knowledge and skills. For example, the essay-writing skills a student learns through his history teacher may be credited to his English teacher, even if she assigns no writing; the math content and skills he learns in his physics class may be credited to his math teacher. Specific skills and topics taught in one year may not be tested until later, if at all. Some students receive tutoring, as well as help from well-educated parents. A teacher who works in a well-resourced school with specialist supports may appear to be more effective than one whose students do not receive these supports. As noted by Henry Braun, an expert in measurement and evaluation:

It is always possible to produce estimates of what the model designates as teacher effects. These estimates, however, capture the contributions of a number of factors, those due to teachers being only one of them. So treating estimated teacher effects as accurate indicators of teacher effectiveness is problematic.¹¹

Initial research on the use of value-added methods to dismiss some teachers and award bonuses to others shows that value-added ratings often do not agree with the ratings teachers receive from skilled observers and are influenced by all of the factors described above.

For example, among several teachers dismissed in Houston as a result of their value-added test scores, one 10-year veteran had been voted “Teacher of the Month” and “Teacher of the Year” and was rated each year as “exceeding expectations” by her supervisor.¹² She showed positive VA scores on 8 of 16 of tests over four years (50% of the total observations), with wide fluctuations from year to year and both across and within subjects. (See Table 2.) It is worth noting that this teacher’s lower value-added in grade 4, when English learners are mainstreamed in Houston, was a pattern for many other teachers as well.

Table 2. EVAAS Scores by Subject, Grade, and Year for One Teacher

EVAAS Scores	2006–2007	2007–2008	2008–2009	2009–2010
	Grade 5	Grade 4	Grade 3	Grade 3
Math	-2.03	+0.68*	+0.16*	+3.46
Reading	-1.15	-0.96*	+2.03	+1.81
Language Arts	+1.12	-0.49*	-1.77	-0.20*
Science	+2.37	-3.45	n/a	n/a
Social Studies	+0.91*	-2.39	n/a	n/a
ASPIRE Bonus	\$3,400	\$700	\$3,700	\$0

Notes: (1) Scores with asterisks () signify that the scores are not detectably different from the reference gain scores of other teachers across HISD.

Source: Amrein-Beardsley & Collins (forthcoming).

The wide variability shown in this teacher's ratings from year to year, like that documented in many other studies, was not unusual for teachers in Houston. Teachers reported that they could not identify a relationship between their instructional practices and their ratings on VA, which appear unpredictable. As one teacher noted:

I do what I do every year. I teach the way I teach every year. [My] first year got me pats on the back. [My] second year got me kicked in the backside. And for year three my scores were off the charts. I got a huge bonus, and now I am in the top quartile of all the English teachers. What did I do differently? I have no clue.¹³

Another teacher classified her past three years as “bonus, bonus, disaster.” And another noted:

We had an 8th grade teacher, a very good teacher, the “real science guy,” [who was a] very good teacher...[but] every year he showed low EVAAS growth. My principal flipped him with the 6th grade science teacher who was getting the highest EVAAS scores on campus. Huge EVAAS scores. [And] now the 6th grade teacher [is showing] no growth, but the 8th grade teacher who was sent down is getting the biggest bonuses on campus.

This example of two teachers whose value-added ratings flip-flopped when they exchanged assignments is an example of a phenomenon found in other studies that document a larger association between the class taught and value-added ratings than the individual teacher effect itself. The notion that there is a stable “teacher effect” that is a function of the teacher's teaching ability or effectiveness is called into question if the specific class or grade-level assignment is a stronger predictor of the value-added rating than the teacher.

Another teacher who was dismissed, also consistently rated as “exceeding expectations” or

“proficient” by her supervisor, and also receiving positive VA scores about 50% of the time, had a noticeable drop in her value-added ratings when she was assigned to teach a large number of English-language learners who were transitioned into her classroom. Overall, the study found that, in this system:

- Teachers teaching in grades in which English-language learners (ELLs) are transitioned into mainstreamed classrooms are the least likely to show “added value.”
- Teachers teaching larger numbers of special education students in mainstreamed classrooms are also found to have lower value-added scores, on average.
- Teachers teaching gifted students show little value-added because their students are already near the top of the test score range.
- Ratings change considerably when teachers change grade levels, often from “ineffective” to “effective” and vice versa.

The following kinds of comments from teachers were typical:

Every year I have the highest test scores, [and] I have fellow teachers that come up to me when they get their bonuses... One recently came up to me [and] literally cried, “I’m so sorry.” ... I’m like, “Don’t be sorry. It’s not your fault.” Here I am ... with the highest test scores and I’m getting \$0 in bonuses. It makes no sense year to year how this works. You know, I don’t know what to do. I don’t know how to get higher than 100%.

I went to a transition classroom, and now there’s a red flag next to my name. I guess now I’m an ineffective teacher? I keep getting letters from the district, saying “You’ve been recognized as an outstanding teacher” ... this, this, and that. But now because I teach English-language learners who “transition in,” my scores drop? And I get a flag next to my name for not teaching them well?

I’m scared to teach in the 4th grade. I’m scared I might lose my job if I teach in an [ELL] transition grade level, because I’m scared my scores are going to drop, and I’m going to get fired because there’s probably going to be no growth.

It is not surprising, given these findings, that teachers in Houston report seeking to boost their scores by avoiding certain subjects and types of students, and by seeking assignments to teach particular subjects/grades, while being increasingly confused and demoralized by the system.

The long-run implications for teacher recruitment and retention in districts that use such measures has yet to be studied empirically. Nor have those implications been studied for schools and classrooms serving students who appear to negatively influence measured gains.

Professional Consensus About the Use of Value-Added Methods in Teacher Evaluation

For all of these reasons, most researchers have concluded that value-added modeling (VAM) is not appropriate as a primary measure for evaluating individual teachers. A major report by the RAND Corporation concluded that:

The research base is currently insufficient to support the use of VAM for high-stakes decisions about individual teachers or schools.¹⁴

Similarly, Henry Braun concluded in his review of research:

VAM results should not serve as the sole or principal basis for making consequential decisions about teachers. There are many pitfalls to making causal attributions of teacher effectiveness on the basis of the kinds of data

available from typical school districts. We still lack sufficient understanding of how seriously the different technical problems threaten the validity of such interpretations.¹⁵

Finally, the National Research Council's Board on Testing and Assessment concluded that:

VAM estimates of teacher effectiveness that are based on data for a single class of students should not be used to make operational decisions because such estimates are far too unstable to be considered fair or reliable.¹⁶

Other Approaches to Teacher Evaluation

While VAMs based on student test scores are problematic for making evaluation decisions for individual teachers, they are useful for looking at groups of teachers for research purposes, for example, to examine how specific teaching practices or measures of teaching influence the learning of large numbers of students. The larger scale of these studies reduces error, and their frequent use of a wider range of outcome measures allows more understanding of the range of effects of particular strategies or interventions.

These kinds of analyses provide other insights for teacher evaluation since there is a large body of evidence over many decades concerning how specific teaching practices influence student learning gains. For example, there is considerable evidence that effective teachers:

- understand subject matter deeply and flexibly;
- connect what is to be learned to students' prior knowledge and experience;
- create effective scaffolds and supports for learning;
- use instructional strategies that help students draw connections, apply what they are learning, practice new skills, and monitor their own learning;
- assess student learning continuously and adapt teaching to student needs;
- provide clear standards, constant feedback, and opportunities for revising work; and
- develop and effectively manage a collaborative classroom.¹⁷

These aspects of effective teaching, supported by research, have been incorporated into professional standards for teaching that offer some useful approaches to teacher evaluation.

Using Professional Standards for Teacher Evaluation

Professional standards defining accomplished teaching were first developed by the National Board for Professional Teaching Standards to guide assessments for veteran teachers. Subsequently, a group of states working together under the auspices of the Council for Chief State School Officers created the Interstate New Teacher Assessment and Support Consortium (INTASC), which translated these into standards for beginning teachers, adopted by over 40 states for initial teacher licensing. A recent revision of the INTASC teaching standards has been aligned with the Common Core Standards in order to reflect the kind of teacher knowledge, skills, and understandings needed to enact the standards.

These standards have become the basis for assessments of teaching that produce ratings which are much more stable than value-added measures. At the same time, they incorporate classroom evidence of student learning, and they have recently been shown in larger-scale studies to predict teachers' value-added effectiveness. So they help ground evaluation in student learning in more stable ways. Typically the performance assessments ask teachers to document their plans and

teaching for a unit of instruction linked to the state standards, adapt them for special education students and English-language learners, videotape and critique lessons, and collect and evaluate evidence of student learning.

A number of studies have found that the National Board Certification assessment process identifies teachers who are more effective in raising student achievement than other teachers.¹⁸ Equally important, studies have found that teachers' participation in the National Board process stimulates improvements in their practice.¹⁹ Similar performance assessments, used with beginning teachers in Connecticut and California, have been found to predict their students' achievement gains on state tests.²⁰ The Performance Assessment for California Teachers (PACT) has also been found to improve beginning teachers' competence and to stimulate improvements in the teacher education programs that use it as a measure.²¹

Professional standards have also been translated into teacher evaluation instruments in use at the local level. In a study of three districts using standards-based evaluation systems, researchers found significant relationships between teachers' ratings and their students' gain scores on standardized tests, and evidence that teachers' practice improved as they were given frequent feedback in relation to the standards.²² In the schools and districts studied, assessments of teachers were based on well-articulated standards of practice evaluated through evidence including observations of teaching along with teacher pre- and post-observation interviews and, sometimes, artifacts such as lesson plans, assignments, and samples of student work.

Finding Additional Measures Related to Teacher Effectiveness

The Gates Foundation has launched a major initiative to find additional tools that are validated against student achievement gains and that can be used in teacher evaluation at the local level. The Measure of Effective Teaching (MET) Project has developed a number of tools, some of them based on the standards-based assessments described above and others taking a new tack. Among these are observations or videotapes of teachers, supplemented with other artifacts of practice (lesson plans, assignments, etc.), that can be scored according to a set of standards which reflect practices associated with effective teaching. Also included are tools such as student surveys about teaching practice, which have been found, in an initial study, to be significantly related to student achievement gains.²³

Countries such as Singapore include a major emphasis on teacher collaboration in their evaluation systems. This kind of measure is supported by studies that have found that stronger value-added gains for students are supported by teachers who work together as teams²⁴ and by higher levels of teacher collaboration for school improvement.²⁵

Some systems ask teachers to assemble evidence of student learning as part of the overall judgment of effectiveness. Such evidence is drawn from classroom- and school-level assessments and documentation, including pre- and post-test measures of student learning in specific courses or curriculum areas, and evidence of student accomplishments in relation to teaching activities. A study of Arizona's career ladder program, which requires the use of various methods of student assessment to complement evaluations of teachers' practice, found that, over time, participating teachers demonstrated an increased ability to create locally developed assessment tools to assess student learning gains in their classrooms; to develop and evaluate pre- and post-tests; to define measurable outcomes in hard-to-quantify areas such as art, music, and physical education; and to

monitor student learning growth. They also showed a greater awareness of the importance of sound curriculum development, more alignment of curriculum with district objectives, and increased focus on higher quality content, skills, and instructional strategies.²⁶ Thus, the development and use of student learning evidence, in combination with examination of teaching performance, can stimulate improvements in practice.

Building Systems for Teacher Evaluation That Support Improvement and Decision Making

Systems that help teachers improve and that support timely and efficient personnel decisions have more than good instruments. Successful systems use multiple classroom observations throughout the year by expert evaluators looking at multiple sources of data that reflect a teacher's instructional practice, and they provide timely and meaningful feedback to the teacher.

For example, the Teacher Advancement Program (TAP), which is based on the standards of the National Board and INTASC, as well as the standards-based assessment rubrics developed in Connecticut,²⁷ ensures that teachers are evaluated four to six times a year by master/mentor teachers or principals who have been trained and certified in a rigorous 4-day training. The indicators of good teaching are practices that have been found to be associated with desired student outcomes. Teachers also study the rubric and its implications for teaching and learning, view and evaluate videotaped teaching episodes using the rubric, and engage in practice evaluations. After each observation, the evaluator and teacher meet to discuss the findings and to make a plan for ongoing growth. Ongoing professional development, mentoring, and classroom support are provided to help teachers meet these standards. Teachers in TAP schools report that this system, along with the intensive professional development offered, is substantially responsible for improvements in their practice and the gains in student achievement that have occurred in many TAP schools.²⁸

In districts that use peer assistance and review (PAR) programs, highly expert mentor teachers conduct some aspects of the evaluation and provide assistance to teachers who need it. Key features of these systems include not only the instruments used for evaluation but also the expertise of the consulting teachers or mentors—skilled teachers in the same subject areas and school levels who have released time to serve as mentors to support their fellow teachers—and the system of due process and review that involve a panel of both teachers and administrators in making recommendations about personnel decisions based on the evidence presented to them from the evaluations. Many systems using this approach have been found not only to improve teaching, but also to successfully identify teachers for continuation and tenure as well as intensive assistance and personnel action.²⁹

Summary and Conclusions

New approaches to teacher evaluation should take advantage of research on teacher effectiveness. While there are considerable challenges in the use of value-added test scores to evaluate individual teachers directly, the use of value-added methods can help to validate measures that are productive for teacher evaluation.

With respect to value-added measures of student achievement tied to individual teachers, current research suggests that high-stakes, individual-level decisions, as well as comparisons across highly dissimilar schools or student populations, should be avoided. Valid interpretations require aggregate-level data and should ensure that background factors, including overall classroom

composition, are as similar as possible across groups being compared. In general, such measures should be used only in a low-stakes fashion when they are part of an integrated analysis of what the teacher is doing and who is being taught.

Other teacher evaluation tools that have been found to be both predictive of student learning gains and productive for teacher learning include *standards-based evaluation processes*. These include systems like National Board Certification and performance assessments for beginning teacher licensing as well as district- and school-level instruments based on professional teaching standards. Effective systems have developed an integrated set of measures that show what teachers do and what happens as a result. These measures may include evidence of student work and learning, as well as evidence of teacher practices derived from observations, videotapes, artifacts, and even student surveys.

These tools are most effective when embedded in systems that support evaluation expertise and well-grounded decisions, by ensuring that evaluators are trained, evaluation and feedback are frequent, mentoring and professional development are available, and processes are in place to support due process and timely decision making by an appropriate body.

With these features in place, evaluation can become a more useful part of a productive human-capital system, supporting accurate information about teachers, helpful feedback, and well-grounded personnel decisions

Endnotes

1. See for example, studies of teaching practices found associated with student learning gains, summarized in L. Darling-Hammond & J. Bransford (2005). *Preparing Teachers for a Changing World: What Teachers Should Learn and be Able to Do*. San Francisco: Jossey-Bass. See also, studies of National Board Certification as a measure of teacher effectiveness summarized in National Research Council (2008). *Assessing Accomplished Teaching: Advanced level Certification Programs*. Washington, DC: National Academy Press.
2. X. Newton, L. Darling-Hammond, E. Haertel, & E. Thomas (2010). Value-Added Modeling of Teacher Effectiveness: An Exploration of Stability Across Models and Contexts. *Educational Policy Analysis Archives* 18(23).
3. Newton et al. (2010). J. Rothstein. (2007). Do Value-Added Models Add Value? Tracking, Fixed Effects, and Causal Inference. CEPS Working Paper No. 159. National Bureau for Economic Research.
4. D. Briggs, & B. Domingue. (2011). Due Diligence and the Evaluation of Teachers: A Review of the Value-Added Analysis Underlying the Effectiveness Rankings of Los Angeles Unified School District Teachers by the Los Angeles Times. Boulder, CO: National Education Policy Center.
5. J. R. Lockwood, D. F. McCaffrey, L. S. Hamilton, B. Stetcher, V. N. Le, & J. F. Martinez. (2007). The Sensitivity of Value-Added Teacher Effect Estimates to Different Mathematics Achievement Measures. *Journal of Educational Measurement* 44 (1):47–67.
6. Bill & Melinda Gates Foundation (2010). *Learning About Teaching: Initial Findings From the Measures of Effective Teaching Project*. Seattle: Author. Rothstein, Jesse (2011). *Review of 'Learning About Teaching: Initial Findings from the Measures of Effective Teaching Project*. Boulder, CO: National Education Policy Center.
7. Lockwood et al. (2007).
8. Sean P. Corcoran, Jennifer L. Jennings, and Andrew A. Beveridge (2011). *Teacher Effectiveness on High- and Low-Stakes Tests*. Working paper. New York: New York University.
9. Briggs & Domingue (2011). Rothstein, Jesse (2010). Teacher Quality in Educational Production: Tracking, Decay, and Student Achievement, *Quarterly Journal of Economics* 125(1): 175–214.
10. Newton et al. (2010).
11. Henry Braun, H. (2005). *Using Student Progress to Evaluate Teachers: A Primer on Value-Added Models*. Princeton, NJ: Educational Testing Service.
12. This and other examples that follow are from Audrey Amrein-Beardsley & C. Collins (forthcoming). *The SAS® Education Value-Added Assessment System (EVAAS®): Its Intended and Unintended Effects in a Major Urban School System*. Arizona State University.
13. Amrein-Beardsley & Collins (forthcoming).
14. Daniel F. McCaffrey, Daniel Koretz, J. R. Lockwood, & Laura S. Hamilton (2005). *Evaluating Value-Added Models for Teacher Accountability*. Santa Monica: RAND Corporation.
15. Henry Braun, *Using Student Progress to Evaluate Teachers: A Primer on Value-Added Models*. Princeton, NJ: ETS, 2005, p. 17.
16. National Research Council, Board on Testing and Assessment (2009). Letter Report to the U.S. Department of Education.
17. For a summary of studies, see Darling-Hammond & J. Bransford (2005).

-
18. See, for example, L. Bond, T. Smith, W. Baker, & J. Hattie (2000). *The Certification System of the National Board for Professional Teaching Standards: A Construct and Consequential Validity Study* (Greensboro, NC: Center for Educational Research and Evaluation); L. Cavaluzzo (2004). *Is National Board Certification an Effective Signal of Teacher Quality?* National Science Foundation No. REC-0107014. Alexandria, VA: The CNA Corporation; D. Goldhaber, & E. Anthony (2005). *Can Teacher Quality Be Effectively Assessed?* Seattle, WA: University of Washington and the Urban Institute; T. Smith, B. Gordon, S. Colby, & J. Wang (2005). *An Examination of the Relationship of the Depth of Student Learning and National Board Certification Status*, Office for Research on Teaching, Appalachian State University. L. G. Vandevoort, A. Amrein-Beardsley, & D. C. Berliner (2004). *National Board Certified Teachers and Their Students' Achievement*. *Education Policy Analysis Archives* 12(46):117.
19. S. Athanases (1994). Teachers' Reports of the Effects of Preparing Portfolios of Literacy Instruction. *Elementary School Journal* 94(4):421-43; M. Sato, R. C. Wei, & L. Darling-Hammond (2008). Improving Teachers' Assessment Practices Through Professional Development: The Case of National Board Certification, *American Educational Research Journal* 45:669-700; S. M. Tracz, S. Sienty, K. Todorov, J. Snyder, B. Takashima, R. Pensabene, B. Olsen, L. Pauls, & J. Sork (1995, April). Improvement in teaching skills: Perspectives from National Board for Professional Teaching Standards field test network candidates. Paper presented at the annual meeting of the American Educational Research Association. San Francisco, CA.
20. M. Wilson & P.J. Hallum (2006). Using Student Achievement Test Scores as Evidence of External Validity for Indicators of Teacher Quality: Connecticut's *Beginning Educator Support and Training Program*. Berkeley, CA: University of California at Berkeley; S. P. Newton (2011). *Predictive Validity of the Performance Assessment for California Teachers*. Stanford, CA: Stanford Center for Opportunity Policy in Education, 2010.
21. R. Chung (2008). Beyond Assessment: Performance Assessments in Teacher Education, *Teacher Education Quarterly* 35(1):7-28; R.C. Wei and R. Pecheone (2010). Teaching Performance Assessments as Summative Events and Educative Tools. In Mary Kennedy (ed.), *Teacher Assessment and Teacher Quality: A Handbook* (New York: Jossey-Bass).
22. A. Milanowski, S.M. Kimball, B. White (2004). *The Relationship Between Standards-Based Teacher Evaluation Scores and Student Achievement*. University of Wisconsin-Madison: Consortium for Policy Research in Education; Anthony Milanowski (2004). The Relationship Between Teacher Performance Evaluation Scores and Student Achievement: Evidence From Cincinnati, *Peabody Journal of Education* 79(4):33-53. Jonah Rockoff & Cecilia Speroni (2010). *Subjective and Objective Evaluations of Teacher Effectiveness*, New York: Columbia University.
23. Bill & Melinda Gates Foundation (2010).
24. C. K. Jackson & E. Bruegmann (2009). *Teaching Students and Teaching Each Other: The Importance of Peer Learning for Teachers*. Washington, DC: National Bureau of Economic Research.
25. Y. Goddard & R. D. Goddard (2007). A theoretical and empirical investigation of teacher collaboration for school improvement and student achievement in public elementary schools. *Teachers College Record* 109(4):877-896.
26. R. Packard & M. Dereshiwsy (1991). *Final Quantitative Assessment of the Arizona Career Ladder Pilot-Test Project*. Flagstaff: Northern Arizona University.
27. The teacher responsibility rubrics were designed based on several teacher accountability systems currently in use, including the Rochester (NY) Career in Teaching Program, Douglas County (CO) *Teacher's Performance Pay Plan*, Vaughn Next Century Charter School (Los Angeles, CA) Plan, and Rolla (MO) Professional Based Teacher Evaluation
28. Lewis Solomon, J. Todd White, Donna Cohen, & Deborah Woo (2007). *The Effectiveness of the Teacher Advancement Program*. National Institute for Excellence in Teaching.
29. National Commission on Teaching and America's Future (1996). *What Matters Most: Teaching for America's Future*. NY: NCTAF; Piet Van Lier (2008). *Learning From Ohio's Best Teachers: A Homegrown Model to Improve Our Schools*. Policy Matters Ohio.