

# Using Science Simulations to Support Powerful Formative Assessments of Complex Science Learning

Edys S. Quellmalz, Michael J. Timms, and Barbara C. Buckley  
WestEd

## Introduction

An integral component of student learning is the assessment of student progress in order to inform instructional decisions. Classroom-based assessments that are used in *formative* ways have been repeatedly shown to significantly benefit student learning (Black & Wiliam, 1998). The Council of Chief State School Officers (CCSSO) state collaborative on formative assessment for students and teachers (FAST SCASS) developed a definition of formative assessment based on the research literature. “Formative assessment is a process used by teachers and students during instruction that provides feedback to adjust ongoing teaching and learning to improve students’ achievement of intended instructional outcomes” (CCSSO FAST SCASS, 2008). Formative assessment practices include not only a variety of informal strategies teachers may use “on the fly”, but also formal assessment tasks and methods systematically designed to help teachers probe and promote student thinking and reasoning. Design principles drawn from research on learning and assessment can shape formal assessment tasks and interactions specifically designed to elicit evidence of what students know and can do in relation to clearly articulated learning targets (Messick, 1994; Pellegrino et al., 2001). Such principled assessment designs are greatly needed given the uneven technical quality of classroom assessments developed by teachers or those that are embedded in published curriculum materials (Wilson & Sloan, 2000; Mislevy & Haertel, 2006).

This presentation describes how SimScientists projects at WestEd are studying ways in which simulation-based assessments support effective formative assessment to promote rich science learning. Science simulations present opportunities for students to actively conduct investigations by manipulating models of complex scientific phenomena involving dynamic, causal, and temporal interactions. Simulations can test if students’ have integrated knowledge of systems in the natural world and also measure students’ use of inquiry skills—science content and processes that are addressed in only limited ways or not at all by paper-based tests (Quellmalz et al., 2005).

As computer-based testing increases in large-scale, summative assessments, more attention is being paid to the capabilities of technology for innovative assessment designs. Interactive computer-based assessments to measure students’ abilities to use inquiry practices are components of the 2009 National Assessment of Educational Progress (NAEP) in science. Several state science tests are initiating innovative, dynamic assessment scenarios. However, high quality assessment tasks that teachers may use for formative assessment are in short supply. Moreover, technically sound, research-based assessments that can provide formative feedback and support *during* instruction are difficult for teachers to develop on their own, and are seldom provided in published curriculum materials. Nor, are tests readily available that could serve as benchmarks of complex science learning tied to curriculum units.

In this symposium, we describe the design and development of assessments in WestEd's SimScientists program that aim to take advantage of simulations to offer a new generation of technology-enhanced assessments to support best classroom formative assessment practices. We describe methods for designing and studying innovative technology-enhanced classroom assessments that will transform what, how, when, and where learning is assessed. We end with a description of how the SimScientists projects are studying the benefits of simulation-based science assessments in balanced state science assessment systems.

## **SimScientists Assessments**

In one of WestEd's SimScientists projects, curriculum-embedded and benchmark science assessments are being developed in an NSF-funded study, Calipers II: Using Simulations to Assess Complex Science Learning. Goals of the Calipers II project include: (1) develop simulation-based assessments that can be embedded in curriculum units for formative purposes and as benchmark assessments that can be administered at the end of a unit for summative information about proficiency on intended standards, (2) document the re-usable designs and processes employed, (3) provide evidence of the technical quality, feasibility, and utility of the new assessments, and (4) study the influence of formative assessments on complex science and inquiry learning.

The simulation-based assessments are designed to present dynamic, engaging interactive tasks of established technical quality that test complex science knowledge and inquiry skills that go well beyond the capabilities of print tests. Benchmark assessments are designed to test end-of-unit achievement of content and inquiry standards addressed in curriculum units on a middle school science topic such as human body systems or climate. Sets of shorter, embedded assessments are designed to be used during the instructional unit. The simulation-based embedded science assessments are intended to function as formative resources in three ways: (1) provide immediate feedback contingent on an individual student's performance, (2) offer graduated levels of coaching in real-time, and (3) provide diagnostic information to guide offline reflection and extension activities. The software and technical infrastructure of these simulation-based assessments can overcome many of the practical constraints that have limited the use and effectiveness of formative assessments. The SimScientists projects couple embedded assessment simulations with follow-on, off-line self-assessment and reflection activities. The embedded assessment and reflection activities incorporate the features of effective formative assessment: *frequent* use of standards-based classroom assessments; feedback that is *timely, individualized*, and *diagnostic*; online *supplementary instruction* that is individualized; and off-line *self-assessment* and reflection activities that help students confront misunderstandings, make new connections, and become more reflective, self-regulating learners (Herman et al., 2005).

**The SimScientists Assessment Designs.** The projects are using evidence-centered design to create and link *student models* (content and inquiry targets) to *task models* that will elicit evidence of the targeted knowledge and inquiry, to *evidence models* of the scoring and reporting that will be used to describe student progress and achievement on the targets (Mislevy & Haertel, 2006). The SimScientists simulation-based assessments are currently being developed for the middle school level, for two systems each within physical, life, and earth science. The design

principles shaping the assessments draw on recommendations from research on methods to promote science learning. They include: (1) a focus on integrated knowledge about the dynamic relationships among structures, behaviors, and mechanisms within models of science systems, (2) use of authentic, problem-driven inquiry practices, (3) scaffolding that generates immediate, individualized feedback, and levels of customized coaching, (4) metacognitive self-assessment and reflection, (5) scientific explanations and arguments, and (5) use of the affordances of simulations to provide multiple representations of dynamic, causal, and temporal phenomena and to offer multiple ways for students to interact and respond (Duschl et al., 2007).

The SimScientists assessment design process begins with examinations of middle school level national science standards set forth in the National Science Education Standards (NSES), *AAAS Benchmarks for Scientific Literacy* and the 2009 National Assessment for Educational Progress (NAEP) Science Framework to identify standards related to understanding systems and using inquiry practices. Project staff also analyze representative middle school science curricula to identify the science standards addressed and the types of activities and representations of science systems employed. SimScientists project science experts specify the components of the science system to be represented in the simulation model, and the types of inquiry that will promote students' understanding and use of the model. Following identification of content and inquiry targets to be addressed in a benchmark, end-of-unit assessment for a particular system (e.g., ecosystem, force and motion), the project staff conduct literature reviews to identify major misconceptions related to the assessment targets and, especially, misconceptions that are likely to reveal incorrect or naïve understandings of the target science system model.

Simulation shells are then developed to provide an overview of the assessment content and structure. Simulation shells include: (1) specification of an authentic, driving problem to structure the assessment, (2) an outline of the sequence and types of assessment tasks and questions students will complete as they conduct inquiry to address the overarching problem, (3) alignment of each of the types of assessment tasks with national content and inquiry standards, (4) identification of key misconceptions, and (5) descriptions of the features of the simulation model students will see and manipulate.

An alignment matrix of intended assessment content and inquiry targets linked to national standards and the 2009 NAEP science framework, along with the simulation shell, are reviewed by project science experts and external reviewers prior to development of a storyboard.

Storyboards sketch screen-by-screen features of the simulation to be presented to the students, specify the interactivity of the student with the simulation, and draft student directions and questions. Internal and external science and assessment experts review the storyboards for alignment of the intended tasks and questions to the standards they are intended to test, quality of the science content and representations, and the appropriateness of the simulations for representing the target model of the science system.

As storyboards are developed, universal design guidelines are also employed to make the assessments accessible to the widest range of diverse student populations. The task features of simulation-based assessments allow multiple representations of the science system models and multi-modal opportunities for students to respond. One aim is for students with diverse language backgrounds and learning styles to have better opportunities to understand the tasks and to demonstrate their learning.

Once the set of targets and misconceptions to be assessed in the end-of-unit benchmark assessments have been identified, the end-of-unit targets are divided into component models or concepts to be addressed in embedded assessments that can be used by the teacher and students at multiple points during the unit. Designs of the embedded assessments specify the immediate feedback to be generated by the computer for responses that can be automatically scored and also present sample appropriate responses for text or graphical constructed-responses. In addition to immediate feedback, the assessment software presents the student with graduated coaching, ranging from directions to try again to fully worked examples. Students are guided through multiple trials to respond appropriately to an assessment task. The system provides sufficient coaching to allow a student to successfully complete each task and move along in the embedded assessment.

Figure 1 presents a screen shot with an example of feedback and coaching provided for an assessment task in which students are asked to draw arrows in a food web to represent the flow of energy among organisms the students have previously observed. The technology enhances the feedback and coaching possible in multiple ways, such as allowing students to observe again the roles of organisms by re-running an animation of the organisms in the environment as they interact. The feedback and coaching system can highlight incorrect arrows a student has drawn between organisms, and provide more scaffolding by running an animation that highlights arrows being drawn correctly, and prompting the student to then draw the correct arrows before the student is allowed to proceed to the next screen.

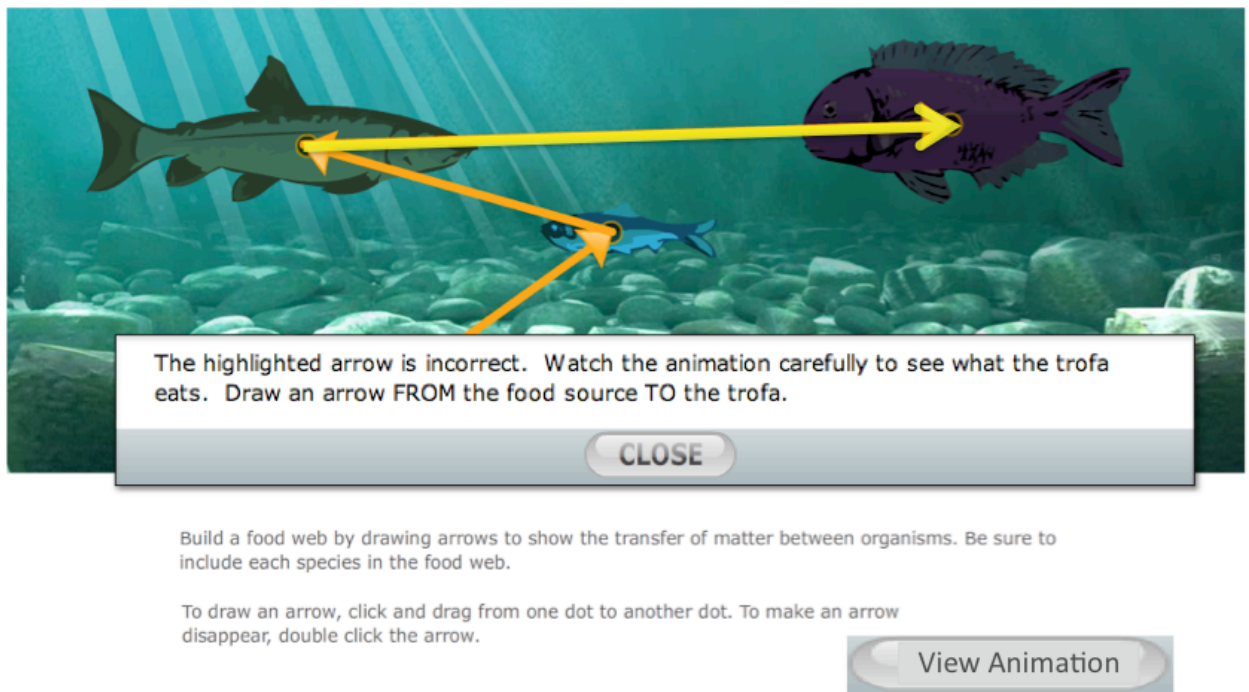


Figure 1. Sample of feedback on task in embedded assessment.

Figure 2 presents a screenshot of an ecosystem simulation in which students can conduct multiple investigations of the effects of varying the number of organisms on a model of the population of organisms in the system.

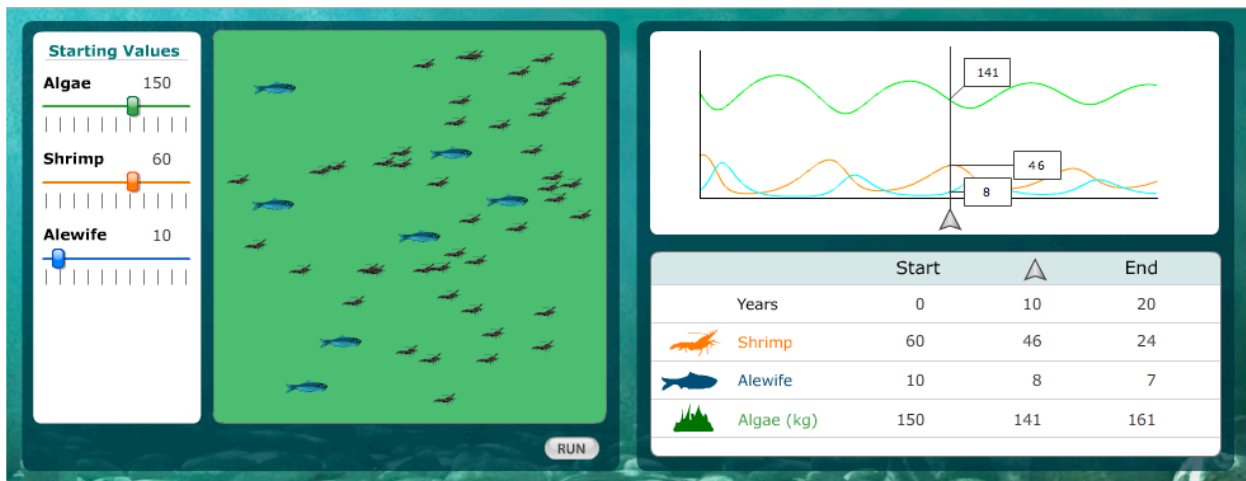


Figure 2. The population simulation for use in both embedded and benchmark assessment tasks.

The model of the population level of an ecosystem is simulated in three ways. A physical representation of changes in fish and algae populations are animated in the left box. On the right, a graph and table represent the changing population levels over time. The simulation allows students to run multiple experiments. A graph inspector arrow allows students to read the graphic data generated from the simulation at different points in time to examine the relationship among the numbers of organisms. Students are asked to predict how the number of an organism will be affected initially and later when another organism is added. Feedback and coaching for these tasks in which students use the model to predict, observe, and test their predictions graduates from informing them that their prediction was incorrect and to try again, to prompting students to use the graph inspector, to highlighting the places on the graph and data in the table and pointing out the changing population numbers.

Each curriculum-embedded assessment is accompanied by guidelines to help the teacher conduct an off-line, follow-up self-assessment and reflection activity. The self-assessments are intended to offer further support for promoting learning of concepts and/or inquiry skills identified in the reports generated by the embedded assessment system as needing to be strengthened. Self-assessment activities would include both students' participation in using criteria to evaluate their own constructed-responses and also guidelines for additional instruction that might involve, for example, group or class discussion of concepts or inquiry skills the teacher judges need further scaffolding and practice. Student self-assessments of constructed-responses are scaffolded by sample student responses for levels of a scoring rubric. The self-assessment and reflection activities are designed to build students' metacognitive skills for planning and self-monitoring. The reflection activities are also designed to engage students in scientific argumentation and discourse.

**Technical Infrastructure for the Embedded Formative Assessments.** Through the use of carefully designed algorithms that assess the level of a student response, the software system is able to give tailored feedback to a student based on the pattern of his/her responses. The system assesses not only student responses to conventional selected response items, but also elements of the students' interaction with the simulation such as latency (time taken), variables manipulated in a simulation, the values assigned to variables in the simulation (e.g., the number of shrimp), and the order and number of trials run in an experiment. These indicators can be used to appraise students' problem solving strategies that cannot be assessed in conventional paper and pencil science tests.

The SimScientists software system uses algorithms to process automated scores and levels of coaching received in order to track the student's understanding and inquiry skills and report them back to the student and to the teacher. The reports to students are intended to show the progress they are making and affirm what they have demonstrated that they know and can do. The format of the formative assessment reports will be, for example, "You have shown that you understand the role of producers in the lake ecosystem. You still need to work on understanding the role of consumers." The reports to teachers are intended to provide information that will allow them to identify individual students, or groups of students, who need further assistance in understanding targeted concepts or applying the inquiry skills. For example, the report might identify that 65% of the students understand the role of producers in the ecosystem, but that 35% are still unclear about this and need additional instruction.

### **Studying the Technical Quality, Utility, and Feasibility of the SimScientists Assessments**

The current generation of SimScientists assessments draw upon findings from a prior NSF-funded project, Calipers I, which aimed to demonstrate the technical quality, feasibility, and utility of simulation-based science summative, benchmark assessments (Quellmalz et al., 2007; Quellmalz et al., in press). The Calipers I project used a range of qualitative and quantitative measures. Major validation methods triangulated alignment of the assessments with intended standards, cognitive labs to confirm the construct validity—e.g., that the tasks and items were eliciting the intended knowledge and skills—with standard empirical analyses of student performance (Quellmalz, et al., 2005). Utility was studied by interviews with teachers who pilot tested the assessments in their classes. Feasibility was studied in early usability trials with individuals and in observations of the pilot testing of the assessments in intact classrooms.

In the Calipers I demonstration project, the intended alignment with national science standards for the content and inquiry abilities that resulted from the systematic development process was confirmed by independent expert and teacher reviews of the assessments. The cognitive labs conducted with individual students also confirmed that the assessment tasks and items were eliciting the intended constructs. Analyses of student performance on the assessments provided further evidence that the simulation-based assessments were reliable and valid. Multidimensional analyses showed that the assessments did detect science content knowledge and inquiry ability constructs as they were designed to, and the assessments were shown to distinguish clearly between different categories of students based on an external measure of science achievement. Students' content and inquiry scores overall were consistent with teachers' rankings of their students' science achievement level (High, Medium, Low). As expected,

students with high rankings performed better than students with medium or low rankings. Pilot testing of the assessments with middle school students showed that the items performed within commonly accepted levels on standard psychometric measures of item difficulty and fit. Student scores were also consistent with teacher and student reports on whether students had received instruction (Opportunity to Learn) related to the specific targets for content and inquiry. Although we were not able to fully investigate how student problem solving actions during the simulations related to student performances, we did learn that some actions could be significant predictors of outcomes. Furthermore, we were able to identify different strategies used by students to solve problems we posed. We also learned many practical lessons about constraining students' interactions with the simulations in order to better assess students' reasoning. The Calipers I classroom pilot testing provided strong evidence of the benchmark assessments' quality as summative assessments of test end-of unit achievement of the targeted complex science learning.

The pilot testing also supported the feasibility and utility of the Calipers I benchmark assessments. Observations of the administration and implementation of the simulation-based assessments revealed few logistical issues. Teachers were interviewed about their perceptions of the Calipers simulation-based assessments. Teachers indicated that they thought using the assessments would be practical given sufficient access to computers. All the pilot teachers were very positive about the Calipers simulation-based assessments. The teachers felt that simulation-based assessments probed depth of understanding rather than rote learning. The teachers appreciated that the students had to use their minds and science skills to solve the problems posed in the assessments. One teacher remarked that what the simulations do, and the paper tests cannot do, is to *animate and show the students the results* of the answer they chose. In addition, teachers saw that the assessments provided information about the processes students were using, not just the answer. Teachers observed that the assessments presented authentic problems that allowed students to see how the science they were studying related to real life and answered the question, "Why do we care?" Several teachers suggested that the simulations would help English language learners and students who didn't do well on traditional tests. All of the teachers said that the most useful information would be score reports on their students' progress and difficulties related to the specific content and inquiry skills. Teachers felt that real time scoring and immediate results would help them decide what to do next with students. Four of the teachers remarked that they would like to use the simulations for instruction and to administer them during a unit as formative assessments.

The results of the pilot testing for the Calipers I demonstration project supported the the technical quality, feasibility, and utility of the benchmark assessments. The results also provided evidence that the particular exemplars and others that could be modeled after them would be likely to provide credible data for summative accountability purposes.

Consequently, the SimScientists program at WestEd launched a set of projects to further study the promise of simulation-based science assessments. The Calipers II project, funded by the National Science Foundation, broadens the use of science simulations to additional topics in life, physical and earth science. Importantly, the Calipers II project is augmenting the designs of the science simulations so that they can serve as formative assessment resources, as illustrated in Figures 1 and 2. To study the technical quality, feasibility, and utility of science simulations for both formative and summative uses, the Calipers II project is employing and extending the

rigorous methods used in Calipers I. Research on science learning and model-based reasoning is guiding the identification of knowledge and inquiry targets and key misconceptions (Buckley, 2000; Buckley, in press; Gobert, 2000; Clement, 2008;). Evidence-centered design aligns the designs of assessment tasks with the knowledge and skills they are intended to elicit, and the evidence of learning the scoring and reporting will provide. Formative reviews of the alignment of the assessments with national standards and also the scientific quality of the assessments have been conducted by AAAS for the ecosystem assessments, and are scheduled for the other topics. Cognitive labs are currently in progress in which approximately a small number of middle school students think aloud as they work through the embedded and benchmark assessments. Cognitive labs are also underway in which middle school teachers currently teaching ecosystems think aloud as they complete the embedded and benchmark assessment tasks. In these initial cognitive labs, we can check on the usability of the software and gauge if the items do provoke the intended thinking about science content and application of inquiry skills. Further validity evidence gathered in subsequent classroom pilot and field testing based will come from IRT analysis of individual items to establish that they are of appropriate ranges of difficulty and are fitting well to the measurement construct. Evidence based on relations to other variables will include correlations to student performance on more traditional items on the same topics, administered as pretests and posttests.

Feasibility of the assessments will be judged by conducting small scale pilot testing within a few classrooms involving a few hundred students, followed by larger field tests studying the effects of the embedded assessments on student learning.

Two additional projects are studying the potential of the embedded and benchmark assessments as components of balanced statewide science assessment systems. In these projects, the SimScientists simulation-based embedded and benchmark assessments will be studied in schools in three US states involving thousands of students. Common accommodations will be designed and studied for their impacts on samples of English language learners and students with disabilities. Teachers will complete weekly logs about their use of the assessments to inform us of how they are being implemented and any difficulties they encounter. In two of the projects, CRESST will serve as an external evaluator, reviewing the assessment development methods and data collections and conducting case studies of their use.

The multi-year SimScientists projects are in their first years of development. In the next year, classroom pilot testing will yield more data on the technical quality, feasibility, and utility of simulation-based science assessments for formative and summative purposes. The SimScientists research and development program will provide models and evidence of the power of science simulations for formative assessment.

## References

- American Association for the Advancement of Science. (1993). *Benchmarks for science literacy*. New York: Oxford University Press.
- Duschl, R.A., Schweingruber, H.A., Shouse, A.W. (2007). *Taking Science to School: Learning and Teaching Science in Grades K-8*. Committee on Science Learning, Kindergarten Through Eighth Grade. Board on Science Education, Center for Education. Division of Behavioral and Social Sciences and Education. Washington, DC: The National Academies Press.



- Black, P., & Wiliam, D. (1998). *Inside the black box: Raising standards through classroom assessment*. London: King's College.
- Duschl, R.A., Schweingruber, H.A., Shouse, A.W. (2007). *Taking Science to School: Learning and Teaching Science in Grades K-8*. Committee on Science Learning, Kindergarten Through Eighth Grade. Board on Science Education, Center for Education. Division of Behavioral and Social Sciences and Education. Washington, DC: The National Academies Press.
- Buckley, B. C. (2000a). Interactive multimedia and model-based learning in biology. *International Journal of Science Education*, 22(9), 895-935.
- Buckley, B. C., & Boulter, C. J. (2000b). Investigating the role of representations and expressed models in building mental models. In J. K. Gilbert & C. J. Boulter (Eds.), *Developing models in science education* (pp. 105-122). Dordrecht, Holland: Kluwer.
- Buckley, B. C., Gobert, J., Horwitz, P., & O'Dwyer, L. (in press, 2008). Looking inside the black box: Assessing model-based learning and inquiry in BioLogica. *International Journal of Learning Technologies*.
- Buckley, B. C., Gobert, J., Kindfield, A. C. H., Horwitz, P., Tinker, B., Gerlits, B., et al. (2004). Model-based Teaching and Learning with Hypermodels: What do they learn? How do they learn? How do we know? *Journal of Science, Education and Technology*, 13(1), 23-41.
- Clement, J. J., & Rea-Ramirez, M. A. (Eds.). (2008). *Model Based Learning and Instruction in Science*. London: Springer.
- Duschl, R.A., Schweingruber, H.A., Shouse, A.W. (2007). *Taking Science to School: Learning and Teaching Science in Grades K-8*. Committee on Science Learning, Kindergarten Through Eighth Grade. Board on Science Education, Center for Education. Division of Behavioral and Social Sciences and Education. Washington, DC: The National Academies Press.
- Gobert, J. D., & Buckley, B. C. (2000). Introduction to model-based teaching and learning in science education. *International Journal of Science Education*, 22(9), 891-894.
- Herman, J. L., Osmundson, E., Ayalya, C., Schneider, S., & Timms, M. (2005). *The nature and impact of teachers' formative assessment practices*. Paper presented at the annual meeting of the American Educational Research Association, Montreal, Quebec, Canada.
- Messick, S. (1994). The interplay of evidence and consequences in the validation of performance assessments. *Educational Researcher*, 32, 13-23.
- Mislevy, R. J., & Haertel, G. D. (2006). Implications of evidence-centered design for educational testing. *Educational Measurement: Issues and Practice*, 25(4), 6-20.
- National Committee on Science Education Standards and Assessment. (1996). *National Science Education Standards: 1996*. Washington, D.C.: National Academy Press.
- Pellegrino, J., Chudowsky, N., & Glaser, R. (2001). *Knowing what students know: The science and design of educational assessment*. Washington, DC: National Academy Press.
- Quellmalz, E.S., Buckley, B., & Timms, M. (in press). Exploring the Role of Technology-Based Simulations in Science Assessment: The Calipers Project. *International Journal of Learning Technology*.
- Quellmalz, E., Kreikemeier, P., DeBarger, A. H., and Haertel, G. (2007). *A study of the alignment of the NAEP, TIMSS, and New Standards Science Assessments with the inquiry abilities in the National Science Education Standards*. Presented at the Annual Meeting of the American Educational Research Association, April 9-13, Chicago, IL.

- Quellmalz, E. S., DeBarger, A., Haertel, G., & Kreikemeier, P. (2005). *Validities of science inquiry assessments: Final report*. Menlo Park, CA: SRI International.
- Quellmalz, E.S., DeBarger, A.H., Haertel, G., Schank, P., Buckley, B., Gobert, J., Horwitz, P., Ayala, C. (2008). Exploring the Role of Technology-Based Simulations in Science Assessment: The Calipers Project. In Coffrey, R.Douglas, and C. Stearns (Eds.) *Assessing Science Learning: Perspectives*
- Quellmalz, E. S., & Haertel, G. (2004, May). *Technology supports for state science assessment Systems*. Paper commissioned by the National Research Council Committee on Test Design for K-12 Science Achievement.
- Wilson, M., & Sloan, K. (2000). From principles to practice: An embedded assessment system. *Applied Measurement in Education*, 13 (2), 181-208.