

# Validation of ELA and Mathematics Assessments: A General Approach<sup>1</sup>

Joan L. Herman & Kilchan Choi

## Overview

Validity refers to the degree to which an assessment actually measures what it claims to measure and well serves intended purposes. From this perspective, assessments themselves are not valid or not; rather evidence of validity must be established in the context of specific interpretations and uses of test scores. A test may be well suited for one use and not for another. For example, a mathematics test appropriate for accountability or summative purposes may not well serve formative assessment purposes for teachers. Moreover, validity is a matter of degree; validation requires the *accumulation of evidence* to support the argument that scores derived from a given test yield accurate inferences to support intended interpretations and uses (AERA, APA, NCME, 1999; Kane, 2001). Finally, our validity definition requires consideration of both (1) what is measured—that an assessment measures what it is intended to measure—and (2) what interpretations and uses the assessment is intended to serve. For both summative and formative assessments, the definition implies that assessments must yield technically sound measures of student learning. Moreover, results should be useful and used for intended purposes and not carry serious unintended or negative consequences.

Modern theory suggests that validation be approached as an argument that needs to be substantiated. A validity argument lays out the claims that an assessment and its scores must satisfy to serve its proposed purpose(s) and/or use(s). Validation efforts then focus on the collection of evidence to document how well the assessment satisfies each claim. Below, we set out these claims as a set of criteria that educational assessments used for summative or accountability purposes should meet and present plans to systematically collect evidence to evaluate these claims throughout our test development and field-testing process.

Validity is a matter of degree; validation requires the accumulation of evidence to support the argument that scores derived from a given test yield accurate inferences to support intended interpretations and uses.

## Validity Criteria

The criteria combine the *Standards for Psychological and Educational Testing* (AERA, APA, & NCME, 1999) with the National Research Council's core dictum that assessments intended to both measure and benefit learning must first and foremost be themselves learning-based (Pellegrino, Chudowsky, & Glaser, 2001). This means that in addition to adhering to the *Standards'* core concepts of validity, accuracy and reliability, comparability, and fairness, assessments should fully reflect the intended learning domain and a model of how such learning develops. In the current national context, the Common Core State Standards define the prevailing domain specification and figure prominently in our validation approach. We draw on CRESST's ontology methodology—representations of major idea and principles that define a subject matter domain, the enduring concepts, knowledge, skills and attendant cognitive demands that constitute understanding of each, and the inter-relationships among them (see Baker et al., 2009; Iselli, 2010)—for both documenting the alignment of standards and assessment and, over time, validating and refining essential paths for learning.



National Center for Research  
on Evaluation, Standards, & Student Testing

UCLA | Graduate School of Education & Information Studies

<sup>1</sup> Copyright © 2012 The Regents of the University of California. The findings and opinions expressed in this report are those of the authors and do not necessarily reflect the positions or policies of the National Center for Research on Education, Standards, and Student Testing (CRESST).

The intent that both accountability and formative assessments will support the improvement of learning also brings with it the need for tests that are instructionally sensitive and have value for teaching and learning. Instructional sensitivity – that test scores reflect the quality of instruction, rather than home background and/or native ability—is a fundamental assumption underlying current accountability policy. If teaching improves, test scores will do likewise – otherwise, it is not reasonable to hold teachers and schools accountable for performance. Similarly, if scores are to be used to assess student progress – for example, the extent to which students are on track to

college readiness—then scores must be comparable across time and vertically scaled.

Validity, in short, resides in evidence of the extent to which an assessment embodies the characteristics that support its intended purpose(s) and the extent to which the scores from an assessment yield meaningful inferences to support intended decisions and uses. These concerns suggest two basic kinds of criteria for judging test quality: one set based on the attributes or the characteristics of the test itself, and the other based on the validity of score interpretations. We offer an initial list in *Table 1* and *Table 2* below:

**Table 1: Criteria Related to Assessment Attributes & Validity of Score Interpretation**

Criteria Related to Assessment Attributes
<p><b>Learning Based:</b></p> <ul style="list-style-type: none"> <li>• Aligned with significant learning goals, as reflected in Common Core State Standards and represented through CRESST ontologies, including expected content and cognitive demand</li> <li>• Comprehensive in representation of target constructs/domain intended for assessment, in both content and cognitive demands</li> <li>• Linked to expected trajectories of learning</li> </ul>
<p><b>Fairness</b></p> <ul style="list-style-type: none"> <li>• Accessible, enabling all students to show what they know</li> <li>• To the extent possible, free of knowledge and skills that are irrelevant to the target of the assessment, (e.g., construct irrelevant language demands)</li> <li>• Sensitive to a range of student abilities; appropriate for students at the range of developmental levels likely in assessed population</li> <li>• Provide accommodations where needed</li> </ul>
<p><b>Learning/Instructional Value</b></p> <ul style="list-style-type: none"> <li>• Incorporates transfer; new, authentic applications</li> <li>• Cognitively expansive (i.e., provides opportunities for students to organize and expand their knowledge, (e.g., through explanation, modeling, etc.)</li> <li>• Engaging</li> <li>• Consequences: models/supports rich curriculum learning opportunities; well communicates important learning goals</li> </ul>
<p><b>Utility</b></p> <ul style="list-style-type: none"> <li>• Timely</li> <li>• Useable/interpretable by teachers and/or students; provides actionable feedback for intended users</li> <li>• Instructionally useful, at the right grain size to guide subsequent, intended decision-making and action</li> </ul>
<p><b>Credibility</b></p> <ul style="list-style-type: none"> <li>• Educators</li> <li>• Students and parents</li> <li>• Public</li> </ul>

**Table 1: Criteria Related to Assessment Attributes & Validity of Score Interpretation (cont.)**

Criteria Related to Validity of Score Interpretation
<b>Technical Soundness</b> <ul style="list-style-type: none"><li>• Score reliability at level of intended use(s) (e.g., if assessment is formative, scores provide reliable diagnostic information)</li><li>• Reliability of scoring</li><li>• Accuracy for intended decision-making purpose(s), including monitoring progress and determining attainment of standards</li></ul>
<b>Fairness/lack of bias</b>
<b>Model/theory-based relationships with other measures</b>
<b>Generalizability, extent of transfer</b>
<b>Instructional sensitivity</b>
<b>Comparable, as necessary: across sites, across time, within and across years</b>
<b>Consequences</b> <ul style="list-style-type: none"><li>• Supports productive changes in teachers pedagogical and assessment practices</li><li>• Supports improvement in student learning</li><li>• Supports equity in student outcomes.</li></ul>

Validation Plan

Assessment validity cannot be an afterthought but rather needs to be built in throughout a systematic test development process. That process starts with clear delineations of the purpose of the assessment, the domains to be assessed, characteristics of the examinee population and intended users. Within existing constraints and resources, test blueprints can then be developed to specify the number and types of items to be sampled for specified content and cognitive demand domains. Task design/templates then may be developed or selected to reflect the blueprint specification and used to generate items to populate it. Once initially developed or selected, test items, tasks and test forms go through multiple levels of review and trial to fine-tune and assure quality and to generate evidence of validity. Just as assessment has been described as a process of “reasoning from evidence (Mislevy et al. 2002; Pellegrino et al, 2001),” so too is assessment development an evidence-based process.

Mirroring the two general types of criteria describe above, our validation plans are designed to triangulate evidence from two basic frames: (1) expert and educator review and (2) empirical study. For example, subject matter, learning and cognition, and language experts can be convened at multiple stages of the test development process to examine the alignment between the assessment to the

Common Core State Standards (as represented in content and cognitive demand ontologies), the assessment’s incorporation of research-based design principles to directly enhance learning, and to identify/analyze construct irrelevant bias that may be introduced by linguistic complexity or cultural differences. Drawing on the teachers involved in the empirical validity studies described below, feedback on issues of utility, usability and feasibility can be solicited. A sensitivity review by multi-cultural experts also is planned. Review strategies figure prominently in the early states of test development; empirical studies are built into pilot, field tests, and continuing study of operational use and consequences. These plans are summarized in Table 3.

Just as assessment has been described as a process of "reasoning from evidence," so too is assessment development an evidence-based process.



**Table 3: Validation Plan**

Test Development Step	Validity Criterion	Design/Data Source
1. Test Specification	<b>Learning-based:</b> <ul style="list-style-type: none"> <li>Alignment with standards/ontology</li> <li>Comprehensiveness</li> </ul>	Expert review
2. Task specifications (Design), including initial scoring rubrics	<b>Learning-based:</b> <ul style="list-style-type: none"> <li>Alignment with standards/ontology</li> <li>Linked to expected trajectory</li> </ul>	Expert review
	<b>Fairness</b>	Pedagogical/educator review
	<b>Learning/Instructional value</b>	Educator review
	<b>Feasibility/credibility</b>	
3. Task Development/Developmental Testing	<b>Learning-based:</b> <ul style="list-style-type: none"> <li>Alignment with standards/ontology</li> <li>Linked to expected trajectory</li> </ul>	
	<b>Fairness</b>	Sensitivity/accessibility reviews
	<b>Learning/Instructional value</b>	
	<b>Feasibility/credibility</b>	Empirical trial: 10-50 students from target grade levels for each task— practicality, timing, clarity
4. Field Test (test forms/sets of items and tasks)	<b>Learning-based:</b> <ul style="list-style-type: none"> <li>Alignment</li> <li>Comprehensiveness</li> <li>Relationships with other measures</li> </ul>	<b>Expert review</b> <b>Empirical study: <math>n</math> = at least 20 classes and 400 students, representing a full range of student ability.</b> <ul style="list-style-type: none"> <li>Oversample depending on intended focal groups for DIF analysis</li> <li>Each student responds to at least two tasks</li> <li>4-6 scorers</li> <li>Supplement with “expert sample,” <math>n</math> = 50-100.</li> <li>Teacher data: OTL, utility, credibility, feasibility</li> <li>Other concurrent student data: grades, standardized test, demographics</li> </ul>
	<b>Technical Soundness</b> <ul style="list-style-type: none"> <li>Score reliability at level of intended use</li> <li>Scorer reliability</li> </ul>	
	<b>Fairness/lack of bias</b>	
	<b>Generalizability</b>	
	<b>Instructional sensitivity (preliminary)</b>	
	<b>Model relationships</b> <ul style="list-style-type: none"> <li>Concurrent validity</li> <li>Expert-novice comparisons</li> <li>Overall vs. subgroup models</li> </ul>	
5. Operational Use	<b>Instructional Sensitivity</b>	<b>Empirical studies: <math>n</math> = 80 classrooms randomly assigned to instructed and non instructed groups; 3 tasks; data as above</b>  <b>Other variations: through course/end course administration</b>  <b>Teacher data: utility, credibility, feasibility</b>  <b>Instructional practice data</b>  <b>Longitudinal sample across three grades (<math>n</math> = 300 students), data as above with external growth and success indicators.</b>
	<b>Technical Soundness</b>	
	<b>Model Relationships as above</b>	
	<b>Teacher interpretation/use</b>	
	<b>Consequences/differential impact</b>	
	<b>Quality of implementation/teacher scoring</b>	
	<b>Validity of standard setting/vertical trajectory across grades</b>	

## References

---

- AERA, APA, NCME. (1999). *Standards for educational and psychological testing*. Washington, DC: APA.
- Bodilly, S., & Beckett, M. K. (2005). *Making out-of-school time matter: Evidence for an action agenda*. Santa Monica, CA: RAND.
- Baker, E., Chung, G., & Herman, J. (2009). *Ontology-based educational design: Seeing is believing*. Los Angeles, CA: CRESST.
- Iseli, M. (2010). *Ontology system overview*. Los Angeles, CA: CRESST.
- Kane, M. (2001). Current concerns in validity theory. *Journal of Educational Measurement*, 38(4), 319–342.
- Mislevy, R.J., Steinberg, L.S., Breyer, F.J., Almond, R.G., & Johnson, L. (2002). Making sense of data from complex assessment. *Applied Measurement in Education*, 15, 363-378
- Pellegrino, J., Chudowsky, N., & Glaser, R. (2001). *Knowing what students know: The science and design of educational assessment*. Washington, DC: National Academy Press.