

01 **Chapter 13**
 02 **Epilogue: Achieving Quality 21st Century**
 03 **Assessment**

AQ1

04
 05
 06
 07 **Betsy Jane Becker and Valerie J. Shute**
 08

AQ2

09
 10
 11
 12
 13 Writing a chapter to summarize all of the contents of this multifaceted book presents
 14 a challenge. Rather than create a laundry list of the conclusions drawn by our
 15 authors, we have tried to draw out three themes that appear across the works herein.
 16 These themes represent ideas that we believe will need careful attention from assess-
 17 ment experts, measurement professionals, teachers, principals, learning scientists
 18 and many others if the field is to move forward to develop better assessments
 19 that promote learning as well as provide fair means of accountability for students,
 20 teachers, and schools. We argue that

- 21
 22
 - Assessment must capitalize on advances in technology,
 - Assessment is a contextualized, social activity, and
 - Assessment must serve teaching and learning.
 23
 24

25
 26 We discuss in turn how each of these themes is raised by the authors of our book,
 27 and also touch on potential areas for research suggested by their work. We do not,
 28 however, mention every instance in which every author touches on these themes.
 29 We apologize if we have omitted points on these themes that are important to our
 30 authors.
 31

32
 33
 34 **13.1 Assessment Must Capitalize on Technology**
 35

36 Technological advances have already clearly affected the world of assessment in
 37 many ways. Even mundane components of assessment, such as the scoring of mul-
 38 tiple choice questions, were long ago made easier by Scantron machines and other
 39 scanning devices (Clarke, Madaus, Horn, & Ramos, 2000). However, while such
 40 devices enabled the rapid increase and widespread use of testing in the schools in the
 41

42

 43 B.J. Becker (✉)
 44 Educational Psychology and Learning Systems Department, Florida State University, Tallahassee,
 45 FL, USA
 e-mail: bbecker@fsu.edu

46 1960s and 1970s, they did not always lead to improvements in what we know about
47 students, or in what students learn (e.g., Epstein, Epstein, & Brosvic, 2001). Indeed,
48 while modern technologies clearly have made an impact on how testing is done, it
49 is clear they present both new challenges and new possibilities for assessment in the
50 twenty-first century (e.g., Naglieri et al., 2004).

51 Hickey, Honeyford, Clinton and McWilliams examine a context where technol-
52 ogy is inherently part of the assessment—the assessment of competence with new
53 media. New media literacies include activities as diverse as social networking, cre-
54 ation of fan fiction, music remixing, and blogging. Because nearly all new media are
55 based in technology, to assess competencies in these domains requires that technol-
56 ogy be embedded in the assessment process. However, it is also true that traditional
57 literacy skills—in writing, reading and spoken communication—are both needed for
58 and enhanced by use of new media (see also Leu, O’Byrne, Zawilinski, McVerry, &
59 Everett-Cacopardo, 2009). While endorsing the idea of assessing new media com-
60 petencies, Hickey and his colleagues raise considerable concerns about whether the
61 advent of accountability in this domain will narrow views of “proficiencies” to what
62 is easily measured. Indeed, since many new media skills are inherently social (see
63 our next theme), but most existing assessment systems are fundamentally individ-
64 ualized, clear tensions and conflicts will play out as assessment of these new skills
65 moves forward. Given the view of many scholars that media literacy is (and must
66 be) entwined with a participatory culture, such tensions will be a key concern for
67 the field. But as with many challenges, also see interesting possibilities for research.
68 For example, how should we best assess media literacy? Is it ever possible to gauge
69 individual contributions to fully participatory activities? Many interesting research
70 questions will emerge in relation to this context.

71 Several authors in our volume attend to the role technology must play in the
72 future of assessment. An argument for an elegant system of evidence centered
73 design (ECD) for assessment is given by Russell Almond. Almond lays out key
74 aspects of the ECD philosophy he has developed with collaborators Robert Mislevy
75 and Linda Steinberg. His chapter makes concrete how a complex mathematical
76 modeling framework can be combined with thoughtful consideration of the skills
77 desired of an examinee population to produce both improved learning and qual-
78 ity assessment all in one comprehensive system. The system relies on technology
79 in its fundamental use of a Bayesian framework for evaluating student capacities.
80 Information based on prior knowledge of examinee capabilities, along with data
81 from observable events associated with a collection of tasks is fed back into a
82 system to create posterior distributions of (hopefully changed) student skill lev-
83 els. Almond also points out that an ECD system could aim at tracking growth in
84 multiple competencies, based on related process or product observables, and can
85 even interface with automated scoring systems (like those described by Shermis
86 for essay scoring). Almond argues that eventually ECD, combined with modern
87 technologies could support “seamless” (text page 30) collection of observables
88 embedded in ongoing work—assessment that would seem so natural students would
89 not realize that it had even occurred (see also Shute, in press for more on this
90 topic).

13 Epilogue: Achieving Quality 21st Century Assessment

91 Clearly such assessment would not cause the stress and disruptions present in so
92 much high-stakes testing (e.g., Cizek & Burg, 2005; Suen & Yu, 2006).

93 While they endorse many of Almond's ideas, commenters Ellington and Verges
94 point out that currently schools do not have the infrastructure to support these inno-
95 vations. For example, problems and incompatibilities in hardware and software
96 could lead to glitches with data collection from diverse sources. Also while they
97 acknowledge the benefits of using complex tasks in assessments, they raise several
98 practical problems (such as cost and extensive field test requirements) that would
99 limit implementation especially in the current economic climate. Clearly such a
100 system for assessment does not currently exist, but the goal of realizing it presents
101 another set of fascinating possibilities for research.

102 A more conventional take on technology and assessment comes from Mark
103 Shermis. His chapter describes how technology can assist with improved student
104 learning of writing skills when the electronic scoring of essays is built into a system
105 of writing improvement, along with revision, feedback, and teacher participation.
106 Shermis describes how automated essay scoring or AES works, from the develop-
107 ment of "proxes" or features used to represent the quality of a writing sample, to
108 the evaluation of the statistical models used in algorithmic scoring. Shermis cites
109 evidence that with careful development and good rubrics, an AES system is at least
110 as reliable as human raters, and can in some cases help avoid biases that human
111 raters cannot seem to eliminate from their rating behavior. One goal of this chapter
112 is to simply describe how AES works, and Shermis further illustrates that by way of
113 a detailed description of the "Intellimetric" system. However Shermis goes farther,
114 and makes the controversial claim that automated essay scoring (and the teaching
115 structures associated with its ongoing use) can replace high-stakes writing assess-
116 ments. Consistent with Almond's arguments for incorporating multiple pieces of
117 evidence in on going assessment system, Shermis argues that an assessment based
118 on multiple instances of writing (such as essays produced throughout the year)
119 would provide a more useful and valid evaluation of student writing than a single-
120 occasion high stakes test. He describes how an integrated writing assessment system
121 could also support instruction by providing feedback aimed at each essay produced
122 by a student. In addition, the release of "used" writing prompts could provide mate-
123 rials for ongoing instruction. Finally he argues that all of this can be accomplished
124 at costs lower than those incurred with human scoring of similar writing products.
125 This is an excellent example of how capitalizing on technology can enhance not
126 only the assessment itself but also student learning. And as with our other chap-
127 ters that touch on technology, various ideas for research arise from Shermis's work.
128 For example, what measurement models best suit this kind of assessment system?
129 How often would one need to take samples of writing throughout the year to get
130 solid information about change in writing competence? Could an AES system have
131 a built-in way to assess the impact of particular kinds or amounts of teacher (or sys-
132 tem) feedback? How would we assess whether the system is providing appropriate
133 feedback to students (i.e., what kinds of "quality control" would be needed)? These
134 and other practical and theoretical questions provide a rich set of ideas for those
135 interested in the future of assessment.

13.2 Assessment is a Contextualized, Social Activity

The strongest advocate of the social view of assessment in our collection of authors is James Gee. Gee lays out the case for assessment of twenty-first century skills in domains as part of social “appreciative systems”. Loosely these are sets of conventions, values—perhaps even rules—for what is acceptable or valuable in a certain domain. He argues that appreciative systems are shared across people, and gives quite a few examples of how such systems develop. He also argues that most assessment goes on as a part of ongoing human interactions and activity, and is not—and does not need to be—formalized. King Beach, who comments on Gee’s chapter, provides several examples of such real-life assessment in out of school settings in western Nepal.

Gee also contends that groups themselves can formalize assessment in a quite natural way. He argues that this often occurs via “Pro-Am” communities, or groups of “... innovative, committed and networked amateurs working to professional standards” (Leadbeater & Miller, 2004, p. 9). A compelling example comes from his research on online communities (Gee & Hayes, in press) where a young girl learned to create virtual clothing for the virtual world *Second Life*. Eventually various discoveries lead her to provide clothes to virtual people on the Internet, first for free then later at a price, after she realized that the “appreciative” community of virtual shoppers highly valued her product. More controversially, Gee goes on to argue that schools could promote twenty-first century skills by encouraging and equipping students to become high-status members of Pro-Am communities that value and promote such skills.

Gee notes that society has in some cases formalized assessment by removing it from such Pro-Am communities. Institutions, including schools, have been created in support of this formalization. But Gee believes this kind of assessment is “backwards” (text page 28) because it occurs at such an abstract, disembedded level, removed from real problems. His arguments for authentic assessments echo in part those of measurement scholars (e.g., Wiggins, 1990) and others (e.g., Darling-Hammond & Snyder, 2000) who called for more realism and context in assessment years ago.

Finally, Gee raises the radical idea that formalized assessments may not be needed if existing communities (like the online buyers of simulated clothing) have already assessed and accepted someone’s skills. He states that, “The job of educators ought to be designing such social organizations and letting them run.” (text page 32) Again we see plenty of possible research avenues in this work, and Beach’s commentary raises one interesting question—what is the appropriate unit of assessment in such circumstances? Could we find a way to examine the developmental relationship between a learner and the domain to be assessed, over time? Could such a complex entity be assessed by the “indigenous workings” (Gee page 32) of the group, as advocated by Gee? It is intriguing to consider how one could study and obtain empirical evidence on such a system of assessment.

Hickey and colleagues also deal with the social context for learning, which is fundamentally a part of examining competencies with “new media”. They begin by

13 Epilogue: Achieving Quality 21st Century Assessment

181 examining the positions of groups like the National Council of Teachers of English
182 and the National Writing Project, which argue that writing is an inherently social
183 activity. This is a launching point for the efforts of Hickey and his colleagues to
184 create a cohesively designed assessment framework for language arts. Their efforts
185 have used not only a classic text (*Moby Dick*) but also new media sources such as
186 related music videos and theatre “re-mixes”. They argue that by using multiple lev-
187 els of assessment (designed after Ruiz-Primo et al. 2002 five levels) the teacher can
188 focus initially and closely on social and interactive aspects of writing, with a quick
189 time frame for feedback and a very informal assessment context. The assessments
190 and classroom activities then move to more individualized, but distant aspects at
191 other levels. Their chapter gives many examples of assessment activities at their five
192 levels, but the contextualized nature of assessment is most evident at their first two
193 levels. Only at the highest level would fully individualized assessments (e.g., exter-
194 nal tests with essay items) be used to tap into an abstract context, and might, for
195 instance, measure performance on content standards such as those at a state level.

196 We consider one example from their second level—“close-level activity-oriented
197 reflections” (page 29 of text). Here the focus is on discussion questions presented
198 to the class, either orally, in written form, or online. The authors found that the
199 informal nature of these assessments and an attendant emphasis on “communal dis-
200 cussion” of ongoing classroom activities enhanced broader student participation.
201 The teacher then could consider the nature of the discussion, and focus further
202 discussion in ways that led students “. . . to create more compelling and creative arti-
203 facts” (page 30 of text). This is an interesting example of how a teacher can assess
204 student understanding in a group context, provide feedback based on questions that
205 are less formal, but very targeted to the activities of the learners.

206 Allan Jeong’s work in examining conceptual causal maps provides a nice exam-
207 ple of how learning can be enhanced by the use of social interactions and by
208 assessments that provide a window on the understandings of experts and of other
209 learners. Jeong describes the use of jMAP, a program that enables learners to create
210 and evaluate causal maps, and also allows for the assessment of changes in maps
211 over time. Students use jMAP to identify critical components of a causal system,
212 to draw interconnections among those components, and to identify the strength of
213 those connections. The connection to “others” is embedded in the structure and use
214 of jMAP. Specifically, after creating their own maps, students can be exposed to the
215 maps of experts, to discussions about the nature of the causal connections (as was
216 the case in Jeong’s studies), and to composite maps made by aggregating the maps
217 of various subsets of learners or all learners in a class. Subsequent rounds of maps
218 can be drawn and changes in the maps of learners can be examined. Jeong shows
219 how comparisons of individual versus aggregate maps can lead to changes in subse-
220 quent maps of individual learners. Also the nature of discussions held (e.g., whether
221 links in the maps are supported or challenged by others, whether explanations are
222 provided, etc.) can impact how learners change their maps on subsequent drawings.

223 For those interested in how learner interactions can impact learning, jMAP pro-
224 vides intriguing tools for analysis of causal maps. We challenge those interested in
225 such phenomena to consider how learning in other domains of understanding (i.e.,

226 other than causal maps) might be represented (using technology) and then exam-
227 ined for changes due to exposure to the knowledge and ideas of others. Domains
228 suggested by the work of our other authors might include writing, where automated
229 scoring systems like those described by Shermis could produce indices of change
230 following peer or teacher feedback of different sorts, or Gee’s “Pro-Am skills”
231 where changes in the products or skill sets of members of Pro-Am communities
232 might be evaluated for evolution as individuals receive feedback from relevant com-
233 munity members. Roehrig and Christesen’s Classroom AIMS is another assessment
234 device where modal (typical) or expert performance ratings could be shared with
235 teachers, and then further measurement instances could be evaluated for change due
236 to those different kinds of feedback.

237
238
239

240 **13.3 Assessment Must Serve Teaching and Learning**

241

242 Mari Pearlman views assessment as integral to the educational enterprise. In her
243 chapter she argues that we must “. . .ally assessment with instruction” (page 20 of
244 text version) and she makes a provocative case for using an assessment based on
245 the vast architecture of the National Assessment of Educational Progress or NAEP
246 as a vehicle to accomplish this goal. Pearlman outlines a variety of problems in our
247 educational system, and among them she lists a need to align all components of
248 the system—curriculum, teacher practice, teacher preparation and assessments—in
249 concert to move student learning ahead. She argues that to date these components
250 have been manipulated by the states, but mainly in efforts to achieve “AYP” or
251 adequate yearly progress, not in efforts to increase learning of content identified
252 as important (See for example Kane, Staiger, & Geppert’s, 2002 views on gam-
253 ing the AYP system so states can look most successful). Pearlman endorses the
254 move towards national standards, and argues that having clear frameworks could
255 lead to national “conversations” about content, curriculum and most importantly
256 equity across states. Equity is an issue conveniently (but sadly) avoided in the cur-
257 rent context, where each state can use its own distinct tests to assess progress, and
258 also can set different goals. In such a context, it is not hard to see that all states
259 could theoretically measure up as “adequate” while in fact being quite different in
260 terms of what their students actually learn. Initial examinations however showed that
261 some states did the opposite—setting such high standards that they were virtually
262 unreachable (Linn, Baker, & Betebenner, 2002).

263 In terms of the theme of assessment in support of learning, Pearlman chal-
264 lenges us to find ways to use the many items developed as part of the
265 NAEP assessments in support of learning. Indeed, thousands of released NAEP
266 items are currently available for public use via the NAEP Questions Tool
267 (see <http://nces.ed.gov/nationsreportcard/itmrlsx/landing.aspx>), and can be used by
268 teachers and others in a variety of ways. Much like the position taken by Almond,
269 Pearlman argues that we must first know what we want students to learn, then what
270 we want to measure, then only last can we design lessons and activities to support

13 Epilogue: Achieving Quality 21st Century Assessment

271 those goals. Pearlman argues that if benchmarks such as those set for NAEP were
272 used as goals, and NAEP’s current plan of sampling students were replaced with
273 every-pupil-testing, we would soon move towards better outcomes and towards
274 equity across states in terms of student learning.

275 In her response to Pearlman, Lynn Wicker emphasizes a point made by Pearlman.
276 That is, both agree that the culture of schools relative to assessment must change
277 so that a stronger link to learning can be achieved. Wicker argues that the use of
278 multiple ongoing assessments tied with targeted interventions “should be viewed as
279 a non-negotiable in the learning process” (page 6 of text version). We also strongly
280 endorse this view.

281 Martineau and Dean give perhaps the most comprehensive proposal in our vol-
282 ume for how assessment, in many forms and at many levels, can serve instruction
283 and learning. They draw on the idea of balanced assessment (Redfield, Roeber,
284 Stiggins, & Philip, 2008) and elaborate it to describe how formative, summative and
285 interim assessments can be used to provide both indices for accountability as well
286 as detailed input to teachers in support of their instruction. Their system begins with
287 clear and focused K-12 content standards, aimed at supporting students’ progres-
288 sion towards specified high-school outcomes. They argue that curriculum materials
289 (“model curriculum units”) can then be developed in support of these standards, and
290 be made available to all teachers (but not mandated as a required curriculum). This
291 set of content standards and materials would be paired with professional develop-
292 ment for teachers aimed at helping them to understand the content standards and
293 how they can be used, but more critically how to use data from assessments—both
294 classroom assessments and more formalized “secure” assessments to modify their
295 instruction.

296 A requirement that teachers receive instruction in classroom assessment is
297 already a part of the Florida Department of Education’s preservice teaching require-
298 ments. Martineau and Dean want continued support in the form of professional
299 development for both teachers and other administrators. They also argue that
300 accountability purposes can also be served by their system, but only with a mul-
301 tifaceted system in which teachers, teacher preparation institutions, administrators
302 and students are all held to account.

303 Working from a very different perspective, Hickey and colleagues also describe a
304 system of assessment that focuses on different kinds and levels of assessment in the
305 context of new media literacy. Based in the technological context discussed earlier,
306 the first three of their five levels of assessment—the immediate, close and proximal
307 levels—are completely entwined with student products (artifacts) and interactions
308 with teachers and other students as they create those products in the classroom. At
309 the third (proximal) level the assessment tasks move towards more individualiza-
310 tion, but even at this level reflection questions and student products together allow
311 for targeted teacher feedback. Their system illustrates how assessment can be tied
312 directly into ongoing instruction.

313 A more targeted approach to both assessing and supporting instruction is
314 described by Roehrig and Christesen. They describe the development and use
315 of the Classroom AIMS instrument, which examines how teachers create a

316 positive classroom atmosphere, implement instruction and classroom management,
317 and engage students in learning. Roehrig and Christesen start with the premise that
318 teachers' behaviors and activities are more likely to predict student outcomes than
319 teacher characteristics such as teacher knowledge. They describe how a set of exem-
320 plary behaviors were identified by observing teachers who had produced students
321 with strong learning gains in reading and writing performance. More importantly for
322 this "theme" of our book, the authors go on to describe how the AIMS instrument
323 can be used not just to observe the current behaviors of teachers, but also to diagnose
324 possible areas for teachers to improve. They describe the use of the AIMS with pre-
325 service teachers as well as practicing teachers working with mentors. Other research
326 questions could be asked about the Classroom AIMS. How does AIMS function as
327 a measure of the effects of professional development for practicing teachers? Can it
328 differentiate between more and less effective interventions? How well does it work
329 across different subject areas? These and other questions may be fertile areas for
330 future research.

331

332

333 **13.4 Conclusion**

334

335 As mentioned in the Prelude, our goal for the symposium and this book was to
336 bring together groups of individuals who normally do not converse, but who we
337 believe should communicate—researchers from different areas, policymakers, and
338 educational professionals. In all chapters, the call for educational reform is clear and
339 there is no shortage of problems to be addressed with innovative thinking and high-
340 quality research. The linchpin for such reform—reform that aims to fully support
341 students' success in the twenty-first Century—is assessment.

342 The chapters in this book present a broad swath of assessment issues and pos-
343 sible solutions, and embrace three main theses: (a) assessment must capitalize on
344 advances in technology, (b) assessment is a contextualized, social activity, and (c)
345 assessment must serve teaching and learning. Each of these alone can move the
346 assessment conversation and ensuing research forward, but we contend that when
347 these issues are considered collectively, important breakthroughs in assessment and
348 educational reform will surely follow.

349

350

351

352 **References**

353

- 354 Cizek, G. J., & Burg, S. S. (2005). *Addressing test anxiety in a high-stakes environment: Strategies*
355 *for classrooms and schools*. Thousand Oaks, CA: Corwin Press.
- 356 Clarke, M. M., Madaus, G. F., Horn, C. L., & Ramos, M. A. (2000). Retrospective on educational
357 testing and assessment in the 20th century. *Journal of Curriculum Studies*, 32(2), 159–181.
- 358 Darling-Hammond, L., & Snyder, J. (2000). Authentic assessment of teaching in context. *Teaching*
359 *and Teacher Education*, 16(5–6), 523–545.
- 360 Epstein, M. L., Epstein, B. B., & Brosvic, G. M. (2001). Immediate feedback during academic
testing. *Psychological Reports*, 88(3), 889–894.

13 Epilogue: Achieving Quality 21st Century Assessment

- 361 Kane, T. J., Staiger, D. O., & Geppert, J. (2002). Randomly accountable. *Education Next*, 2(1),
 362 Downloaded 12/16/09 from <http://educationnext.org/randomly-accountable/>
- 363 Leadbeater, C., & Miller, P. (2004). *The Pro-Am Revolution: How enthusiasts are chang-*
 364 *ing our economy and society*. Demos. Downloaded on January 12, 2010 from <http://www.demos.co.uk/files/proamrevolutionfinal.pdf?1240939425>
- 365 Leu, D. J., O'Bryne, I., Zawilinski, J., McVerry, J. G., & Everett-Cacopardo, H. (2009). Expanding
 366 the new literacies conversation. *Educational Researcher*, 38(4), 264–269.
- 367 Linn, R. L., Baker, E. L., & Betebenner, D. W. (2002). Accountability systems: Implications of
 368 requirements of the No Child Left Behind Act of 2001. *Educational Researcher*, 31(6), 3–16.
- 369 Naglieri, J. A., Drasgow, F., Schmit, M., Handler, L., Prifitera, A., Margolis, A., et al. (2004).
 370 Psychological testing on the Internet: New problems, old issues. *American Psychologist*, 59(3),
 150–162.
- 371 Redfield, D., Roeber, E., Stiggins, R., & Philip, F. (2008). *Building balanced assessment systems*
 372 *to guide educational improvement*. A background paper for the keynote panel presentation at
 373 the National Conference on Student Assessment, June 15, 2008, Orlando, FL. Downloaded on
 374 December 15, 2009, from <http://www.ccsso.org/content/PDFs/OpeningSessionPaper-Final.pdf>
- 375 Ruiz-Primo, M. A., Shavelson, R. J., Hamilton, L., & Klein, S. (2002). On the evaluation of sys-
 376 temic science education reform: Searching for instructional sensitivity. *Journal of Research in*
Science Teaching, 39(5), 369–393.
- 377 Suen, H. K., & Yu, L. (2006). Chronic consequences of high-stakes testing? Lessons from the
 378 Chinese Civil Service Exam. *Comparative Education Review*, 50(1), 46–65.
- 379 Wiggins, G. (1990). *The name assigned to the document by the author. This field may also con-*
 380 *tain sub-titles, series names, and report numbers. The case for authentic assessment*. ERIC
 381 Digest. ERIC Document Reproduction Service Number ED328611. Washington, DC: ERIC
 Clearinghouse on Tests Measurement and Evaluation.
- 382 Gee, J. P., & Hayes, E. R. (in press). *Women and gaming: The Sims and 21st century learning*.
 383 New York: Palgrave/McMillan.
- 384 Shute, V. J. (in press). Stealth assessment in computer-based games to support learning. To
 385 appear In S. Tobias & J. D. Fletcher (Eds.), *Computer games and instruction*. Charlotte, NC:
 386 Information Age Publishers.

AQ3

AQ4

386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405

Chapter 13

406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450

Q. No.	Query
AQ1	There is mismatch in the chapter title between the chapter opening page and contents. We have followed as per Chapter opening page. Please confirm
AQ2	Please provide Abstract and Keywords for this chapter in order to maintain book consistency
AQ3	Please update the year for the reference “Gee and Hayes, in press”
AQ4	Please update the year for the reference “Shute, in press”

UNCORRECTED PROOF