

Evaluating the Content and Quality of Next Generation Assessments

NANCY DOOREY AND MORGAN POLIKOFF

Foreword by Amber M. Northern and Michael J. Petrilli



The Thomas B. Fordham Institute is the nation's leader in advancing educational excellence for every child through quality research, analysis, and commentary, as well as on-the-ground action and advocacy in Ohio. It is affiliated with the Thomas B. Fordham Foundation, and this publication is a joint project of the Foundation and the Institute. For further information, please visit our website at www.edexcellence.net or write to the Institute at 1016 16th St. NW, 8th Floor, Washington, D.C. 20036. The Institute is neither connected with nor sponsored by Fordham University.

Contents

<i>Foreword</i>	4
<i>Executive Summary</i>	11

Introduction **25**

Overview	25
The Four Assessments	26
Study Design and Key Questions	26
Organization of the Report	27
Approach to Alignment and Quality	27
Test Program Differences	28

Section I: Assessments of English Language Arts/Literacy and Mathematics **31**

Overview of the Methodology	31
Study Phases	31
Selection of Test Forms	33
Selection of Review Panels and Assignment to Forms	34
The Study Criteria	36
Methodology Modifications	38
Findings	40
ELA/Literacy Content Criteria	41
ELA/Literacy Depth Criteria	47
Mathematics Content Criteria	54
Mathematics Depth Criteria	57
Criterion Level Ratings	63
Summary Ratings	64
Program Strengths and Areas for Improvement	66

Section II: Recommendations **70**

For State Policymakers	70
For Test Developers	72

Section III: Suggestions for Methodological Improvement **74**

English Language Arts/Literacy	74
Mathematics	77
Across Subjects	79
Summary	79

Appendices **80**

Appendix A: Depth of Knowledge (DOK) of the Four Assessment Programs as Compared to the CCSS and Other Policy-Relevant Assessments	80
Appendix B: Key Terminology	82
Appendix C: The Methodology as Written	85
Appendix D: Author Biographies	93
Appendix E: Review Panelist Biographies	94
Appendix F: Full ELA/Literacy and Math Ratings and Summary Statements (Grades 5 and 8)	100
Appendix G: Testing Program Responses to Study and Descriptions of Test Changes for 2015–2016	113

Foreword

By Amber M. Northern and Michael J. Petrilli

We at the Thomas B. Fordham Institute have been evaluating the quality of state academic standards for nearly twenty years. **Our very first study**, published in the summer of 1997, was an appraisal of state English standards by Sandra Stotsky. Over the last two decades, we’ve regularly reviewed and reported on the quality of state K–12 standards for **mathematics, science, U.S. history, world history, English language arts, and geography**, as well as the **Common Core, International Baccalaureate, Advanced Placement** and other influential standards and frameworks (such as those used by **PISA, TIMSS, and NAEP**). In fact, evaluating academic standards is probably what we’re best known for.

For most of those two decades, we’ve also dreamed of evaluating the tests linked to those standards—mindful, of course, that in most places the tests are the real standards. They’re what schools (and sometimes teachers and students) are held accountable to and they tend to drive actual curricula and instruction. (That’s probably the reason we and other analysts have never been able to demonstrate a close relationship between the quality of standards per se and changes in student achievement.) We wanted to know how well aligned the assessments were to the standards, whether they were of high quality, and what type of cognitive demands they placed on students.

But with fifty-one different sets of tests, such an evaluation was out of reach—particularly since any bona fide evaluation of assessments must get under the hood (and behind the curtain) to look at a sizable chunk of actual test items. Getting dozens of states—and their test vendors—to allow us to take a peek was nigh impossible.

So when the opportunity came along to conduct a groundbreaking evaluation of Common Core-aligned tests, we were captivated. We were daunted too, both by the enormity of the task and by the knowledge that our unabashed advocacy of the standards would likely cause any number of doubters and critics to sneer at such an evaluation coming from us, regardless of its quality or impartiality.

So let’s address that first. It’s true that we continue to believe that children in most states are better off with the Common Core standards than without them. If you don’t care for the standards (or even the concept of “common” standards), or perhaps you come from a state that never adopted these standards or has since repudiated them, you should probably ignore this study. Our purpose here is not to re-litigate the Common Core debate. Rather, we want to know, for states that are sticking with the common standards, whether the “next generation assessments” that have been developed to accompany the standards deliver what they promised by way of strong content, quality, and rigor.

It is also true that the study was funded by a number of foundations that care about assessment quality and the Common Core (seven, in fact, including the Fordham Institute’s own foundation). If you think that big private foundations are ruining public education, this study is also not for you.

Now is an especially opportune time to look closely at assessments, since the national testing landscape is in a state of flux. According to the Education Commission of the States, as of October 2015, six states and the District of Columbia planned to administer the Partnership for Assessment of Readiness for College and Careers (PARCC)

test in 2015–16 and fifteen states will deploy the Smarter Balanced Assessment Consortium (Smarter Balanced) test.¹ At least twenty-five others will administer state-specific assessments in math and English language arts. Some (Florida, Ohio, and Utah) will use tests developed by the American Institutes for Research (AIR); others (Indiana, Kentucky, and Virginia) are using Pearson-developed products; still others are choosing “blended” versions of consortium and state-developed items (Michigan and Massachusetts). A handful are undecided and currently in the midst of evaluating test vendors through their RFP process (Maine, Louisiana, and South Carolina). About half the states also require an additional assessment for college admissions, such as the ACT or SAT, which is generally administered in grade 11 (and sometimes statewide). And let’s not forget that the new SAT will be unveiled in **March 2016**.

Hence there’s no way any single study could come close to evaluating all of the products in use and under development in today’s busy and fluid testing marketplace. But what we were able to do was to provide an in-depth appraisal of the content and quality of three “next generation” assessments—ACT Aspire, PARCC, and Smarter Balanced—and one best-in-class state test, the Massachusetts Comprehensive Assessment System (MCAS, 2014). In total, over thirteen million children (about 40 percent of the country’s students in grades 3–11) took one of these four tests in spring 2015. Of course it would be good to encompass even more. Nevertheless, our study ranks as a major accomplishment—as well as possibly the most complex and ambitious single project ever undertaken by Fordham.

After we agreed to myriad terms and conditions, we and our team of nearly forty reviewers (more about them below) were granted secure access to operational items and test forms for grades 5 and 8 (the elementary and middle school capstone grades that are this study’s focus).²

This was an achievement in its own right. It’s no small thing to receive access to examine operational test forms. This is especially true in a divisive political climate where anti-testing advocates are looking for *any* reason to throw the baby out with the bathwater and where market pressure gives test developers ample reason to be wary of leaks, spies, and competitors. Each of the four testing programs is to be commended for allowing this external scrutiny of their “live” tests—tests that cost them much by way of blood, sweat, tears, and cash to develop and bring to market. They could have easily said “thanks, but no thanks.” But they didn’t.

Part of the reason they said yes was the care we took in recruiting smart, respected individuals to help with this project. Our two lead investigators, Nancy Doorey and Morgan Polikoff, together bring a wealth of experience in educational assessment and policy, test alignment, academic standards, and accountability. Nancy has authored reports for several national organizations on advances in educational assessment and she co-piloted the Center for K–12 Assessment and Performance Management at ETS. Morgan is assistant professor of education at the University of Southern California and a well-regarded analyst of the implementation of college and career readiness standards and the influence of curricular materials and tests on that implementation. He is an associate editor of the *American Educational Research Journal*, serves on the editorial board for *Educational Administration Quarterly*, and is the top finisher in the RHSU 2015 Edu-Scholar rankings for junior faculty.³

Nancy and Morgan were joined by two well-respected content experts who facilitated and reviewed the work of the ELA/Literacy and math review panels. Charles Perfetti, Distinguished University Professor of Psychology at University of Pittsburgh, served as the ELA/Literacy content lead, and Roger Howe, Professor of Mathematics at Yale, served as the math content lead.

1. J. Woods, “State Summative Assessments: 2015–2016 School Year” (Denver, CO: Education Commission of the States, 2015), <http://www.ecs.org/ec-content/uploads/12141.pdf>. According to ECS, fifteen states are members of the Smarter Balanced Assessment Consortium, and all but one plan to administer the full assessment in grades 3–8 math and English language arts.

2. The study targets “summative,” not “formative,” assessments, though most of these same test developers also make the latter available.

3. R. Hess, “2016 RHSU Edu-Scholar Public Influence: Top Tens,” *Education Week* (blog), January 7, 2016, http://blogs.edweek.org/edweek/rick_hess_straight_up/2016/01/2016_rhsu_edu-scholar_public_influence_top_tens.html.

Given the importance and sensitivity of the task at hand, we spent months recruiting and vetting the individuals who would eventually comprise the panels led by Dr. Perfetti and Dr. Howe. We began by soliciting recommendations from each participating testing program and other sources, including content and assessment experts, individuals with experience in prior alignment studies, and several national and state organizations. Finalists were asked to submit CVs and detailed responses to a questionnaire regarding their familiarity with the Common Core, their prior experience in conducting alignment evaluations, and any potential conflicts of interest. Individuals currently or previously employed by participating testing organizations and writers of the Common Core were not considered. Given that most card-carrying experts in content and assessment have earned their experience by working on prior alignment or assessment-development studies, and that it's nearly impossible to find experts with zero conflicts, we prioritized balance and fairness. In the end, we recruited at least one reviewer recommended by each testing program to serve on each panel; this strategy helped to ensure fairness by equally balancing reviewer familiarity with the various assessments. (Their bios can be found in Appendix E.)

Which brings us to the matter at hand: How did our meticulously assembled panels go about evaluating the tests—and what did they find? You can read plenty on both questions in the Executive Summary and report itself, which includes ample detail about the study design, testing programs, criteria, and selection of test forms, and review procedures, among other topics.

But the short version is this: we deployed a brand new methodology developed by the Center for Assessment to evaluate the four tests—a methodology that was itself based on the Council of Chief State School Officers' 2014 **“Criteria for Procuring and Evaluating High-Quality Assessments.”** Those criteria, say their authors, are “intended to be a useful resource for any state procuring and/or evaluating assessments aligned to their college and career readiness standards.” This includes, of course, tests meant to accompany the Common Core standards.

About Those Criteria...

The CCSSO Criteria address the “content” and “depth” of state tests in both English language arts and mathematics. For ELA, “content” spans topics such as whether students are required to use evidence from texts; for math, they are concerned with whether the assessments focus strongly on the content most needed for success in later mathematics. The “depth” criteria for both subjects include whether the tests required a range of “cognitively demanding,” high-quality items that make use of various item types (e.g., multiple choice, constructed response, etc.), among other things.

The Center for Assessment took these criteria and transformed them into a number of measurable elements that reviewers addressed. In the end, the newly minted methodology wasn't perfect. Our rock-star reviewers improved upon it and wanted others following in their footsteps to benefit from their learned experience. So we made adjustments along the way (see Section I, *Methodology Modifications* for more).

The panels essentially evaluated the extent of the match between the assessment and a key element of the CCSSO document. They assigned one of four ratings to each ELA and math-specific criterion, such that tests received one of four “match” ratings: Excellent, Good, Limited/Uneven, or Weak Match. To generate these marks, each panel reviewed the ratings from the grade 5 and grade 8 test forms, considered the results from the analysis of the program's documentation (which preceded the item review), and came to consensus on the rating.

What did they ultimately find? The summary findings appear below.

TABLE F 1

Overall Content and Depth Ratings for ELA/Literacy and Mathematics

	ACT Aspire	MCAS	PARCC	Smarter Balanced
	L	L	E	E
	G	G	E	G
	L	L	G	G
	G	E	G	G

LEGEND E Excellent Match G Good Match L Limited/Uneven Match W Weak Match

As shown, the PARCC and Smarter Balanced assessments earned an Excellent or Good Match to the subject-area CCSSO Criteria for both ELA/Literacy and mathematics. This was the case with both Content and Depth.

ACT Aspire and MCAS (along with the others) also did well regarding the quality of their items and the depth of knowledge assessed (both of which are part of the Depth rating). But the panelists also found that they did not adequately assess—or in some cases did not really assess at all—some of the priority content in both ELA/Literacy and mathematics at one or both grade levels in the study (Content).

What do we make of these bottom-line results? Simply put, developing a test—like all major decisions and projects in life—is full of trade-offs. PARCC and Smarter Balanced are a better match to the CCSSO criteria, which is not surprising, given that they were both developed with the Common Core in mind. ACT Aspire, on the other hand, was not developed for that explicit purpose. In a paper on their website, ACT officials Sara Clough and Scott Montgomery explain that ACT Aspire was

under development prior to the release of the Common Core State Standards [and] not designed to directly measure progress toward those standards. However, since ACT data, empirical research, and subject matter expertise about what constitutes college and career readiness was lent to the Common Core development effort, significant overlap exists between the Common Core State Standards and the college and career readiness constructs that ACT Aspire and the ACT measure.⁴

Our reviewers also found some “overlap” in MCAS given that the state had added new Common Core items to its 2014 test. Yet the Bay State’s intention was not a full redesign, particularly since it was then in the midst of deciding between MCAS and PARCC as its test of choice (the state ultimately decided on a hybrid).⁵ To the extent that states want their tests to reflect the grade-level content in the new standards, they should choose accordingly.

The CCSSO Criteria do not consider testing time, cost, or comparability. But those are nonetheless key considerations for states as they make assessment decisions. Although PARCC and Smarter Balanced are a better match to the Criteria, they also take longer to administer and are more expensive. The estimated testing time for

4. S. Clough and S. Montgomery, “How ACT Assessments Align with State College and Career Readiness Standards” (Iowa City, IA: ACT, 2015), http://www.discoveractaspire.org/pdf/ACT_Alignment-White-Paper.pdf.

5. J. Fox, “Education Board Votes to Adopt Hybrid MCAS-PARCC Test,” *Boston Globe*, November 17, 2015, <https://www.bostonglobe.com/metro/2015/11/17/state-education-board-vote-whether-replace-mcas/aex1nGyBYZW2sucEW2o82L/story.html>. To the extent that states want their tests to reflect the grade-level content in the new standards, they should choose accordingly.

students in grades 5 and 8, on average, to complete both the ELA/Literacy and mathematics assessments for all four programs is as follows:

- ◆ ACT Aspire: three to three and one-quarter hours for all four tests (English, reading, writing, and mathematics)
- ◆ MCAS 2014: three and a half hours
- ◆ PARCC: seven to seven and a half hours⁶
- ◆ Smarter Balanced: five and a half hours

The longer testing times for PARCC and Smarter Balanced are due primarily to their inclusion of extended performance tasks. Both programs use these tasks to assess high-priority skills within the CCSS, such as the development of written compositions in which a claim is supported with evidence drawn from sources; research skills; and solving complex multi-step problems in mathematics. In addition to requiring more time than multiple-choice items, these tasks are also typically costlier to develop and score.⁷

Another trade-off pertains to inter-state comparability. Some states want the autonomy and uniqueness that come with having their own state test developed by their own educators. Other states prioritize the ability to compare their students with those in other states via a multi-state test. We think the extra time and money,⁸ plus the comparability advantage, are trade-offs worth making, but we can't pretend that they're not tough decisions in a time of tight budgets and widespread anxiety about testing burden.

Of course we're mindful—as anyone in this field would be—of the recent backlash to testing and the so-called “opt-out movement.” We understand that some local and state officials are wary of adopting longer tests. We also suspect that most of the concerns that parents have isn't with the length of one test in May, but with the pressure that educators feel to teach to the test and narrow the curriculum.

If we're right and that's the real problem, the answer is stronger tests, which encourage better, broader, richer instruction, and which make traditional “test prep” ineffective. Tests that allow students of all abilities, including both at-risk and high-achieving youngsters, to demonstrate what they know and can do. More rigorous tests that challenge students more than they've been challenged in the past. But, again, those tests tend to take a bit longer (say, five hours rather than two and a half hours) and cost a bit more. Our point is not to advocate for any particular tests but to root for those that have qualities that enhance, rather than constrict, a child's education and give her the best opportunity to show what she's learned.

A discussion of such qualities, and the types of trade-offs involved in obtaining them, are precisely the kinds of conversations that merit honest debate in states and districts.

We at Fordham don't plan to stay in the test-evaluation business. The challenge of doing this well is simply too overwhelming for a small think tank like ours. But we sincerely hope that others will pick up the baton, learn from

6. The 2015–16 PARCC revisions will reduce this time by an estimated one and a half hours.

7. That said, Matthew Chingos, in a 2012 study on state assessment spending, found that “collaborating with other states achieves cost savings simply by spreading fixed costs over more students...” (page 22). See M. Chingos, “Strength in Numbers: State Spending on K–12 Assessment Systems” (Washington, D.C.: Brookings Institution, November 29, 2012), <http://www.brookings.edu/research/reports/2012/11/29-cost-of-ed-assessment-chingos>.

8. Note that the per-pupil costs for PARCC, Smarter Balanced, and ACT Aspire are in the same ballpark, ranging from roughly \$22 to \$25 depending on the tested subjects. The MCAS, typically viewed as a higher-quality state test, costs \$42 per student. The costs associated with many of the prior state tests were considerably lower than these figures so changing tests represented an increase for them. See M. Chingos, “Strength in Numbers.” Cost estimates for PARCC and Smarter Balanced can be found here: <http://www.parcconline.org/news-and-video/press-releases/248-states-select-contractor-to-help-develop-and-implement-parcc-tests>; <http://www.smarterbalanced.org/faq/7-what-does-it-cost-for-a-state-to-participate-in-smarter-balanced-assessment-consortium/>. Per MCAS, “The approximate cost of the legacy MCAS assessment is \$42 per student for ELA and mathematics per estimates presented to the Massachusetts State Board of Elementary and Secondary Education in fall 2015” (personal email communication with Michol Stapel, January 22, 2016). Per ACT Aspire, “The estimated price for 2016 is \$25 per student and includes English, Mathematics, Reading, Writing, and Science subject tests” (personal email communication with Elizabeth Sullivan, January 21, 2016).

our experience, and provide independent evaluations of the assessments in use in the states that have moved away from PARCC, Smarter Balanced, and ACT Aspire.

Not only will such reviews provide critical information for state and local policymakers, as well as educators, curriculum developers and others, they might also deter the U.S. Department of Education from pursuing a dubious plan to make states put their new assessments through a federal evaluation system. In October 2015, the Department issued procedures for the “peer review” process that had been on hold for the last three years. The guidelines specify that states must produce evidence that they “used sound procedures in design and development to state tests aligned to academic standards, and for test administration and security.” Count us among those who think renewed federal vetting of state tests invites unwanted meddling from Uncle Sam (and could spark another round of backlash akin to what happened to the Common Core itself a few years back.) Besides, twelve years during which the Department already had such guidance in place did little to improve the quality of state tests—hence the recent moves to improve them.

Parting Thoughts

We are living in a time of political upheaval, divisiveness, and vitriol. The public’s faith in government and other large institutions is at an all-time low. So we’re glad to be the bearers of good news for a change. All four tests we evaluated boasted items of high technical quality. Further, the next generation assessments that were developed with the Common Core in mind have largely delivered on their promises. Yes, they have improvements to make (you’ll see that our reviewers weren’t shy in spelling those out). But they tend to reflect the content deemed essential in the Common Core standards and demand much from students cognitively. They are, in fact, the kind of tests that many teachers have asked state officials to build for years.

Now they have them.

Acknowledgments

This research was made possible through the generous support of the Louis Calder Foundation, the High-Quality Assessment Project (including the Bill & Melinda Gates Foundation, the Lumina Foundation, the Charles and Lynn Schusterman Family Foundation, the William and Flora Hewlett Foundation, and the Helmsley Trust), and our sister organization, the Thomas B. Fordham Foundation.

We owe a debt of gratitude to Nancy Doorey, project manager and report coauthor, and Morgan Polikoff, alignment expert and report coauthor, for their invaluable contributions to this project. This study had its share of difficulties, and through their tireless efforts Nancy and Morgan proved themselves highly equipped to handle all of them. We also extend our thanks to Dr. Roger Howe and Dr. Charles Perfetti, who served as math and ELA/Literacy content leads for the study and assisted with everything from creating initial reviewer training materials to overseeing the review process and synthesizing final study findings. Thanks also to Melisa Howey and Lynne Olmos for their special assistance and to the rest of our esteemed panelists for the thoughtfulness and care with which they conducted their reviews.

We also extend sincere thanks to each of the testing organizations who participated in the study (ACT Aspire, MCAS, PARCC, and Smarter Balanced) and to the many members of their staff who conducted initial reviewer trainings for their respective programs, responded to the panelists' questions, and reviewed the final report drafts for accuracy. In particular, we thank Elizabeth Sullivan, Carrie Conaway, Judy Hickman, and Francine Markowitz for facilitating this work.

We also appreciate the contributions of our colleagues at the Human Resources Research Organization (HumRRO), who led a similar evaluation at the high school level (reported separately) and with whom we conducted several joint reviews; the National Center for the Improvement of Educational Assessment (NCIEA) for producing the study's methodology; Judy Wurtzel and Joanne Weiss for facilitating the study on behalf of funders; and Student Achievement Partners (SAP) for helping to design and deliver, along with Jami-Jon Pearson, Morgan Polikoff, and HumRRO staff, portions of the reviewer training.

Fordham Research Manager Victoria Sears skillfully helped manage all aspects of the project, led recruitment of panelists, provided input on drafts, and shepherded the project across the finish line. Chester E. Finn, Jr. provided valuable feedback and edits to drafts, Alyssa Schwenk handled funder and media relations, and Shep Ranbom assisted with managing the report's dissemination. We also thank Fordham interns Megan Lail, Damien Schuster, and Stephan Shehy for their assistance throughout the project, and Jonathan Lutton, who ushered the report through production. Finally, we thank Shannon Last, who served as our copy editor; Edward Alton, who designed the report's layout; and Thinkstock.com, from which our cover design originated.

Executive Summary

Approximately one-third of American freshmen at two-year and four-year colleges require remedial coursework and over 40 percent of employers rate new hires with a high school diploma as “deficient” in their overall preparation for entry-level jobs.^{9, 10} Yet, over the past decade, as these students marched through America’s public education system, officials repeatedly told them, and their parents, that they were on track for success. They passed their courses, got good grades, and aced state annual tests. To put it plainly, it was all a lie. Imagine being told year after year that you’re doing just fine—only to find out when you apply for college or a job that you’re simply not as prepared as you need to be.

Thankfully, states have taken courageous steps to address this preparedness gap. Over the past five years, every state has upgraded its K–12 academic standards to align with the demands of college and career readiness (CCR), either by adopting the Common Core State Standards (CCSS) or working with their own higher education and career training providers to strengthen or develop standards. New assessments intended to align to these more-rigorous standards made their debut in the past year or two, and, as was widely expected (and, indeed, inevitable), student proficiency rates are lower than on previous tests—often significantly lower. State and local officials must decide whether to forge ahead with the new tests and higher expectations or back down in order to cast more schools and students in a positive (if, perhaps, illusory) light.

Of course, test scores that more accurately predict students’ readiness for entry-level coursework or training are not enough. The content of state assessments, too, is an important predictor of the impact of those tests on what is taught and learned. For instance, low-quality assessments poorly aligned with the standards will undermine the content messages of the standards; given the tests’ role in accountability under the newly reauthorized Elementary and Secondary Education Act, it is only logical that such tests might contribute to poor-quality instruction.

In short, good tests matter. Of critical importance to this conversation, therefore, is whether the new tests are indeed good and worth fighting for. That’s the central question this study seeks to answer.

The Tests

In the pages that follow, we evaluate the quality of four standardized assessments—three new, multi-state assessments and a well-regarded existing state assessment—to determine whether they meet new criteria developed by the Council of Chief State School Officers (CCSSO) for test quality. These new criteria, as explained in the following pages, ask that evaluators take a deep look at whether the assessments target and reliably measure the essential skills and knowledge needed at each grade level to achieve college and career readiness by the end of high school.

9. National Center for Education Statistics (NCES), Digest of Education Statistics, *Percentage of First-Year Undergraduate Students Who Took Remedial Education Courses, by Selected Characteristics: 2003–04 and 2007–08*, Table 241 (Washington, D.C.: NCES, 2010), https://nces.ed.gov/programs/digest/d10/tables/dt10_241.asp.

10. Conference Board et al., “Are They Really Ready To Work? Employers’ Perspectives on the Basic Knowledge and Applied Skills of New Entrants to the 21st Century U.S. Workforce” (New York, NY: Conference Board, 2006), http://www.p21.org/storage/documents/FINAL_REPORT_PDF09-29-06.pdf.

We evaluate English language arts/literacy and mathematics assessments for grades 5 and 8 for this quartet of testing programs:

- ◆ ACT Aspire
- ◆ The Partnership for Assessment of Readiness for College and Careers (PARCC)
- ◆ The Smarter Balanced Assessment Consortium (Smarter Balanced)
- ◆ The Massachusetts Comprehensive Assessment System (MCAS, 2014)

The Study Design

The analysis that follows was designed to answer three questions:

- 1 Do the assessments place strong emphasis on the most important content for college and career readiness (CCR), as called for by the Common Core State Standards and other CCR standards? (**Content**)
- 2 Do they require all students to demonstrate the range of thinking skills, including higher-order skills, called for by those standards? (**Depth**)
- 3 What are the overall strengths and weaknesses of each assessment relative to the examined criteria for ELA/Literacy and mathematics? (**Overall Strengths and Weaknesses**)

To answer these questions, we use a new methodology based on the CCSSO's 2014 "Criteria for Procuring and Evaluating High-Quality Assessments."¹¹ Developed by experts at the National Center for the Improvement of Educational Assessment (NCIEA), this methodology evaluates the degree to which test items and supporting program documentation (e.g., test blueprints and documents describing the item creation process) measure the critical competencies reflected in college and career readiness standards, thereby sending clear signals about the instructional priorities for each grade.¹²

The evaluation was conducted by review panels composed of practitioners, content experts, and specialists in assessment. Following reviewer training and a calibration exercise, the panels evaluated test items across various dimensions, with three to four experts reviewing each test form. Results were aggregated for each test form, discussed among the panel members, combined with results from a review of program documentation, and turned into group ratings and summary statements about each program.

The quality and credibility of an evaluation of this type rests largely on the expertise and judgment of the individuals serving on the review panels. To recruit highly qualified yet impartial reviewers, the study team requested recommendations from each of the four testing programs; from other respected content, assessment, and alignment experts; and from several national and state organizations. Reviewers were carefully vetted for their familiarity with the CCSS, their experience with developing or evaluating assessment items, and potential conflicts of interest. Individuals currently or previously employed by participating testing organizations and writers of the CCSS were not considered. (For more information, see Section I, *Selection of Review Panels*.) To ensure fairness and a balance of reviewer familiarity with each assessment, each of the panels included at least one reviewer recommended by each testing program.

Two university-affiliated content leads facilitated and reviewed the work of the ELA/Literacy and math review panels. Dr. Charles Perfetti, Distinguished University Professor of Psychology at University of Pittsburgh, served as the ELA/Literacy content lead, and Dr. Roger Howe, Professor of Mathematics at Yale University, served as the mathematics content lead. The names and biographical summaries of all panelists appear in Appendix E.

11. Council of Chief State School Officers (CCSSO), "Criteria for Procuring and Evaluating High-Quality Assessments" (Washington, D.C.: CCSSO, 2014).

12. The National Center for the Improvement of Educational Assessment, Inc. (NCIEA), "Guide to Evaluating Assessments Using the CCSSO Criteria for High Quality Assessments: Focus on Test Content" (Dover, NH: NCIEA, February 2016): http://www.nciea.org/publication_PDFs/Guide%20to%20Evaluating%20CCSSO%20Criteria%20Test%20Content%20020316.pdf.

This study evaluates English language arts and math assessments at grades 5 and 8, while a parallel study led by the Human Resources Research organization (HumRRO) evaluates the high school assessments from the same four testing programs (see Table ES-1). Because both organizations used the same methodology, it made sense to conduct two portions of the review jointly and across all grades: the documentation review and the accessibility review. Documentation results specific to grades 5 and 8 are addressed in this report. Please see HumRRO’s report for the results from their evaluation of the high school assessments, as well as results from the joint accessibility review (all grades).¹³

TABLE ES 1

Overview of the Parallel Fordham and HumRRO Studies

	ELA/Literacy Review	Math Review	Documentation Review	Accessibility Review
Fordham Study	Grades 5 and 8	Grades 5 and 8	Joint Panel	Joint Panel
HumRRO Study	High School	High School	(grades 5 and 8 findings presented in this report; high school findings presented in HumRRO report)	(presented in HumRRO report)

The CCSSO Criteria for High-Quality Assessments

To evaluate assessments intended to measure student mastery of the Common Core State Standards, we needed a new methodology that would capture their key dimensions. Traditional alignment methodologies offer the advantage of having been studied extensively, but treat each of the grade-level standards with equal importance, creating an inadvertent incentive for tests—and instruction—to be “a mile wide and an inch deep.”

The CCSSO’s “Criteria for Procuring and Evaluating High-Quality Assessments” was the basis of the new methodology. Specifically designed to address tests of college and career readiness, these criteria focus the evaluation on the highest priority skills and knowledge at each grade in the CCSS, addressing foundational as well as complex skills. By using the CCSSO Criteria as the basis of the methodology, the evaluation rewards those tests that focus on the essential skills and give clear signals about the instructional priorities for each grade.

The CCSSO Criteria address six domains, but only two pertain to the research questions addressed in this study: those for the assessment of ELA/Literacy standards and the assessment of mathematics standards (see Table ES-2).

In addition, CCSSO defined ratings for test content and depth, each of which is based on a subset of ratings. The Content rating reflects the degree to which each test assesses the material most needed for college and career readiness, and the Depth rating reflects the degree to which each test assesses the depth and complexity of the college and career readiness standards.

13. This study also originally included an evaluation of test program transparency, or the extent to which programs provide sufficient information to the public regarding assessment design and expectations (CCSSO criterion A.6). Due to several challenges associated with this review, however, we ultimately decided to drop this criterion from our study. Review panelists were not able to review all relevant documentation for each program, due to the vast volume of materials provided and publicly available. In addition, many test programs continued to release additional information (such as sample items) since our review occurred, rendering this panel’s findings somewhat outdated.

TABLE ES 2

CCSSO Criteria Evaluated in This Study

Assessment of ELA/Literacy Standards

Test Content Criteria

- B.3 Requiring students to read closely and use evidence from texts
- B.5 Assessing writing from sources
- B.6 Emphasizing vocabulary and language skills
- B.7 Assessing research and inquiry
- B.8 Assessing speaking and listening

Test Depth Criteria

- B.1 Using a balance of high-quality literary and informational texts
- B.2 Focusing on the increasing complexity of texts across grades
- B.4 Requiring a range of cognitive demand
- B.9 Ensuring high-quality items and a variety of item types

Assessment of Mathematics Standards

Test Content Criteria

- C.1 Focusing strongly on the content most needed for success in later mathematics (i.e., the major work of the grade)
- C.2 Assessing a balance of concepts, procedures, and applications

Test Depth Criteria

- C.3 Connecting mathematics practices to mathematical content
- C.4 Requiring a range of cognitive demand
- C.5 Ensuring high-quality items and a variety of item types

Findings

Results are organized around the key research questions above.

✓ RESULTS FOR QUESTIONS #1 AND #2

Do the assessments place strong emphasis on the most important content for college and career readiness (CCR) as called for by the Common Core State Standards and other CCR standards? **(Content)**

Do they require all students to demonstrate the range of thinking skills, including higher-order skills, called for by those standards? **(Depth)**

The panels assigned one of four ratings to each ELA/Literacy and math criterion: Excellent Match, Good Match, Limited/Uneven Match, or Weak Match. To generate these, each panel reviewed the ratings from the grade 5 and grade 8 test forms, considered the results of the documentation review, and came to consensus on the criterion rating.

Table ES-3 shows the ratings for test content and depth in ELA/Literacy and mathematics across the four programs.

The PARCC and Smarter Balanced assessments earned an Excellent or Good Match to the CCSSO Criteria for both ELA/Literacy and mathematics. While ACT Aspire and MCAS did well regarding the quality of items (see Section I, *Results*) and the Depth of Knowledge assessed (Depth), the panelists found that these two programs do not adequately assess—or may not assess at all—some of the priority content in both ELA/Literacy and mathematics at one or both grades in the study (Content).

TABLE ES 3

Overall Content and Depth Ratings for ELA/Literacy and Mathematics

	ACT Aspire	MCAS	PARCC	Smarter Balanced
ELA/Literacy CONTENT	L	L	E	E
ELA/Literacy DEPTH	G	G	E	G
Mathematics CONTENT	L	L	G	G
Mathematics DEPTH	G	E	G	G

LEGEND E Excellent Match G Good Match L Limited/Uneven Match W Weak Match

Criterion Level Results for ELA/Literacy and Mathematics

The Content and Depth ratings are based on the results of subsets of the CCSSO Criteria, as described above. NCIEA also recommended that certain criteria be “emphasized,” meaning awarded greater weight in the final determinations (though precise weightings were not specified). The panels, however, sometimes chose not to adhere to the weighting based on their level of confidence in reviewing each criterion (see Section I, *Methodology Modifications*).

Tables ES-4A and 4B show the distribution of the ELA/Literacy and math criteria ratings. Immediately striking in ELA is that the two consortia assessments (PARCC and Smarter Balanced, which received development grants from the U.S. Department of Education) earned twice as many ratings of Good and Excellent Match as the other two programs, earning eight high ratings to the four of ACT Aspire and MCAS. PARCC earned the most Excellent Match ratings (six), while Smarter Balanced was the only assessment with no ratings of Weak Match (partly because it was also the only program to test listening on the summative assessment).

TABLE ES-4A

ELA/Literacy Ratings Tally by Program

ACT Aspire	E	G	G	G	L	L	L	W	W
MCAS	E	G	G	G	L	L	W	W	W
PARCC	E	E	E	E	E	E	G	G	W
Smarter Balanced	E	E	E	E	G	G	G	G	L

TABLE ES-4BMathematics Ratings Tally by Program¹⁴

ACT Aspire	E	E	L	W
MCAS	E	E	E	L
PARCC	E	G	G	G
Smarter Balanced	E	G	G	L

LEGEND E Excellent Match G Good Match L Limited/Uneven Match W Weak Match

14. Although all four programs require the assessment of conceptual understanding, procedural skill/fluency, and applications (criterion C.2), final ratings could not be determined with confidence due to variations in how reviewers understood and implemented this criterion.

The ratings for mathematics (Table ES-4B) were more similar between programs, with PARCC earning four Excellent or Good Match ratings, Smarter Balanced and MCAS three each, and ACT Aspire two. MCAS scored particularly well on the three Depth criteria in mathematics, while PARCC is the only assessment that earned all Good Match or better scores.

Tables ES-5A and ES-5B on the following pages provide the final criterion ratings for each program, organized by Content and Depth. They also provide the specific attributes required to fully meet each criterion as indicated in the methodology.¹⁵ Those criteria followed by an asterisk were awarded greater emphasis during development of the Content and Depth ratings.

TABLE ES-5A

























Criterion Ratings for ELA/Literacy

CONTENT	ACT Aspire	MCAS	PARCC	Smarter Balanced
B.3* Reading: Items require close reading and use of direct textual evidence, and focus on central ideas and important particulars. To Meet the Criterion: 1) Nearly all reading items require close reading and analysis of text, rather than skimming, recall, or simple recognition of paraphrased text. 2) More than half of the reading score points are based on items that require direct use of textual evidence. 3) Nearly all items are aligned to the specifics of the standards. 4) More than half of the reading score points are based on items that require direct use of textual evidence.	L	G	E	E
B.5* Writing: Test programs assess a variety of types and formats of writing and the use of writing prompts that require students to confront and use evidence from texts or other stimuli directly. To Meet the Criterion: 1) All three writing types (expository, narrative, and persuasive/argument) are approximately equally represented across all forms in the grade band (K–5 and 6–12), allowing blended types (those that combine types) to contribute to the distribution. 2) All writing prompts require writing to sources (are text-based).	L	W	E	E
B.6 Vocabulary and Language Skills: Test forms place adequate emphasis on language and vocabulary items on the assessment, assess vocabulary that reflect requirements for college and career readiness, and focus on common student errors in language questions. To Meet the Criterion: 1) The large majority of vocabulary items (i.e., three-quarters or more) focuses on Tier 2 words and requires use of context, and more than half assess words important to central ideas. 2) A large majority (i.e., three-quarters or more) of the items in the language skills component and/or scored with a writing rubric mirror real-world activities, focus on common errors, and emphasize the conventions most important for readiness. 3) Vocabulary is reported as a sub-score or at least 13 percent of score points are devoted to assessing vocabulary/language. 4) Same as #3 for language skills.	G	L	E	G

* Criterion awarded greater weight in determination of Content and Depth rating.

LEGEND  Excellent Match  Good Match  Limited/Uneven Match  Weak Match

15. Note: As first implementers of the methodology, the reviewers made a number of modifications they deemed important for improvement. See Section I, *Methodology Modifications*.

CONTENT	ACT Aspire	MCAS	PARCC	Smarter Balanced
B.7 Research and Inquiry: Test forms include research items/tasks requiring students to analyze, synthesize, organize, and use information from multiple sources. To Meet the Criterion: The large majority (i.e., three-quarters or more) of the research items require analysis, synthesis, and/or organization of information.				
B.8 Speaking and Listening: <i>(Not yet required by the criteria, so not included in the Content rating. Listening requirements are listed here because one program assesses listening.)</i> Items assess students' listening skills using passages with adequate complexity and assess students' speaking skills through oral performance tasks. To Meet the Criterion: 1) The large majority (i.e., at least three-quarters) of listening items meet the requirements outlined in B.1 and B.2 and evaluate active listening skills.				
DEPTH	ACT Aspire	MCAS	PARCC	Smarter Balanced
B.1* Text Quality and Types: Test forms include a variety of text types (narrative and informational) that are of high quality, with an increasing focus on diverse informational texts across grades. To Meet the Criterion: 1) Approximately half of the texts at grades 3–8 and two-thirds at high school are informational, and the remainder literary. 2) Nearly all passages are high quality (previously published or of publishable quality). 3) Nearly all informational passages are expository in structure. 4) For grades 6–12, the informational texts are split nearly evenly between literary nonfiction, history/social science, and science/technical texts.				
B.2 Complexity of Texts: <i>(based on documentation review only)</i> Assessments include texts that have appropriate levels of text complexity for the grade or grade band (grade bands identified in the CCSS are K–5 and 6–12). To Meet the Criterion: 1) The documentation clearly explains how quantitative data are used to determine grade band placement. 2) Texts are then placed at the grade level recommended by qualitative review. 3) Text complexity rating process results in nearly all passages being placed at a grade band and grade level justified by complexity data.				
B.4 Matching the Complexity of the Standards: Each test form contains an appropriate range of cognitive demand that adequately represents the cognitive demand of the standards. To Meet the Criterion: 1) The distribution of cognitive demand of the assessment matches the distribution of cognitive demand of the standards as a whole and matches the higher cognitive demand (DOK 3+) of the standards. <i>(Note: This is not a rating of test difficulty. Assessments that do not match the DOK distribution of the standards, even if there are too many high DOK items, may receive a rating less than Excellent Match. See Appendix A for more information.)</i>				
B.9 High-Quality Items and Variety of Item Types: Test items are of high quality, lacking technical or editorial flaws and each test form contains multiple item types including at least one type in which students construct, rather than select, a response. To Meet the Criterion: 1) All or nearly all operational items reviewed reflect technical quality and editorial accuracy. 2) At least two item formats are used, including one that requires students to generate, rather than select, a response.				

* Criterion awarded greater weight in determination of Content or Depth rating.






LEGEND  Excellent Match  Good Match  Limited/Uneven Match  Weak Match
 Cells for which the ratings are not used in determining Content and Depth ratings
(See Section I, Weighting of Criteria for Content and Depth Ratings.)

TABLE ES-5B

Criterion Ratings for Mathematics

CONTENT	ACT Aspire	MCAS	PARCC	Smarter Balanced
C.1* Focus: Each test form contains a strong focus on the content most crucial for success in later mathematics (i.e., the major work of the grade). To Meet the Criterion: The vast majority (i.e., at least three-quarters in elementary grades, at least two-thirds in middle school grades, and at least half in high school) of score points in each assessment focus on the content that is most important for students to master in that grade band in order to reach college and career readiness (the major work of the grade).				
C.2 Concepts, Procedures, and Applications: Each test form contains items that assess conceptual understanding, procedural skill/fluency, and application in approximately equal proportions. To Meet the Criterion: The distribution of score points reflects a balance of mathematical concepts, procedures/fluency, and applications.	Due to variations in how reviewers understood and implemented this criterion, final ratings could not be determined with confidence.			
DEPTH	ACT Aspire	MCAS	PARCC	Smarter Balanced
C.3 Connecting Practice to Content: Assessments test students' use of mathematical practices through test items that connect these practices with grade-level content standards. To Meet the Criterion: All or nearly all items that assess mathematical practices also align to one or more content standards.				
C.4* Matching the Complexity of the Standards: Each test form contains an appropriate range of cognitive demand that adequately represents the cognitive demand of the standards. To Meet the Criterion: 1) The distribution of cognitive demand of the assessment matches the distribution of cognitive demand of the standards as a whole and matches the higher cognitive demand (DOK 3+) of the standards. (Note: This is not a rating of test difficulty. Assessments that do not match the DOK distribution of the standards, even if there are too many high DOK items, may receive a rating less than Excellent Match. See Appendix A for more information.)				
C.5* High-Quality Items and Variety of Item Types: Test items are of high quality, lacking technical or editorial flaws, and each test form contains multiple item types, including at least one type in which students construct, rather than select, a response. To Meet the Criterion: 1) All or nearly all operational items reviewed reflect technical quality and editorial accuracy. 2) At least two item formats are used, including one that requires students to generate, rather than select, a response.				

* Criterion awarded greater weight in determination of Content or Depth rating.

LEGEND Excellent Match Good Match Limited/Uneven Match Weak Match

In the ELA/Literacy assessments, all four programs receive high ratings for the quality of items and variety of item types. In addition, all pay close attention to the use of high-quality informational and literary texts and increasing the complexity of tests across grades, which are significant advances over many previous state ELA assessments. Significant differences exist across the testing programs, however, in the degree to which their writing tasks require students to use evidence from sources and the extent to which research skills are assessed. In these areas, PARCC and Smarter Balanced perform well, receiving higher ratings than either ACT Aspire, which receives a rating of Limited/Uneven Match on these criteria, or MCAS, which receives a rating of Weak Match. PARCC and Smarter Balanced assessments also contain a distribution of cognitive demand that better reflects that of the standards, when compared to ACT Aspire and MCAS.

In mathematics, PARCC and Smarter Balanced receive a rating of Good Match for the degree to which their tests focus on the most important content of the grade. ACT Aspire test forms receive a rating of Weak Match on this prioritized criterion, due to their test design choice, in which off-grade standards are assessed in order to monitor mastery across grades. MCAS receives a rating of Limited/Uneven because its grade 5 forms do not contain

Supplemental Analysis: Assessment of Higher-Order Thinking Skills

CCSSO criteria B.4 and C.4 capture the degree to which the range of cognitive demand on the test forms match that of the CCSS. We used Webb’s Depth of Knowledge (DOK) taxonomy to assess cognitive demand, as it is by far the most widely used approach to categorizing cognitive demand. Webb’s DOK is composed of four levels. Level 1 is the lowest level (recall), Level 2 requires use of a skill or concept, and Levels 3 and 4 are higher-order thinking skills. We compared the DOK of the assessments to those of the Common Core State Standards, which were coded by content experts (see Section I, *Selection of Review Panels and Assignment to Forms*). We also compared the tests’ DOK distributions to those of fourteen highly regarded previous state assessments, as well as the distribution reflected in several national and international assessments—including Advanced Placement (AP), the National Assessment of Education Progress (NAEP), and the Program for International Student Assessment (PISA).^{16, 17}

We found that the CCSS call for greater emphasis on higher-order skills than fourteen highly regarded previous state assessments in ELA/Literacy at both grades 5 and 8 as well as in math at grade 8 (they are similar at grade 5). In addition, the grade 8 CCSS in both ELA/Literacy and math call for greater emphasis on higher-order thinking skills than either NAEP or PISA, both of which are considered to be high-quality, challenging assessments.

Overwhelmingly, the assessments included in our study were found to be more challenging—placing greater emphasis on higher-order skills—than prior state assessments, especially in mathematics (where prior assessments rarely included items at DOK 3 or 4 at all). In some cases, the increase was dramatic: PARCC’s DOK in grade 8 exceeds even that of AP and PISA in both subjects. See Appendix A for more.

However, the panels found significant variability in the degree to which the four assessments match the distribution of DOK in the CCSS. In some cases, the panels found significant variability between the grade 5 and grade 8 assessments for a given program. PARCC tests generally have the highest DOK in ELA/Literacy, while ACT Aspire had the highest in mathematics. See Section I, Tables 14 and 22 for the DOK distribution of each program.

16. L. Yuan and V. Le, *Estimating the Percentage of Students who were Tested on Cognitively Demanding Items through the State Achievement Tests* (Santa Monica, CA: RAND Corporation, 2012).

17. *Ibid.*

sufficient focus on the critical content for the grade. With respect to item quality, ACT Aspire and MCAS receive the highest rating of Excellent Match, whereas PARCC receives a rating of Good Match and Smarter Balanced a rating of Limited/Uneven Match.¹⁸

✓ RESULTS FOR QUESTION #3

What are the overall strengths and weaknesses of each assessment relative to the examined criteria for ELA/Literacy and mathematics? **(Overall Strengths and Weaknesses)**

Each of the review panels developed summary statements for each assessment program, detailing their strengths and areas of improvement in ELA/Literacy and mathematics. In addition, they created summary statements for each test's Content and Depth ratings based on the prioritization of criteria recommended in the study methodology (see Appendix F). They also generated final statements summarizing the observed strengths and areas of improvement for each program.

ACT Aspire

English Language Arts:

In ELA/Literacy, ACT Aspire receives a Limited/Uneven to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The combined set of ELA/Literacy tests (reading, writing, and English) requires close reading and adequately evaluates language skills. More emphasis on assessment of writing to sources, vocabulary, and research and inquiry, as well as increasing the cognitive demands of test items, will move the assessment closer to fully meeting the criteria. Over time, the program would also benefit by developing the capacity to assess speaking and listening skills.

Content: ACT Aspire receives a Limited/Uneven match to the CCSSO Criteria for Content in ELA/Literacy. The assessment program includes an emphasis on close reading and language skills. However, the reading items fall short on requiring students to cite specific textual information in support of a conclusion, generalization, or inference and in requiring analysis of what has been read. In order to meet the criteria, assessing writing to sources, vocabulary, as well as research and inquiry need to be strengthened.

Depth: ACT Aspire receives a rating of Good Match for Depth in ELA/Literacy. The program's assessments are built on high-quality test items and texts that are suitably complex. To fully meet the CCSSO Criteria, more cognitively demanding test items are needed at both grade levels, as is additional literary narrative text—as opposed to literary informational texts.¹⁹

Mathematics:

In mathematics, ACT Aspire receives a Limited/Uneven to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. Some of the mismatch with the criteria is likely due to intentional program design, which requires that items be included from previous and later grades.

18. The nature and timing of this review required Smarter Balanced to make the test items and forms available to reviewers through an alternate test interface that was more limited than the actual student interface used for the summative assessments, particularly with regard to how items appeared on the screen and how erroneous responses were handled. Though reviewers were not able to determine the extent to which these interface limitations impacted their findings, the study team worked with Smarter Balanced to ascertain which item issues were caused by interface differences and which were not. All item-relevant statements in the report reflect data not prone to interface differences.

19. ACT Aspire does not classify literary nonfiction texts that are primarily narrative in structure as “informational.” See Appendix G for more information about ACT Aspire's interpretation of CCSSO criterion B.1.

The items are generally high quality and test forms at grades 5 and 8 have a range of cognitive demand, but in each case the distribution contains significantly greater emphasis at DOK 3 than reflected in the standards. Thus, students who score well on the assessments will have demonstrated a strong understanding of the standards' more complex skills. However, the grade 8 test may not fully assess standards at the lowest level of cognitive demand. The tests would better meet the CCSSO Criteria with an increase in the number of items focused on the major work of the grade and the addition of more items at grade 8 that assess standards at DOK 1.

Content: ACT Aspire receives a Limited/Uneven Match to the CCSSO Criteria for Content in Mathematics. The program does not focus exclusively on the major work of the grade, but rather, by design, assesses material from previous and later grades. This results in a weaker match to the criteria. The tests could better meet the criteria at both grades 5 and 8 by increasing the number of items that assess the major work of the grade.

Depth: ACT Aspire receives a good match to the CCSSO Criteria for Depth in Mathematics. The items are well crafted and clear, with only rare instances of minor editorial issues. The ACT Aspire tests include proportionately more items at high levels of cognitive demand (DOK 3) than the standards reflect and proportionately fewer at the lowest level (DOK 1). This finding is both a strength, in terms of promoting strong skills, and a weakness, in terms of ensuring adequate assessment of the full range of cognitive demand within the standards. While technically meeting the criterion for use of multiple item types, the range is nonetheless limited, with the majority comprising multiple-choice items. The program would better meet the criteria for Depth by including a wider variety of item types and relying less on traditional multiple-choice items.

MCAS

English Language Arts:

In ELA/Literacy, MCAS receives a Limited/Uneven to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The test requires students to closely read high-quality texts and a variety of high-quality item types. However, MCAS does not adequately assess several critical skills—including reading informational texts, writing to sources, language skills, and research and inquiry; further, too few items assess higher-order skills. Addressing these limitations would enhance the ability of the test to signal whether students are demonstrating the skills called for in the standards. Over time, the program would also benefit by developing the capacity to assess speaking and listening skills.

Content: MCAS receives a Limited/Uneven Match to the CCSSO Criteria for Content in ELA/Literacy. The assessment requires students to read closely well-chosen texts and presents test questions of high technical quality. However, the program would be strengthened by assessing writing annually, assessing the three types of writing called for across each grade band, requiring writing to sources, and placing greater emphasis on assessing research and language skills.

Depth: MCAS receives a rating of Good Match for Depth in ELA/Literacy. The assessments do an excellent job in presenting a range of complex reading texts. To fully meet the demands of the CCSSO Criteria, however, the test needs more items at higher levels of cognitive demand, a greater variety of items to test writing to sources and research, and more informational texts—particularly those of an expository nature.

Mathematics:

In mathematics, MCAS receives a Limited/Uneven Match to the CCSSO Criteria for Content and an Excellent Match for Depth relative to assessing whether students are on track to meet college and career readiness standards. The MCAS mathematics test items are of high technical and editorial quality. Additionally, the content is distributed well across the breadth of the grade level standards, and test forms closely reflect the range of cognitive demand of the standards. Yet the grade 5 tests have an insufficient degree of focus on the major work of the grade.

While mathematical practices are required to solve items, MCAS does not specify the assessed practices(s) within each item or their connections to content standards. The tests would better meet the criteria through increased focus on major work at grade 5 and identification of the mathematical practices that are assessed—and their connections to content.

Content: MCAS receives a Limited/Uneven Match to the CCSSO Criteria for Content in Mathematics. While the grade 8 assessment focuses strongly on the major work of the grade, the grade 5 assessment does not, as it samples more broadly from the full range of standards for the grade. The tests could better meet the Criteria through increased focus on the major work of the grade on the grade 5 test.

Depth: MCAS receives an Excellent Match to the CCSSO Criteria for Depth in Mathematics. The assessment uses high-quality items and a variety of item types. The range of cognitive demand reflects that of the standards of the grade. While the program does not code test items to math practices, mathematical practices are nonetheless incorporated within items. The program might consider coding items to the mathematical practices and making explicit the connections between specific practices and content standards.

PARCC

English Language Arts:

In ELA/Literacy, PARCC receives an Excellent Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The tests include suitably complex texts, require a range of cognitive demand, and demonstrate variety in item types. The assessments require close reading, assess writing to sources, research, and inquiry, and emphasize vocabulary and language skills. The program would benefit from the use of more research tasks requiring students to use multiple sources and, over time, developing the capacity to assess speaking and listening skills.

Content: PARCC receives an Excellent Match to the CCSSO Criteria for Content in ELA/Literacy. The program demonstrates excellence in the assessment of close reading, vocabulary, writing to sources, and language, providing a high-quality measure of ELA/Literacy content as reflected in college and career readiness standards. The tests could be strengthened by the addition of research tasks that require students to use two or more sources and, as technologies allow, a listening and speaking component.

Depth: PARCC receives a rating of Excellent Match for Depth in ELA/Literacy. The PARCC assessments meet or exceed the depth and complexity required by the Criteria through a variety of item types that are generally of high quality. A better balance between literary and informational texts would strengthen the assessments in addressing the Criteria.

Mathematics:

In mathematics, PARCC receives a Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The assessment is reasonably well aligned to the major work of each grade. At grade 5, the test includes a distribution of cognitive demand that is similar to that of the standards. At grade 8, the test has greater percentages of higher-demand items (DOK 3 and 4) than reflected by the standards, such that a student who scores well on the grade 8 PARCC assessment will have demonstrated strong understanding of the standards' more complex skills. However, the grade 8 test may not fully assess standards at the lowest level (DOK 1) of cognitive demand.

The test would better meet the CCSSO Criteria through additional focus on the major work of the grade, the addition of more items at grade 8 that assess standards at DOK 1, and increased attention to accuracy of the items—primarily editorial, but in some instances mathematical.

Content: PARCC receives a Good Match to the CCSSO Criteria for Content in Mathematics. The test could better meet the criteria by increasing the focus on the major work at grade 5.

Depth: PARCC receives a Good Match to the CCSSO Criteria for Depth in Mathematics. The tests include items with a range of cognitive demand, but at grade 8, that distribution contains a higher percentage of items at the higher levels (DOK 2 and 3) and significantly fewer items at the lowest level (DOK 1). This finding is both a strength, in terms of promoting strong skills, and a weakness, in terms of ensuring adequate assessment of the full range of cognitive demand within the standards. The tests include a variety of item types that are largely of high quality. However, a range of problems (from minor to severe) surfaced relative to editorial accuracy and, to a lesser degree, technical quality. The program could better meet the Depth criteria by ensuring that all items meet high editorial and technical standards and by ensuring that the distribution of cognitive demand on the assessments receives sufficient information across the range.

Smarter Balanced

English Language Arts:

In ELA/Literacy, Smarter Balanced receives a Good to Excellent Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The tests assess the most important ELA/Literacy skills of the CCSS, using technology in ways that both mirror real-world uses and provide quality measurement of targeted skills. The program is most successful in its assessment of writing and research and inquiry. It also assesses listening with high-quality items that require active listening, which is unique among the four programs. The program would benefit by improving its vocabulary items, increasing the cognitive demand in grade 5 items, and, over time, developing the capacity to assess speaking skills.

Content: Smarter Balanced receives an Excellent Match to the CCSSO Criteria for Content in ELA/Literacy. The program demonstrates excellence in the areas of close reading, writing to sources, research, and language. The listening component represents an important step toward adequately measuring speaking and listening skills—a goal specifically reflected in the standards. Overall, Smarter Balanced is a high-quality measure of the content required in ELA/Literacy, as reflected in college and career readiness standards. A greater emphasis on Tier 2 vocabulary would further strengthen these assessments relative to the criteria.

Depth: Smarter Balanced receives a rating of Good Match for Depth in ELA/Literacy. The assessments use a variety of item types to assess student reading and writing to source. The program could better meet the depth criteria by increasing the cognitive demands of the grade 5 assessment and ensuring that all items meet high editorial and technical quality standards.

Mathematics:

In mathematics, Smarter Balanced has a Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The test provides adequate focus on the major work of the grade, although it could be strengthened at grade 5.

The tests would better meet the CCSSO Criteria through increased focus on the major work at grade 5 and an increase in the number of items on the grade 8 tests that assess standards at the lowest level of cognitive demand. In addition, removal of serious mathematical and/or editorial flaws, found in approximately one item per form, should be a priority.²⁰

Content: Smarter Balanced receives a Good Match to the CCSSO Criteria for Content in Mathematics. The tests could better meet the criteria by increasing the focus on the major work for grade 5.

Depth: Smarter Balanced receives a Good Match to the CCSSO Criteria for Depth in Mathematics. The exam includes a range of cognitive demand that fairly represents the standards at each grade level. The tests have a strong variety of item types including those that make effective use of technology. However, a range of problems (from minor to severe) surfaced relative to editorial accuracy and, to a lesser degree, technical

20. See footnote 18 for more on Smarter Balanced test interface.

quality. A wide variety of item types appear on each form, and important skills are assessed with multiple items, as is sound practice. The program could better meet the Depth criteria by ensuring that all items meet high editorial and technical standards and that a given student is not presented with two or more virtually identical problems.

For too many years, state assessments have generally focused on low-level skills and have given parents and the public false signals about students' readiness for postsecondary education and the workforce. They often weren't very helpful to educators or policymakers either. States' adoption of college and career readiness standards has been a bold step in the right direction. Using high-quality assessments of these standards will require courage: these tests are tougher, sometimes cost more, and require more testing time than the previous generation of state tests. Will states be willing to make the tradeoffs?

Introduction

Overview

Approximately one-third of American freshmen at two-year and four-year colleges require remedial coursework, and over 40 percent of employers rate new hires with a high school diploma as “deficient” in their overall preparation for entry-level jobs.^{21, 22} Yet over the past decade, as these students marched through America’s public education system, officials repeatedly told them, and their parents, that they were on track for success. They passed their courses, got good grades, and aced state annual tests. To put it plainly, it was all a lie. Imagine being told year after year that you’re doing just fine—only to find out when you apply for college or a job that you’re simply not as prepared as you need to be.

Thankfully, states have taken courageous steps to address this preparedness gap. Over the past five years, every state has upgraded its K–12 academic standards to align with college and career readiness, either by adopting the Common Core State Standards or by working with its own higher education and career training providers to strengthen or develop standards. But whether or not these improved standards will be faithfully implemented in the classroom depends a great deal on whether enough states adopt rigorous, high-quality assessments, and whether those assessments truly measure the complex skills and knowledge called for by the standards. If the tests are weak measures that, like many of their predecessors, focus too much on low-level and easy-to-measure skills, they make it less likely that the standards will achieve their ultimate goal of ensuring that high school graduates are indeed prepared for successful transition into the job market or higher education.

The selection of state assessments has recently been a topic of great debate in policy and political circles. As new assessments intended to align to the new college and career readiness standards made their debut in the past year or two, student proficiency rates have (predictably and inevitably) been lower—sometimes significantly so—than on previous tests, forcing state and local officials nationwide to make critical decisions. Should they forge ahead with the new tests and higher expectations or back down in order to cast more schools and students in a positive (if, perhaps, illusory) light?

Of course, test scores that more accurately predict students’ readiness for entry-level coursework or training are not enough. The content of state assessments, too, is an important predictor of the impact of those tests on what is taught and learned. For instance, low-quality assessments with poorly aligned tasks undermine the content messages of the standards; given the tests’ role in accountability, it is only logical that such tests might contribute to poor-quality instruction.

In short, good tests matter. Of critical importance to this conversation, however, is whether the new tests are indeed good and worth fighting for. That’s the central question this study seeks to answer.

This report, then, provides much-needed information to policymakers, practitioners, and researchers about the quality of four current assessments and their potential to support effective implementation of college and career readiness standards. It takes an unprecedented, in-depth look at three new multi-state assessments and an

21. NCES, *Percentage of First-Year Undergraduate Students who Took Remedial Education Courses*.

22. Conference Board et al., “Are They Really Ready To Work?”

existing best-in-class state assessment. All four assert that they measure college and career readiness standards generally and the CCSS specifically. The study results are particularly relevant to states that have yet to make final decisions on which assessments they'll use in future years to measure student learning in English language arts/literacy and mathematics.

The Four Assessments

Three of the assessments evaluated in this study are currently in use across multiple states: the Partnership for the Assessment of Readiness for College and Careers (PARCC), the Smarter Balanced Assessment System (Smarter Balanced), and ACT Aspire. The first two are membership-based organizations that are governed by participating states;²³ the latter is developed by the maker of the widely known ACT college admissions test.

The fourth assessment, the 2014 version of the Massachusetts Comprehensive Assessment System (MCAS), is a highly regarded state assessment program from a state that has consistently outscored the rest of the nation on a host of achievement metrics.^{24, 25} In 2011, Massachusetts adopted new curriculum frameworks in mathematics and ELA/Literacy, which incorporated the CCSS with Massachusetts-specific standards, and began to transition the MCAS assessments accordingly. By the spring of 2014, the transition was complete: MCAS included items that covered the CCSS and additional Massachusetts-specific standards.²⁶ Within this study, the 2014 MCAS serves as a comparison point or “best-case” example for the solo state option.

Study Design and Key Questions

The study utilizes a new methodology developed by the National Center for the Improvement of Educational Assessment (NCIEA) to determine how well new tests measure the requisite content, knowledge, and critical skills at key grade levels and, in doing so, whether they sufficiently tap higher-order thinking skills.

This study evaluates the summative (end-of-year) assessments administered in grades 5 and 8, the capstone grades for the elementary and middle school levels.²⁷ A parallel study led by the Human Resources Research Organization (HumRRO) evaluates the capstone assessment in grades 9–12 (typically administered at grade 10 or 11, depending on the program). Because both organizations used the same methodology, it made sense to conduct two portions of the review jointly—the documentation review and the accessibility review. The composition of these joint panels is described below. (HumRRO has published a separate report with the results of the evaluation of the high school assessments, which will also include findings from our joint accessibility review.)²⁸

23. Congress and the U.S. Department of Education allocated \$350 million in federal grants to groups of states seeking to develop new assessments aligned to college and career readiness standards. Two consortia of states were awarded Race to the Top assessment grants: the Partnership for Assessment of Readiness for College and Careers (PARCC) and the Smarter Balanced Assessment Consortium. By 2012, forty-five states were participating either as governing or participating (advisory) members in either or both of the consortia. That number is now 21 states, however, as of January 2016.

24. V. Bandeira de Mello et al., “Mapping State Proficiency Standards onto NAEP Scales: 2005–2007” (Washington, D.C.: National Center for Education Statistics, Institute of Education Sciences, 2009).

25. National Assessment Governing Board, “Nation's Report Card,” 2015, http://www.nationsreportcard.gov/reading_math_2015/.

26. MCAS items assessing standards unique to Massachusetts—meaning not part of the CCSS—were removed from the test forms prior to use in this study. In 2014, school districts in Massachusetts were allowed to choose either MCAS or PARCC as their summative assessment.

27. Some of these testing programs provide aligned formative and/or benchmark/interim assessments, which are not used for consequential purposes. Those assessments were not part of this study.

28. This study also originally included an evaluation of test program transparency, or the extent to which programs provide sufficient information to the public regarding assessment design and expectations (CCSSO criterion A.6). Due to several challenges associated with this review, however, we ultimately decided to drop this criterion from our study. Review panelists were not able to review all relevant documentation for each program, due to the vast volume of materials provided and publicly available. In addition, many test programs continued to release additional information (such as sample items) since our review occurred, rendering this panel's findings somewhat outdated.

TABLE 1

Parallel Makeup of the Fordham and HumRRO Studies

	ELA/Literacy Review	Math Review	Documentation Review	Accessibility Review
Fordham Study	Grades 5 and 8	Grades 5 and 8	Joint Panel	Joint Panel
HumRRO Study	High School	High School	(grades 5 and 8 findings presented in this report; high school findings presented in HumRRO report)	(presented in HumRRO report)

Together, these studies provide the public with the first in-depth look at this new generation of highly touted assessments.

Our study was designed to address this trio of questions:

- 1 Do the assessments place strong emphasis on the most important content for college and career readiness (CCR), as called for by the Common Core State Standards and other CCR standards? (**Content**)
- 2 Do they require all students to demonstrate the range of thinking skills, including higher-order skills, called for by those standards? (**Depth**)
- 3 What are the overall strengths and weaknesses of each assessment relative to the examined criteria for ELA/Literacy and mathematics? (**Overall Strengths and Weaknesses**)

Questions 1 and 2, regarding the Content and Depth of the assessments, are the major emphases of this evaluation, and Question 3 highlights the major findings of the panels.

Organization of the Report

The report is organized as follows. First, we take a high-level look at the ways in which this study's approach to the evaluation of test quality differs significantly from previous approaches and describe major design differences within and among the four participating test programs.

Next, we explain the methodology developed by the NCIEA and how it was operationalized for this study (Section I). The findings from our mathematics and English language arts/literacy review follow this (see Section I, *Results*), then the section concludes with a summary of the various strengths and weaknesses of each testing program relative to the CCSSO Criteria. Section II includes our recommendations for policymakers and test developers, and we conclude in Section III by offering suggestions for future improvements to the methodology. The multiple appendices include, among other topics, a discussion of Depth of Knowledge (DOK), key terminology, the biographies of the authors and panelists, and responses to the study from the testing programs.

Approach to Alignment and Quality

To evaluate assessments intended to measure student mastery of the Common Core State Standards (CCSS), we needed a methodology that would capture their key dimensions. Traditional alignment methodologies, such as Webb's alignment tool²⁹ and the Surveys of Enacted Curriculum,³⁰ offer the advantage of having been studied extensively, but they measure only a fairly narrow range of test-standards alignment issues. For instance,

29. N. L. Webb, *Alignment of Science and Mathematics Standards and Assessments in Four States*, Research Monograph No. 18 (Madison, WI: University of Wisconsin–Madison, National Institute for Science Education, 1999).

30. A. C. Porter, "Measuring the Content of Instruction: Uses in Research and Practice," *Educational Researcher* 31, no. 7 (2002): 3–14.

neither of these methodologies could be used to report on issues such as a math test's coverage of the CCSS for Mathematical Practice, or whether an ELA/Literacy test contains passages with appropriate text complexity for the grade level.

Because of the limitations of existing approaches, a new methodology was created based on the CCSSO's "Criteria for Procuring and Evaluating High-Quality Assessments."³¹ Specifically designed to address tests of college and career readiness, these criteria focus the evaluation on the highest-priority skills and knowledge at each grade in the CCSS, addressing foundational as well as complex skills.

The methods were developed by experts at the NCIEA and refined in several cycles of revision. They represent a very different approach to gauging test quality and alignment than earlier generations of such studies. Here we discuss two of these fundamental differences, both of which are the result of critical design attributes of the CCSS and similar college and career readiness standards.

First, most prior alignment approaches assume a one-to-one match between standards and test items. The CCSS and kindred standards of college and career readiness, however, describe numerous complex competencies that cannot be assessed with individual test items, such as the ability to craft clearly written arguments supported by evidence drawn from multiple sources or the ability to apply multiple mathematical skills to solve a complex problem. Such competencies are essential for postsecondary education, training, and citizenship, and states need to know whether their assessments are measuring them. This new approach to alignment and quality asks evaluators to determine the degree to which the tests measure these more complex competencies.

The second key difference of this approach is that it focuses the evaluation on the highest-priority skills and knowledge at each grade. Prior approaches treated each of the grade-level standards with equal importance, creating an inadvertent incentive for tests—and instruction—to be “a mile wide and an inch deep.” The CCSS clearly define the prioritized skills and competencies at each grade level and call upon both instruction and assessment to focus strongly, although not exclusively, on these priorities. By focusing the review on the critical competencies, this methodology rewards those tests that give clear signals about the instructional priorities for each grade.

Test Program Differences

Although each of the four tests in this study was developed to assess college and career readiness standards, they differ in significant ways. These differences provide important context for interpreting the results of our study.

- 1 Test Design:** Content standards are not written to guide the creation of assessments; rather, they are written to define the specific content and skills students are to master at each grade level. Thus, in creating assessments from standards, programs must first interpret and bundle standards into assessment targets that define what is to be measured. Even though each testing program was tasked with assessing student learning against college and career readiness standards, and the CCSS in particular, the level of emphasis on specific skills and knowledge—such as explaining mathematical reasoning—varies. Simply put, some programs, by design, follow the contours of the CCSS more closely than others, and the impact of these test-design choices is reflected in the study results. As another example, ACT Aspire's ELA/Literacy assessment is broken into three components—reading, writing, and English—while the other programs do not have this split.

31. Council of Chief State School Officers (CCSSO), “Criteria for Procuring and Evaluating High-Quality Assessments” (Washington, D.C.: CCSSO, 2014).

- 2 **Test Delivery:** The 2014 MCAS is a paper-based assessment, whereas the other three are designed and primarily delivered as computer-based tests (although each currently offers a paper version). Administering the test via computer offers a wider array of options for the presentation of test items as well as how student responses are entered. For example, audio, video, and animation can be used in the presentation of computerized test items, and student responses may involve manipulation of data to identify patterns or live editing of texts with grammatical errors. Paper-based versus computer delivery also impacts the degree to which a test program can assess certain standards, such as listening and modeling skills.³²
- 3 **Test Forms and Computer Adaptivity:** The set of items presented to a student in a test form is predetermined for the PARCC, MCAS, and ACT Aspire assessments. Smarter Balanced assessments, however, are computer adaptive, meaning that items or small groups of items are selected as the student proceeds through the test, based at least in part on all of the student's prior responses. Their summative tests have an adaptive and a performance task (PT) component. The PT is a fixed set of items that are not chosen adaptively. The adaptive portions of the test adjust after each item is administered and scored, except in cases in which multiple questions are based on the same reading or listening passage(s) or scenario. Such groups of items are administered as a unit.³³ The majority of students see only on-grade-level items, but those performing near the ends of the performance spectrum may be given off-grade-level items near the end of the test to better assess their performance and increase the precision of their score.

This performance-based adaptation allows the test to produce scores with smaller margins of error for students near the ends of the performance spectrum and/or to shorten the overall length of the test for most students.³⁴ It also, however, creates the potential for greater variation in content across student forms, which could impact the degree to which individual forms meet the CCSSO Criteria. This study followed the processes recommended by the NCIEA for the creation of simulated Smarter Balanced test forms and, in order to evaluate the degree to which other forms would vary, also included the results of a simulation study of 1,000 forms per grade and content area.

- 4 **Testing Time:** Test developers are forced to strike a balance between accurately measuring student learning and the time and cost required to do so. Our four testing programs show different solutions to this balancing act. The estimated testing time for students in grades 5 and 8, on average, to complete both the ELA/Literacy assessments and the mathematics assessments for the respective programs are as follows:

- ♦ ACT Aspire: three to three-and-a-quarter hours for all four tests (English, reading, writing, and mathematics)
- ♦ MCAS: three-and-a-half hours
- ♦ PARCC: seven to seven-and-a-half hours (the 2015–16 revisions will reduce this by an estimated one-and-a-half hours)³⁵
- ♦ Smarter Balanced: five-and-a-half hours

The longer testing times for PARCC and Smarter Balanced are primarily due to the inclusion of the extended performance tasks. Both programs use these tasks to assess high-priority skills within the CCSS, such as the development of written compositions in which a claim is supported with evidence drawn from

32. The CCSS define modeling as the process of choosing and using appropriate mathematics and statistics to analyze empirical situations, to understand them better, and to improve decisions.

33. The Smarter Balanced adaptive engine, which selects the items to be presented to the student, is programmed to first ensure that the requirements in the test blueprint for content and cognitive demand are met. Within these constraints, it also adapts to improve score precision.

34. In CAT, the algorithm is typically programmed to continue asking a student questions until a certain level of precision is achieved OR until other conditions are met, such as a time limit or the assigned item pool is exhausted. See, for example, H. Wainer et al., *Computerized Adaptive Testing: A Primer* (2000).

35. See Appendix G for a description of the changes made to the 2015–16 versions of PARCC and the other three assessments, as relevant.

sources; research skills; and solving complex multi-step problems in mathematics. In addition to requiring more time than selected-response items, these tasks are typically more costly to develop and score.

- 5 **Ownership of the Assessments and Setting Proficiency Cut Scores:** PARCC and Smarter Balanced are governed by their respective groups of member states, which collectively set, as required by the initial grants, the proficiency cut scores to be used by all members for federal reporting of student performance. Each member state, however, may determine the cut scores to be used within its own state accountability systems—cut scores that may impact grade-to-grade promotions, the awarding of high school diplomas, and/or educator evaluations.³⁶

In contrast, the ACT Aspire assessments are the property of ACT Aspire. Individual states contract with ACT Aspire for their use and may set their own proficiency cut scores to be used for both federal and state accountability purposes. For grade 10, ACT Aspire is able to inform states' cut-score decisions with historical data regarding the minimum scores associated with a high likelihood of success in credit-bearing first-year college courses³⁷ (i.e., the ACT Aspire College Readiness Benchmarks).

The MCAS is a custom assessment developed and owned by the state of Massachusetts. State policymakers are charged with setting the proficiency cut scores used for both federal and state accountability purposes.

Next we turn our attention to the methodology.

36. Both Smarter Balanced and PARCC plan to analyze the relationship between student performance on the high school assessments and subsequent performance in entry-level post-secondary courses to inform future state decision making.

37. "Success" is defined as a 50 percent chance of obtaining a B or higher, or about a 75 percent chance of obtaining a C or higher, in corresponding credit-bearing first-year college courses.

Section I: Assessments of English Language Arts/Literacy and Mathematics

Overview of the Methodology

In 2014, the Council of Chief State School Officers (CCSSO) released a set of criteria approved by the chiefs for evaluating the new generation of assessments designed to measure college and career readiness standards, the “Criteria for Procuring and Evaluating High-Quality Assessments.”³⁸ CCSSO’s stated goal was to provide guidance and support to states seeking to ensure that new assessments aligned to college and career readiness standards are not only valid, reliable, and fair, but also “match the depth, breadth, and rigor of the standards; accurately measure student progress toward college and career readiness; and provide valid data to inform teaching and learning.”

Under the leadership of Brian Gong and Scott Marion, the NCIEA utilized these criteria as the basis for a new methodology for evaluating assessments based on college and career readiness standards, including the CCSS. Our study and the parallel high school study conducted by HumRRO are the first to implement this new methodology.

Below we provide an overview of its key elements. We first describe the various phases of the evaluation—from the review of individual test items to the creation of final program-level summary statements. Next, we present details on the selection of reviewers and test forms as well as the assignment of reviewers to forms. Third, we review the CCSSO Criteria and describe how they were used to address our key questions. Finally, we describe several modifications made to the NCIEA-written methodology, due to issues that arose during implementation.

Study Phases

The study comprises two major phases. The first is a review of the test items (the “item review”) and documentation (the “generalizability review”). The second is the development of program-level (overall) ratings for each CCSSO criterion.

The heart of the evaluation is the review of actual test items and forms. The results of the item reviews are “rolled up” to the form levels, which are then rolled up to the grade level, then to the overall-testing-program level—the latter of which comprise the ratings.

Phase 1: Item and Documentation Review

The item review requires individual evaluation of actual student test items against each of the subject area sub-criteria outlined in the CCSSO Criteria. These sub-criteria are described alongside the criterion-level results in the *Results* section.

38. Council of Chief State School Officers (CCSSO), “Criteria for Procuring and Evaluating High-Quality Assessments” (Washington, D.C.: CCSSO, 2014).

The programs were first asked to provide test item metadata. Metadata refer to descriptive data about test items and passages, such as their alignment to standards, their level of text complexity (i.e., qualitative and quantitative measures of the complexity of text passages), and item type (e.g., multiple choice, constructed-response, and technology-enhanced). These metadata were pre-populated into customized electronic coding sheets and used throughout the review process. Reviewers used these coding sheets to move through each test form one item at a time, rating each item on each of the applicable dimensions. Some of these ratings relied on the metadata provided by the programs, and others were strictly based on expert judgment.³⁹

The methodology also requires that panels evaluate each program's documentation to determine whether the results from reviewing one or two test forms per grade are likely generalizable to the test program as a whole (see *Selection of Forms* for more). In other words, would item-review ratings likely remain the same, improve, or decline if all possible test forms built from the same blueprints and other test specifications had been reviewed?⁴⁰

For the documentation review, testing programs were provided with the criteria to be evaluated in the study and asked to provide documentation that reviewers could use to assess them. Each program provided a large volume of materials they deemed to support the criteria, such as test blueprints, item specifications, and cognitive-demand definitions. These materials detailed their approach to building and validating their assessments to demonstrate that they adequately assessed college and career readiness.

Phase 2: Development of Program Ratings

In accordance with the methodology, reviewers next aggregated the ratings of individual panelists across test forms, grade levels, and programs.

To begin this process, panelists first met to discuss their individual scores and comments for each test form reviewed.⁴¹ The panels then decided on a final “group match score” for each form, which was simply their agreed-upon score after talking through rater differences and settling on a score that represented the collective wisdom of the panel.⁴²

Next, review panels developed final sub-criterion ratings and comments for each grade overall, taking into account the final scores for the two forms at that grade level. Reviewers then inspected each program's scores across the grade 5 and grade 8 forms, considered the results of the documentation review, and came to consensus on the degree to which each program “matched” the respective CCSSO criterion. They issued final ratings on the following scale: Excellent, Good, Limited/Uneven, or Weak Match to the criterion.⁴³ They also developed summary statements with a rationale for the ratings. Finally, review panels developed ratings for the Content and Depth of each assessment, based on the prioritization of criteria recommended in the study methodology (see Section I, Table 3). They also generated final statements summarizing the observed strengths and areas of improvement for each program.

In short, each consecutive step of the evaluation built upon the prior step, which helped to foster shared understanding among reviewers and strengthen the internal consistency of the results. Figure 1 illustrates the entire process for arriving at the final Content and Depth scores, as well as the final summary statements for each assessment program.

39. Field-tested items were not included in the formal review.

40. This portion of the study was conducted jointly between HumRRO and Fordham.

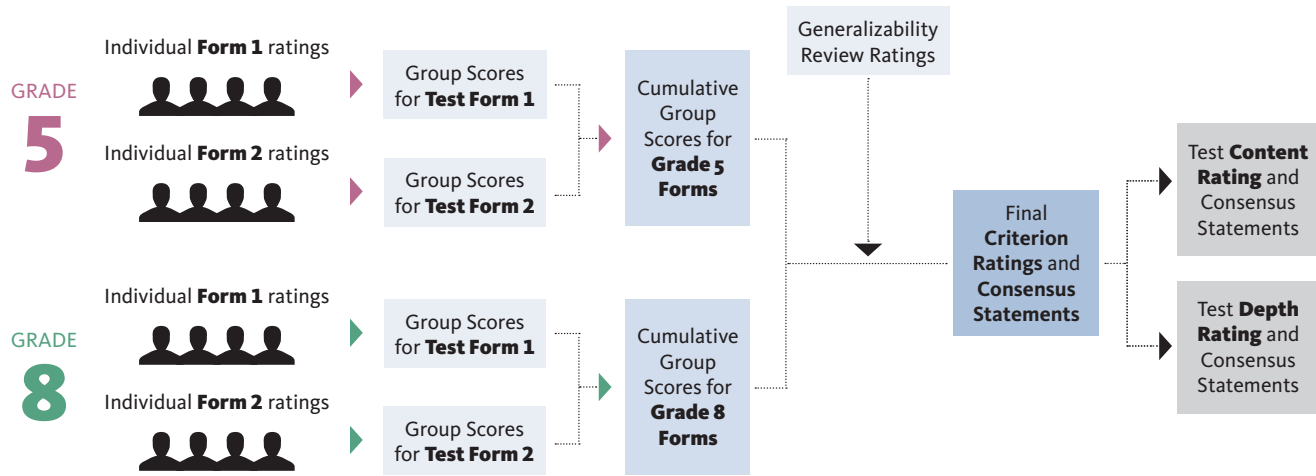
41. Per the methodology, two forms were reviewed for each program per grade and subject area (except for MCAS, which had only one form per grade).

42. In instances where they were unable to reach consensus, minority opinions are reflected in the final statements at the criterion and program level.

43. If the program documentation supported reviewers' form-level findings, suggesting that the results would hold regardless of the number of forms reviewed (i.e., forms developed in different years that are based on the same test blueprints and specifications), the final criterion rating was based solely on the aggregation of the grade 5 and 8 form-level ratings, as explained in the methodology. If, however, the documentation did not demonstrate that the rating would remain the same if additional forms had been reviewed, the panel determined whether or not to adjust the final criterion rating and, if so, stated the rationale.

FIGURE 1

The Process for Generating Final Program Ratings for ELA/Literacy and Mathematics



Selection of Test Forms

The NCIEA's methodology has as its primary goal that evaluators obtain an accurate measure of each test program's content and quality within reasonable time and cost parameters. Thus, for each assessed program, the methodology recommends reviewing two forms, which "should be sufficient basis for an evaluation when coupled with a review of generalizability documentation." This recommendation allows implementers to test for consistency of test quality across forms and helps ensure that results are not idiosyncratic to a particular form. Further, in the case of ELA/Literacy, since each test form can accommodate only a handful of text passages and writing prompts, evaluating two forms gives reviewers a broader look at quality. We followed this two-form recommendation per content area and grade level.⁴⁴

Programs were asked to select two operational forms of their choosing for each subject at each grade level.⁴⁵ They could not be "special forms" (e.g., used only for students with particular special needs).

As indicated, Smarter Balanced is a computer-adaptive test (CAT), meaning each student taking the test can receive a different set of questions based on their response to prior questions (see *Test Program Differences* in the Introduction). Thus there exist hundreds or more possible forms to consider. Methods for analyzing the alignment of computer adaptive assessments are in their infancy.⁴⁶ One approach would be to simply analyze a representative sample of the item pool, ignoring how items were placed on student test forms. Yet doing so would violate a primary study objective, which is to determine whether complete test forms administered to students meet the CCSSO Criteria.

Instead, the NCIEA methodology recommends analyzing two test forms for computer-adaptive tests, one for a student at the 40th percentile of achievement and the other at the 60th percentile, in order to represent a range of possible forms that "typical students" are likely to see. To augment the two forms, the NCIEA also recommends that the program provide simulation data that specifically summarize the characteristics of 1,000

44. The one exception was for MCAS, which has only one operational test form per subject and grade available for review.

45. An operational form refers to a form that was actually used by students.

46. S. L. Wise et al., "Evaluating Content Alignment in Computerized Adaptive Testing," *Educational Measurement: Issues and Practice* (2015; published online ahead of print).

simulated forms per grade and content area. Smarter Balanced provided such information for this study, and the NCIEA summarized those results for our panelists and answered their related questions. The summary explained whether the simulation data provided confidence that the results observed on two forms were likely to apply across the full range of Smarter Balanced forms.

We applied the methodology as written for computer adaptive tests and are confident in the accuracy of our ratings for the forms that panelists reviewed. Still, if we had chosen a different approach to evaluating CAT, such as analyzing larger numbers of forms and at additional student percentiles, our results for Smarter Balanced may have been different.

Selection of Review Panels and Assignment to Forms

Selection and Recruitment of Panelists

The quality and credibility of a complex evaluation of this type rests largely on the expertise and objectivity of the individuals serving on its review panels, such that it is critical to recruit highly qualified yet impartial reviewers with diverse experience and backgrounds. For each panel, we sought to recruit for each panel a mix of practitioners, content experts, and assessment experts.

We began by soliciting reviewer recommendations from each participating testing program and other sources, including content and assessment experts, individuals with experience in prior alignment studies, and several national and state organizations. Finalists were asked to submit CV as well as detailed responses to a questionnaire about their familiarity with the Common Core State Standards, prior experience in conducting alignment evaluations, and potential conflicts of interest. Follow-up phone calls were conducted as necessary. Individuals currently or previously employed by participating testing organizations and writers of the CCSS were not considered. Given that most content and assessment experts have become such by working on prior alignment or assessment-development studies, and that it appeared impossible to find individuals with zero conflicts who are also experts, we prioritized balance and fairness. We recruited at least one reviewer recommended by each testing program to serve on each panel; this strategy helped to ensure fairness by equally balancing reviewer familiarity with the various assessments.⁴⁷

In addition, two university-affiliated content leads facilitated the work of the ELA/Literacy and math review panels. Dr. Charles Perfetti, Distinguished University Professor of Psychology at University of Pittsburgh, served as the ELA/Literacy content lead and Dr. Roger Howe, Professor of Mathematics at Yale University, served as the mathematics content lead. The names and biographical summaries of all panelists appear in Appendix E.

Reviewers received several days of training, both online and in-person, prior to conducting the analysis. See *How Were Reviewers Trained?* for more.

47. One reviewer in ELA/Literacy grade 5 was forced to drop out of the study just before the item review due to medical issues; we were unable to replace her, so some ELA/Literacy grade 5 forms were reviewed by three (and not four) individuals.

How Were Reviewers Trained?

Training for our item review was delivered via several online and in-person sessions. Online training began with a two-hour webinar that described the broader review process and briefly introduced the CCSSO Criteria. The second element of the online training was a series of subject-specific, one-hour presentations by the test vendors describing important features of their assessments, such as their approach to alignment and their form specifications. Panelists also individually reviewed overviews of the methodology and the CCSS for their assigned subject and grade.⁴⁸

In-person training was conducted in Washington, D.C. in the summer of 2015. The three-day training included a series of PowerPoint training modules ranging from one to three hours, each detailing a relevant criterion. Modules laid out the intention of the criterion, the procedures for rating each test item and reading passage, and included a guided-practice activity using sample items. (They also addressed the meaning of key terms used in the methodology.) These modules were developed by Student Achievement Partners (SAP) with input and editing from the study team. The training was delivered by SAP, HumRRO, Dr. Morgan Polikoff, and several consultants, including alignment and assessment experts. Reviewers also received training on how to access the online systems of the various assessments and to correctly complete the scoring worksheets.

Following this training, reviewers independently completed a calibration exercise (led by Dr. Polikoff) using released test items from the 2014 New York Common Core-aligned exam. Individual ratings were collated and examined for disagreement, which was defined as less than 80 percent agreement across raters on any item or passage rating. They were subsequently discussed in-depth at the start of day three, in order to improve consensus, increase inter-rater reliability, and attempt to clarify misunderstandings of the criteria. At this point, the rating of actual test forms began.

The subset of reviewers who participated in the documentation review were also provided training and support by assessment consultant Dr. Jami-Jon Pearson.

Assignment to Forms

There were thirty-two total reviewers across subjects and grade levels. Of these, sixteen were practitioners (practicing teachers in the content area), eight were content experts (from higher education or content area consultants), and eight were assessment experts. Thus, for each grade level and subject there were eight total reviewers, four of whom were practitioners, two who were content experts, and two who were assessment experts (though some had expertise in multiple areas).

For the item review, reviewers were stratified by expertise (practitioners, content experts, and assessment experts) and randomly assigned to test forms. See Table 2 below, which illustrates this jigsaw arrangement.

For each content area and grade level, a total of eight experts reviewed six total forms from ACT Aspire (two forms), PARCC (two forms), and Smarter Balanced (two forms), and eight experts reviewed the form from MCAS. This approach ensured that there were four reviewers for each of the seven test forms at each grade level—see in Table 2 that each column has four X's indicating the four reviewers. These four reviewers included two practicing teachers in the content area, one content expert (either from higher education or a consultant), and one assessment expert. This approach also ensured that there were six unique reviewers across the two test forms for each program at each grade level—except for MCAS, which had only one test form.⁴⁹

48. Separate trainings were provided for the documentation analysis, which were conducted online.

49. The difference in the number of experts is due to the number of forms reviewed for each program. Two forms per grade level (grades 5 and 8) and content area were reviewed for ACT Aspire, PARCC, and Smarter Balanced, and only one form per grade and content area for MCAS.

After the item reviews were completed, a subset of reviewers convened to develop final program ratings. This group comprised seven reviewers in each subject drawn from the sixteen who participated in the item review. As was the case with the item review, each subject-area panel included practitioners, assessment experts, and content experts, plus reviewers recommended by each of the four programs.⁵⁰ Their deliberations occurred over three days via phone and online meetings.

TABLE 2

Reviewer Assignment Jigsaw for Grade 5 Tests in a Given Subject*

Reviewers	ACT Aspire		PARCC		Smarter Balanced		MCAS
	Form 1	Form 2	Form 1	Form 2	Form 1	Form 2	Form 1
Practitioner 1	×			×	×	×	
Practitioner 2	×	×				×	×
Practitioner 3		×	×		×		
Practitioner 4			×	×			×
Content Expert 1	×			×		×	×
Content Expert 2		×	×		×		
Assessment Expert 1	×	×			×	×	
Assessment Expert 2			×	×			×

*Note: The same configuration was implemented for Grade 8 in both subjects.

The Study Criteria

The study methodology evaluates assessment programs against a subset of the CCSSO “Criteria for Procuring and Evaluating High-Quality Assessments.”⁵¹ The Criteria “focus on the critical characteristics that should be met by high-quality assessments aligned to college and career readiness standards” and are based in part on the seminal Standards for Educational and Psychological Testing.⁵² The methodology was created primarily by the NCIEA and refined iteratively over several months in the winter and spring of 2014–15.

The CCSSO Criteria are the backbone of the methodology. These are briefly presented in Tables 3 and 4 (see Appendices A and C for more).

50. Furthermore, all but one of the twenty-eight total test forms (math plus ELA/Literacy) had at least one reviewer on the Phase 2 panel.

51. Council of Chief State School Officers (CCSSO), “Criteria for Procuring and Evaluating High-Quality Assessments” (Washington, D.C.: CCSSO, 2014).

52. American Educational Research Association (AERA), the American Psychological Association (APA), and the National Council on Measurement in Education (NCME) Joint Committee on Standards for Educational and Psychological Testing, “Standards for Educational and Psychological Testing” (Washington, D.C.: AERA, 1999).

TABLE 3

The CCSSO Criteria for the Assessment of ELA/Literacy

	Content	Depth
B.1* focuses on text type (narrative vs. informational) and text quality. It calls for a) a variety of text types, b) an increasing focus on diverse informational texts across grades, and c) the use of high-quality text passages.		×
B.2* focuses on text complexity (qualitative and quantitative). It requires that assessments include texts that have appropriate levels of text complexity for the grade or grade band (grade bands identified in the CCSS are K–5 and 6–12).		×
B.3* focuses on key dimensions of reading items. It includes test items that require close reading and direct textual evidence and that focus on central ideas and important particulars.	×	
B.4 focuses on the cognitive demand of the assessment (i.e., the type of student performance that is required to complete each task on the assessment). It requires an appropriate range of cognitive demand that adequately represents the cognitive demand of the standards.		×
B.5* focuses on key dimensions of writing items. It includes assessing a variety of types and formats of writing and the use of writing prompts that require students to confront and use evidence from texts or other stimuli directly.	×	
B.6 focuses on key dimensions of language and vocabulary items. It includes placing adequate emphasis on language and vocabulary items on the assessment, assessing vocabulary that reflect requirements for college and career readiness, and focusing on common student errors in language questions.	×	
B.7 focuses on key dimensions of research items. It expects that items will require students to analyze, synthesize, organize, and use information from multiple sources.	×	
B.8 focuses on key dimensions of listening and speaking items. While acknowledging that the technology does not yet exist for a complete assessment of these skills, this criterion assesses students' listening skills using passages with adequate complexity and assesses students' speaking skills through oral performance tasks.	×	
B.9 focuses on high-quality items and a variety of item types. It calls for multiple item types, including at least one type where students construct, rather than select, a response as well as high-quality items lacking technical or editorial flaws.		×

Notes:

Nine ELA/Literacy criteria are evaluated. These criteria have between two and eight sub-criteria each, which are the actual targets of the review.

*The methodology suggests that criteria B.3 and B.5 receive greater emphasis in developing the Content rating, and that criteria B.1 and B.2 receive greater emphasis in developing a Depth rating.⁵³ These emphases are indicated by an asterisk.

53. Though the methodology called for B.2 to be emphasized in determining the ELA/Literacy Depth rating, the reviewers ultimately chose not to do so for reasons described in Section III.

TABLE 4

The CCSSO Criteria for the Assessment of Mathematics

	Content	Depth
C.1* focuses on mathematics item content. It requires a heavy focus on the content most crucial for success in later mathematics (the major work of the grade). ⁵⁵	×	
C.2 focuses on the variety of skills assessed. It requires measurement of conceptual understanding, procedural skill, and application in approximately equal proportions.	×	
C.3* focuses on mathematical practices and their connection to content. The criterion calls for students' use of mathematical practices through test items that connect these practices with grade-level content standards.		×
C.4 focuses on the cognitive demand of the assessment. As in ELA/Literacy, this criterion requires an appropriate range of cognitive demand that adequately represents the cognitive demand of the standards.		×
C.5 focuses on high-quality items and a variety of item types. As in ELA/Literacy, this criterion requires multiple item types, including at least one constructed-response type, as well as high-quality items lacking technical or editorial flaws.		×

Notes:

Five mathematics criteria are evaluated. As in ELA/Literacy, each of these criteria has at least two sub-criteria each of which are the targets of the review.

*The methodology suggests that criterion C.1 receives greater emphasis in developing a Content rating and that criterion C.3 receives greater emphasis in developing a Depth rating.⁵⁴ These emphases are indicated by an asterisk.

Methodology Modifications

Fordham and HumRRO were the first organizations to implement the new methodology developed by NCIEA. Thus, while the goal of this initial implementation was to conduct a review as faithful as possible to the original methodology, several logistical or other issues arose that limited our ability to do so. Some pertained to both subjects and others were unique to programs or the content area. We describe seven modifications below.

1 Challenges Associated with Particular Testing Programs

Because the methodology was designed to be neutral with respect to any particular test, it will not perfectly fit each test's design or specifications.

One issue that affected PARCC and Smarter Balanced pertains to the assignment of test items to individual Common Core State Standards. These programs use evidence-centered design (ECD),⁵⁶ which is an approach that bases test-design decisions on the key inferences intended to be made about student performance—also called “claims.” Neither of these programs creates items through a one-to-one mapping of items to standards; instead they map items to the specific evidence statements or targets for their claims. Because our methodology requires metadata that assume a one-to-one or one-to-many mapping of items to specific standards, both programs had to produce metadata they would not normally produce, and that do not fully correspond to a test-construction philosophy based on claims. As a result, both programs' tests conflict in some way with criteria in the methodology. For instance, ELA/Literacy reviewers were asked to evaluate the alignment ratings for individual reading items, but mapping individual items to specific standards is not how either PARCC or Smarter Balanced designed their alignment.

54. The methodology called for C.3 to be emphasized in the determination of the Depth rating but the panel found this criterion to be poorly operationalized and used their professional judgment to reduce its weight. See Section III, *Suggestions for Methodological Improvement*.

55. A list of the major work standards is available at http://achievethecore.org/content/upload/Focus%20in%20Math_091013_FINAL.pdf.

56. ACT Aspire also uses evidence-centered design, but this issue did not affect them because they had item alignment metadata as well.

A unique issue arose with ACT Aspire regarding CCSSO criterion B1. This criterion places narrative informational texts into the “informational text” category along with expository texts. ACT Aspire includes literary fiction and literary nonfiction texts in its Literary Narrative category, but does not classify literary nonfiction texts that are primarily narrative in structure as informational. Review panelists followed the methodology’s guidance when evaluating this criterion.⁵⁷

There was also an issue unique to the MCAS program. For MCAS mathematics, the metadata provided did not include information on mathematical practices. Thus, according to the letter of the C.3 criterion, MCAS did not meet it. However, in reviewing the items, panelists were confident that MCAS items did indeed assess mathematical practices—the program just did not code them in their documentation. Thus panelists took this into account in reaching their final rating for MCAS on this criterion.

2 Item Alignment and Quality

Sub-criteria under B.9 and C.5 call for reviewers to evaluate item quality (including the evaluation of item alignment to standards). In particular, reviewers are asked to verify each of the metadata alignment codes provided by the testing program. Per the methodology, reviewers must agree with all of the alignment codes for 90 percent of test items, in order to receive the top rating for item quality. For multiple reasons,⁵⁸ this procedure would have resulted in vast swaths of items (especially in ELA/Literacy) being rated as poorly aligned, leaving test forms rated weakly on this sub-criterion, even though reviewers believed that items were high quality.

The only approach that resulted in conclusions that had face validity to reviewers was to remove item alignment from B.9 and C.5. Thus, these criteria are based only on item quality and not on judgments of the degree to which reviewers agreed with alignment metadata. Alignment to the standards is still considered in the methodology in the criteria under Content in both subjects, but reviewers did not evaluate alignment to the metadata in their review.

3 Cognitive Demand

The procedures for evaluating criteria B.4 and C.4 (which we evaluated using Webb’s depth of knowledge, or DOK) were intended to ensure that a) the cognitive demand distribution on the assessments adequately matched the depth of the cognitive demand in the standards and b) that the higher-level cognitive demands of the standards, in particular, were not under-assessed. During implementation, reviewers discovered that the distribution of DOK on some assessments was actually higher on average than in the standards. These exams sometimes received low ratings because the DOK distribution did not match the standards, but some reviewers believed that assessments with higher average cognitive demand than the standards should not necessarily be penalized. In calculating final ratings, reviewers used their professional judgment to adjust these cases higher when needed.

4 Test Complexity Metadata

The B.2 sub-criteria focus on text complexity (quantitative and qualitative measures of the difficulty of a reading passage). In order to evaluate these criteria, reviewers need access to the metadata. However, we were unable to include text complexity data in our analysis for several reasons: 1) vendors often used different methods for evaluating text complexity, 2) qualitative text complexity data varied dramatically across programs, and 3) they were often too voluminous to be displayed in the coding worksheet in any readable format. However, as the documentation review did consider text complexity, these results are used in final scoring.

57. For more on ACT Aspire’s interpretation of B1, see Appendix G.

58. For example, some programs indicated that each item aligned to multiple standards; hence, even if reviewers agreed with most of these alignments for a given item, the item would nonetheless be deemed misaligned since the requirement was complete agreement.

5 Major Work of the Grade in Mathematics

The C.1 criterion in mathematics uses the language of focusing “exclusively” on the major work of the grade, which would penalize items that mostly focus on major work but also include some non-major work content. We instructed reviewers to rate an item favorably if it focused “primarily” on major work standards.

6 Application, Conceptual Understanding, and Procedural Skill

Criterion C.2 assesses the extent to which the assessment is adequately balanced across items that assess application, conceptual understanding, and procedural fluency. There were a number of difficulties with the implementation of this criterion; these issues are described in more detail in Section III. The end result is that we chose not to report a criterion-level score, instead providing qualitative feedback only.

7 Weighting of Criteria for Content and Depth Ratings

The methodology recommends that certain criteria be emphasized more heavily as reviewers were rolling up scores from the individual-criterion level to the Content and Depth ratings. In two cases, our reviewers modified these recommendations during implementation. For ELA/Literacy Depth, the recommendation was that B.1 and B.2 be most emphasized in determining the Depth rating; however, because of the issues mentioned above, reviewers did not emphasize the latter in determining the Depth rating. For Mathematics Depth, the methodology recommended that C.3 receive the predominant emphasis. However, as discussed in the Results section, reviewers did not believe the rating for C.3 should drive the Depth rating because of the way it was operationalized. Instead, C.4 and C.5 were equally emphasized.

Overall, we believe the methodology represents a rigorous and detailed approach to evaluating the quality of new assessments against the CCSS and other college and career readiness standards—and were pleased with how our panelists implemented it for this study. As first implementers, we also have a number of recommendations for improving the design and implementation of the methodology in future studies (Section III).

Findings

Findings are organized around three key questions:

- 1 Do the assessments place strong emphasis on the most important content for college and career readiness (CCR), as called for by the Common Core State Standards and other CCR standards? (**Content**)
- 2 Do they require all students to demonstrate the range of thinking skills, including higher-order skills, called for by those standards? (**Depth**)
- 3 What are the overall strengths and weaknesses of each assessment relative to the examined criteria for ELA/Literacy and mathematics? (**Overall Strengths and Weaknesses**)

Results for English language arts appear first, followed by math. For each criterion, we first delineate the requirements for assessments to earn a top score on the evaluation. In cases where it is relevant, we present the tentative scoring guidance proposed by NCIEA (e.g., percentages of items required to score a “2: Meets the criterion”).⁵⁹ These are tentative guidelines that may be revised for future implementations of the methodology.⁶⁰

We also provide additional detail in the form of illustrative quotes and fine-grained analyses to help contextualize the findings. Illustrative quotes are gleaned from reviewers’ comments and responses to open-ended questions, which they were asked to complete after their reviews.⁶¹ Although submitted by individual reviewers, we include

59. Depending on relevance, tentative scoring guidance appears in the report body or the Appendix.

60. The tentative cut-offs were provided to support reviewers in interpreting the more general language of the criteria—but because they were “tentative,” reviewers could (and did) use their professional judgment as they interpreted and applied the criteria.

61. Reviewers were asked to provide feedback on: a) the greatest strength and weakness of the methodology; b) ways in which the methodology could be improved; c) the strengths and weaknesses of each program they evaluated; and d) suggestions for future improvements to each testing program.

only those quotes that represent consensus opinion. In some cases, we provide more fine-grained results than those that appear in the criterion results. (For instance, we present the average proportion of reading items requiring close reading.) Unless otherwise specified, all descriptive statistics are calculated by averaging form-level reviewer ratings across test forms and grades.⁶²

English Language Arts/Literacy Content Criteria

The four test programs varied significantly in the degree to which they emphasize the most important content of the CCSS for the grade level. Summary ratings and panel statements are presented in Table 5.

TABLE 5

Overall ELA/Literacy Content Ratings



ACT Aspire

LIMITED/UNEVEN MATCH

The assessment program includes an emphasis on close reading and language skills.

However, the reading items fall short on requiring students to cite specific textual information in support of a conclusion, generalization, or inference and in requiring analysis of what has been read. In order to meet the criteria, ACT Aspire should strengthen assessing writing to sources, vocabulary, and research and inquiry.



MCAS

LIMITED/UNEVEN MATCH

The assessment requires students to read closely well-chosen texts and presents test questions of high technical quality.

However, the program would be strengthened by assessing writing annually, assessing the three types of writing called for across each grade band, requiring writing to sources, and placing greater emphasis on assessing research and language skills.



PARCC

EXCELLENT MATCH

The program demonstrates excellence in the assessment of close reading, vocabulary, writing to sources, and language, providing a high-quality measure of ELA/Literacy content, as reflected in college and career readiness standards.

The tests could be strengthened by the addition of research tasks that require students to use two or more sources and, as technologies allow, a listening and speaking component.



Smarter Balanced

EXCELLENT MATCH

The program demonstrates excellence in the areas of close reading, writing to sources, research, and language. The listening component represents an important step toward adequately measuring speaking and listening skills—a goal specifically reflected in the standards. Overall, Smarter Balanced is a high quality measure of the content required in ELA/Literacy and literacy, as reflected in college and career readiness standards.

A greater emphasis on Tier 2 vocabulary would further strengthen these assessments relative to the criteria.

Note: These ratings are based on five criteria (B.3, B.5, B.6, B.7, and B.8); the first two receive greater emphasis.

62. Proportions presented in the text and used in determining form-level ratings were typically based on percentages of score points associated with items. Calculating proportions based on item counts instead may have resulted in somewhat different ratings. This may be especially true for computer-adaptive tests such as Smarter Balanced, where the number of score points does not necessarily correspond to the actual weight in determining the final score. This note applies to all criteria for which score points were used.

Criterion B.3

Do the tests require students to read closely and use evidence from texts to obtain and defend responses?

This criterion reflects the high priority in the CCSS that students be able to read and understand increasingly complex texts, both literary and informational. The following were required to fully meet this criterion:

- 1 *Nearly all reading items require close reading and analysis of text, rather than skimming, recall, or simple recognition of paraphrased text.*
- 2 *Nearly all reading items focus on central ideas and important particulars.*
- 3 *Nearly all items are aligned to the specifics of the standards.*
- 4 *More than half of the reading score points are based on items that require direct use of textual evidence.*

As shown in Table 6, both PARCC and Smarter Balanced received the highest rating (Excellent Match) on this criterion. On average across forms and grades, 87 percent of PARCC reading items and 69 percent of Smarter Balanced reading items were scored as requiring direct textual evidence. The 2014 MCAS missed the Excellent rating and fell to Good Match largely due to an insufficient number of items that require direct citing of evidence from texts—just 29 percent. As one MCAS ELA/Literacy grade 8 reviewer explained, “While ... students likely had to use the information in the texts in order to answer the questions, the items, as a whole, did not require students to provide direct evidence from the text.”

TABLE 6

Reading Summary (Criterion B.3)

L	ACT Aspire LIMITED/UNEVEN MATCH
Although most reading items require close reading of some kind, too many can be answered without analysis of what was read. Items that purport to require specific evidence from text often require only recall of information from text. To meet this criterion, the test items should require students to cite specific text information in support of some conclusion, generalization, or inference drawn from the text.	
G	MCAS GOOD MATCH
Most reading items require close reading and focus on central ideas and important particulars. Some questions, however, do not require the students to provide direct textual evidence to support their responses. In addition, too many items do not align closely to the specifics of the standards.	
E	PARCC EXCELLENT MATCH
Nearly all reading items require close reading, the understanding of central ideas and the use of direct textual evidence.	
E	Smarter Balanced EXCELLENT MATCH
Nearly all reading items align to the reading standards requiring close reading, the understanding of central ideas, and use of direct textual evidence in support of a conclusion, generalization, or inference.	

The ACT Aspire test was rated as Limited/Uneven Match, falling short both in the percentage of items requiring direct textual evidence (51 percent) and in the percentage of items requiring analysis of text (75 percent, as compared to 89 percent or more for the other three programs). Multiple ELA/Literacy reviewers commented on ACT Aspire’s lack of questions requiring students to use direct textual evidence or reference multiple texts.

Criterion B.5

Do the tests require students to write narrative, expository, and persuasive/ argumentation essays (across each grade band, if not in each grade) in which they use evidence from sources to support their claims?

The following were required to fully meet this criterion:

- 1 *All three writing types are approximately equally represented across all forms in the grade band (K–5; 6–12), allowing blended types (i.e., writing types that blend two or more of narrative, expository, and persuasive/ argumentation) to contribute to the distribution.*
- 2 *All writing prompts require writing to sources (meaning they are text-based).*

The grade bands identified in the CCSS are K–5 and 6–12. Because this study evaluated one grade from each band (grades 5 and 8), panelists also reviewed documentation for each program to check their writing requirements across each grade band. As shown in Table 7, PARCC and Smarter Balanced each received a rating of Excellent Match, as each writing type was assessed at least once across each grade band and all extended writing tasks required students to use sources. As one ELA/Literacy grade 8 reviewer commented for PARCC, “writing tasks were highly analytical and required close reading and thoughtful analysis.”

TABLE 7

Writing Summary (Criterion B.5)

L	<p>ACT Aspire LIMITED/UNEVEN MATCH</p> <p>Although the program documentation shows that a balance of all three writing types is required across each grade band, the writing prompts do not require writing to sources. As a result, the program insufficiently assesses the types of writing required by college and career readiness standards.</p>
W	<p>MCAS WEAK MATCH</p> <p>Writing is assessed at only one grade level per band, and there is insufficient opportunity to assess writing of multiple types. In addition, the writing assessments do not require students to use sources. As a result, the program inadequately assesses the types of writing required by college and career readiness standards.</p>
E	<p>PARCC EXCELLENT MATCH</p> <p>The assessment meets the writing criterion, which requires writing to sources. Program documentation shows that a balance of all three writing types is required across each grade band.</p>
E	<p>Smarter Balanced EXCELLENT MATCH</p> <p>The writing items are of high quality, and the writing prompts all require the use of textual evidence. Program documentation shows that a balance of all three writing types is required across each grade band.</p>

ACT Aspire’s writing test, available annually at grades 3–8 and early high school, does include the three types of writing called for by the criterion, but does not require students to write using evidence from sources, resulting in a rating of Limited/Uneven Match. Multiple ELA/Literacy reviewers commented on how infrequently students were asked to write to and compare multiple passages.

Given that MCAS does not evaluate writing at grades 5 and 8, it received a score of Limited/Uneven Match.

MCAS does, however, evaluate writing at grades 4 and 7, so we asked our panelists to review those writing items and evaluate them based on the associated CCSS grade-level standards. This exercise was conducted only to give the program information about its treatment of writing; it was not a part of the final scoring. Ultimately, both the writing prompts on the fourth- and seventh-grade assessments did not require writing to sources, which likely would have also resulted in a Weak Match had those prompts been a part of the actual review. As one ELA/Literacy grade 8 reviewer remarked, “the biggest flaw in showing college and career readiness is that writing is not assessed at every grade; therefore, there is very little evidence [relative] to students’ abilities in the various writing modes.”

Criterion B.6

Do the tests require students to demonstrate proficiency in the use of language, including academic vocabulary and language conventions, through tasks that mirror real-world activities?

This criterion calls for assessments to include vocabulary questions that focus on words important to the central ideas of a text and that require students to use context to determine meaning. In addition, the tests should focus on what are referred to as Tier 2 words. According to the standards, these are “general academic” words that are far more likely to appear in written text than in speech. They “appear in all sorts of texts: informational texts (words such as *relative*, *vary*, *formulate*, *specificity*, and *accumulate*), technical texts (*calibrate*, *itemize*, *periphery*), and literary texts (*misfortune*, *dignified*, *faltered*, *unabashedly*).”⁶³

This same criterion also addresses language conventions. As opposed to asking a student to select the correctly punctuated sentence from a list of four or five options, this criterion calls for items, whether multiple-choice or technology-enhanced, that mirror real-world activities (i.e., place errors in context and require students to edit or revise for clarity or correctness).

The following were required to fully meet this criterion:

- 1 *The large majority of vocabulary items (i.e., three-quarters or more) focus on Tier 2 words and require the use of context, and more than half assess words important to central ideas.*
- 2 *A large majority (i.e., three-quarters or more) of the items in the language skills component and/or scored with a writing rubric (i.e., points in writing tasks that are allocated toward a language sub-score), mirror real-world activities, focus on common errors, and emphasize the conventions most important for readiness.*
- 3 *Vocabulary is reported as a sub-score or at least 13 percent of score points are devoted to assessing vocabulary/language.*
- 4 *Language skills are reported as a sub-score or at least 13 percent of score points are devoted to assessing language skills (language skills items plus score points).*

As shown in Table 8, PARCC excelled on this criterion, receiving high marks for both the vocabulary (85 percent of vocabulary items were on Tier 2 words and phrases) and the language items, leading to a score of Excellent Match. As one grade 8 reviewer commented, “[PARCC] items were generally great, especially the technology-enhanced vocabulary items.” ACT Aspire and Smarter Balanced each received a rating of Good Match, and in each case the deficiency was in the number of vocabulary items that assess Tier 2 words (47 percent and 74 percent of

63. Common Core State Standards, “ELA/Literacy, Appendix A,” http://www.corestandards.org/assets/Appendix_A.pdf, 33–36.

vocabulary items, respectively), whereas these assessments met the criteria for language and for the percent of score points devoted to language and vocabulary. MCAS's vocabulary items were generally adequate at grade 8 but insufficient at grade 5 (71 percent of vocabulary items included Tier 2 words and phrases, falling just below the 75 percent target), and the few language conventions items at grade 5 failed to require students to edit or revise for clarity or correctness, resulting in a rating of Limited/Uneven Match.

TABLE 8

Vocabulary and Language Skills Summary (Criterion B.6)



ACT Aspire

GOOD MATCH

Language items meet the criterion for being tested within writing activities, though more items are needed that are embedded in real world tasks such as editing. The vocabulary items do not meet the criterion because there are too few of them and not enough assess Tier 2 words.



MCAS

LIMITED/UNEVEN MATCH

Vocabulary items are sufficient and generally aligned to the criterion; however, the grade 5 items need more words at the Tier 2 level. Furthermore, a lack of program documentation means that the quality of vocabulary assessments cannot be substantiated across forms. MCAS does not meet the criterion for assessing language skills, which call for them to be assessed within writing assessments that mirror real-world activities including editing and revision.



PARCC

EXCELLENT MATCH

The test contains an adequate number of high-quality items for both language use and Tier 2 vocabulary and awards sufficient score points, according to the program's documentation, to both of these areas.



Smarter Balanced

GOOD MATCH

Language skill items are contained in a sub-score and meet the criterion for being assessed within writing and mirroring real-world activities such as editing and revision. The number of items that test vocabulary is a bit low; further, items coded as vocabulary too often did not test Tier 2 vocabulary words.

Criterion B.7

Do the tests require students to demonstrate research skills, including the ability to analyze, synthesize organize, and use information from sources?

This criterion presents a new challenge to test developers in that it requires development of tasks that contain multiple sources and assess whether students have accurately gleaned, analyzed, and synthesized evidence from those sources in their response. The criterion also requires that the tasks mirror real-world research activities, which is difficult when such activities typically involve online searches conducted over days, weeks, or months.

The following was required to fully meet this criterion:

- 1 *Three-quarters or more of the research items on each test form require analysis, synthesis, and/or organization of information.*

PARCC and Smarter Balanced included research tasks within their performance tasks. (The latter are extended tasks with multiple questions that ask students to analyze and use evidence from at least two sources, which may include audio, video, and text.) Both programs earned a rating of Excellent Match (Table 9).

ACT Aspire included a single item per test form that required analysis and organization of information, but the panels found this to be inadequate to earn more than a rating of Limited/Uneven Match. The 2014 MCAS did not contain any items that assessed research skills, resulting in a rating of Weak Match.

Another CCSSO criterion is to be evaluated “over time, as assessment advances allow.” Criterion B.8 reads, “Do the tests measure students’ speaking and listening communication skills?” Only one program, Smarter Balanced, has incorporated listening items, which the panel commended; none of the programs assess speaking skills at this time. Because the criterion calls for speaking and listening to be assessed over time, ratings for this criterion were not included in the overall ELA/Literacy content ratings.

TABLE 9

Research and Inquiry Summary (Criterion B.7)



ACT Aspire

LIMITED/UNEVEN MATCH

Although the one item at each grade level involving research and inquiry did indeed require analysis and organization of information, this single item is insufficient to provide a quality measure of research and inquiry.



MCAS

WEAK MATCH

The assessment has no test questions devoted to research.



PARCC

EXCELLENT MATCH

The research items require analysis, synthesis, and/or organization, as well as the use of multiple sources, therefore meeting the criterion for Excellent.



Smarter Balanced

EXCELLENT MATCH

The research items require analysis, synthesis, and/or organization, as well as the use of multiple sources, therefore meeting the criterion for Excellent.

English Language Arts/Literacy Depth Criteria

All four programs fared very well on Depth, which required students to demonstrate the higher-order thinking skills called for by the CCSS, via high-quality items that reflect a variety of item types (Table 10).

PARCC received the highest rating of Excellent Match and the remaining three programs received ratings of Good Match.

TABLE 10

Overall ELA/Literacy Depth Ratings



ACT Aspire

GOOD MATCH

The program's assessments are built on high-quality test items and texts that are suitably complex.

To fully meet the CCSSO Criteria, more cognitively demanding test items are needed at both grade levels, as well as additional literary narrative text, as opposed to literary informational texts.



MCAS

GOOD MATCH

The assessments do an excellent job in presenting a range of complex reading texts.

To fully meet the demands of the CCSSO Criteria, however, the test needs more items at higher levels of cognitive demand, a greater variety of items to test writing to sources and research, and more informational texts, particularly those of an expository nature.



PARCC

EXCELLENT MATCH

The PARCC assessments meet or exceed the depth and complexity required by the criteria through a variety of item types that are generally high quality.

A better balance between literary and informational texts would further strengthen the assessments in addressing the criteria.



Smarter Balanced

GOOD MATCH

The assessments use a variety of item types to assess student reading and writing to sources.

The program could better meet the depth criteria by increasing the cognitive demands of the grade 5 assessment and ensuring that all items meet high editorial and technical quality standards.

Note: These ratings are based on four criteria (B.1, B.2, B.4, B.9); the first two receive greater emphasis.

Criterion B.1

Do the tests require a balance of high-quality literary and informational texts?

The CCSS place great emphasis on students being able to read literary texts (e.g., classic novels), informational texts (e.g., newspapers, scientific and historical texts), as well as technical documents, in order to build their content knowledge.

The following were required to fully meet this criterion:

- 1 *Approximately half of the texts at grades 3–8 and two-thirds at high school are informational, and the remainder literary.*
- 2 *Nearly all passages are high quality (previously published or of publishable quality).*
- 3 *Nearly all informational passages are expository in structure.*
- 4 *For grades 6–12, the informational texts are split nearly evenly for literary nonfiction, history/social science, and science/technical.*

The texts used across all four programs were found to be of high quality (previously published or of publishable quality). Smarter Balanced tests contained both the requisite balance and text quality, earning a rating of Excellent Match. However, ACT Aspire, MCAS, and PARCC were all found to have a slight imbalance across literary and informational texts, resulting in a rating of Good Match (Table 11).⁶⁴

TABLE 11

Text Quality and Types Summary (Criterion B.1)



ACT Aspire

GOOD MATCH

The texts are of high quality, and the proportion of informational texts meets the criterion.

The assessment would better align to the criterion, however, with additional literary narrative text, as opposed to literary informational text.



MCAS

GOOD MATCH

The quality of the texts is very high.

Regarding the balance of text types, some forms had too few informational texts.



PARCC

GOOD MATCH

Although the passages are consistently of high quality, the tests would have better reflected the criterion with additional literary nonfiction passages.



Smarter Balanced

EXCELLENT MATCH

Overall text quality is high, and among informational texts there is a high proportion of expository text types.

64. ACT Aspire does not classify literary nonfiction texts that are primarily narrative in structure as “informational.” See Appendix G for more on ACT Aspire’s interpretation of CCSSO criterion B.1.

Criterion B.2

Do the tests require appropriate levels of text complexity, increasing the level each year so that students are ready for the demands of college and career by the end of high school?

As explained in the methodology, the intent of this criterion was to gather quantitative and qualitative data from the testing programs regarding text complexity to determine whether reading passages had been assigned to grade bands and grade levels appropriately. Unfortunately, these data could not be used (see Section I, *Methodology Modifications*), so the evaluation of this criterion was based solely on the requirements recorded within each program's test documentation and specifications.

The following were required to fully meet this criterion:

- 1 *Documentation clearly explains how quantitative data are used to determine grade band placement.*
- 2 *Texts are placed at the grade level recommended by the qualitative review.*

In every case, the documentation met the criterion (Table 12). However, because the documentation is not a guarantee of what will appear on actual test forms, the panel decided on a maximum rating of Good Match. All four programs received this rating.

TABLE 12

Complexity of Texts Summary (Criterion B.2)

ACT Aspire GOOD MATCH

It is based solely on the review of program documentation, which is determined to have met the criterion. The test blueprints and other documents clearly and explicitly require texts to increase in complexity grade-by-grade and for texts to be placed in grade bands and grade levels based on appropriate quantitative and qualitative data.

MCAS GOOD MATCH

It is based solely on the review of program documentation, which is determined to have met the criterion. The test blueprints and other documents clearly and explicitly require texts to increase in complexity grade-by-grade and for texts to be placed in grade bands and grade levels based on appropriate quantitative and qualitative data.

PARCC GOOD MATCH

It is based solely on the review of program documentation, which is determined to have met the criterion. The test blueprints and other documents clearly and explicitly require texts to increase in complexity grade-by-grade and for texts to be placed in grade bands and grade levels based on appropriate quantitative and qualitative data.

Smarter Balanced GOOD MATCH

It is based solely on the review of program documentation, which is determined to have met the criterion. The test blueprints and other documents clearly and explicitly require texts to increase in complexity grade-by-grade and for texts to be placed in grade bands and grade levels based on appropriate quantitative and qualitative data.

Criterion B.4

Are all students required to demonstrate a range of high order, analytical thinking skills in reading and writing based on the depth and complexity of the standards?

Research over the last decade has shown that existing state tests were not adequately testing—and therefore not requiring the teaching of—higher-order skills such as analysis, synthesis, the development of a logical argument, and use of concepts to solve non-routine problems.⁶⁵ The CCSSO Criteria recognize these shortcomings of prior tests and recommend that “all students demonstrate a range of higher-order, analytical thinking skills in reading and writing based on the depth and complexity of college and career ready standards....” Webb’s Depth of Knowledge (DOK) taxonomy⁶⁶ was used in this study to classify thinking skills (see Appendix A). In this taxonomy, level 1 is the lowest level (recall), level 2 requires use of a skill or concept, and levels 3 and 4 are considered the higher-order thinking skills.

To receive the highest rating on this criterion, the distribution of cognitive demand on test forms had to match the distribution of cognitive demand of the standards as a whole and match the higher cognitive demand (DOK 3+) of the standards. Note that criterion B.4 is not a rating of test difficulty. Assessments that do not match the distribution of complexity of the standards, including if they have too many high DOK items, may receive a rating of less than Excellent Match.

TABLE 13

Matching the Complexity of the Standards (Criterion B.4)



ACT Aspire

WEAK MATCH

To better reflect the depth and complexity of the standards, both grade-level tests should require more items with higher cognitive demands, although this problem is greater at grade 8.



MCAS

LIMITED/UNEVEN MATCH

More items that measure the higher levels of cognitive demand are needed to sufficiently assess the depth and complexity of the standards.



PARCC

EXCELLENT MATCH

The test is challenging overall; indeed, the cognitive demand of the grade 8 test exceeds that of the CCSS.



Smarter Balanced

GOOD MATCH

The cognitive demand of items cover a sufficient range and, in grade 8, the percentage of more demanding items (DOK 3 and 4) correspond well to the demand of the standards. However, the grade 5 test needs more items at higher levels of cognitive demand to reflect fully the depth and complexity of the standards.

65. Yuan and Le, 2012.

66. We chose the Webb DOK approach because it is widely used and familiar; however, results might have been different if another approach to cognitive complexity had been used. Future iterations of the methodology could use newer or more innovative approaches to cognitive complexity if desired.

To determine whether the tests reflected the depth and complexity of the standards, we first needed to measure those aspects of the CCSS. As described in Appendix A, the study team contracted with content experts in advance of the study to determine the distribution of the cognitive demand of the grade 5 and grade 8 CCSS standards.





The results show a great deal of variation across programs in their emphasis on higher-level skills and match to the standards. PARCC scored an Excellent Match, as the DOK of PARCC at both grade levels exceeded that of the standards.⁶⁷ Smarter Balanced scored Good Match for meeting the DOK of the standards at grade 8 but underemphasizing DOK 3–4 at grade 5. MCAS scored Limited/Uneven Match for underemphasizing DOK 3–4 by roughly 10–20 percent at each grade. ACT Aspire scored Weak Match for underemphasizing DOK 3+ by roughly 30 percent at grade 8.

The results of the DOK analysis for the CCSS and the four programs are shown in Table 14. Interestingly, at each grade, 46 percent of standards content require the use of higher-order skills (level 3 or 4).

More detail on this analysis, including the distributions of DOK on other national and international tests, appears in Appendix A.

TABLE 14

The Distribution of Cognitive Demand in ELA/Literacy: The CCSS vs. Tests

		Grade 5			Grade 8		
		LEVEL 1	LEVEL 2	LEVELS 3 & 4	LEVEL 1	LEVEL 2	LEVELS 3 & 4
CCSS		18%	37%	46%	10%	44%	46%
ACT Aspire		34%	38%	28%	46%	36%	18%
MCAS		10%	63%	27%	5%	59%	37%
PARCC		5%	45%	50%	2%	29%	69%
Smarter Balanced		19%	59%	22%	15%	41%	44%

LEGEND  Excellent Match  Good Match  Limited/Uneven Match  Weak Match

Note: Percentages in the table represent percentages of score points at each DOK level. Results for a particular grade and program were generated by averaging across all raters and forms for that grade and program (e.g., averaging the four raters of ACT Aspire's form 1 and the four raters of ACT Aspire's form 2 at grade 5).

67. The review focused primarily on match to the DOK of the standards, but reviewers also used professional judgment in determining final ratings.

Criterion B.9

Are a variety of item types used, including at least one that requires students to generate, rather than select, a response? Are the test items of high quality?

One aspect of item quality that the CCSSO Criteria emphasize is a diversity of item types. Specifically, the criteria recognize the widespread dissatisfaction with the ability of multiple-choice-only tests to fully measure complex student skills.

The CCSSO Criteria also address the editorial and technical quality of items. On all standardized assessments of student learning, it is imperative that items be of high editorial and technical quality and accuracy. The tests evaluated in this study are used for accountability purposes, so item quality must be above reproach. Each of these programs requires multiple rounds of item review and field-testing; even so, items with quality issues can sometimes make their way into the tests.

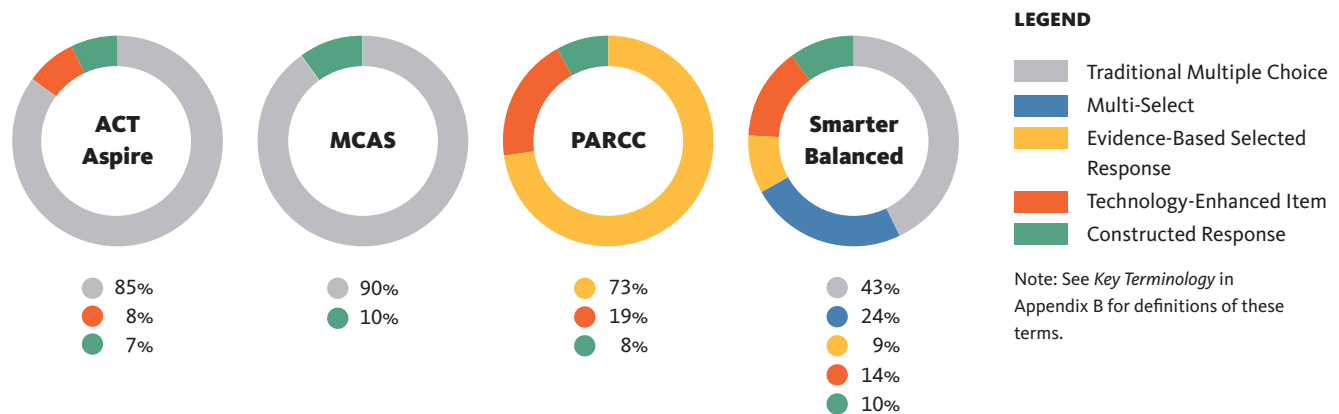
The following were required to fully meet this criterion:

- 1 *At least two item formats are used, including one that requires students to generate, rather than select, a response.*
- 2 *All or nearly all operational items reviewed reflect high editorial and technical quality and accuracy.*

All four programs use multiple item types, including at least one student-constructed response type, so they all met this portion of the criterion. The actual variety of item types, however, differed much more significantly, as shown in Table 15. Reviewers commended the innovative and appropriate use of technologies within Smarter Balanced and PARCC items, such as audio files in listening items and text editing. The ACT Aspire assessments, while also computer-based, had a much more limited set of item types, with heavy reliance on traditional multiple-choice items. The 2014 MCAS is a paper-and-pencil test that has a near-complete focus on traditional multiple-choice items, cited by multiple reviewers as a serious limitation.

TABLE 15

Distribution of Item Types in the ELA/Literacy Tests



The results for item quality are very strong (Table 16). Three programs—ACT Aspire, MCAS and PARCC—were rated as Excellent Match. These programs were each seen as having exceptionally high editorial accuracy and technical quality (less than 5 percent of items with quality issues). On Smarter Balanced items, some items (approximately one item per form) contained quality issues, ranging from editorial errors to readability or technical quality concerns, resulting in a score of Good Match.⁶⁸ For instance, reviewers sometimes felt that multiple answers could be considered correct even on items that were scored as having just one correct answer; identified items in which students had to define vocabulary from context as having insufficient context to determine the definition; and noted minor typographical issues, such as punctuation and spelling.⁶⁹

TABLE 16

High-Quality Items and Variety of Item Types (Criterion B.9)



ACT Aspire

EXCELLENT MATCH

The test includes items that exhibit high technical quality and editorial accuracy. Multiple item formats are used, including student-constructed responses.



MCAS

EXCELLENT MATCH

Multiple item formats are used, including student-generated response items. The items exhibit high technical quality and editorial accuracy. The paper-and-pencil format precludes the use of technology-enhanced items, but the criterion for multiple item types is met.



PARCC

EXCELLENT MATCH

The tests use multiple item formats, including student-constructed responses.



Smarter Balanced

GOOD MATCH

The tests use multiple formats and technology-enhanced items including constructed responses. However, editorial or technical issues, including readability, were noted in a number of items.

68. The nature and timing of this review required Smarter Balanced to make the test items and forms available to reviewers through an alternate test interface that was more limited than the actual student interface used for the summative assessments, particularly with regard to how items appeared on the screen and how erroneous responses were handled. Though reviewers were not able to determine the extent to which these interface limitations impacted their findings, the study team worked with Smarter Balanced to ascertain which item issues were caused by interface differences and which were not. All item-relevant statements in the report reflect data not prone to interface differences.

69. Some Smarter Balanced “multi-select” items have more than one answer by design.

Mathematics Content Criteria

The mathematics review panel found that the four test programs varied somewhat in the degree to which they emphasize the most important content of the CCSS at each grade level (Table 17). PARCC and Smarter Balanced earned ratings of Good Match to the CCSSO Criteria. The ACT Aspire tests and the 2014 MCAS received ratings of Limited/Uneven Match.

TABLE 17

Overall Mathematics Content Ratings

L	<p>ACT Aspire LIMITED/UNEVEN MATCH</p> <p>The program does not focus exclusively on the major work of the grade, but rather, by design, assesses material from previous and later grade(s). This results in a weaker match to the criteria.</p> <p>The tests could better meet the criteria at both grades 5 and 8 by increasing the number of items that assess the major work.</p>
L	<p>MCAS LIMITED/UNEVEN MATCH</p> <p>While the grade 8 assessment focuses strongly on the major work of the grade, the grade 5 assessment does not, as it samples more broadly from the full range of standards for the grade.</p> <p>The tests could better meet the criteria through increased focus on the major work of the grade on the grade 5 test.</p>
G	<p>PARCC GOOD MATCH</p> <p>The test could better meet the criteria by increasing the focus on the major work at grade 5.</p>
G	<p>Smarter Balanced GOOD MATCH</p> <p>The tests could better meet the criteria by increasing the focus on the major work in grade 5.</p>

Note: Because of methodological challenges pertaining to criterion C.2 (see Section III), the Content ratings for mathematics are based primarily on the ratings for C.1.

Criterion C.1

Do the tests focus strongly on the content most needed for success in later mathematics?

Certain CCSS standards (termed “major work of the grade”) have been deemed most important for ensuring that students are on track toward readiness for college or career. This criterion reflects that expectation. To score well on it, assessments must devote adequate attention toward these major work standards and also assess the full breadth of the major work (i.e., must assess, with at least one item, a large proportion of the major work “clusters”).

The following was required to fully meet this criterion:

- 1 *The vast majority (i.e., at least three-quarters at elementary grades, at least two-thirds in middle school grades, and at least half in high school) of score points in each assessment focuses on the content that is most important for students to master in that grade in order to reach college and career readiness (also called the major work of the grade), and at least 90 percent of the major work clusters must be assessed by at least one item.*

As shown in Table 18, both PARCC and Smarter Balanced received a rating of Good Match on this criterion, falling short of the top rating because the grade 5 assessment was insufficiently focused (less than 75 percent of score points for each program) on the major work. The MCAS was rated Limited/Uneven Match, again due to insufficient focus in fifth grade. Only about 53 percent of MCAS score points focused on the major work, far below the 75 percent threshold, and also lower than PARCC (that just missed it at 72 percent) and Smarter Balanced (66 percent).

The ACT Aspire received the lowest score of Weak Match on this criterion because both grades' tests were far below the major work tentative cutoff (32 percent of score points on major work for grade 5 and 44 percent for grade 8, as compared to a cutoff of 75 percent). Though reviewers recognized that ACT Aspire's design intentionally included lower grade items, even those lower grade items were not concentrated on the major work of the respective grade. Reviewers for both grades 5 and 8 commented that ACT Aspire's test questions did not focus adequately on the major work of the grade.

TABLE 18

Focus Summary (Criterion C.1)



ACT Aspire

WEAK MATCH

ACT Aspire forms do not consistently place sufficient emphasis on the major work of the given grade, due in part to intentional test design, which requires inclusion of selected content from earlier and later grades. Still, many of the items coded to standards from lower grades do not address the major work of the relevant grade.



MCAS

LIMITED/UNEVEN MATCH

The grade 8 assessment is focused on the major work of the grade. The grade 5 assessment is significantly less focused on the major work of the grade than called for by the criterion, as it samples content across the full set of standards for the grade.



PARCC

GOOD MATCH

While the grade 8 tests focus strongly on the major work of the grade, the grade 5 tests fall short of the threshold required for the top rating.



Smarter Balanced

GOOD MATCH

While the grade 8 tests focus strongly on the major work of the grade, the grade 5 tests fall short of the threshold required for the top rating.

Criterion C.2

Do the tests assess a balance of concepts, skills, and applications?

The CCSS require students to demonstrate conceptual understanding and procedural skill/fluency, and apply this knowledge.

The following were required to fully meet this criterion:

- 1 *On each test form, at least 25 percent and no more than 50 percent of score points are allocated to each of the three categories: mathematical concepts, procedures/fluency, and applications.*

As described in Section III, there were a number of methodological challenges in implementing this criterion.⁷⁰ As a consequence, our reviewers agreed to provide only qualitative statements, rather than the ratings awarded to the other criteria. This criterion, therefore, was not used in the determination of the overall Content rating.

In general, the test forms from all four programs showed attention to conceptual understanding, procedural skill, and application. However, each program fell short of the goal of balance (which was operationalized as an even distribution) in one way or another. For ACT Aspire at both grades, reviewers noted that items directly assessing procedural skill were underrepresented. For MCAS at grade 5, reviewers found few items assessing conceptual understanding and an overabundance of application items. The grade 5 PARCC exam similarly had an overabundance of application items, some of which reviewers noted had shallow contexts. Finally, the Smarter Balanced exams at both grade levels had a slight wealth of application items, and reviewers also noticed that some forms were more heavily focused on applications than others.

70. To wit: All four programs require, in their program documentation, the assessment of conceptual understanding, procedural skill/ fluency, and application, although most do not clearly distinguish between procedural skill/fluency and conceptual understanding. Also, specific balance across these three types is not required. Due to variation across reviewers in how this criterion was understood and implemented, final ratings could not be determined with confidence.

Mathematics Depth Criteria

All four programs fared well on Depth in mathematics (Table 19), requiring all students to demonstrate the higher-order thinking skills called for by the CCSS. MCAS received a rating of Excellent Match, and the other three programs received a rating of Good Match.

TABLE 19

Overall Mathematics Depth Ratings



ACT Aspire

GOOD MATCH

The items are well crafted and clear, with only rare instances of minor editorial issues.

The ACT Aspire tests include proportionately more items at high levels of cognitive demand (DOK 3) than the standards reflect, and proportionately fewer at the lowest level (DOK 1). This finding is both a strength, in terms of promoting strong skills, and a weakness, in terms of ensuring adequate assessment of the full range of cognitive demand within the standards.

While technically meeting the criterion for use of multiple item types, the range is nonetheless limited, with the large majority comprising multiple-choice items.

The program would better meet the criteria for Depth by including a wider variety of item types and relying less on traditional multiple-choice items.



MCAS

EXCELLENT MATCH

The assessment uses high-quality items and a variety of item types. The range of cognitive demand reflects that of the standards of the grade. While the program does not code test items to math practices, mathematical practices are nonetheless incorporated within items.

The program might consider coding items to the mathematical practices and making explicit the connections between specific practices and specific content standards.



PARCC

GOOD MATCH

The tests include items with a range of cognitive demand, but at grade 8 that distribution contains a higher percentage of items at the higher levels (DOK 2 and 3) and significantly fewer items at the lowest level (DOK 1). This finding is both a strength, in terms of promoting strong skills, and a weakness, in terms of ensuring adequate assessment of the full range of cognitive demand within the standards.

The tests include a variety of item types that are largely of high quality. However, a range of problems (from minor to severe) surfaced relative to editorial accuracy and, to a lesser degree, technical quality.

The program could better meet the Depth criteria by ensuring that all items meet high editorial and technical standards and by ensuring that the distribution of cognitive demand on the assessments provides sufficient information across the range.



Smarter Balanced

GOOD MATCH

The exam includes a range of cognitive demand that fairly represents the standards at each grade level.

The tests have a strong variety of item types, including those that make effective use of technology. However, a range of problems (from minor to severe) surfaced relative to editorial accuracy and, to a lesser degree, technical quality. A wide variety of item types appear on each form, and important skills are assessed with multiple items, as is sound practice. Yet, individual forms sometimes contained two or three items measuring the same skill that were nearly identical, with only the numerical values changed in the item stem and a different set of answer choices. Such near-duplication may not impact the accuracy of the score, but a greater variety of question stems/scenarios is desirable.

The program could better meet the Depth criteria by ensuring that all items meet high editorial and technical standards and that a given student is not presented with two or more virtually identical problems.

Note: The mathematics Depth ratings are an aggregation of ratings from three criteria (C.3–C.5). The methodology recommended that C.3 receive greater emphasis but reviewers chose instead to consider equally C.4 and C.5.

Criterion C.3

Do the tests connect mathematical practices to content?

The CCSS have a focus on mathematics practices through their Standards for Mathematical Practice (SMPs). The SMPs describe a variety of types of mathematical expertise that students are expected to develop, such as reasoning, modeling, and attending to precision.

The following was required to fully meet this criterion:

- 1 *All or nearly all items that assess mathematical practices also align to one or more content standards.*

All four of the programs earned a score of Excellent Match on this criterion (Table 20). All ACT Aspire, PARCC, and Smarter Balanced items that had been coded in item metadata to SMPs also assessed at least one content standard. Since reviewers were not verifying the SMP coding provided by the programs—just checking for alignment to standards—these three programs met this criterion with 100 percent of items. MCAS did not code items to mathematical practices, but reviewers agreed that items assessed these practices nonetheless, and it also earned a score of Excellent Match.

TABLE 20

Connecting Practice to Content (Criterion C.3)



ACT Aspire

EXCELLENT MATCH

All items that are coded to mathematics practices are also coded to one or more content standard.



MCAS

EXCELLENT MATCH

Although no items are coded to mathematical practices, the practices were nonetheless assessed within items that also assessed content.



PARCC

EXCELLENT MATCH

All items that are coded to mathematics practices are also coded to one or more content standard.



Smarter Balanced

EXCELLENT MATCH

All items that are coded to mathematics practices are also coded to one or more content standard.

Criterion C.4

Are all students required to demonstrate a range of high-order, analytical thinking skills in mathematics based on the depth and complexity of the standards?

The approach to studying DOK in mathematics was analogous to ELA/Literacy. Webb's Depth of Knowledge (DOK) taxonomy was again used: level 1 is the lowest level (recall), level 2 requires use of a skill or concept, and levels 3 and 4 are considered the higher-order thinking skills. As in ELA/Literacy, the study team contracted with content experts in advance of the study to determine the distribution of the cognitive demand of the grade 5 and grade 8 standards. The experts coded 7 percent to 9 percent of CCSS content as being on DOK levels 3 or 4 in mathematics, depending on the grade (see Table 22).

To receive the highest rating on this criterion, the distribution of cognitive demand on test forms had to match the distribution of cognitive demand of the standards as a whole and match the higher cognitive demand (DOK 3+) of the standards. (As was the case in the ELA/Literacy review of cognitive demand, this is not a rating of test difficulty.) Assessments that do not match the distribution of complexity of the standards, including if they have too many high DOK items, may receive a rating of less than Excellent Match.

As shown in Table 21, the highest score for this criterion was awarded to MCAS. At both grades, the distribution of cognitive demand in the standards closely mirrored that of the CCSS, which was the goal of this criterion.

TABLE 21

Matching the Complexity of the Standards (Criterion C.4)



ACT Aspire

LIMITED/UNEVEN MATCH

At both grades 5 and 8, the test forms include significantly more items of high cognitive demand (DOK 3) than reflected in the standards, and proportionately fewer at the lowest level (DOK 1). While these items increase the challenge of the tests, standards that call for the lowest level of cognitive demand (DOK 1) may be under-assessed.



MCAS

EXCELLENT MATCH

At each grade level, the distribution of cognitive demand closely reflects that of the standards.



PARCC

GOOD MATCH

The distribution of cognitive demand of items reflects that of the standards very well at grade 5, while the grade 8 test includes proportionately more items at the higher levels of cognitive demand (DOK 2 and 3). As a result, grade 8 standards that call for the lowest level of cognitive demand may be under-assessed.



Smarter Balanced

GOOD MATCH

The distribution of cognitive demand of items reflects that of the standards very well at grade 5. At grade 8, the test includes proportionately fewer items at the lowest levels of cognitive demand (DOK 1) than in the standards, and proportionately more items at the mid-level of cognitive demand (DOK 2). As a result, grade 8 standards that call for the lowest level of cognitive demand may be under-assessed.

PARCC and Smarter Balanced both earned a score of Good Match for similar reasons—while their grade 5 assessments were fairly closely aligned to the DOK of the standards, their grade 8 assessments overemphasized DOK 2 and underemphasized DOK 1 relative to the standards (just 13 percent to 16 percent at DOK 1 versus 51 percent in the standards).⁷¹ (See Table 22 for each program’s distribution of cognitive demand, as compared to the CCSS.) ACT Aspire earned the lowest score of Limited/Uneven Match for this criterion. The ACT Aspire exam was seen as too heavily concentrated at the higher DOK levels relative to the standards (e.g., 35 percent at DOK 3 or 4 in grade 8, versus 9 percent CCSS) and thereby under-assessing the lower-level skills. Reviewers remarked that the test is “tipped heavily” toward higher DOK.

TABLE 22

The Distribution of Cognitive Demand in Mathematics: CCSS vs. Tests

		Grade 5			Grade 8		
		LEVEL 1	LEVEL 2	LEVELS 3 & 4	LEVEL 1	LEVEL 2	LEVELS 3 & 4
CCSS		43%	50%	7%	51%	40%	9%
ACT Aspire	L	23%	40%	37%	20%	45%	35%
MCAS	E	40%	58%	2%	40%	46%	14%
PARCC	G	34%	55%	11%	13%	62%	24%
Smarter Balanced	G	46%	36%	18%	16%	75%	9%

LEGEND E Excellent Match G Good Match L Limited/Uneven Match W Weak Match

Note: Percentages in the table represent percentages of score points at each DOK level. Results for a particular grade and program were created by averaging across all raters and forms for that grade and program (e.g., averaging the four raters of ACT Aspire’s form 1 and the four raters of ACT Aspire’s form 2 at grade 5).

Criterion C.5

Are a variety of item types used, including at least one that requires students to generate, rather than select, a response? Are the test items of high quality?

One aspect of item quality that the CCSSO Criteria emphasize is a diversity of item types. Specifically, the criteria recognize the widespread dissatisfaction with the ability of multiple-choice-only tests to fully measure complex student skills.

The CCSSO Criteria also address the editorial and technical quality of items. The tests evaluated in this study are used for accountability purposes, so item quality must be above reproach. Each of these programs requires multiple rounds of item review and field-testing, but issues with quality can nonetheless make their way in.

The following were required to fully meet this criterion:

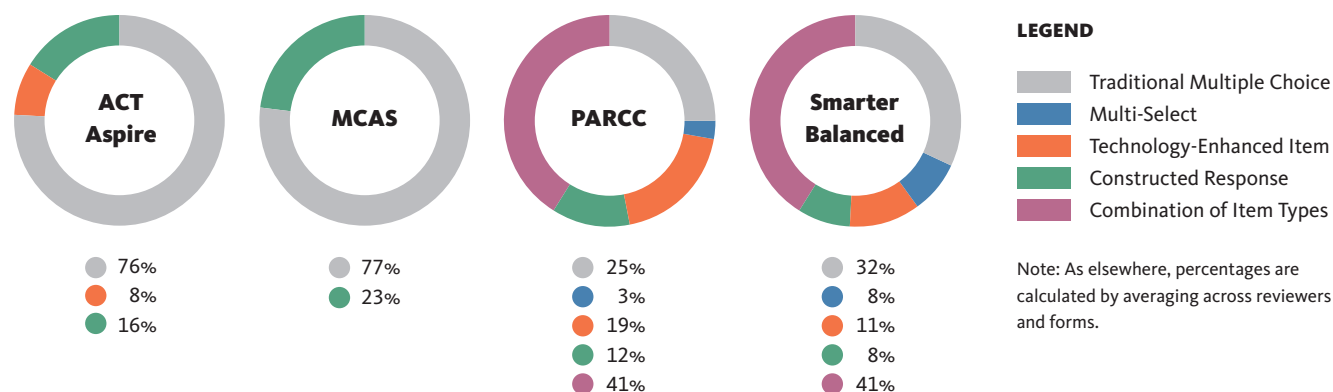
- 1 At least two item formats are used, including one that requires students to generate, rather than select, a response.
- 2 All or nearly all operational items reviewed reflect both high technical quality and high editorial accuracy.

71. PARCC uses a partial credit scoring model for multi-point items, such that students can sometimes earn partial credit for demonstrating DOK 1 or 2 skills on an item that might require DOK 3 for the maximum score. Since reviewers did not use scoring guides in their analysis of item DOK, PARCC’s emphasis on higher-level skills might be overstated.

All four programs use multiple item types, including at least one student-constructed response type, so they all met this portion of the criterion. The actual variety of item types, however, differed much more significantly, as shown in Table 23. PARCC and Smarter Balanced typically had a wide range of item types, including multiple choice, constructed response, multi-select, and technology-enhanced. Both of these assessments had less than 50 percent traditional multiple-choice items. In contrast, both ACT Aspire and MCAS had 75 percent or more of items as traditional multiple-choice items.

TABLE 23

Distribution of Item Types in Mathematics Tests



As shown in Table 24, the results for item quality are more varied. MCAS and ACT Aspire items were rated very highly in terms of quality, with only minor editorial issues. (Though reviewers did note that some ACT Aspire grade 8 items were susceptible to simplification via calculator use.) Thus, both of these assessments earned a score of Excellent Match on C.5.

Reviewers noted more issues related to item quality on the PARCC and Smarter Balanced forms than on the ACT Aspire or MCAS forms. Although this concern applies to a small percentage of items, the review panels expressed the need that a very high bar be set on the quality of items used on consequential tests. Indeed, most of the problems noted on PARCC and Smarter Balanced forms were editorial; some concerned the layout of the item on the screen,⁷² but others were mathematical. Reviewers of both PARCC and Smarter Balanced forms noted that the technology enhancements did not always improve item quality.

Reviewers found quality issues more frequently on the PARCC forms than on the Smarter Balanced forms, albeit those instances were rare (one reviewer noted “very few ‘broken’ problems or major issues” on PARCC items). Still, the panel believed that these issues were significant enough to potentially impact the accuracy of the score. Thus, PARCC earned a score of Good Match on C.5.

Although less frequent, some quality concerns on Smarter Balanced items were viewed as more serious by the panel, at times interfering with the assessed construct. Some reviewers noted mathematical errors or imprecision in certain items (approximately one item per form).⁷³ The most common of these issues pertained to excessive precision in numbers used in application problems (e.g., more digits after the decimal point than would be appropriate given the problem context) and a lack of precise language in stem wording that made multiple

72. For example, there was one instance in which a mathematical formula appeared “broken” across lines on the screen.

73. See footnote 68 for more on Smarter Balanced’s test interface.

answers plausible on items where there was only one right answer allowed. Several reviewers also noted that some Smarter Balanced forms had multiple items measuring the same skill that were nearly identical, with only the numbers in the item stem and the set of answer choices changed.⁷⁴ In light of these issues, Smarter Balanced earned a score of Limited/Uneven Match for C.5.

TABLE 24

High-Quality Items and Variety of Item Types (Criterion C.5)



ACT Aspire

EXCELLENT MATCH

The program uses multiple item types, including constructed response on the Extended Task items. These items, although they carry high point values, are limited in number; the rest of the items are predominantly multiple-choice.

The large majority of items are of high technical and editorial quality, with only very minor issues of editing, language, or accuracy. At the grade 8 level, some items appear to be susceptible to simplification by use of calculators, which are allowed on all items at grade 8, in contrast to the other programs that allow them on a restricted set of items.



MCAS

EXCELLENT MATCH

Both grade 5 and grade 8 forms include multiple item types, including constructed-response. The items are of high technical and editorial quality, with very minor issues of editing, language, and accuracy at grade 8.



PARCC

GOOD MATCH

The program includes a wide variety of item types, including several that require student-constructed responses. However, there are a number of items with quality issues, mostly minor editorial but sometimes mathematical.



Smarter Balanced

LIMITED/UNEVEN MATCH

The program includes a wide variety of item types, many of which make effective use of technology.

The program could be improved by ensuring that virtually identical items are not presented to individual students. Further, a good deal of variability across forms and grades is observed, with some forms fully meeting the item quality criterion and others only partially meeting it. Issues exist with the editorial quality and mathematical accuracy of individual items, most of which are minor but some of which could impact assessment of the targeted skill, resulting in a rating of Limited/Uneven.

74. Smarter Balanced staff indicate that the use of near-duplicate items is intentional in the assessment of mathematical fluency in Smarter Balanced assessments.

Criterion Level Ratings

Tables 25A and 25B show the final tally of the ELA/Literacy and Math criteria ratings. Immediately striking in ELA/Literacy is that the two consortia assessments earned twice as many ratings of Good and Excellent Match as the other two programs, earning eight high ratings to the four for ACT Aspire and MCAS. PARCC earned the most Excellent Match ratings (six), while Smarter Balanced was the only assessment with no ratings of Weak Match (partly because it was also the only program to test listening on the summative assessment).

The ratings for mathematics (Table 25B) were more similar between programs, with PARCC earning four Excellent or Good Match ratings, Smarter Balanced and MCAS three, and ACT Aspire two. MCAS scored particularly well on the three Depth criteria in mathematics, while PARCC is the only assessment that earned all Good Match or better scores.

TABLE 25A

ELA/Literacy Ratings Tally by Program

ACT Aspire	E	G	G	G	L	L	L	W	W
MCAS	E	G	G	G	L	L	W	W	W
PARCC	E	E	E	E	E	E	G	G	W
Smarter Balanced	E	E	E	E	G	G	G	G	L

TABLE 25B

Mathematics Ratings Tally by Program⁷⁵

ACT Aspire	E	E	L	W
MCAS	E	E	E	L
PARCC	E	G	G	G
Smarter Balanced	E	G	G	L

LEGEND E Excellent Match G Good Match L Limited/Uneven Match W Weak Match













































75. Although all four programs require the assessment of conceptual understanding, procedural skill/fluency, and applications (criterion C.2), final ratings could not be determined with confidence due to variations in how reviewers understood and implemented this criterion.

Summary Ratings

Tables 26A and 26B summarize the results across all tests, criteria, and subject areas. Appendix F includes these ratings plus the panel's summary statements, which are shown in separate tables throughout the *Results* section.

TABLE 26 A

ELA/Literacy Ratings Summary

Criteria	ACT Aspire	MCAS	PARCC	Smarter Balanced
I. CONTENT: Assesses the content most needed for College and Career Readiness				
<u>B.3 Reading:</u> * Tests require students to read closely and use specific evidence from texts to obtain and defend correct responses.				
<u>B.5 Writing:</u> * Tasks require students to engage in close reading and analysis of texts. Across each grade band, tests include a balance of expository, persuasive/argument, and narrative writing.				
B.6 Vocabulary and language skills: Tests place sufficient emphasis on academic vocabulary and language conventions as used in real-world activities.				
B.7 Research and inquiry: Assessments require students to demonstrate the ability to find, process, synthesize, and organize information from multiple sources.				
B.8 Speaking and listening: Over time, and as assessment advances allow, the assessments measure speaking and listening communication skills.**				
II. DEPTH: Assesses the depth that reflects the demands of College and Career Readiness				
<u>B.1 Text quality and types:</u> * Tests include an aligned balance of high-quality literary and informational texts.				
<u>B.2 Complexity of texts:</u> * Test passages are at appropriate levels of text complexity, increasing through the grades, and multiple forms of authentic, high-quality texts are used.***				
B.4 Cognitive demand: The distribution of cognitive demand for each grade level is sufficient to assess the depth and complexity of the standards.				
B.9 High-quality items and variety of item types: Items are of high technical and editorial quality and test forms include at least two item types with at least one that requires students to generate a response.				

* The criteria recommended to be more heavily emphasized are underlined.

** The methodology indicates that criterion B.8 (speaking and listening) should be included "over time, and as assessment advances allow." Thus B.8 ratings are not included in the overall rating for Content.

*** The criterion B.2 rating is based solely on program documentation, as reviewers were not able to rate the extent to which quantitative measures are used to place each text in a grade band. Thus, reviewers did not consider the criterion B.2 rating as heavily when deciding the overall depth rating.






























LEGEND  Excellent Match  Good Match  Limited/Uneven Match  Weak Match
 Cells for which the ratings are not used in determining Content and Depth ratings
(See Section I, *Weighting of Criteria for Content and Depth Ratings*.)





TABLE 26 B

Mathematics Ratings Summary

Criteria	ACT Aspire	MCAS	PARCC	Smarter Balanced
I. CONTENT: Assesses the content most needed for College and Career Readiness				
C.1 Focus: * Tests focus strongly on the content most needed in each grade or course for success in later mathematics (i.e., major work).				
C.2: Concepts, procedures, and applications: Assessments place balanced emphasis on the measurement of conceptual understanding, fluency and procedural skill, and the application of mathematics.**	—	—	—	—
II. DEPTH: Assesses the depth that reflects the demands of College and Career Readiness				
C.3 Connecting practice to content: * Test questions meaningfully connect mathematical practices and processes with mathematical content.				
C.4 Cognitive demand: The distribution of cognitive demand for each grade level is sufficient to assess the depth and complexity of the standards.				
C.5 High-quality items and variety of item types: Items are of high technical and editorial quality and test forms include at least two item types, at least one that requires students to generate a response.				

* The criteria recommended to be more heavily emphasized are underlined.

** Both programs require, in their program documentation, the assessment of conceptual understanding, procedural skill/fluency, and application, although most do not clearly distinguish between procedural skill/fluency and conceptual understanding. Also, specific balance across these three types is not required. Due to variation across reviewers in how this criterion was understood and implemented, final ratings could not be determined with confidence. Therefore, for criterion C.2, only qualitative observations are provided for grades 5 and 8. (See Section I, *Findings* for more information.)

LEGEND  Excellent Match  Good Match  Limited/Uneven Match  Weak Match
— Cells for which no quantitative rating could be determined

Program Strengths and Areas for Improvement

After completing the ratings, each subject area panel reviewed their findings and developed summary statements regarding the overall strengths and areas for improvement for each program. These comments pertain to how the programs performed against the CCSSO Criteria as operationalized. Feedback from the panels concerning future improvements to the methodology can be found in Section III.

ACT Aspire

English Language Arts:

In ELA/Literacy, ACT Aspire receives a Limited/Uneven to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The combined set of ELA/Literacy tests (reading, writing, and English) requires close reading and adequately evaluates language skills. More emphasis on assessment of writing to sources, vocabulary, and research and inquiry, as well as increasing the cognitive demands of test items, will move the assessment closer to fully meeting the criteria. Over time, the program would also benefit by developing the capacity to assess speaking and listening skills.

Content: ACT Aspire receives a Limited/Uneven Match to the CCSSO Criteria for Content in ELA/Literacy. The assessment program includes an emphasis on close reading and language skills. However, the reading items fall short on requiring students to cite specific textual information in support of a conclusion, generalization, or inference and in requiring analysis of what has been read. In order to meet the criteria, assessing writing to sources, vocabulary, and research and inquiry need to be strengthened.

Depth: ACT Aspire receives a rating of Good Match for Depth in ELA/Literacy. The program's assessments are built on high-quality test items and texts that are suitably complex. To fully meet the CCSSO Criteria, more cognitively demanding test items are needed at both grade levels, as is additional literary narrative text, as opposed to literary informational texts.⁷⁶

Mathematics:

In mathematics, ACT Aspire receives a Limited/Uneven to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. Some of the mismatch with the criteria is likely due to intentional program design, which requires that items be included from previous and later grade(s).

The items are generally high quality and test forms at grades 5 and 8 have a range of cognitive demand, but in each case the distribution contains significantly greater emphasis at DOK 3 than reflected in the standards. Thus, students who score well on the assessments will have demonstrated strong understanding of the standard's more complex skills. However, the grade 8 test may not fully assess standards at the lowest level of cognitive demand.⁷⁷ The tests would better meet the CCSSO Criteria with an increase in the number of items focused on the major work of the grade and the addition of more items at grade 8 that assess standards at DOK 1.

Content: ACT Aspire receives a Limited/Uneven Match to the CCSSO Criteria for Content in Mathematics. The program does not focus exclusively on the major work of the grade, but rather, by design, assesses material from previous and later grade(s). This results in a weaker match to the criteria. The tests could better meet the criteria at both grades 5 and 8 by increasing the number of items that assess the major work.

76. As discussed previously, ACT Aspire does not classify literary nonfiction texts that are primarily narrative in structure as "informational." See Appendix G for more information about ACT Aspire's interpretation of CCSSO criterion B.1.

77. As noted previously, reviewers did not account for PARCC's partial credit scoring model, which may lower the average DOK of the test.

Depth: ACT Aspire receives a Good Match to the CCSSO Criteria for Depth in Mathematics. The items are well crafted and clear, with only rare instances of minor editorial issues. The ACT Aspire tests include proportionately more items at high levels of cognitive demand (DOK 3) than the standards reflect, and proportionately fewer at the lowest level (DOK 1). This finding is both a strength, in terms of promoting strong skills, and a weakness in terms of ensuring adequate assessment of the full range of cognitive demand within the standards. While technically meeting the criterion for use of multiple item types, the range is nonetheless limited, with the large majority comprising multiple-choice items. The program would better meet the criteria for Depth by including a wider variety of item types and relying less on traditional multiple-choice items.

MCAS

English Language Arts:

In ELA/Literacy, MCAS receives a Limited to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The test requires students to closely read high-quality texts and a variety of high-quality item types. However, MCAS does not adequately assess several critical skills, including reading informational texts, writing to sources, language skills, and research and inquiry; further, too few items assess higher-order skills. Addressing these limitations would enhance the ability of the test to signal whether students are demonstrating the skills called for in the standards. Over time, the program would also benefit by developing the capacity to assess speaking and listening skills.

Content: MCAS receives a Limited/Uneven Match to the CCSSO Criteria for Content in ELA/Literacy. The assessment requires students to read closely well-chosen texts and presents test questions of high technical quality. However, the program would be strengthened by assessing writing annually, assessing the three types of writing (narrative, expository, and persuasive/argumentation) called for across each grade band, requiring writing to sources, and placing greater emphasis on assessing research and language skills.

Depth: MCAS receives a rating of Good Match for Depth in ELA/Literacy. The assessments do an excellent job in presenting a range of complex reading texts. To fully meet the demands of the CCSSO Criteria, however, the test needs more items at higher levels of cognitive demand, a greater variety of items to test writing to sources and research, and more informational texts, particularly those of an expository nature.

Mathematics:

In mathematics, MCAS receives a Limited/Uneven Match to the CCSSO Criteria for Content and an Excellent Match for Depth relative to assessing whether students are on track to meet college and career readiness standards. The MCAS mathematics test items are of high technical and editorial quality. Additionally, the content is distributed well across the breadth of the grade level standards, and test forms closely reflect the range of cognitive demand of the standards. Yet the grade 5 tests have an insufficient degree of focus on the major work of the grade.

While mathematical practices are required to solve items, MCAS does not specify the assessed practices(s) within each item or their connections to content standards. The tests would better meet the criteria through increased focus on the major work at grade 5 and identification of the mathematical practices that are assessed—and their connections to content.

Content: MCAS receives a Limited/Uneven Match to the CCSSO Criteria for Content in Mathematics. While the grade 8 assessment focuses strongly on the major work of the grade, the grade 5 assessment does not, as it samples more broadly from the full range of standards for the grade. The tests could better meet the criteria through increased focus on the major work of the grade on the grade 5 test.

Depth: MCAS receives an Excellent Match to the CCSSO Criteria for Depth in Mathematics. The assessment uses high-quality items and a variety of item types. The range of cognitive demand reflects that of the standards of the grade. While the program does not code test items to math practices, mathematical practices are nonetheless incorporated within items. The program might consider coding items to the mathematical practices and making explicit the connections between specific practices and specific content standards.

PARCC

English Language Arts:

In ELA/Literacy, PARCC receives an Excellent Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The tests include suitably complex texts, require a range of cognitive demand, and demonstrate variety in item types. The assessments require close reading, assess writing to sources, research, and inquiry, and emphasize vocabulary and language skills. The program would benefit from the use of more research tasks requiring students to use multiple sources and, over time, developing the capacity to assess speaking and listening skills.

Content: PARCC receives an Excellent Match to the CCSSO Criteria for Content in ELA/Literacy. The program demonstrates excellence in the assessment of close reading, vocabulary, writing to sources, and language, providing a high-quality measure of ELA/Literacy content, as reflected in college and career readiness standards. The tests could be strengthened by the addition of research tasks that require students to use two or more sources and, over time, a listening and speaking component.

Depth: PARCC receives a rating of Excellent Match for Depth in ELA/Literacy. The PARCC assessments meet or exceed the depth and complexity required by the criteria through a variety of item types that are generally of high quality. A better balance between literary and informational texts would strengthen the assessments in addressing the criteria.

Mathematics:

In mathematics, PARCC receives a Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The assessment is reasonably well aligned to the major work of each grade. At grade 5, the test includes a distribution of cognitive demand that is similar to that of the standards. At grade 8, the test has greater percentages of higher-demand items (DOK 3 and 4) than reflected by the standards, such that a student who scores well on the grade 8 PARCC assessment will have demonstrated strong understanding of the standard's more complex skills. However, the grade 8 test may not fully assess standards at the lowest level (DOK 1) of cognitive demand.

The test would better meet the CCSSO Criteria through additional focus on the major work of the grade at grade 5, the addition of more items at grade 8 that assess standards at DOK 1, and increased attention to accuracy of the items—primarily editorial, but in some instances mathematical.

Content: PARCC receives a Good Match to the CCSSO Criteria for Content in Mathematics. The test could better meet the criteria by increasing the focus on the major work at grade 5.

Depth: PARCC receives a Good Match to the CCSSO Criteria for Depth in Mathematics. The tests include items with a range of cognitive demand, but at grade 8, that distribution contains a higher percentage of items at the higher levels (DOK 2 and 3) and significantly fewer items at the lowest level (DOK 1). This finding is both a strength, in terms of promoting strong skills, and a weakness, in terms of ensuring adequate assessment of the full range of cognitive demand within the standards. The tests include a variety of item types that are largely of high quality. However, a range of problems (from minor to severe) surfaced relative to editorial accuracy and, to a lesser degree, technical quality. The program could better meet the Depth criteria by ensuring that all items meet high editorial and technical standards and by ensuring that the distribution of cognitive demand on the assessments provides sufficient information across the range.

Smarter Balanced

English Language Arts:

In ELA/Literacy, Smarter Balanced receives a Good to Excellent Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The tests assess the most important ELA/Literacy skills of the CCSS, using technology in ways that both mirror real-world uses and provide quality measurement of targeted skills. The program is most successful in its assessment of writing and research and inquiry. It also assesses listening with high quality items that require active listening, which is unique among the four programs. The program would benefit by improving its vocabulary items, increasing the cognitive demand in grade 5 items, and, over time, developing the capacity to assess speaking skills.

Content: Smarter Balanced receives an Excellent Match to the CCSSO Criteria for Content in ELA/Literacy. The program demonstrates excellence in the areas of close reading, writing to sources, research, and language. The listening component represents an important step toward adequately measuring speaking and listening skills—a goal specifically reflected in the standards. Overall, Smarter Balanced is a high-quality measure of the content required in ELA/Literacy, as reflected in college and career readiness standards. A greater emphasis on Tier 2 vocabulary would further strengthen these assessments relative to the criteria.

Depth: Smarter Balanced receives a rating of Good Match for Depth in ELA/Literacy. The assessments use a variety of item types to assess student reading and writing to source. The program could better meet the depth criteria by increasing the cognitive demands of the grade 5 assessment and ensuring that all items meet high editorial and technical quality standards.

Mathematics:

In mathematics, Smarter Balanced is a Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The test provides adequate focus on the major work of the grade, although it could be strengthened at grade 5.

The tests would better meet the CCSSO Criteria through increased focus on the major work at grade 5 and an increase in the number of items on the grade 8 tests that assess standards at the lowest level of cognitive demand. In addition, removal of serious mathematical and editorial flaws, found in approximately one item per form, should be a priority.

Content: Smarter Balanced receives a Good Match to the CCSSO Criteria for Content in Mathematics. The tests could better meet the criteria by increasing the focus on the major work in grade 5.

Depth: Smarter Balanced receives a Good Match to the CCSSO Criteria for Depth in Mathematics. The exam includes a range of cognitive demand that fairly represents the standards at each grade level. The tests have a strong variety of item types, including those that make effective use of technology. However, a range of problems (from minor to severe) surfaced relative to editorial accuracy and, to a lesser degree, technical quality. A wide variety of item types appear on each form, and important skills are assessed with multiple items, as is sound practice. The program could better meet the Depth criteria by ensuring that all items meet high editorial and technical standards.

Section II: Recommendations

The size and complexity of this report make it easy to miss its most important finding: these new, computer-based tests are major improvements over the previous generation of state tests. There is no better evidence than the fact that a best-in-class state assessment, the 2014 Massachusetts Comprehensive Assessment System, does not measure many of the important competencies that are part of today's college and career readiness standards. To be fair, MCAS is quite strong in the quality of its items and the degree to which its math test matches the depth of knowledge in the standards, but significant enhancements would be needed for it to meet the CCSSO Criteria for a high-quality assessment. In contrast, both PARCC and Smarter Balanced, while not perfect, measured a much wider array of skills and content called for in the Common Core State Standards.

But those two new assessments remain far from ubiquitous. Three years ago, forty-five states were members of either PARCC or Smarter Balanced for ELA/Literacy and mathematics. As of January 2016, nearly thirty states are planning to use either a customized state assessment or a vendor-developed option, such as ACT Aspire.⁷⁸ Based on the results of this evaluation, and with an eye to the future, we offer the following recommendations for state policymakers and test developers.

For State Policymakers

1 **Make quality non-negotiable.**

State assessments serve as a yardstick for gauging the quality and progress of public schools and the students in them. Whether used for high-stakes accountability purposes (such as teacher evaluation) or lower stakes (such as public reporting of school results), weak assessments leave state leaders (as well as educators and parents) with blinders on, misleading them or showing them only a portion of the truth. In the worst cases, weak assessments paint a rosier-than-justified picture of student performance, despite serious skill and knowledge deficits. This sets students up for unsuccessful transitions into college or the workplace, leaves states without the warning signals needed to protect the integrity of their future workforce, and deprives parents of an accurate picture of their children's educational progress. Quality assessments are also important for formative purposes for educators in that they reinforce the expectations of the standards, send teachers clear messages about what they should be teaching, and provide useful feedback on student progress.

The recently revived peer-review process for assessments at the U.S. Department of Education is explicitly focused on assessment quality along multiple dimensions; we view this as a promising development. That said, peer-review guidance was in place throughout the No Child Left Behind (NCLB) era, yet states still adopted assessments of questionable quality and alignment. We hope the peer-review process is taken seriously and that the Department helps to ensure that, whatever tests states adopt, that they are high quality and support effective standards implementation. The approach we used in this report provides relevant evidence about assessment quality that could feed into a peer-review evaluation.

⁷⁸ National Council of State Legislators, "Hot Topics in Higher Education: Reforming Remedial Education," <http://www.ncsl.org/research/education/improving-college-completion-reforming-remedial.aspx>.

2 When developing or revising assessments, carefully prioritize the set of skills and knowledge at each grade that should serve as the focus of instruction, building public understanding and support as you do so.

State tests typically do not assess every grade level standard, just as none of the four tests in this evaluation assess every Common Core State Standard. The selection of the standards to be tested, then, is incredibly important, as it sends powerful signals throughout the educational system about the priorities for instructional time.

States that have joined a consortium have made these selections, delineated in their testing blueprints, and will likely be revised from time to time by the consortia. Yet the thirty or so states that are embarking solo on assessment development should pause to clarify the essential subset of skills at each grade to be assessed. The CCSSO's "Criteria for Procuring and Evaluating High Quality Assessments," the foundation of this evaluation, provides helpful guidance. This document lays out what educational leaders from a number of states, informed by research from higher education and employers, have determined to be the essential building blocks of college and career readiness. Most, if not all, of this "short list" of skills resonates with policymakers and the public as common-sense expectations. For example, the ability to read several texts about a topic and craft a well-written argument supported by evidence is an important skill for all high school graduates today. The state's employers and higher education institutions may want to augment CCSSO's short list with a few skills that are particularly important to them.

We also recommend active public involvement during this process because the assessment of college and career readiness is neither quick nor inexpensive—nor lacking controversy. Some of the skills most highly valued by colleges, such as the writing of clear, well-supported arguments or the development, testing, and refinement of mathematical models, cannot be assessed with short, low-cost, multiple-choice questions. Instead, they require a willingness to invest in extended-response items with human scoring or in technology-enhanced items. It is vital to build public support for both the cost of developing such an assessment and the time needed for students to complete it. Pressure to roll back testing decisions, such as recent opt-out initiatives, are due in part to the failure to build public understanding and buy-in for high-quality tests that measure high-priority skills and knowledge.

3 Ensure quality is maintained while addressing concerns about testing time and costs.

Even though high-quality tests tend to require more testing time and come with larger price tags, there are constructive steps that states can take to minimize both testing time and cost. A study by the Council of Great City Schools found great variability in the number of hours per year that students spend taking mandatory tests, ranging from ten hours to nearly thirty.⁷⁹ Many of the tests turned out to be redundant or not aligned to the state standards and assessments. The 2015–16 versions of the tests evaluated in this study will require no more than three hours per content area, which we find reasonable given the set of complex skills they are assessing. So a logical first step in minimizing testing time is to eliminate redundant and misaligned tests at the state and local levels. CCSSO has released a framework that states can use to identify unnecessary tests and to "ensure [that] every student is taking assessments that are necessary, of high quality and provide meaningful information."⁸⁰

States can also ease testing costs while maintaining high-quality assessments. Three of the tests in this study serve as examples: states can share costs by participating in a consortium of states (PARCC or Smarter Balanced) or can purchase a test that has already been developed by a test provider (ACT Aspire). Both options, however, require that the state yield some control over a) the specification of the set of skills and

79. Council of Great City Schools, press release, "Student Assessments in Public Schools Not Strategic, Often Redundant," October 24, 2015, <http://www.cgcs.org/cms/lib/DC00001581/Centricity/Domain/4/Testing%20Report.pdf>.

80. CCSSO, "Comprehensive Statewide Assessment Systems: A Framework for the Role of the State Education Agency in Improving Quality and Reducing Burden," June 2015, http://www.ccsso.org/Resources/Publications/Comprehensive_Statewide_Assessment_Systems_A_Framework_for_the_Role_of_the_State_Education_Agency_in_Improving_Quality_and_Reducing_Burden.html.

knowledge assessed and b) the setting of cut scores. A third option is to secure high-quality test items from a credible source—such as one of the consortia or a bank of secure items shared across multiple states—while maintaining full state control over the specific items used and how the state will define proficiency.⁸¹

Policymakers may also encounter critics who demand the elimination of all state tests, some of whom assert that the goal of improving education could be better served by using those dollars at the classroom level. However, this is almost certainly false, considering how small testing budgets are compared to the overall costs of education. Matthew Chingos of the Urban Institute determined that the elimination of all federally and state-mandated testing would save an average of just \$34 per student per year, an amount that would reduce pupil-teacher ratios by just 0.1 student.⁸² Taking steps to minimize testing costs while maintaining quality is sound governance, but eliminating the yardstick altogether would be foolhardy.

4 Work with other state leaders to press the assessment industry and researchers for improvements in test item types and scoring engines to better measure key constructs in a cost-effective way.

Great advances have been made in testing over the past five years, as the majority of states have transitioned to computer-based delivery and the two multi-state consortia have used their unprecedented buying power to fund rapid research and development. Nevertheless, further advances are needed to assess accurately and efficiently all of the skills that states value and to improve the assessment of students with disabilities. The fact that forty states use the CCSS or a variant of them makes it possible for these states to use their collective influence to press for further improvements.⁸³ For example, the mathematical practices call for high school students to adeptly use various common technologies, such as spreadsheets and statistical software, and to “make sound decisions about when each of these tools might be helpful, recognizing both the insights to be gained and their limitations.”⁸⁴ But advances in both technology-enhanced items and schools’ bandwidth are needed to incorporate such technological tools and assess student’s ability to use them strategically. Another example: the ELA/Literacy standards call for students to express orally well-supported ideas and probe the ideas of others, yet no current state test includes the assessment of speaking skills. By working through the CCSSO or other national organizations, all states with college and career readiness standards can prioritize areas for improvement and collectively strategize about how best to meet them.

For Test Developers

1 Ensure that every item meets the highest standards for editorial accuracy and technical quality.

Simple errors, the review panels noted, can reduce test validity and sometimes impact the accuracy of scores. When test scores are used for so many consequential purposes, these tests must maintain a very high bar for item quality. It’s true that testing programs already have extensive review processes in place, but our reviewers still noted editorial and substantive issues that undermined item integrity. Even inconsequential errors can also spread rapidly through social media and news stories, exacerbating public understanding of the quality of state assessment items.

81. The Massachusetts State Board of Education voted on November 17, 2015 to award a new MCAS contract to include a next-generation assessment for English language arts and math using both PARCC items and items specific to Massachusetts. See Massachusetts Department of Elementary and Secondary Education, press release, “Board of Elementary and Secondary Education Approves Path to Next-Generation MCAS,” November 17, 2015, <http://www.doe.mass.edu/news/news.aspx?id=21314>.

82. M. Chingos, “Testing Costs a Drop in the Bucket” (Washington, D.C.: Brookings Institution, February 2, 2015), <http://www.brookings.edu/blogs/up-front/posts/2015/02/02-standardized-tests-chingos>.

83. National Conference of State Legislatures, “Common Core Status Map,” updated November 23, 2015, <http://www.ccrslegislation.info/CCR-State-Policy-Resources/common-core-status-map>.

84. Common Core State Standards, “Standards for Mathematical Practice,” www.corestandards.org/Math/Practice.

2 Use technology-enhanced items (TEIs) strategically to improve test quality and enhance student effort.

Used well, TEIs can expand the depth and breadth of what can be measured and help maintain student engagement and effort during the assessment, which are important to obtaining meaningful results. Panelists were concerned, however, that in some cases TEIs (i.e., drag-and-drop items, equation editors, and line plotters) were used seemingly to no advantage—though information about their impact on student engagement was not part of the review. Because TEIs tend to be more expensive to develop than simpler item types, using them only when they enhance the quality of measurement will help to minimize costs.

3 Focus research and development on areas of targeted importance relative to measuring student performance on CCR standards.

Most of the skills and abilities the CCSS and other college and career readiness standards call for can currently be measured well. A few, however, require advances in the field of measurement and/or test development and scoring costs, if they are to be used at scale. These include:

- ◆ *Research skills:* The ability to formulate an inquiry into a topic, locate multiple sources of information, evaluate their quality, and cull evidence to develop and support a claim is a skill needed by all adults—whether it's put to use in researching jobs, choosing an insurance plan, or deciding how to vote. Current tests do not evaluate this integrated set of skills, but focus on some of the skill components such as the ability to organize information or to use evidence from sources to develop and support a claim. Test developers should have as their goal to measure robustly this integrated set of skills within reasonable cost and test-time constraints.
- ◆ *Use of technological math tools, such as spreadsheets and software designed for use in statistics, geometry, and modeling:* These tools are widely used in colleges and the workplace, which means that high school graduates need to be able to use them in solving complex problems. Incorporating such tools within state assessments is a new expectation, but one that warrants focused effort.
- ◆ *Speaking and listening skills:* Of the four assessments evaluated in this study, only Smarter Balanced appraised listening skills, and none gauged mastery of speaking skills. How to reliably measure these skills is a challenge for those in the assessment field, but advances here would obviously enhance the value of future tests.

We close with a note of optimism. For too many years, state assessments have generally focused on low-level skills and have given students and parents false signals about the readiness of their children for postsecondary education and the workforce. (They weren't very helpful to educators and policymakers, either.) Many students were forced to waste time and money on remedial coursework. Thus, the choices concerning state assessments are some of the most important decisions that state education leaders face. Students deserve carefully crafted tests that will measure the skills and knowledge they need to transition successfully into postsecondary education or the workplace. Educators need and deserve good tests that honor their hard work and give useful feedback, which enables them to improve their craft and boost their students' success. And policymakers, parents, and the public need tests that will tell them whether their students are developing the skills and knowledge that today's high school graduates need.

States' adoption of college and career readiness standards has been a bold step in the right direction. The fact that multiple states are already using ACT Aspire, PARCC, and Smarter Balanced, each of which is as good or better than the previous best-in-class Massachusetts state assessment, fuels our optimism. Yet, adopting and sticking with high-quality assessments requires courage. These tests are tougher, sometimes cost more, and can require more testing time than the previous generation of state tests.

Are we up to the task?

Section III: Suggestions for Methodological Improvement

The methodology used in this study is highly comprehensive and includes evaluating individual test items, passages, prompts, and documentation along many dimensions. Given the breadth and scope of what reviewers were asked to do, it's notable that they deemed the final results an accurate analysis of each assessment against the CCSSO Criteria.

For instance, one ELA/Literacy reviewer noted:

“This methodology’s greatest strength is its comprehensiveness of approach. It takes all tests’ supporting documents into consideration, relies on quality training of evaluators, and sound research approaches to arrive at useful and worthy metrics on how state/national assessments are doing.”

A mathematics reviewer also remarked:

“The development from scratch of a procedure that obtains such an in-depth review of a variety of tests is impressive. I am sure there are refinements to the process that can give a better evaluation of these tests, but the one we implemented gave, in my opinion, a fair and in depth evaluation of the tests.”

As compared to prior alignment methodologies (e.g., Webb’s alignment tool and the Surveys of Enacted Curriculum), the task for reviewers implementing this methodology is great. While the methodology offers some advances over existing methodologies—most notably its close reflection of the principles of new college and career readiness standards—there are ways in which it could be improved to provide a more complete picture of the quality of new assessments.

As the first implementers of these new methods, it is not surprising that we identified a number of improvements that could be addressed in subsequent iterations. In this chapter, we lay out the challenges, and, where appropriate, suggestions for revising the methodology to address these limitations—many of which would make it even more comprehensive (as well as time- and labor-intensive). The section is organized by content area and criterion.

English Language Arts/Literacy

Methodological Concerns

- B.1** This criterion includes both text quality and text type/balance. Several reviewers noted that these are different dimensions, and would have preferred that they not be combined. For future implementations of the methodology, we recommend that text quality and text type/balance be reported separately.
- B.1.1** Balance of text types was determined simply by counting the number of texts of each type that students encountered on a test form. While this is a reasonable approach, reviewers noted that some passages were longer or included many more items than others, and that these might also be appropriate ways to weight text passages. We recommend considering these alternative ways of weighting text passages, either in the study design or in the roll-up from items to forms, to determine whether the results are sensitive to the choice of weighting scheme.

- B.1.2** The purpose of this criterion is to evaluate the quality of texts—whether they are previously published or of publishable quality. In practice, almost all texts used on these assessments were previously published, so tests generally received full credit for this criterion. However, reviewers felt that text quality varied among texts, even though they were all previously published, but the methodology did not offer them the opportunity to render a judgment on text quality.

We recommend that future iterations of the methodology consider how to provide a measure of text quality that includes reviewer judgments.

- B.1.3** The purpose of this criterion is to evaluate the structure of informational texts, with the goal being that “nearly all informational texts are expository in structure” (as opposed to narrative in structure). To operationalize this goal, the suggested cutoffs are 90 percent or above of informational text being expository for a score of 2, 75 percent to 89 percent for a score of 1, and 74 percent or below for a score of 0. Unfortunately, these suggested cutoffs don’t work well with tests that have a small number of informational texts. For example, if an assessment has just two informational texts, it can only earn a score of 2 (fully met) or 0 (not met) on this criterion.

This is one of several of the criteria that are challenging to implement for a single test form. We recommend that certain criteria might be more appropriately evaluated across multiple forms, and this may be one. For instance, it is not necessary to go through an entire test form item by item to evaluate text passages on their structure; it may be more appropriate to evaluate this criterion across many forms. (Of course, this approach is not possible for tests, such as MCAS, that have just one or two forms.) Evaluating passages on multiple forms will result in a more precise determination of text passage quality and structure. Alternatively, the suggested cutoffs in the methodology might be made more flexible to account for the fact that there are typically small numbers of passages (and, as indicated, can sometimes swing the score from a 2 to a 0).

- B.3.1 and B.3.4** These two criteria focus on “close reading” and “direct textual evidence.” In general, reviewers understood the latter to be a subset of the former, but even after extensive training and discussion, there were often disagreements about how these two terms should be understood and implemented. The variability was often around whether direct textual evidence was truly required when the student only had to use the text to find an answer, as opposed to having to show the actual text evidence that supported an answer. It is not clear whether the misunderstanding was a problem with the methodology or with the training provided to reviewers in this study.

We recommend strengthening and clarifying definitions of these two terms and providing reviewers with more detailed training on them.

- B.4.1** The methodology calls for the highest score when the distribution of DOK on a form matches that of the CCSS (a DOK index of .80 or above), and for lower scores to be assigned as the degree of match declines. However, this guideline does not differentiate between forms where the DOK misalignment is due to the test having too low DOK versus too high. There was broad agreement among reviewers that tests with a greater emphasis on DOK 3 and 4 items should be rated more positively than tests with a greater emphasis on DOK 1 and 2.

There are many possible solutions to this problem. The best approach would likely be to add guidance in the scoring criteria to account for this possibility. For instance, one way to determine whether the DOK misalignment occurred because the test is “too high” versus “too low” would be to create an aggregate DOK value for each test by treating the DOK index as ordered.⁸⁵ Revised guidance could specify that a lower DOK index threshold would be acceptable if the average DOK of the test exceeded that of the standards. Another approach would be to specify a priori DOK distributions that would be acceptable—for instance, one-third each at levels 1, 2, and 3 to 4.

- B.6.3** On these criteria, assessments get top scores when they report a sub-score for language or vocabulary, even when those assessments include very few items in those sub-scales. In contrast, some assessments
and
B.6.4 that have many items on vocabulary or language but do not report a sub-score receive no or limited credit. While it is important to provide parents and teachers with direct feedback in the form of sub-scores about student knowledge of language and vocabulary, reviewers felt that assessments should not receive full credit if those sub-scales were based on an unreliably low number of items.

We recommend that the scoring guidance be revised so that assessments that have language or vocabulary sub-scales cannot receive the maximum score, unless they have an adequate number of test items in those areas.

- B.7.1** This criterion evaluates test items assessing research. Reviewers felt the definition of “research item” was too broad and allowed credit for items that did not truly mirror research activities.⁸⁶

We recommend that the definition of a research item be enhanced to ensure that the focus is on the use of two or more discrete sources (e.g., texts, audio recording, or videos) and of research skills applied in an authentic way. Also, the sufficiency of research items (e.g., the percentage of the test devoted to research items) should be addressed, because a test that has only a single research item comprising just a small proportion of the total score points could score a 2 (met the criterion) under the current methodology.

- B.9.1** The current scoring guidance for this criterion requires only a single non-multiple choice test item to earn a score of 2. However, reviewers noted that assessments varied a great deal in the extent to which they employed a variety of closed and open-ended questions. Reviewers believed that assessments with a wider variety of item types measuring more authentic skills should receive more credit.

We recommend that the scoring guidance be revised to set suggested score point cutoffs for constructed-response items. For example, the scoring guidance might indicate that the assessment must have 25 percent of score points from constructed-response items in order to receive a score of 2.

Implementation Concerns

- B.5.2** This criterion is focused on the proportion of writing prompts that require writing to a source. However, the denominator used to calculate this proportion is the number of items coded to a writing standard. Thus, depending on how vendors coded writing items, some assessments had many items that counted toward this criterion that were not actual writing items. This would have resulted in them failing the criterion if we did not manually correct the issue as we were implementing the study.

We recommend that future versions of the coding document use the correct denominator for this calculation, which is the number of writing prompts.

85. This approach was used with a somewhat different DOK scale in Polikoff & Struthers, 2013. As an example, suppose the standards were 25 percent at DOK 1 and 2 and 50 percent at DOK 3, while the assessment was 10 percent at DOK 1 and 2 and 80 percent at DOK 3. The DOK index for this comparison would be .70, well below the cutoff for a 2 (which is .80). But treating the scale as ordered the standards would have an average DOK of 2.25 ($25\% \times 1 + 25\% \times 2 + 50\% \times 3$), while the test would have an average DOK of 2.7.

86. That definition is as follows: These tasks require students to use multiple (minimum of two) informational texts about one topic (may be written, audio, or visual), analyze these multiple texts, and synthesize and/or organize information across texts. Research tasks may include multiple-choice and technology-enhanced items and must include at least one constructed-response item for which the student writes a response. See Appendix B, *Key Terminology*.

- B.9.2** This criterion is focused on item quality. According to the methodology, two measures are used to determine the rating: the proportion of items with technical or editorial issues and the proportion of items for which reviewers agree with the vendors' alignment determinations (as provided in the metadata). The reviewers felt strongly that it was inappropriate to use alignment in evaluating this criterion, especially because it had been used in other criteria (e.g., B.3.1). Further, using alignment to evaluate B.9.2 is not fair to vendors who take different approaches, such as aligning items to claims rather than to individual content standards. Finally, for an item to meet the alignment criterion for B.9.2, the reviewers had to agree with all of the alignment ratings offered by the vendor, which may differentially affect vendors that list more standards compared to those who list fewer.

We recommend the approach used in our implementation—pulling apart alignment from item quality. For this sub-criterion, we believe reviewers should only evaluate items for technical and editorial quality. If other implementers also want to evaluate the extent to which the vendors' alignment ratings are accurate, that could either be presented as a separate index, or as its own sub-criterion with scoring guidance.

Mathematics

Methodological Concerns

- C.1.1** This criterion concerns the extent to which the test focuses on the major work of the grade. While it is required that the large majority of the major work clusters be assessed,⁸⁷ the methodology does not require that each cluster be adequately addressed. Thus, a test that has just one item per cluster receives full credit, even though the test could not possibly provide a reliable score for that cluster. As a related point, the methodology does not call for a review of the degree to which the items assessing a given cluster address different standards or skills within that cluster or are redundant. As a result, the perception is that tests covering all the clusters could meet this criterion, even if they are quite unbalanced and only measure a small slice of the standards.

We recommend revising this criterion to include a measure of balance of coverage across the major work clusters. For instance, Webb's methodology includes measures of balance that could be applied at the cluster level. This would help ensure that assessments are not overly focused on a small subset of the standards in the major work.

- C.2.1** This criterion assesses the extent to which the assessment is adequately balanced across application, conceptual understanding, and procedural fluency. There were a number of difficulties with the implementation of this criterion. First, the requirement of categorizing items by their predominant focus led to a failure to recognize and give credit for items that address two or more categories. Second, and related, items that were coded as measuring “combined” skills were not counted in any way, so assessments with more “combined” items were penalized in accordance with the tentative scoring guidance. Third, the broad definition of “application” (i.e., any item that includes a context) resulted in many items that also assessed conceptual understanding and/or procedural skill/fluency to be categorized as only application because they included use of a context (even if trivial). This resulted in a lowered rating and a failure to recognize the other competencies being addressed.

To address these issues, we recommend that the methodology be revised to allow reviewers to indicate 1, 2, or 3 of the available skill categories, rather than encouraging them to select just one and penalizing the assessment if “combined” is selected. Another approach would allow raters to allocate emphasis across each category (e.g., rate an item as two-thirds procedural and one-third application). If either of these changes were made, the scoring criteria need not be changed—a goal of equal balance of application procedures, and conceptual understanding is still appropriate and could be calculated easily.

87. The major work clusters are the groups of content standards that correspond to the major work of the grade. The clusters are listed here, “Key Instructional Shifts of the Common Core State Standards for Mathematics,” http://achievethecore.org/content/upload/Focus%20in%20Math_091013_FINAL.pdf.

- C.3.1** This criterion is intended to evaluate the test's measurement of the standards for mathematical practice (SMP). However, as currently implemented, it is inadequate in this regard. The methodology simply requires that items coded to a SMP also be coded to align with *any* content standard—not even grade level content standards. The result is that all assessments either earn a perfect score for this criterion (because they code to the SMPs) or a zero (because they do not). A related concern is that the methodology requires reviewers to accept the program's designations of alignment to SMPs, and these were, at times, found to be generous to aggressive. Finally, the methodology does not require coverage of the SMPs, allowing a program that assesses only one or two SMPs across all items to be awarded a high rating.

We have several recommendations for this criterion. First, reviewers might be asked to evaluate vendors' claims of SMP alignment, rather than accepting them. This will require additional training and time to complete the review. Second, the methodology might allow reviewers to evaluate the extent to which the mathematical practices are adequately covered by the assessment. Again, this could be done using indices of coverage or balance, such as those employed in the Webb alignment methodology. Third, implementers might remove the requirement of content standard coverage from this criterion; every item on every test was reported as aligned to a content standard.

- C.4.1** The same issues and recommendations apply in mathematics as for criterion B.4.1 in ELA/Literacy.

Another issue in mathematics was that reviewers were dissatisfied that the difficulty of the item was not captured in the methodology. Some reviewers felt that difficulty was a de facto component of the mathematics DOK levels, with DOK 1 and 2 seen as focusing on the number of steps and DOK 3 and 4 focusing on the level of reasoning. We recommend considering item difficulty as a separate dimension.

A final issue in mathematics (although an issue in both subjects, it was only raised by the mathematics panel) is that some reviewers were uncomfortable with the methodology's requirement that the DOK distribution of the test match the DOK distribution of the standards, since the standards themselves do not call for such a match. As mentioned above in the ELA/Literacy section, one solution to this problem would be to specify an a priori distribution of DOK that is desirable.

- C.5.1** The same issues and recommendations apply in mathematics as for criterion B.9.1 in ELA/Literacy.

- C.5.2** Although the methodology distinguishes between technical quality and editorial accuracy, it does not distinguish between major quality issues that would impact score accuracy or meaningfulness, and trivial quality issues that would not have such impacts. For example, panelists found numerous instances in which multiple strategies could be used to solve a given test item. If the purpose of the item is to determine whether the student understands and can apply one specific strategy, having multiple could weaken the meaning of the score and, consequently test validity. While this is not a problem for constructed-response items where work is evaluated, it can be for multiple-choice items if the student's work is not provided.

We recommend that additional options be given to allow reviewers to indicate the severity of the item quality problem. Furthermore, we recommend reviewers be allowed to indicate cases in which the manner in which a student solves a problem does not invalidate an item, but does seriously obscure what competencies the item is testing. The panel also recommends that reviewers be given access to the vendor's worked solutions for all test items, so that the intention for student-solution methods can be directly evaluated.

Across Subjects

Finally, we summarize several concerns that apply to both mathematics and ELA/Literacy.

First, reviewers noted that the methodology does not allow them to express enthusiasm for items, only to indicate problems. In both subjects reviewers identified a number of items they viewed as especially high quality and they had no way to differentiate them. We recommend that the methodology allow reviewers to note and describe especially high-quality items based on their professional judgment.

Second, the quality of the discussions during Phase 2 of the methodology (when data are rolled up from test forms to programs) depended in large part on the quality of the comments left by reviewers throughout the process. Because only certain comments populated forward in the Excel review sheets designed for the study, it was not always possible for the Phase 2 reviewers to understand the rationales behind other reviewers' scoring decisions. We recommend that reviewers be required to provide substantive comments at key locations in the coding documents, so that these comments can be carried forward and inform subsequent discussions.

Third, the methodology, as written, required reviewers to examine answer keys to ensure technical accuracy. However, several test vendors had privacy policies that prohibited their releasing answer keys to reviewers, or that limited answer-key release to specific times and places (i.e., in locked rooms using paper answer keys). In practice, this made the review of answer keys impractical. For the purposes of this review we made two changes to the methodology, which we recommend in cases where the answer key is not fully available across all assessments under review:

- ◆ Remove the option “item incorrectly keyed” from the list of possible quality issues in B.9.2 and C.5.2; and
- ◆ Instruct reviewers to use the option “unintended correct answer” to identify selected-response items where there appear to be more correct answers than allowed (e.g., two correct answers exist though the item calls for just one). While these solutions limit the scope of the item quality review, we see no practical alternative if vendors have strict answer key privacy policies.

Fourth and finally, the methodology relies heavily on test vendors' metadata, especially their standards alignment data. But several programs do not take a one-to-one standards-to-item approach to test construction. While reviewers found a way to work around this problem, the quality of the item alignment determinations may have been less accurate from those vendors that used an evidence-centered design, resulting in weaker scores. Considering how to make the methodology more useful for all programs, regardless of their approach to alignment, may be beneficial.

Summary

Taken together, these suggestions provide a wealth of good counsel from our reviewers about how the methodology could be improved for future implementation. Though there are many recommendations, reviewers were nonetheless quite satisfied overall with the methodology as written and implemented and broadly agreed with the conclusions reached.

We are confident that the methodology, as implemented in this study, represents an important first step in rethinking techniques to evaluate test quality, which should inform future generations of assessments in U.S. schools. With some or all of the proposed revisions incorporated, these methods will provide an even stronger, clearer message about the kind of assessment we need to ensure valid measurement of student mastery of college and career readiness standards.

Appendix A

Depth of Knowledge (DOK) of the Four Assessment Programs as Compared to the CCSS and Other Policy-Relevant Assessments

We also conducted a supplementary analysis of the Depth of Knowledge (DOK) of the four assessment programs against the Common Core State Standards (CCSS). HumRRO had previously conducted a DOK analysis of the CCSS, so we started with those results. Next, we contracted with one subject-area expert in each subject to independently code the standards in terms of DOK (the content standards only, not the Standards for Mathematical Practice). Finally, we contracted with another content-area expert in each subject to adjudicate in any instance where our reviewers disagreed with the HumRRO reviewers. In each instance, reviewers were allowed to place each standard into one or more DOK levels, and the standards were equally weighted to arrive at the final DOK distribution for the standards, as is common in alignment studies.⁸⁸ To help contextualize the findings, we compared the DOK of our four assessments against CCSS and against two analyses of a) a set of state tests in the NCLB era regarded as having the highest likelihood of assessing deeper knowledge and b) national and international assessments previously conducted by the RAND Corporation in 2012 and 2014, respectively.⁸⁹ The results of this analysis are shown below.

In broad terms, DOK refers to the cognitive demands required by a task to produce an acceptable response. In each subject, there are four levels of DOK. Generally, level 1 refers to rote or reproductive skills (e.g., identifying an obvious detail in a text, conducting a straightforward one-step operation in mathematics). Level 2 refers to skills and concepts such as multi-step operations in mathematics and comprehension across one or more sentences. Level 3 refers to strategic thinking, such as solving a mathematics problem that has multiple possible approaches or identifying complex themes across an entire passage. Level 4 refers to extended thinking, such as extended mathematical investigations or synthesis and analysis across multiple texts.

The results in ELA/Literacy, shown in Table A.1, show that the DOK of the PARCC assessments is the highest of the four studied at both grade levels. This difference is especially large in eighth grade, where nearly 70 percent of PARCC score points are on DOK 3 or 4. In fifth grade, the other three assessments fall short of the CCSS's emphasis on DOK 3+, with Smarter Balanced the lowest as 22 percent. In eighth grade, ACT Aspire has just 19 percent of score points at DOK 3+, as compared to 46 percent in the standards and 36 percent or more for the other assessments. With few exceptions, the assessments studied here have far higher proportions of DOK 3+ score points than was found on the fourteen NCLB-era state tests studied by Yuan and Le.

The results in mathematics, shown in Table A.2, show that ACT Aspire has a much larger proportion of DOK 3+ score points than the other assessments studied in this or the two Yuan and Le studies. The difference is especially notable at grade 5. With the exception of MCAS at grade 5, the assessments studied here meet or exceed the DOK of the standards. They also typically exceed the higher DOK emphasis of the fourteen NCLB-era assessments, and they match up favorably to assessments such as AP and PISA.

88. A.C. Porter, "Measuring the Content of Instruction: Uses in Research and Practice," *Educational Researcher* 31, no. 7: 2002, 3–14.

89. Yuan and Le, 2012.

TABLE A 1

Depth of Knowledge of ELA/Literacy Assessments

Reading	DOK 1	DOK 2	DOK 3	DOK 4
CCSS grade 5	17.6%	36.8%	37.5%	8.1%
ACT Aspire grade 5	34.3%	37.5%	28.1%	0.0%
MCAS grade 5	9.6%	63.5%	26.9%	0.0%
PARCC grade 5	4.5%	45.0%	50.5%	0.0%
Smarter Balanced grade 5	19.0%	58.9%	20.8%	1.2%
CCSS grade 8	10.1%	44.2%	41.8%	3.8%
ACT Aspire grade 8	44.3%	36.8%	15.0%	3.8%
MCAS grade 8	4.8%	58.7%	33.7%	2.9%
PARCC grade 8	1.6%	29.1%	46.4%	22.9%
Smarter Balanced grade 8	15.0%	40.8%	36.7%	7.5%
14 states	31.5%	43.8%	23.3%	1.6%
AP	11.0%	33.0%	56.0%	0.0%
NAEP	20.4%	45.2%	33.3%	1.1%
PISA	37.3%	26.2%	36.5%	0.0%
TIMSS	55.2%	17.8%	26.1%	0.9%

Note: Results for AP, NAEP, PISA, PIRLS and fourteen states from Yuan and Le (2014). IB results not provided because IB does not have a reading assessment. AP, NAEP, PISA, and PIRLS results are for reading assessments, whereas MCAS, PARCC, Smarter Balanced, ACT Aspire, and CCSS results are for combined ELA/Literacy assessments. The programs evaluated in this study are based on percentages of score points, whereas other assessment results are based on percentages of items.

TABLE A 2

Depth of Knowledge of Mathematics Assessments

Math	DOK 1	DOK 2	DOK 3	DOK 4
CCSS grade 5	43.3%	49.6%	7.1%	0.0%
ACT Aspire grade 5	23.0%	40.2%	36.8%	0.0%
MCAS grade 5	40.3%	57.7%	2.0%	0.0%
PARCC grade 5	34.4%	54.5%	11.1%	0.0%
Smarter Balanced grade 5	46.6%	35.6%	15.2%	2.6%
CCSS grade 8	50.9%	39.8%	9.3%	0.0%
ACT Aspire grade 8	19.9%	45.1%	34.3%	0.7%
MCAS grade 8	40.2%	45.8%	14.0%	0.0%
PARCC grade 8	13.3%	62.0%	24.2%	0.5%
Smarter Balanced grade 8	16.1%	74.5%	8.6%	0.8%
14 states	41.2%	51.0%	7.9%	0.0%
AP	30.7%	54.2%	14.0%	1.0%
IB	21.0%	50.0%	28.0%	1.0%
NAEP	38.6%	54.0%	6.7%	0.6%
PISA	32.8%	50.4%	16.8%	0.0%
TIMSS	52.0%	46.5%	1.5%	0.0%

Note: Results for AP, IB, NAEP, PISA, TIMSS and fourteen states from Yuan and Le. The programs evaluated in this study are based on percentages of score points, whereas other assessment results are based on percentages of items.

Appendix B

Key Terminology

Accessibility: The degree to which an assessment allows all test takers, including English language learners and those with disabilities, to demonstrate their mastery of the knowledge, skills, and abilities being tested.

Alignment: The degree of agreement, overlap, or intersection among standards, instruction, and assessments.⁹⁰

Applications: A type of mathematics item that includes a context and requires students to use that context to decide which mathematical skills or concepts to use to solve the problem.⁹¹

Blueprints: A series of documents that together describe the content and structure of a test.

Calibration: A process used to make sure that evaluators or reviewers apply the scoring standards correctly and uniformly.⁹²

Claims: Statements about the knowledge, skills, or abilities of test takers who have attained a specified level of performance on the test.

Close reading: Close reading of text involves an investigation of a short piece of text, with multiple readings completed over multiple instructional lessons. Through text-based questions and discussion, students are guided to deeply analyze and appreciate various aspects of the text, such as key vocabulary and how its meaning is shaped by context; attention to form, tone, imagery and/or rhetorical devices; the significance of word choice and syntax; and the discovery of different levels of meaning as passages are read multiple times.⁹³

Cognitive demand: The type of thinking required to solve a task. (See DOK for additional information.)

College and career readiness: The level of academic preparation (in this case in ELA/Literacy and mathematics) needed by a student to both enroll in and successfully complete entry-level postsecondary collegiate or vocational programs without remedial academic work or assistance.

Computer-adaptive testing (CAT): A type of testing in which the questions presented to the test taker are selected on the basis of the test taker's previous responses. Correct answers by the test taker lead to harder questions; incorrect answers lead to easier questions. The purpose of adaptive testing is to use testing time

90. N. L. Webb, "Criteria for Alignment of Expectations and Assessments in Mathematics and Science Education," Council of Chief State School Officers and National Institute for Science Education Research Monograph No. 6 (Madison, WI: University of Wisconsin–Madison, Wisconsin Center for Education Research, 1997).

91. The National Center for the Improvement of Educational Assessment, Inc. (NCIEA), "Guide to Evaluating Assessments Using the CCSSO Criteria for High Quality Assessments: Focus on Test Content" (Dover, NH: NCIEA, February 2016): http://www.nciea.org/publication_PDFs/Guide%20to%20Evaluating%20CCSSO%20Criteria%20Test%20Content%2020316.pdf.

92. See Educational Testing Service, "Glossary of Standardized Testing Terms," https://www.ets.org/understanding_testing/glossary/.

93. S. Brown and L. Kappes, *Implementing the Common Core State Standards: A Primer on "Close Reading of Text"* (Washington, D.C.: The Aspen Institute, 2012).

more efficiently and to produce more precise scores for students toward the extremes of the performance distribution.⁹⁴

Conceptual Understanding: In mathematics, problems that assess conceptual understanding require students to use their understanding of mathematical concepts, operations and relations, as opposed to memorized facts, to solve problems.⁹⁵

Constructed-response item: A test question that requires the test taker to supply the answer, instead of choosing it from a list of possibilities. These items can take the form of a single word, number, or symbol, to a sentence, equation, or full paragraph.

Cut score: A specific test score used for classifying the test takers into groups on the basis of performance. Scores at or above that point are interpreted to mean something different from scores below that point, such as grade-level proficiency or lack thereof.

Depth: In the context of this study, Depth is a rating of the degree to which a test or testing program assesses the depth and complexity of skills and knowledge called for by college and career readiness standards.

Depth of Knowledge (DOK): The complexity or depth of understanding required to answer or explain an assessment related item, developed by Norman Webb. Level 1 includes basic recall of facts, concepts, information, or procedures; Level 2 includes skills and concepts, such as the use of information (graphs) or requires two or more steps with decision points along the way; Level 3 includes short-term strategic thinking; Level 4 includes extended thinking and, often, the application of concepts. Levels 3 and 4 are also referred to as “higher-order thinking skills.”⁹⁶

Evidence-based selected response items: These items feature two parts. The first (Part A) resembles a traditional multiple-choice item, and the second (Part B) calls for students to select/provide evidence to support the answer they chose in the first part. Part B can be multiple-choice, short-answer constructed response, multi-select, or technology-enhanced.

Evidence-centered design: A systematic approach to test development that involves a) the development of claims (see above), b) statements that describe the evidence needed from the student to prove the claims, and c) tasks designed to provide that evidence.⁹⁷

Form: A set of test questions given to a student.

Generalizability: The degree to which the inferences drawn from a sample are representative of the whole population.

High-quality texts: Texts that are content rich, exhibit exceptional craft and thought, and/or provide useful information.⁹⁸

Higher-order thinking skills: Complex thinking skills, including analysis, evaluation, creation, logical reasoning, judgment and critical thinking, as well as problem solving.

94. See Glossary of Education Reform, “Computer-Adaptive Test,” updated August 29, 2013, <http://edglossary.org/computer-adaptive-test>.

95. The National Center for the Improvement of Educational Assessment, Inc. (NCIEA), “Guide to Evaluating Assessments Using the CCSSO Criteria for High Quality Assessments: Focus on Test Content” (Dover, NH: NCIEA, February 2016): http://www.nciea.org/publication_PDFs/Guide%20to%20Evaluating%20CCSSO%20Criteria%20Test%20Content%20020316.pdf.

96. N. L. Webb, “Criteria for Alignment of Expectations and Assessments in Mathematics and Science Education.”

97. R. J. Mislevy, et al., “A Brief Introduction to Evidence-Centered Design,” CSE Technical Report 632, (Los Angeles, CA: University of California–Los Angeles, 2004).

98. *Ibid.*

Item: A test question, including the question itself, any stimulus material provided with the question, and the answer choices (for a multiple-choice item) or the scoring rules (for a constructed-response item).⁹⁹

Multi-select items: Similar to multiple-choice items, but students can (or must) select more than one correct answer.

Operational items: Test items administered to students as part of their official state assessment.

Performance task: Small clusters of items that address a specific scenario and ask students to perform a task to demonstrate their knowledge, understanding, and proficiency. They can be used to assess several standards or outcomes, and typically require students to apply knowledge and skills to novel situations.

Procedural fluency: As defined in the CCSS, “skill in carrying out procedures flexibly, accurately, efficiently, and appropriately.”¹⁰⁰

Proficiency: The level of performance needed to meet or exceed grade level standards.

Research tasks: These tasks require students to use multiple (minimum of two) informational texts about one topic (may be written, audio, or visual), analyze these multiple texts, and synthesize and/or organize information across texts. Research tasks may include multiple-choice and technology-enhanced items and must include at least one constructed-response item for which the student writes a response.¹⁰¹

Rubric: A set of rules for scoring the responses on a constructed-response item. Sometimes called a “scoring guide.”

Summative assessment: An assessment administered at or near the conclusion of a grade level or course to monitor student learning and/or to meet accountability requirements.

Task: Also referred to as a “performance task,” an activity that requires students to select and apply a set of skills to respond. Tasks typically involve a scenario and a series of test items, one or more of which are constructed-response items.

Technology-enhanced items: Items that require students to perform some action that cannot be executed on a traditional paper-and-pencil test, such as drag-and-drop, highlighting, or sequencing of objects.

Tier 2 vocabulary: According to the CCSS, these are “general academic” words that are far more likely to appear in written text than in speech. They “appear in all sorts of texts: informational texts (words such as *relative*, *vary*, *formulate*, *specificity*, and *accumulate*), technical texts (*calibrate*, *itemize*, *periphery*), and literary texts (*misfortune*, *dignified*, *faltered*, *unabashedly*).”¹⁰²

Traditional multiple-choice items (also referred to as selected response): Items that require the student to select one answer from a set of provided options.

Writing to sources: A writing activity that requires students to read or listen to one or more sources (texts, videos, audio files, etc.) and to draw evidence from the source(s) to support a conclusion, generalization, or inference in the written response.¹⁰³

99. See Educational Testing Service, “Glossary of Standardized Testing Terms,” https://www.ets.org/understanding_testing/glossary/.

100. CCSS, “Standards for Mathematical Practice,” <http://www.corestandards.org/Math/Practice/>.

101. The National Center for the Improvement of Educational Assessment, Inc. (NCIEA), “Guide to Evaluating Assessments Using the CCSSO Criteria for High Quality Assessments: Focus on Test Content” (Dover, NH: NCIEA, February 2016): http://www.nciea.org/publication_PDFs/Guide%20to%20Evaluating%20CCSSO%20Criteria%20Test%20Content%20020316.pdf.

102. CCSS, “Common Core State Standards for English/Language Arts and Literacy in History/Social Studies, Science, and Technical Subjects,” Appendix A, http://www.corestandards.org/assets/Appendix_A.pdf, 33–36.

103. *Ibid.*

Appendix C

The Methodology as Written

ELA/Literacy Criteria

There are nine criteria in English language arts (see Section I, *The Study Criteria*). Underneath these nine criteria are a total of thirty-nine sub-criteria. Of these thirty-nine sub-criteria, twenty-three are outcome sub-criteria (based on test items or forms) and sixteen are documentation sub-criteria.

Criterion B.1

The assessments are English language arts and literacy tests that are based on an aligned balance of high-quality literary and informational texts.

Item review. There are six sub-criteria under B.1. These sub-criteria focus on the texts that students are required to read as part of the tests. In particular, criteria B.1.1 and B.1.4 focus on the balance of informational and literary texts, criteria B.1.2 and B.1.5 focus on the quality of texts, and criteria B.1.3 and B.1.6 focus on the type of informational texts. These correspond to major focuses of the standards—that students should read both informational and literary texts (transitioning to a greater proportion of informational texts in higher grades), that these texts should be high quality and authentic, and that informational texts should be balanced across types (science/technical, history/social science, and literary nonfiction).

For the item review, raters read each text passage on the test and make four judgments about the text passages. First, they decide whether the text passage is literary (fiction) or informational (non-fiction). Second, for informational texts, they decide whether the text is primarily narrative (to tell a story a series of events) or expository (to inform or explain). Third, for informational texts on grades 6–12 assessments, they decide whether the text is science/technical, history/social science, or literary nonfiction (selecting more than one option is acceptable on this criterion). Fourth, they review the test metadata (the details provided by the test vendor to describe the characteristics of test items and passages), indicating whether the text has been previously published, and make a judgment as to whether the text is of publishable quality.

Form roll-up. To reach a rating for each of the criteria for a test form, the ratings from individual items must be rolled up to the test form (i.e., to give a complete picture of each test form on each criterion). For all test form roll-ups, individual reviewers use the tentative scoring guidance offered in the methodology. For instance, middle school criterion B.1.1 reads:

Texts are balanced across literary and informational text types and across genres, with more informational than literary texts used as the assessments move up in the grade bands. Goals include: In grades 3–8, approximately half of the texts are literature and half are informational.

This is operationalized in the tentative scoring guidance as 2 – Meets: 45–55 percent informational texts; 1 – Partially Meets: 33–44 percent or 56–84 percent; and 0 – Does Not Meet: 0–32 percent or 85–100 percent. Reviewers classify each passage on the test form as literary or informational, then calculate the proportion of

text passages that are informational on the entire test form, which is compared against the scoring guidance to create a tentative score. The reviewer is then given the opportunity to consider whether the tentative score matches her assessment of the entire test form. If not, he or she may change the score and provide comments. The reviewer's final score is passed on to Phase 2 of the review. The other two sub-criteria are described in the full scoring guidance in the appendix.

Document review. The documentation review for B.1 criterion is quite similar, as seen in the full scoring template. The proposed scoring guidelines are the same and the three documentation sub-criteria parallel the item review sub-criteria one-to-one.

Criterion B.2

The assessments require appropriate levels of text complexity; they raise the bar for text complexity each year so students are ready for the demands of college- and career-level reading no later than the end of high school. Multiple forms of authentic, previously published texts are assessed, including written, audio, visual, and graphic, as technology and assessment constraints permit.

Criterion B.2 focuses on text complexity. The goal of this criterion is to ensure that students are presented with grade-appropriate texts, as defined by both qualitative and quantitative measures of text complexity. The complexity of assessment passages should increase across grades, just as the complexity of readings in ELA/Literacy classes, such that the student is ready for college- or career-level reading by the end of high school.

Item review. To evaluate sub-criterion B.2.1 for the item review, raters again focus on the individual passages. Here, there are three ratings. First, reviewers indicate whether there is evidence that both quantitative and qualitative measures of text complexity have been used to select the passage. Second, reviewers indicate whether there is evidence that the text has been placed in an appropriate grade band (the CCSS grade bands are 2–3, 4–5, 6–8, 9–10, and 11–12) based on quantitative text complexity data. Finally, reviewers indicate whether there is evidence that the text has been placed in an appropriate grade level based on qualitative text complexity. The scoring at the form level is straightforward: forms are evaluated on the proportion of texts that fall into appropriate grade bands and levels based on the quantitative and qualitative evidence.

Document review. The documentation review for B.2 involves just one criterion, which asks reviewers to evaluate the process that is used to select texts and the extent to which that process adequately uses quantitative and qualitative text complexity data.

Criterion B.4

The assessments require all students to demonstrate a range of higher-order, analytical thinking skills in reading and writing based on the depth and complexity of college- and career-ready standards, allowing robust information to be gathered for students with varied levels of achievement.

While the first two criteria focus on reading passages, the remaining seven focus on test items. Criteria B.4 and B.9 are recommended for reviewers to evaluate first, because these are the only criteria that are evaluated for all items. The other criteria focus on only subsets of items depending on what is being tested (e.g., reading, writing, and language).

Outcome review. Criterion B.4 focuses on the cognitive demand of the test items relative to the standards. To evaluate the test items, reviewers use the Webb Depth of Knowledge (DOK) framework specific to the content area being evaluated.¹⁰⁴ Each item is rated on the DOK framework; items may be placed into one or two of the four available levels.

¹⁰⁴ N. L. Webb, "Depth of Knowledge Levels for Four Content Areas," 2002.

Form roll-up. Based on the item-level ratings, the DOK distribution of the entire test is calculated by averaging across items. To ensure an accurate calculation, test items are weighted by the number of score points associated with each item (e.g., on a two-item test where the first item is a one-point item placed at DOK 1 and the second item is a two-point item placed at DOK 3, the DOK distribution for the test would be 33 percent DOK 1 and 67 percent DOK 3).

To reach a 0–1–2 score, the DOK distribution for the test form is compared to the recommended DOK distribution for the grade-level standards. The standards' DOK distribution is obtained in a parallel fashion—each standard is rated for DOK and the results are averaged across all grade-level standards to determine the proportion at each DOK level. In this case each standard is equally weighted.

The comparison between the DOK distribution of the test and that of the standards is based on two measures. First, the DOK distributions are compared by creating a DOK index, which is based on the proportional agreement between the test form and the standards in DOK distribution. Mathematically, this is calculated as the sum of the cell-by-cell minima between the two documents. For example, suppose the standards were coded as being 25 percent at each of the four DOK levels, and the test was coded as being 40 percent at DOK 1, 40 percent at DOK 2, and 20 percent at DOK 3. The DOK index would be .70 (25 percent from DOK 1, 25 percent from DOK 2, and 20 percent from DOK 3). (See Table C-1.)

Second, because a key problem of prior-generation assessments was their low overall DOK, the test is compared with the standards specifically on coverage of higher-level DOK (3+), with the goal of ensuring that the proportion of the test on DOK 3+ is not markedly lower than that of the standards.

Document review. The documentation review for B.4 focuses on the extent to which DOK is an explicit part of the test documentation. That is, there is a research-based definition of cognitive demand, a way of operationalizing cognitive demand at the item level, and a rationale for and specification of distribution of cognitive demand for each test form.

Criterion B.9

High-quality items and a variety of types are strategically used to appropriately assess the standard(s).

The next criterion, which is also evaluated for every item on each test form, focuses on overall item and test quality. Reviewers are asked to consider quality along multiple dimensions. For test form quality, the main focus is that a variety of item types are used (and that at least one of these item types requires students to construct, rather than select, a response). For item quality, reviewers are asked to consider issues, such as alignment to standards (whether the vendor's rating of alignment is accurate), editorial accuracy, and whether there are multiple correct answers.

Outcome review. For the item review, the first step is to determine the item type for each item. (In practice, the item type is almost always provided as part of test metadata, so no reviewer judgment is needed in this step.)

Second, reviewers are presented with all of the CCSS listed in the test's metadata as being measured by the item (as many as twelve standards, depending on the vendor). For each item, the reviewers indicate whether they agree with each of the alignment ratings (yes or no). They also are asked to list any additional on- or off-grade standards to which they believe the item is aligned.

TABLE C 1

Example DOK Index Calculation

	Standards	Test	Minimum
DOK 1	.25	.40	.25
DOK 2	.25	.40	.25
DOK 3	.25	.20	.20
DOK 4	.25	0	0
			sum = .70

Third, reviewers rate each item on six possible quality issues, indicating as many as apply. These are:

- 1 Item may not yield valid evidence of targeted skill.
- 2 Item has issues with readability.
- 3 Item is incorrectly keyed.¹⁰⁵
- 4 Content is inaccurate.
- 5 Item has unintended correct answer.
- 6 Item has issues with editorial accuracy.

Reviewers are also given the option to indicate there is another problem not included in the six possible issues and provided a text box in which they can provide additional notes.

Form roll-up. These item-level ratings are rolled up to the whole test form for criteria B.9.1 and B.9.2. For B.9.1, the scoring is as follows:

- 2 – Meets:** At least two item formats are used, including one that requires students to generate, rather than select, a response (i.e., constructed response, extended writing);
- 1 – Partially Meets:** At least two formats (but not including constructed response) are used, including technology-based formats and/or two-part selected response formats; and
- 0 – Does Not Meet:** Only a traditional multiple choice format is used. For B.9.2, the scoring is based on the proportion of items that exhibit high quality (none of the quality issues) and alignment to standards, with the cutoffs for receiving a 2 being 95 percent of items having technical quality and 90 percent reflecting alignment to standards.

Document review. The documentation criteria for B.9 mirror the item review criteria, with B.9.3 focusing on item types and B.9.4 focusing on technical quality.

Criterion B.3

Reading assessments consist of test questions or tasks, as appropriate, which demand that students read carefully and deeply and use specific evidence from increasingly complex texts to obtain and defend correct responses.

Criterion B.3 focuses on important dimensions of reading as called for in the CCSS, such as close reading, focusing on central ideas, and citing direct textual evidence. This criterion includes six sub-criteria (four relative to the item review, two relative to documentation).

Outcome review. To evaluate criterion B.3, reviewers are asked to make four judgments about each relevant item:

- ◆ First, whether the reading item aligns to the specifics of the reading standard. This rating uses the ratings reviewers made for each individual standard that they did under B.9. In particular, the judgment here is whether the item aligns to one or more reading standards other than standard R.1, which underlies all reading items. (In other words, if an item aligns to only standard R.1, it does not receive credit here. If an item aligns to any other reading standard than R.1, it does receive credit.)
- ◆ Second, whether the item requires close reading and analysis—that is, whether it requires the student to carefully read the text, rather than drawing on prior or background knowledge.

¹⁰⁵ Some testing programs would not make their answer keys available to all reviewers due to confidentiality reasons. Thus, we removed from the evaluation the determination of “Item incorrectly keyed” in B.9.2 and C.5.2; our review does not verify the accuracy of the answer keys used by the testing programs.

- ◆ Third, whether the item focuses on central ideas and important particulars, rather than superficial or peripheral content.
- ◆ Fourth, whether the item requires direct use of textual evidence—that is, whether it requires students to provide or cite specific quotes or close paraphrases in order to answer the question.

Form roll-up. These four ratings are simply rolled up to the form level and compared against the tentative guidance, in order to reach score ratings. Sub-criterion B.3.1 is based on the proportion of test items that require close reading. Sub-criterion B.3.2 is based on the proportion of test items that focus on central ideas and important particulars. Sub-criterion B.3.3 is based on the proportion of test items that are aligned to the specifics of the standards. Each of these three sub-criteria are scores on the same scale:

- 2 – Meets:** 90–100 percent;
- 1 – Partially Meets:** 75–89 percent; and
- 0 – Does Not Meet:** 0–74 percent.

The fourth sub-criterion, B.3.4, is based on the proportion of test items that require direct use of textual evidence. The thresholds for this sub-criterion are slightly lower: 51–100 percent, 33–50 percent, and 0–32 percent, respectively.

Document review. The two documentation sub-criteria for B.3 mirror the outcome sub-criteria. Sub-criterion B.3.5 considers the extent to which documentation indicates expectations for close reading, central ideas and important particulars, as well as standards alignment. Sub-criterion B.3.6 considers the extent to which documentation requires text dependency.

Criterion B.5

Assessments emphasize writing tasks that require students to engage in close reading and analysis of texts so that students can demonstrate college- and career-ready abilities.

This criterion focuses on writing. There are four sub-criteria. Sub-criteria B.5.1 and B.5.3 focus on the distribution of writing tasks across exposition, argument, and narrative types. Sub-criteria B.5.2 and B.5.4 focus on the extent to which writing tasks are text-based (that is, they focus on writing to texts).

Item review. The task of evaluating writing items is straightforward for reviewers. First, reviewers determine the type of writing that is called for by the item, choosing from among persuasive/argumentative, narrative, expository, and blended (i.e., a combination of two or more of these types). Second, reviewers indicate whether the writing task requires writing to a text (that is, to confront text or other stimuli directly, to draw on textual evidence, and to support valid inferences from text or stimuli).

Form roll-up. To roll-up to the whole form, reviewers evaluate whether the writing prompts are approximately balanced across the three types in order to rate criterion B.5.1. To evaluate criterion B.5.2, reviewers simply indicate the proportion of writing tasks that require writing to texts.

Document review. The documentation analysis for writing is completely analogous and simply requires evaluating the extent to which the documentation lays out appropriate expectations for the distribution of writing types and the proportion of writing tasks that require writing to texts.

Criterion B.6

The assessments require students to demonstrate proficiency in the use of language, including vocabulary and conventions.

The B.6 criterion focuses on language. This criterion has the largest number of sub-criteria at eight. These comprise four sets of two sub-criteria, each of which has one outcome criterion and one documentation sub-criterion. Criteria B.6.1 and B.6.5 consider vocabulary items and the extent to which they focus on Tier 2

words and phrases, require the use of context, and assess words important to central ideas. Criteria B.6.2 and B.6.6 consider the extent to which language items mirror real-world activities, focus on common errors, and emphasize the conventions most important for readiness. Criteria B.6.3 and B.6.7 consider the emphasis placed on vocabulary on the test (i.e., the proportion of test points), whereas criteria B.6.4 and B.6.8 consider the emphasis placed on language on the test.

Item review. Ratings on criterion B.6 depend on whether the item is a vocabulary item or a language item. If it is a vocabulary item, the reviewers rate it on three dimensions including whether: 1) the item tests a Tier 2 word or phrase; 2) the item tests a word central to understanding the text; and 3) the item requires the use of context. If the item is a language item, reviewers rate it on three separate dimensions: whether it mirrors real-world activities; whether it covers skills in the CCSS language progression skills chart (i.e., the skills considered most important for readiness); whether the item focuses on common student errors.

Form roll-up. The roll-up to the 2/1/0 score for the four outcome criteria is straightforward. For B.6.1, reviewers evaluate the joint distribution of the proportion of Tier 2 words and phrases and the proportion of items assessing words important to central ideas against the thresholds laid out in the tentative cutoffs. For B.6.2, the reviewers consider the proportion of items assessing language meet the three language dimensions. For B.6.3 and B.6.4, reviewers consider the proportion of total score points allocated to vocabulary and language, respectively, where the tentative cutoffs for each are as follows:

- 2 – Meets:** Vocabulary (B.6.3)/language (B.6.4) is reported as a sub-score or ≥ 13 percent of score points;
- 1 – Partially Meets:** 10–12 percent of score points; and
- 0 – Does Not Meet:** 0–9 percent of score points.

Document review. The documentation criteria in B.6 exactly parallel the four outcome criteria just described.

Criterion B.7

The assessments require students to demonstrate research and inquiry skills, demonstrated by the ability to find, process, synthesize, organize, and use information from sources.

This criterion asks whether the assessments require students to demonstrate research and inquiry skills, demonstrated by the ability to find, process, synthesize, organize, and use information from sources.

Item review. In order to qualify as a research item, it must reference two or more reading passages. For those items that do reference two or more reading passages, there is only one rating to be made under this criterion: whether the item requires analysis, synthesis, and/or organization of information (mirroring real-world activities).

Form roll-up. To roll this sub-criterion up to the form-level for B.7.1, the index is the proportion of total research points that indeed require analysis, synthesis, and/or organization of information (mirroring real-world activities).

Document review. There are two documentation criteria for B.7. Documentation criterion B.7.3 mirrors B.7.1 and asks the proportion of research items that mirror real-world activities in the ways described. Documentation criterion B.7.2 simply asks whether a research score is reported or whether there is some other indication that the test vendor indicates that research is an important skill.

Criterion B.8

Over time, and as assessment advances allow, the assessments measure the speaking and listening communication skills students need for college and career readiness.

The final ELA/Literacy criterion—measuring speaking and listening—does not apply to the majority of existing assessments. It has four total sub-criteria—two for item review and two for documentation, and these are in matched pairs. Criteria B.8.1 and B.8.3 assess whether the listening “passages” reflect the criteria for passage quality from sections B.1 and B.2, as well as whether the items require active listening skills. Criteria B.8.2 and B.8.4 assess the quality of speaking tasks against several dimensions shown in the full scoring template.

Of the four tests, only Smarter Balanced currently assesses listening items, and none of the tests assess speaking. For these criteria, reviewers were allowed to provide an “insufficient evidence” rating, though as noted previously, these were ultimately scored as a 0 if the test did not offer items that meet them.

Item review. The B.8 criterion requires reviewers to make three ratings. First, raters indicate whether the listening passage meets the definition for high-quality reading passages laid out in criterion B.1. Second, they indicate whether the passage meets the definition for reading passage complexity laid out in criterion B.2. Finally, they rate whether the passage requires active listening skills.

Form roll-up. The roll-up to the form for B.8.1 considers the proportion of items that meet the criteria for high-quality items and active listening skills. The roll-up to the form for B.8.2 considers the proportion of speaking items that meet the criteria for high-quality speaking items.

Document review. The documentation review for B.8.3 and B.8.4 consider the extent to which the documentation indicates attention to the issues raised in B.8.1 and B.8.2.

Mathematics Criteria

There are five mathematics criteria, comprising six item review sub-criteria and seven documentation sub-criteria.

Criterion C.1

The assessments help educators keep students on track to readiness by focusing strongly on the content most needed in each grade or course for later mathematics.

The first mathematics criterion addresses the extent to which the assessment focuses on the content most needed for success in later mathematics (the “major work of the grade”). These are a set of content clusters in each grade’s standards that the standards’ creators have determined are the most central to success in future mathematics.

Criterion C.1 includes rating assessments’ mathematical progressions—descriptions of the development of student understanding of a particular topic organized across grade levels. Because we were only evaluating assessments at grades 5 and 8, we could not consider the coherence of assessments’ mathematical progressions across grades. This is an important area for future research using adjacent-grade assessments.

Item review. The task for a reviewer was in three steps:

- 1 verify the standards listed for each item by the test vendor;
- 2 add any additional standards that were tapped by the item; and
- 3 answer yes or no whether the item exclusively covered major work standards.

The index for this criterion is therefore the proportion of items that exclusively cover major work standards.

Document review. The documentation sub-criterion for C.1 focuses on the same issue, but using test blueprints. Again, the question is what proportion of the test content focuses on the major work of the grade.

Criterion C.2

The assessments measure conceptual understanding, fluency and procedural skill, and application of mathematics, as set out in college and career readiness standards.

The second mathematics criterion is based on the claim that the best mathematics assessments will measure a balance of conceptual understanding, procedural fluency, and application. Procedural skill items generally assess the extent to which test-takers can solve a purely mathematical item that has no or purely superficial context. Application items assess the extent to which test-takers can use problem context to choose skills or techniques to solve mathematical problems. Conceptual understanding items generally do not require procedures or

applications, though they may be made easier if students have those skills. Rather, conceptual understanding items typically test a sort of generalization from a particular example to a broader rule or class of problems.

Item review. The task for reviewers is to rate each mathematics item as covering one of these three types of knowledge (or a combination). Then, the form-level averages are compared to pre-established targets (with the goal of approximately an equal distribution across the three item types).

Document review. The two documentation criteria underneath C.2 also relate to the equal distribution of items across these three types. The first (C.2.2) simply asks whether item blueprints indicate there is an approximately equal distribution of items across these types. The second (C.2.3) asks whether all or nearly all students are likely to receive test forms that have an equal balance across the item types. This C.2.3 sub-criterion is designed to ensure that test forms are designed in such a way as to ensure that no form overly focuses on one item type at the expense of others. This may be especially relevant for computer-adaptive tests, where each student may take a unique form.

Criterion C.3

The assessments include brief questions and also longer questions that connect the most important mathematical content of the grade or course to mathematical practices, for example, modeling and making mathematical arguments.

The third mathematics criterion relates to the standards for mathematical practice. Specifically, this criterion seeks to ensure that each item that assesses a standard for mathematical practice (SMP) also assesses a content standard.

Item review. Reviewers simply rate each item that, according to the metadata, assesses an SMP for whether it also assesses a content standard. If the test vendor does not indicate standards for mathematical practice in the metadata, then this criterion cannot be addressed. The scoring guidelines suggest that 90 percent or more of items that assess a standard for mathematical practice also assess a content standard in order to award the highest score of 2.

Document review. The documentation criterion is the same, focusing instead on the extent to which the documentation indicates that items assessing mathematical practices also assess content standards.

Criterion C.4

The assessments require all students to demonstrate a range of higher-order analytical thinking skills in reading and writing based on the depth and complexity of college and career readiness standards, allowing robust information to be gathered for students with varied levels of achievement. Assessments include questions, tasks, and prompts about the basic content of the grade or course as well as questions that reflect the complex challenge of college- and career-ready standards.

As with ELA/Literacy criterion B.4, criterion C.4 focuses on the cognitive demand of the test items relative to the standards. The task for mathematics reviewers is exactly the same as for ELA/Literacy reviewers, with the only difference being that the DOK language is subject-specific for mathematics.

Criterion C.5

High-quality items and a variety of item types are strategically used to appropriately assess the standard(s).

The final four criteria for mathematics are under C.5: ensuring high-quality items and a variety of item types. Again, the procedures for C.5 are exactly parallel to those for B.9.

Appendix D

Author Biographies

Nancy Doorey, project manager and report co-author, has been deeply involved in educational reform for more than thirty years, serving as a teacher, state and district policymaker, program director, and consultant in the areas of assessment, teacher quality, and leadership. She has authored reports for several national organizations regarding advances in educational assessment, the six federally funded assessment consortia, and education governance. In 2009, Nancy co-led the creation of the Center for K–12 Assessment & Performance Management at the Educational Testing Service (ETS), which was charged with serving as a catalyst for advances in K–12 testing to support student learning, and authored its widely utilized series “Coming Together to Raise Achievement: New Assessments for the Common Core State Standards.” As the director of programs, she formulated the agendas and managed five national research symposia and six webcasts regarding advances and challenges in K–12 assessment. Nancy received a master’s degree in elementary education from Simmons College and a Certificate of Advanced Studies in computer science from Harvard University Extension; she also completed doctoral studies in educational leadership at Columbia University (ABD).

Morgan Polikoff, alignment expert and report co-author, is an assistant professor of education at the University of Southern California’s Rossier School of Education. His areas of expertise include K–12 education policy; college and career readiness standards; assessment policy; alignment; and the measurement of classroom instruction. Recent work has investigated teachers’ instructional responses to content standards and critiqued the design of school and teacher accountability systems. Ongoing work focuses on the implementation of college and career readiness standards and the influence of curriculum materials and assessments on implementation. He is an associate editor of *American Educational Research Journal* and serves on the editorial boards for *AERA Open* and *Educational Administration Quarterly*. His research is currently supported by the National Science Foundation, the WT Grant Foundation, and the Institute of Education Sciences, among other sources. Polikoff received his doctorate from the University of Pennsylvania’s Graduate School of Education in 2010 with a focus on education policy and his bachelor’s in mathematics and secondary education from the University of Illinois at Urbana-Champaign in 2006.

Appendix E

Review Panelist Biographies

ELA/Literacy and Mathematics Review Panelists

Philip M. Anderson is a Professor in the Department of Secondary Education at Queens College and also in the PhD program in Urban Education at the Graduate Center, City University of New York. Since earning a PhD from the University of Wisconsin–Madison’s Department of Curriculum and Instruction, he has published widely on developmental aesthetics, literacy and the literature curriculum, and cultural theory in curriculum design. Dr. Anderson has also served as an expert reviewer for the Massachusetts Comprehensive Assessment System (MCAS) English Language Arts Assessment, grades 3–10, for more than a decade.

Elizabeth Angell is an English language arts teacher at Atlantic Middle School in Quincy, Massachusetts. She has been a middle school language arts teacher for the last fourteen years and also served as a reading interventionist. Eight years ago, Ms. Angell began working with the Massachusetts Department of Elementary and Secondary Education (in collaboration with Measured Progress) as a member of the Assessment Development Committee for Grade 8 ELA. Ms. Angell also served on the Grade 8 ELA Assessment Development Committee for the MCAS for the Massachusetts Department of Elementary and Secondary Education.

Melody Arabo is a third-grade teacher at Keith Elementary in Walled Lake Consolidated School District. She is a National Education Association (NEA) Master Teacher, a Michigan Educator Voice Fellow, Oakland County’s Elite 40 under 40 Winner, and the 2015 Michigan Teacher of the Year. Mrs. Arabo is also the author of *Chaldean for Kids* and *Diary of a Real Bully* picture books, and served as an External Reviewer for ACT Aspire in 2014.

Jonathan Budd is K–12 Director of Curriculum, Instruction, and Assessments for Connecticut’s Trumbull Public School District. His particular expertise is literacy, with a focus on text complexity. A classroom teacher of nineteen years prior to becoming a district administrator, Dr. Budd holds a BA in English and Music from Connecticut College, an MA in English from Trinity College (Hartford), a Sixth-Year Diploma in Educational Leadership from the University of Connecticut, and a PhD in English Education from Teachers College, Columbia University. Dr. Budd also served as a content-area expert on SBAC’s Range-Finding Committee and Pre-Range-Finding Committee for ELA in 2013 and as a content-area expert on two Item Review Committees for Measured Progress in 2013.

Amber Chandler is a National Board English language arts teacher and has taught at the Frontier Central School District in New York for the last thirteen years. A 2014 New York Educator Voice Fellow, Ms. Chandler is an active educational blogger, certified school building leader, and serves in a staff developer role through her local teacher center on topics such as “Danielson’s Domains,” “Preventing Plagiarism Using the I-Search Paper,” and “Differentiation.” Ms. Chandler has also served as director of Frontier Central Summer School. Ms. Chandler has also worked with the New York Department of Education as a final eyes reviewer and test item developer.

Wally Etterbeek is an emeritus professor of mathematics at California State University Sacramento, where he was a faculty member for thirty-five years and department chair for twelve. During his early tenure, Dr. Etterbeek taught grades 2–6 for Project SEED, and more recently, taught gifted high school students in the San Juan District and Higher Level International Baccalaureate Mathematics. For four decades, Dr. Etterbeek has also been a

participating member and statistician for the Mathematics Diagnostic Texting Program (MDTP), a joint University of California/CSU project that develops readiness tests for grade 7 through first-year college mathematics courses. Dr. Etterbeek has also served as a member of UCLA's National Center for Research on Evaluation, Standards, and Student Testing (CRESST) panel that reviewed the alignment of the ACT Aspire Quality Core Algebra I and Geometry I courses and the Cambridge mathematics courses for grades 9–10 to the CCSS.

Shannon Garrison has been a fourth- and fifth-grade teacher at Solano Avenue Elementary School in downtown Los Angeles for the past eighteen years. She is a National Board Certified Teacher and was honored as a 2008 Milken Educator in California. In 2010, she was appointed to the National Assessment Governing Board, which sets policy for the National Assessment of Educational Progress (also known as the Nation's Report Card). Ms. Garrison also serves as the Chair of the Assessment Development Committee and was just reappointed for a second term by Secretary Duncan.

Jocelyn Gukeisen has been a middle school math educator in independent schools for fifteen years. Ms. Gukeisen currently serves as the Middle School Director and Director of Studies at The McGillis School in Salt Lake City where she coordinates and facilitates curriculum development across school divisions and across disciplines, and oversees the development of various programs and pedagogical strategies, as well as implementation of standardized testing. She has also served as an external item writer and consultant for ACT Aspire since June 2004.

Joey Hawkins was a middle school humanities public school teacher for over twenty years. A founding member of the Vermont Writing Collaborative and co-author of *Writing for Understanding: Using Backward Design to Help All Students Write Effectively*, she now provides professional development for K–12 teachers, focusing on the integration of literacy into content-based thinking and writing.

Roger Howe is the William Kenan Jr. Professor of Mathematics at Yale University. He has been teaching and conducting research in the Mathematics Department at Yale University for over forty years. His research interests focus on applications of symmetry, including representation theory, automorphic forms, harmonic analysis, and invariant theory. Dr. Howe has held visiting positions at many universities and research institutes in the United States, Europe, and Asia. In 1997–98, he served as a Phi Beta Kappa Visiting Scholar. He has served on many committees, including the NRC Study Committee for *Adding It Up*. He is currently a member of the U.S. National Commission on Mathematics Instruction and the Executive Committee of the International Commission on Mathematics Instruction. He is also a member of the American Academy of Arts and Sciences and the National Academy of Sciences.

Melisa R. Howey is an Education Specialist for Area Cooperative Educational Services (ACES) in Hamden, CT. She graduated from the University of Connecticut with a degree in Mathematics, received a Master's in Education from Quinnipiac University, and a sixth-year degree in Administration from Sacred Heart University. Ms. Howey has previously taught middle school math in both West Hartford and East Hartford, CT, and served as a summer school principal, interim math department head, and the K–6 district math specialist in East Hartford, CT. She also was part of a team who created grade 5 performance tasks for the Smarter Balanced Assessment.

Patricia Ford Kavanagh has been a mathematics teacher at Manchester-Essex Regional School District, Massachusetts since 1982. From 2008–2014, she developed and reviewed test items on the eighth-grade Math Assessment Development Team for the Massachusetts Department of Elementary and Secondary Education.

Bowen Kerins is a senior curriculum designer at Education Development Center (EDC). He develops current techniques and procedures used in the design, development, and implementation of curriculum, curriculum-based professional development, instruction, and assessments. Before joining EDC, Mr. Kerins was a high school mathematics teacher in Massachusetts. He has a BS in mathematics from Stanford University and an MA in teaching secondary mathematics from Boston University. Mr. Kerins also reviewed publicly released items for PARCC from 2012–2014.

David Kirshner is a Professor of Mathematics Education and co-Director of the Gordon A. Cain Center for STEM Literacy at Louisiana State University. He holds a doctoral degree in mathematics education and a master's degree in mathematics, both from University of British Columbia. Dr. Kirshner has served as president of the

North American Chapter of the International Group for the Psychology of Mathematics Education (PME-NA) and Chair of the Editorial Board of the Journal for Research in Mathematics Education. He recently co-edited the third edition of the *Handbook of International Research in Mathematics Education*.

Shelli Klein is a freelance consultant in educational assessment. She has more than twenty-five years of experience in education and educational publishing. As a former Publishing Director for an assessment company, she led large-scale assessment projects, developed and analyzed test blueprints, and trained the English language arts staff in interpreting and implementing the Common Core State Standards. In addition, Ms. Klein has experience teaching multiple grade levels, from elementary school through community college. She has also performed senior reviews of ELA field-test items for both the PARCC and the Smarter Balanced consortia, focusing on alignment of items to the specifications and standards in 2013–14.

Lianne Markus is a fourth- and fifth-grade ELA/Math teacher at Hope Township School, Hope, New Jersey. She received her undergraduate degree in journalism from the University of South Carolina and her NJ Teaching Certificate from Centenary College. As a member of the PARCC NJ Educator Leader Cadre, she represents the consortium at local, state, and national events. Ms. Markus participated in the PARCC Communications Boot Camp and served as a voiceover for Advance Illinois PARCC 101 Explainer Video, and also worked on the PARCC bias and sensitivity state educator review panel. She was nominated as NJ Governor's Teacher/Educational Professional of the Year in 2011, and nominated as Teacher of the Year 2010, 2012, 2013, 2014, and 2015.

Joe McGonegal has taught high school English and college composition for fifteen years in public and private settings in Wisconsin, Florida, New Hampshire, and Massachusetts before returning to school to study journalism. Since 2013, he has been a writer and strategist for MIT Technology Review and the MIT Alumni Association, where he currently serves as director of alumni education. Mr. McGonegal has also consulted on the ACT Aspire reading assessment since 2007.

Jennifer McPartland is a reading specialist at Gordon W. Mitchell Middle School in East Bridgewater, Massachusetts. She has worked in public education for seventeen years. She began her career teaching fifth grade, and then moved to serve as the school's reading specialist, where she has spent the last ten years. In Massachusetts, Ms. McPartland holds certifications in Elementary Education, Reading, and Administration. Since 2007, she has served as a member of the Massachusetts Department of Elementary and Secondary Education's MCAS Assessment Development Committee (ADC) for grade 6 ELA.

Robert Noreen is a professor and former Chair of the Department of English at California State University (CSU), Northridge. Previously, he was the Chair of the California English Placement Test Development Committee (1989–2003) and CSU English Assessment Program. Dr. Noreen was a member of the U.S. Department of Education committee that reviewed early drafts of the Common Core assessments in English. He has also served as a scoring leader and/or as a member of the item development team for the Graduate Management Admission Test (GMAT), National Teacher Examination (NTE), AP-English exam, Test of Written English/Test of English as a Foreign Language (TWE/TOEFL), and essay portion of the American Medical Colleges Admission Test (AMCAT).

Lynne Olmos is a National Board-certified teacher currently teaching English and drama in rural Mossyrock, Washington. Mrs. Olmos specializes in assessment preparation and creation, having served on district, state, and national committees to create, review, range-find and score assessment materials, for both students and teachers. Beginning in 2012, she worked on Washington State-managed committees for item writing and passage selection for Smarter Balanced. In the fall of 2014, Mrs. Olmos served as an item reviewer during Smarter Balanced's online achievement level setting. Her most recent projects involve art and the Common Core State Standards, including membership in Washington's Arts Cadre for arts educators concerned with standards and overall student achievement.

Charles Perfetti is the Distinguished University Professor of Psychology and Director of the Learning Research and Development Center at the University of Pittsburgh. Dr. Perfetti's research has addressed components of reading skill—reading comprehension, word identification, word learning, learning from multiple texts—and comparisons of reading across writing systems. His work has been published in several books and over 200

articles and book chapters. He has also received the Distinguished Scientific Contribution Award from both the Society for the Scientific Study of Reading and the Society for Text and Discourse.

Michael Roach is a visiting clinical assistant professor in mathematics education at Indiana University-Purdue University Indianapolis. He has worked as a high school mathematics teacher, a mathematics specialist at the Indiana Department of Education (IDOE), and a mathematics coach for Indianapolis Public Schools. Dr. Roach has a wide variety of experiences with large-scale assessment, including coordinating the development of an online testing program while at the IDOE and conducting secondary analyses of National Assessment of Educational Progress data.

Wilfried Schmid is the Dwight Parker Robinson Professor of Mathematics at Harvard University. He received his BA from Princeton University in 1964 and his PhD in 1967 from the University of California, Berkeley. Prior to joining the Harvard faculty in 1978, Dr. Schmid taught at Berkeley and at Columbia University. He helped to revise the Mathematics Curriculum Framework for Massachusetts; served as Mathematics Advisor to the Massachusetts Department of Elementary and Secondary Education; as member of the Steering Committee of Mathematics NAEP; as member of the Program Committee of the International Congress of Mathematics Education 2004; and as member of the National Mathematics Advisory Panel. Dr. Schmid assisted with the final review of MCAS mathematics test items from 2003–2004, and participated in the final review of mathematics SAT test items in 2008–09.

Colin Starr is a Professor of Mathematics at Willamette University, where he has taught since 2003. Dr. Starr has been involved in numerous projects that combine mathematics and teaching, including the Algebraic Thinking Project and serving as a co-principal investigator on two National Science Foundation Research Experiences for Undergraduates grants. He has reviewed MCAS items since 2005 as an expert reviewer under work facilitated by Measured Progress and worked with various other groups examining the teaching of mathematics.

Rose Switalski is currently an eighth-grade math teacher at an art integrated charter school in Las Vegas. She completed her bachelor's degree in secondary mathematics education at the University of Nevada. Over the past eight years, she has taught grades 6–12.

Erin Thompson is a secondary English language arts assessment specialist at the Indiana Department of Education (IDOE), where she has worked for the past five years. Prior to working at the IDOE, Ms. Thompson taught English/language arts at grades 6–9 for nine years. She has also participated in PARCC's Operational Work Group for English language arts, where she was involved in passage and assessment review, policy decisions, and creation of performance level descriptors.

Howard Tuttmann is a K–8 Instructional Math Coach in Everett, Massachusetts, an urban district just north of Boston. For several years, he has served on the Massachusetts Department of Elementary and Secondary Education's Assessment Development Committee (ADC), reviewing fourth-grade math items for the Massachusetts Comprehensive Assessment System (MCAS). He has also served on the PARCC's Math Item Review Committee, reviewing questions for the grades 3–5, and participated in PARCC's Mathematics Performance Level Setting Committee for grades 5–6.

Tad Watanabe is a Professor of Mathematics Education at Kennesaw State University (KSU), where he has taught mathematics content courses for prospective elementary school teachers since 2006. Prior to coming to KSU, he taught four years at the Pennsylvania State University and ten years at Towson University in Towson, Maryland. Dr. Watanabe is interested in various mathematics education practices in Japan, and has written several articles on Japanese lesson study and elementary school mathematics curriculum. He has also worked with numerous lesson study groups throughout the United States.

Rosalynn Wolfe is a graduate student at the University of Iowa, where she is studying Educational Measurement and Statistics. Prior to this Rosalynn was a math teacher for ten years, mostly grades 6–12, in Alaska, Arizona, and Iowa. Ms. Wolfe has worked on curriculum development for Tucson Unified School District and several charter schools in Arizona. She has also reviewed and written items for alignment with academic standards for ACT Aspire for the past three years.

Christopher Yakes began a position as Associate Professor of Mathematics at Salem State University in September 2015. For the previous nine years, he was a faculty member at California State University (Chico) where he taught mathematics courses to undergraduate mathematics majors and pre-service teachers. While at Chico State, he also worked extensively with in-service teachers through the California Mathematics Project and other grant-funded projects, which included holding professional development workshops, facilitating lesson study groups, and serving as a math coach. For the past three years, Dr. Yakes has also worked with the California Department of Education on rewriting the California Mathematics Framework for Public Schools.

Amy Youngblood is the founder of EduOptimus, which provides professional development to school districts in Missouri and across the Midwest. She began her career as a kindergarten teacher in Southwest Missouri, and has been an elementary teacher, a K–8 gifted teacher, a 9–12 gifted teacher, federal programs director, Missouri School Improvement Program (MSIP) Coordinator, and a K–12 curriculum director. Additionally, Ms. Youngblood has served as the President for the Missouri Affiliate of the Association for Supervision and Curriculum Development (ASCD), a member of the ASCD Leadership Council, a member of the Missouri Expert Curriculum Review Team for MSIP, and a member of the Staff Development Leadership Council. Currently, she is a Lead Facilitator for EdReports and a member of the Educators Evaluating the Quality of Instructional Products (EQulP) Peer Review Panel, which was created by Achieve to identify high-quality instructional materials aligned to the Common Core State Standards.

Generalizability Review Panelists

Heather Goodwin-Nelson is a Special Education Instructional Facilitator for Utah Virtual Academy. She received her BS in Elementary Education from Brigham Young University-Hawaii in 1994. During her undergraduate work, Ms. Goodwin-Nelson taught various grade levels before receiving a master's degree in Special Education from Brigham Young University.

Roger Howe is the William Kenan Jr. Professor of Mathematics at Yale University. He has been teaching and conducting research in the Mathematics Department at Yale University for over forty years. His research interests focus on applications of symmetry, including representation theory, automorphic forms, harmonic analysis, and invariant theory. Dr. Howe has held visiting positions at many universities and research institutes in the U.S., Europe, and Asia. In 1997–98, he served as a Phi Beta Kappa Visiting Scholar. He has served on many committees, including the NRC Study Committee for *Adding It Up*. He is currently a member of the U.S. National Commission on Mathematics Instruction and the Executive Committee of the International Commission on Mathematics Instruction. He is also a member of the American Academy of Arts and Sciences and the National Academy of Sciences.

Shelli Klein is a freelance consultant in educational assessment. She has more than twenty-five years of experience in education and educational publishing. As a former publishing director for an assessment company, she led large-scale assessment projects, developed and analyzed test blueprints, and trained the English language arts staff in interpreting and implementing the Common Core State Standards. In addition, Ms. Klein has experience teaching multiple grade levels, from elementary school through community college. She has also performed senior reviews of ELA field-test items for both the PARCC and the Smarter Balanced consortia, focusing on alignment of items to the specifications and standards in 2013–14.

Tim LaVan is a math teacher in the Oil City Area School District (Pennsylvania) in his twenty-fifth year of teaching. He is highly experienced in current techniques and procedures used in the design, development, and implementation of curriculum, instruction, and assessments, being a member of the Mathematics Advisory Committee for Pennsylvania. Dr. LaVan has been a member of PARCC advisory committees representing Pennsylvania as well. He holds a BS in mathematics education from Clarion University of PA, an MS in Computational Science from Michigan State University, and completed doctoral work in Curriculum and Instruction at the University of Pittsburgh.

Robert Noreen is a professor and former Chair of the Department of English at California State University (CSU), Northridge. Previously, he was the Chair of the California English Placement Test Development Committee (1989–2003) and CSU English Assessment Program. Dr. Noreen was a member of the U.S. Department of

Education committee that reviewed early drafts of the Common Core assessments in English. He has also served as a scoring leader and/or as a member of the item development team for the Graduate Management Admission Test (GMAT), National Teacher Examination (NTE), AP-English exam, Test of Written English/Test of English as a Foreign Language (TWE/TOEFL), and essay portion of the American Medical Colleges Admission Test (AMCAT).

Tabitha Pacheco is a National Board-certified teacher in exceptional needs, a 2015 National Teaching Fellow for the Hope Street Group, and serves on the Practitioners Advisory Group for The Centers on Great Teachers and Leaders. For the past nine years, Ms. Pacheco has worked with students with disabilities and is an expert on accommodating and scaffolding the CCSS to meet the needs of all learners. She has a bachelor's degree from Brigham Young University in Family Life and a post-baccalaureate degree in special education from Brigham Young University.

Michael Roach is a visiting clinical assistant professor in mathematics education at Indiana University-Purdue University Indianapolis. He has worked as a high school mathematics teacher, a mathematics specialist at the Indiana Department of Education (IDOE), and a mathematics coach for Indianapolis Public Schools. Dr. Roach has a wide variety of experiences with large-scale assessment, including coordinating the development of an online testing program while at the IDOE and conducting secondary analyses of National Assessment of Educational Progress data.





Charlie Wayne has been an assessment specialist in mathematics with the Pennsylvania Department of Education (PDE) for over seventeen years. He has been Pennsylvania's mathematics representative with PARCC (Pennsylvania is a participating state in PARCC and Smarter Balanced). He has also participated in various alignment studies with Dr. Norman Webb, Achieve, the American Association for the Advancement of Science (AAAS), and HumRRO. Mr. Wayne has a BS in Economics and an MS in Mathematics along with a graduate certificate in Large-Scaled Assessment Education. Prior to coming to PDE, he taught in elementary and middle schools, at the post-secondary level, and at a business institute.

Appendix F





Full ELA/Literacy and Math Ratings and Summary Statements (Grades 5 and 8)









OVERALL SUMMARY STATEMENTS FOR ELA/LITERACY	
Program	Summary Statements
ACT Aspire	In ELA/Literacy, ACT Aspire receives a Limited to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The combined set of ELA/Literacy tests (reading, writing, and English) requires close reading and adequately evaluates language skills. More emphasis on assessment of writing to sources, vocabulary, and research and inquiry, as well as increasing the cognitive demands of test items, will move the assessment closer to fully meeting the criteria. Over time, the program would also benefit by developing the capacity to assess speaking and listening skills.
MCAS	In ELA/Literacy, MCAS receives a Limited to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The test requires students to closely read high-quality texts and a variety of high-quality item types. However, MCAS does not adequately assess several critical skills, including reading informational texts, writing to sources, language skills, and research and inquiry; further, too few items assess higher-order skills. Addressing these limitations would enhance the ability of the test to signal whether students are demonstrating the skills called for in the standards. Over time, the program would also benefit by developing the capacity to assess speaking and listening skills.
PARCC	In ELA/Literacy, PARCC receives an Excellent Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The tests include suitably complex texts, require a range of cognitive demand, and demonstrate variety in item types. The assessments require close reading, assess writing to sources, research, and inquiry, and emphasize vocabulary and language skills. The program would benefit from the use of more research tasks requiring students to use multiple sources, and, over time, developing the capacity to assess speaking and listening skills.
Smarter Balanced	In ELA/Literacy, Smarter Balanced receives a Good to Excellent Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The tests assess the most important ELA/Literacy skills of the CCSS, using technology in ways that both mirror real-world uses and provide quality measurement of targeted skills. The program is most successful in its assessment of writing and research and inquiry. It also assesses listening with high-quality items that require active listening, which is unique among the four programs. The program would benefit by improving its vocabulary items, increasing the cognitive demand in grade 5 items, and, over time, developing the capacity to assess speaking skills.









ELA/LITERACY CONTENT CRITERIA OVERVIEW

Criteria	Program and Rating	Summary Statements
I. Assesses the content most needed for College and Career Readiness.	 ACT Aspire LIMITED/UNEVEN MATCH	<p>ACT Aspire receives a Limited/Uneven Match to the CCSSO Criteria for Content in ELA/Literacy. The assessment program includes an emphasis on close reading and language skills.</p> <p>However, the reading items fall short on requiring students to cite specific textual information in support of a conclusion, generalization, or inference and in requiring analysis of what has been read. In order to meet the criteria, assessing writing to sources, vocabulary, and research and inquiry need to be strengthened.</p>
	 MCAS LIMITED/UNEVEN MATCH	<p>MCAS receives a Limited/Uneven Match to the CCSSO Criteria for Content in ELA/Literacy. The assessment requires students to read closely well-chosen texts and presents test questions of high technical quality.</p> <p>However, the program would be strengthened by assessing writing annually, assessing the three types of writing called for across each grade band, requiring writing to sources, and placing greater emphasis on assessing research and language skills.</p>
	 PARCC EXCELLENT MATCH	<p>PARCC receives an Excellent Match to the CCSSO Criteria for Content in ELA/Literacy. The program demonstrates excellence in the assessment of close reading, vocabulary, writing to sources, and language, providing a high-quality measure of ELA content, as reflected in college and career ready standards.</p> <p>The tests could be strengthened by the addition of research tasks that require students to use two or more sources and, as technologies allow, a listening and speaking component.</p>
	 Smarter Balanced EXCELLENT MATCH	<p>Smarter Balanced receives an Excellent Match to the CCSSO Criteria for Content in ELA/Literacy. The program demonstrates excellence in the areas of close reading, writing to sources, research, and language. The listening component represents an important step toward adequately measuring speaking and listening skills—a goal specifically reflected in the standards. Overall, Smarter Balanced is a high-quality measure of the content required in ELA and literacy, as reflected in college and career readiness standards.</p> <p>A greater emphasis on Tier 2 vocabulary would further strengthen these assessments relative to the criteria.</p>





ELA/LITERACY CONTENT CRITERIA

Criteria	Program and Rating	Summary Statements
B.3 Reading: Tests require students to read closely and use specific evidence from texts to obtain and defend correct responses.	 ACT Aspire LIMITED/UNEVEN MATCH	<p>On “requiring students to read closely and use evidence from texts,” the rating is Limited/Uneven Match. Although most reading items require close reading of some kind, too many can be answered without analysis of what was read. Items that purport to require specific evidence from text often require only recall of information from text. To meet this criterion, the test items should require students to cite specific text information in support of some conclusion, generalization, or inference drawn from the text.</p>
	 MCAS GOOD MATCH	<p>On “requiring students to read closely and use evidence from texts,” the rating is Good Match. Most reading items require close reading and focus on central ideas and important particulars. Some questions, however, do not require the students to provide direct textual evidence to support their responses. In addition, too many items do not align closely to the specifics of the standards.</p>
	 PARCC EXCELLENT MATCH	<p>On “requiring students to read closely and use evidence from texts,” the rating is Excellent Match. Nearly all reading items require close reading, the understanding of central ideas, and the use of direct textual evidence.</p>
	 Smarter Balanced EXCELLENT MATCH	<p>On “requiring students to read closely and use evidence from texts,” the rating is Excellent Match. Nearly all reading items align to the reading standards requiring close reading, the understanding of central ideas, and use of direct textual evidence in support of a conclusion, generalization, or inference.</p>





ELA/LITERACY CONTENT CRITERIA		
Criteria	Program and Rating	Summary Statements
B.5 Writing: Tasks require students to engage in close reading and analysis of texts. Across each grade band, tests include a balance of expository, persuasive/argument, and narrative writing.	 ACT Aspire LIMITED/UNEVEN MATCH	On “assessing writing,” the rating is Limited/Uneven. Although the program documentation shows that a balance of all three writing types is required across each grade band, the writing prompts do not require writing to sources. As a result, the program insufficiently assesses the types of writing required by college and career readiness standards.
	 MCAS WEAK MATCH	On “assessing writing,” the rating is Weak Match. Writing is assessed at only one grade level per band, and there is insufficient opportunity to assess writing of multiple types. In addition, the writing assessments do not require students to use sources. As a result, the program inadequately assesses the types of writing required by college and career readiness standards.
	 PARCC EXCELLENT MATCH	On “assessing writing,” the rating is Excellent Match. The assessment meets the writing criterion, which requires writing to sources. Program documentation shows that a balance of all three writing types is required across each grade band.
	 Smarter Balanced EXCELLENT MATCH	On “assessing writing,” the rating is Excellent Match. The writing items are of high quality, and the writing prompts all require the use of textual evidence. Program documentation shows that a balance of all three writing types is required across each grade band.
B.6 Vocabulary and language skills: Tests place sufficient emphasis on academic vocabulary and language conventions as used in real-world activities.	 ACT Aspire GOOD MATCH	On “emphasizing vocabulary and language skills,” the rating is Good Match. Language items meet the criterion for being tested within writing activities, though more items are needed that are embedded in real world tasks such as editing. The vocabulary items do not meet the criterion because there are too few of them and not enough assess Tier 2 words.
	 MCAS LIMITED/UNEVEN MATCH	On “emphasizing vocabulary and language skills,” the rating is Limited/Uneven Match. Vocabulary items are sufficient and generally aligned to the criterion; however, the grade 5 items need more words at the Tier 2 level. Furthermore, a lack of program documentation means that the quality of vocabulary assessments cannot be substantiated across forms. MCAS does not meet the criterion for assessing language skills, which call for them to be assessed within writing assessments that mirror real-world activities including editing and revision.
	 PARCC EXCELLENT MATCH	On “emphasizing vocabulary and language skills,” the rating is Excellent Match. The test contains an adequate number of high-quality items for both language use and Tier 2 vocabulary and awards sufficient score points, according to the program's documentation, to both of these areas.
	 Smarter Balanced GOOD MATCH	On “emphasizing vocabulary and language skills,” the rating is Good Match. Language skill items are contained in a sub-score and meet the criterion for being assessed within writing and mirroring real-world activities such as editing and revision. The number of items that test vocabulary is a bit low; further, items coded as vocabulary too often did not test Tier 2 vocabulary words.









ELA/LITERACY CONTENT CRITERIA		
Criteria	Program and Rating	Summary Statements
B.7 Research and inquiry: Assessments require students to demonstrate the ability to find, process, synthesize, and organize information from multiple sources.	 ACT Aspire LIMITED/UNEVEN MATCH	On “assessing research and inquiry,” the rating is Limited/Uneven Match. Although the one item at each grade level involving research and inquiry did indeed require analysis and organization of information, this single item is insufficient to provide a quality measure of research and inquiry.
	 MCAS WEAK MATCH	On “assessing research and inquiry,” the rating is Weak Match. The assessment has no test questions devoted to research.
	 PARCC EXCELLENT MATCH	On “assessing research and inquiry,” the rating is Excellent Match. The research items require analysis, synthesis, and/or organization and the use of multiple sources, therefore meeting the criterion for Excellent.
	 Smarter Balanced EXCELLENT MATCH	On “assessing research and inquiry,” the rating is Excellent Match. The research items require analysis, synthesis, and/or organization, and the use of multiple sources.
B.8 Speaking and listening: Over time, and as assessment advances allow, the assessments measure speaking and listening communication skills.	 ACT Aspire WEAK MATCH	On assessing “speaking and listening,” the rating is Weak Match. The program does not assess speaking or listening at this time. Because this criterion is to be met “over time, as assessment advances allow,” this rating is not included in the overall rating for Content.
	 MCAS WEAK MATCH	On assessing “speaking and listening,” the rating is Weak Match. The program does not assess speaking or listening at this time. Because this criterion is to be met “over time, as assessment advances allow,” this rating is not included in the overall rating for Content.
	 PARCC WEAK MATCH	On assessing “speaking and listening,” the rating is Weak Match. The program does not assess speaking or listening at this time. Because this criterion is to be met “over time, as assessment advances allow,” this rating is not included in the overall rating for Content.
	 Smarter Balanced LIMITED/UNEVEN MATCH	On assessing “speaking and listening,” the rating is Limited/Uneven Match. Listening is tested with high-quality items that assess active listening skills. Speaking is not assessed. Because this criterion is to be met “over time, as assessment advances allow,” this rating is not included in the overall rating for Content.





ELA/LITERACY DEPTH CRITERIA OVERVIEW

Criteria	Program and Rating	Summary Statements
II. Assesses the depth that reflects the demands of College and Career Readiness.	 ACT Aspire GOOD MATCH	<p>ACT Aspire receives a rating of Good Match for Depth in ELA/Literacy. The program's assessments are built on high-quality test items and texts that are suitably complex.</p> <p>To fully meet the CCSSO Criteria, at both grade levels more cognitively demanding test items are needed at both grade levels, as well as additional literary narrative text, as opposed to literary informational texts.</p>
	 MCAS GOOD MATCH	<p>MCAS receives a rating of Good Match for Depth. The assessments do an excellent job in presenting a range of complex reading texts.</p> <p>To fully meet the demands of the CCSSO Criteria, however, the test needs more items at higher levels of cognitive demand, a greater variety of items to test writing to sources and research, and more informational texts, particularly those of an expository nature.</p>
	 PARCC EXCELLENT MATCH	<p>PARCC receives a rating of Excellent Match for Depth in ELA/Literacy. The PARCC assessments meet or exceed the depth and complexity required by the criteria through a variety of item types that are generally high quality.</p> <p>A better balance between literary and informational texts would further strengthen the assessments in addressing the criteria.</p>
	 Smarter Balanced GOOD MATCH	<p>Smarter Balanced receives a rating of Good Match for Depth in ELA/Literacy. The assessments use a variety of item types to assess student reading and writing to sources.</p> <p>The program could better meet the depth criteria by increasing the cognitive demands of the grade 5 assessment and ensuring that all items meet high editorial and technical quality standards.</p>

ELA/LITERACY DEPTH CRITERIA

Criteria	Program and Rating	Summary Statements
B.1 Text quality and types: Tests include an aligned balance of high-quality literary and informational texts.	 ACT Aspire GOOD MATCH	<p>On “the balance of high-quality literary and informational texts,” the rating is Good Match. The texts are of high quality, and the proportion of informational texts meets the criterion.</p> <p>The assessment would better align to the criterion, however, with additional literary narrative text, as opposed to literary informational text.</p>
	 MCAS GOOD MATCH	<p>On “the balance of high-quality literary and informational texts,” the rating is Good Match. The quality of the texts is very high.</p> <p>Regarding the balance of text types, some forms had too few informational texts.</p>
	 PARCC GOOD MATCH	<p>On “the balance of high-quality literary and informational texts,” the rating is Good Match.</p> <p>Although the passages are consistently of high quality, the tests would have better reflected the criterion with additional literary nonfiction passages.</p>
	 Smarter Balanced EXCELLENT MATCH	<p>On “the balance of high-quality literary and informational texts,” the rating is Excellent Match. Overall text quality is high, and among informational texts there is a high proportion of expository text types.</p>





ELA/LITERACY DEPTH CRITERIA		
Criteria	Program and Rating	Summary Statements
B.2 Complexity of texts: Test passages are at appropriate levels of text complexity, increasing through the grades, and multiple forms of authentic, high-quality texts are used.	 ACT Aspire GOOD MATCH	The rating for “text complexity” is Good Match. It is based solely on the review of program documentation, which is determined to have met the criterion. The test blueprints and other documents clearly and explicitly require texts to increase in complexity grade-by-grade and for texts to be placed in grade bands and grade levels based on appropriate quantitative and qualitative data.
	 MCAS GOOD MATCH	The rating for “use of appropriate levels of text complexity” is Good Match. It is based solely on the review of program documentation, which is determined to have met the criterion. The test blueprints and other documents clearly and explicitly require texts to increase in complexity grade-by-grade and for texts to be placed in grade bands and grade levels based on appropriate quantitative and qualitative data.
	 PARCC GOOD MATCH	The rating for “use of appropriate levels of text complexity” is Good Match. It is based solely on the review of program documentation, which is determined to have met the criterion. The test blueprints and other documents clearly and explicitly require texts to increase in complexity grade-by-grade and for texts to be placed in grade bands and grade levels based on appropriate quantitative and qualitative data.
	 Smarter Balanced GOOD MATCH	The rating for “use of appropriate levels of text complexity” is Good Match. It is based solely on the review of program documentation, which is determined to have met the criterion. The test blueprints and other documents clearly and explicitly require texts to increase in complexity grade-by-grade and for texts to be placed in grade bands and grade levels based on appropriate quantitative and qualitative data.
B.4 Cognitive demand: The distribution of cognitive demand for each grade level is sufficient to assess the depth and complexity of the standards.	 ACT Aspire WEAK MATCH	On “requiring a range of cognitive demand,” the rating is Weak Match. To better reflect the depth and complexity of the standards, both grade-level tests should require more items with higher cognitive demands, although this problem is greater at grade 8.
	 MCAS LIMITED/UNEVEN MATCH	On “requiring a range of cognitive demand,” the rating is Limited/Uneven Match. More items that measure the higher levels of cognitive demand are needed to sufficiently assess the depth and complexity of the standards.
	 PARCC EXCELLENT MATCH	On “requiring a range of cognitive demand,” the rating is Excellent Match. The test is challenging overall; indeed the cognitive demand of the grade 8 test exceeds that of the CCSS.
	 Smarter Balanced GOOD MATCH	On “requiring a range of cognitive demand,” the rating is Good Match. The cognitive demand of items cover a sufficient range and, in grade 8, the percentage of more demanding items (DOK 3 and 4) correspond well to the demand of the standards. However, the grade 5 test needs more items at higher levels of cognitive demand to reflect fully the depth and complexity of the standards.

ELA/LITERACY DEPTH CRITERIA		
Criteria	Program and Rating	Summary Statements
B.9 High-quality items and variety of item types: Items are of high technical and editorial quality and each test form includes at least two items types, including at least one that requires students to generate rather than select a response.	 ACT Aspire EXCELLENT MATCH	On “ensuring high-quality items and a variety of item types,” the rating is Excellent Match. The test includes items that exhibit high technical quality and editorial accuracy. Multiple item formats are used, including student-constructed responses.
	 MCAS EXCELLENT MATCH	On “ensuring high-quality items and a variety of item types,” the rating is Excellent Match. Multiple item formats are used, including student-generated response items. The items exhibit high technical quality and editorial accuracy. The paper-and-pencil format precludes the use of technology-enhanced items, but the criterion for multiple item types is met.
	 PARCC EXCELLENT MATCH	On “ensuring high-quality items and a variety of item types,” the rating is Excellent Match. The tests use multiple item formats, including student-constructed responses.
	 Smarter Balanced GOOD MATCH	On “ensuring high-quality items and a variety of item types,” the rating is Good Match. The tests use multiple formats and technology-enhanced items including constructed responses. However, editorial or technical issues, including readability, were noted in a number of items.





OVERALL SUMMARY STATEMENTS FOR MATHEMATICS

Program	Summary Statements
ACT Aspire	<p>In mathematics, ACT Aspire receives a Limited/Uneven to Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. Some of the mismatch with the criteria is likely due to intentional program design, which requires that items be included from previous and later grade(s).</p> <p>The items are generally high quality and test forms at grades 5 and 8 have a range of cognitive demand, but in each case the distribution contains significantly greater emphasis at DOK 3 as reflected by the standards. Thus, students who score well on the assessments will have demonstrated strong understanding of the standard's more complex skills. However, the grade 8 test may not fully assess standards at the lowest level of cognitive demand.</p> <p>The tests would better meet the CCSSO Criteria with an increase in the number of items focused on the major work of the grade and the addition of more items at grade 8 that assess standards at DOK 1.</p>
MCAS	<p>In mathematics, MCAS receives a Limited Match to the CCSSO Criteria for content and an Excellent Match for depth relative to assessing whether students are on track to meet college and career readiness standards. The MCAS mathematics test Items are of high technical and editorial quality. Additionally, the content is distributed well across the breadth of the grade level standards, and test forms closely reflect the range of cognitive demand of the standards.</p> <p>Yet, the grade 5 tests have an insufficient degree of focus on the major work of the grade.</p> <p>While mathematical practices are required to solve items, MCAS does not specify the assessed practices(s) within each item or their connections to content standards.</p> <p>The tests would better meet the criteria through increased focus on major work at grade 5 and identification of the mathematical practices that are assessed—and their connections to content.</p>
PARCC	<p>In mathematics, PARCC receives a Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The assessment is reasonably well aligned to the major work of each grade. At grade 5, the test includes a distribution of cognitive demand that is similar to that of the standards. At grade 8, the test has greater percentages of higher-demand items (DOK 3 and 4) than reflected by the standards, such that a student who scores well on the grade 8 PARCC assessment will have demonstrated strong understanding of the standard's more complex skills. However, the grade 8 test may not fully assess standards at the lowest level (DOK 1) of cognitive demand.</p> <p>The test would better meet the CCSSO Criteria through additional focus on the major work of the grade, the addition of more items at grade 8 that assess standards at DOK 1 and increased attention to accuracy of the items, primarily editorial, but in some instances mathematical.</p>
Smarter Balanced	<p>In mathematics, Smarter Balanced has a Good Match to the CCSSO Criteria relative to assessing whether students are on track to meet college and career readiness standards. The test provides adequate focus on the major work of the grade, although it could be strengthened at grade 5.</p> <p>The tests would better meet the CCSSO Criteria through increased focus on the major work at grade 5, an increase in the number of items on the grade 8 tests assessing standards at DOK 1, and attention to serious mathematical or editorial flaws in some items.</p>

MATHEMATICS CONTENT CRITERIA OVERVIEW





Criteria	Program and Rating	Summary Statements
I. Assesses the content most needed for College and Career Readiness.	 ACT Aspire LIMITED/UNEVEN MATCH	<p>ACT Aspire provides a Limited/Uneven Match to the CCSSO Criteria for Content in Mathematics. The program does not focus exclusively on the major work of the grade, but rather, by design, assesses material from previous and later grade(s). This results in a weaker match to the criteria.</p> <p>The tests could better meet the criteria at both grades 5 and 8 by increasing the number of items that assess the major work of the grade.</p>
	 MCAS LIMITED/UNEVEN MATCH	<p>MCAS provides a Limited/Uneven Match to the CCSSO Criteria for Content in Mathematics. While the grade 8 assessment focuses strongly on the major work of the grade, the grade 5 assessment does not, as it samples more broadly from the full range of standards for the grade.</p> <p>The tests could better meet the criteria through increased focus on the major work of the grade on the grade 5 test.</p>
	 PARCC GOOD MATCH	<p>PARCC provides a Good Match to the CCSSO Criteria for Content in Mathematics.</p> <p>The test could better meet the criteria by increasing the focus on the major work at grade 5.</p>
	 Smarter Balanced GOOD MATCH	<p>Smarter Balanced provides a Good Match to the CCSSO Criteria for Content in Mathematics:</p> <p>The tests could better meet the criteria by increasing the focus on the major work in grade 5.</p>









MATHEMATICS CONTENT CRITERIA





Criteria	Program and Rating	Summary Statements
C.1 Focus: Tests focus strongly on the content most needed in each grade or course for success in later mathematics (i.e., Major Work).	 ACT Aspire WEAK MATCH	<p>On “focusing strongly on the content most needed for success in later mathematics,” the rating is Weak Match.</p> <p>ACT Aspire forms do not consistently place sufficient emphasis on the major work of the given grade, due in part to intentional test design, which requires inclusion of selected content from earlier and later grades. Still, many of the items coded to standards from lower grades do not address the major work of the relevant grade.</p>
	 MCAS LIMITED/UNEVEN MATCH	<p>On “focusing strongly on the content most needed for success in later mathematics,” the rating is Limited/Uneven Match.</p> <p>The grade 8 assessment is focused on the major work of the grade. The grade 5 assessment is significantly less focused on the major work of the grade than called for by the criterion, as it samples content across the full set of standards for the grade.</p>
	 PARCC GOOD MATCH	<p>On “focusing strongly on the content most needed for success in later mathematics,” the rating is Good Match.</p> <p>While the grade 8 tests focus strongly on the major work of the grade, the grade 5 tests fall short of the threshold required for the top rating.</p>
	 Smarter Balanced GOOD MATCH	<p>On “focusing strongly on the content most needed for success in later mathematics,” the rating is Good Match.</p> <p>While the grade 8 tests focus strongly on the major work of the grade, the grade 5 tests fall short of the threshold required for the top rating.</p>

MATHEMATICS CONTENT CRITERIA		
Criteria	Program and Rating	Summary Statements
*C.2 Concepts, procedures, and applications: Assessments place balanced emphasis on the measurement of conceptual understanding, fluency and procedural skill, and the application of mathematics.	ACT Aspire N/A	<p>The test forms contain items that assess conceptual understanding, procedural fluency, and application.</p> <p>However, forms contain few test problems in which procedural skill/fluency is the predominant emphasis and therefore may under-assess these skills.</p>
	MCAS N/A	<p>The test forms contain items that assess conceptual understanding, procedural fluency, and application.</p> <p>Particularly in fifth grade, however, test forms are overly focused on procedural skill/fluency and application relative to conceptual understanding.</p>
	PARCC N/A	<p>The test forms contain items that assess conceptual understanding, procedural fluency, and application. Some of the application problems, however, have shallow contexts that are not necessary or important to the problem.</p>
	Smarter Balanced N/A	<p>The test forms contain items that assess conceptual understanding, procedural fluency, and application. However, the percentage of each problem type varies across forms, with some forms having a relative excess of application problems and other forms an excess of procedural skill/fluency problems.</p>

*All four programs require, in their program documentation, the assessment of conceptual understanding, procedural skill/fluency, and application, although most do not clearly distinguish between procedural skill/fluency and conceptual understanding. Also, specific balance across these three types is not required. Due to variation across reviewers in how this criterion was understood and implemented, final ratings could not be determined with confidence. Therefore, only qualitative observations are provided.

MATHEMATICS DEPTH CRITERIA OVERVIEW		
Criteria	Program and Rating	Summary Statements
II. Assesses the depth that reflects the demands of College and Career Readiness.	 ACT Aspire GOOD MATCH	<p>ACT Aspire provides a Good Match to the CCSSO Criteria for Depth in Mathematics. The items are well crafted and clear, with only rare instances of minor editorial issues.</p> <p>The ACT Aspire tests include proportionately more items at high levels of cognitive demand (DOK 3) than the standards reflect, and proportionately fewer at the lowest level (DOK 1). This finding is both a strength, in terms of promoting strong skills, and a weakness, in terms of ensuring adequate assessment of the full range of cognitive demand within the standards.</p> <p>While technically meeting the criterion for use of multiple item types, the range is nonetheless limited, with the large majority comprising multiple-choice items.</p> <p>The program would better meet the criteria for Depth by including a wider variety of item types and relying less on traditional multiple-choice items.</p>
	 MCAS EXCELLENT MATCH	<p>MCAS provides an Excellent Match to the CCSSO Criteria for Depth in Mathematics. The assessment uses high-quality items and a variety of item types. The range of cognitive demand reflects that of the standards of the grade. While the program does not code test items to math practices, mathematical practices are nonetheless incorporated within items.</p> <p>The program might consider coding items to the mathematical practices and making explicit the connections between specific practices and specific content standards.</p>
	 PARCC GOOD MATCH	<p>PARCC provides a Good Match to the CCSSO Criteria for Depth in Mathematics. The tests include items with a range of cognitive demand, but at grade 8 that distribution contains a higher percentage of items at the higher levels (DOK 2 and 3) and significantly fewer items at the lowest level (DOK 1). This finding is both a strength, in terms of promoting strong skills, and a weakness, in terms of ensuring adequate assessment of the full range of cognitive demand within the standards.</p> <p>The tests include a variety of item types that are largely of high quality. However, a range of problems (from minor to severe) surfaced relative to editorial accuracy and, to a lesser degree, technical quality.</p> <p>The program could better meet the Depth criteria by ensuring that all items meet high editorial and technical standards and by ensuring that the distribution of cognitive demand on the assessments provides sufficient information across the range.</p>
	 Smarter Balanced GOOD MATCH	<p>Smarter Balanced provides a Good Match to the CCSSO Criteria for Depth in Mathematics. The exam includes a range of cognitive demand that fairly represents the standards at each grade level.</p> <p>The tests have a strong variety of item types including those that make effective use of technology. However, a range of problems (from minor to severe) surfaced relative to editorial accuracy and, to a lesser degree, technical quality. A wide variety of item types appear on each form, and important skills are assessed with multiple items, as is sound practice. Yet, individual forms sometimes contained two or three items measuring the same skill that were nearly identical, with only the numerical values changed in the item stem and a different set of answer choices. Such near-duplication may not impact the accuracy of the score, but a greater variety of question stems/scenarios is desirable.</p> <p>The program could better meet the Depth criteria by ensuring that all items meet high editorial and technical standards and that a given student is not presented with two or more virtually identical problems.</p>

MATHEMATICS DEPTH CRITERIA		
Criteria	Program and Rating	Summary Statements
C.3 Connecting practice to content: Test questions meaningfully connect mathematical practices and processes with mathematical content.	 ACT Aspire EXCELLENT MATCH	<p>On “connecting practice to content,” the rating is Excellent Match.</p> <p>All items that are coded to mathematics practices are also coded to one or more content standard.</p>
	 MCAS EXCELLENT MATCH	<p>On “connecting practice to content,” the rating is Excellent Match.</p> <p>Although no items are coded to mathematical practices, the practices were nonetheless assessed within items that also assessed content.</p>
	 PARCC EXCELLENT MATCH	<p>On “connecting practice to content,” the rating is Excellent Match.</p> <p>All items that are coded to mathematics practices are also coded to one or more content standard.</p>
	 Smarter Balanced EXCELLENT MATCH	<p>On “connecting practice to content,” the rating is Excellent Match.</p> <p>All items that are coded to mathematics practices are also coded to one or more content standard.</p>
C.4 Cognitive demand: The distribution of cognitive demand for each grade level is sufficient to assess the depth and complexity of the standards.	 ACT Aspire LIMITED/UNEVEN MATCH	<p>On “requiring a range of cognitive demand,” the rating is Limited/Uneven Match.</p> <p>At both grades 5 and grade 8, the test forms include significantly more items of high cognitive demand (DOK 3) than reflected in the standards, and proportionately fewer at the lowest level (DOK 1). While these items increase the challenge of the tests, standards that call for the lowest level of cognitive demand (DOK 1) may be under-assessed.</p>
	 MCAS EXCELLENT MATCH	<p>On “requiring a range of cognitive demand,” the rating is Excellent Match.</p> <p>At each grade level, the distribution of cognitive demand closely reflects that of the standards.</p>
	 PARCC GOOD MATCH	<p>On “requiring a range of cognitive demand,” the rating is Good Match.</p> <p>The distribution of cognitive demand of items reflects that of the standards very well at grade 5, while the grade 8 test includes proportionately more items at the higher levels of cognitive demand (DOK 2 and 3). As a result, Grade 8 standards that call for the lowest level of cognitive demand may be under-assessed.</p>
	 Smarter Balanced GOOD MATCH	<p>On “requiring a range of cognitive demand,” the rating is Good Match.</p> <p>The distribution of cognitive demand of items reflects that of the standards very well at grade 5. At grade 8, the test includes proportionately fewer items at the lowest levels of cognitive demand (DOK 1) than in the standards, and proportionately more items at the mid-level of cognitive demand (DOK 2). As a result, grade 8 standards that call for the lowest level of cognitive demand may be under-assessed.</p>

MATHEMATICS DEPTH CRITERIA		
Criteria	Program and Rating	Summary Statements
C.5 High-quality items and variety of item types: Items are of high technical and editorial quality, are aligned to the standards, and each test form includes at least two items types including at least one that requires students to generate rather than select a response.	 ACT Aspire EXCELLENT MATCH	<p>On “ensuring high-quality items and a variety of item types,” the rating is Excellent Match.</p> <p>The program uses multiple item types, including constructed response on the Extended Task items. These items, although they carry high point values, are limited in number; the rest of the items are predominantly multiple-choice.</p> <p>The large majority of items are of high technical and editorial quality, with only very minor issues of editing, language, or accuracy. At the grade 8 level, some items appear to be susceptible to simplification by use of calculators, which are allowed on all items at grade 8, in contrast to the other programs that allow them on a restricted set of items.</p>
	 MCAS EXCELLENT MATCH	<p>On “ensuring high-quality items and a variety of item types,” the rating is Excellent Match.</p> <p>Both grade 5 and grade 8 forms include multiple item types, including constructed-response. The items are of high technical and editorial quality, with very minor issues of editing, language, and accuracy at grade 8.</p>
	 PARCC GOOD MATCH	<p>On “ensuring high-quality items and a variety of item types,” the rating is Good Match.</p> <p>The program includes a wide variety of item types, including several that require student-constructed responses. However, there are a number of items with quality issues, mostly minor editorial but sometimes mathematical.</p>
	 Smarter Balanced LIMITED/UNEVEN MATCH	<p>On “ensuring high-quality items and a variety of item types,” the rating is Limited/Uneven Match.</p> <p>The program includes a wide variety of item types, many of which make effective use of technology.</p> <p>The program could be improved by ensuring that virtually identical items are not presented to individual students. Further, a good deal of variability across forms and grades is observed, with some forms fully meeting the item quality criterion and others only partially meeting it. Issues exist with the editorial quality and mathematical accuracy of individual items, most of which are minor, but some of which could impact assessment of the targeted skill, resulting in a rating of Limited/Uneven.</p>

Appendix G

Testing Program Responses to Study and Descriptions of Test Changes for 2015–2016

ACT Aspire

Response to Report

One primary purpose of this study was to identify areas for improvement in each of the four evaluated testing programs. The study's findings include insights that promise to advance the industry standards for assessment quality. ACT Aspire is taking this opportunity to listen to the findings and implement adjustments to the assessment.

There are also aspects of the assessment that we have identified as areas for change through our own internal analyses, and we have already made design adjustments that will improve ACT Aspire in 2015–16. Both ACT Aspire's response to the study and the 2015–16 design adjustments will be discussed here.

Response to Study Findings about ACT Aspire

ACT Aspire is planning changes to two key elements for which the study found limited alignment with the CCSSO Criteria in English Language Arts:

1 Writing

- ◆ Although the ACT Aspire Writing test was intentionally designed to have writing tasks that do not contain the heavy reading load of “writing to sources” tasks, we are currently exploring updated designs that would supplement the current items with tasks that measure these valuable literacy skills. These tasks would also improve coverage of the “Assessing research and inquiry” criterion in the CCSSO framework.

2 Reading

- ◆ In response to the findings about distribution of Depth of Knowledge (DOK), ACT Aspire has already increased the percentage of upper-level DOK items. This effort will build on changes already in effect for the 2016 assessments (DOK 3 items in grades 5 and 8 will increase from 31 percent in 2015 to 38 percent in 2016).
- ◆ ACT Aspire is adding new technology-enhanced item designs that emphasize selecting evidence directly from the passage to support claims and interpretations. While some of these new TE items will be operational in 2016, ACT Aspire is continuing to explore new ways to assess student use of evidence from texts.

ACT Aspire would also like to make a clarification about terminology used to classify texts on the Reading test:

- ◆ The study findings indicate that while ACT Aspire is a Good Match in Depth, the tests should have “additional literary narrative text, as opposed to literary informational texts.” The study’s ELA/Literacy panel has made a different interpretation than ACT Aspire of CCSSO criterion B.1 that this finding refers to. Differences of interpretation around genre definitions are understandable, but it is important to note the effects on the study outcomes. The CCSSO criteria B.1 refers to texts that are “balanced across literary and informational text types and genres” and does not specify a balance of fiction and nonfiction in the literary category. The study’s panel interprets “literary” text types as only including literary narrative *fiction*. ACT Aspire, however, interprets “literary” to include both literary fiction and literary nonfiction passages that have a narrative structure. In accordance with B.1 (“In all grades, informational texts are primarily expository rather than narrative in structure”), ACT Aspire does not include texts that have a primarily narrative structure in the informational category. Aspire’s interpretation of criterion B.1 results in a stronger match to the specified balance of text types.

In math, ACT Aspire is enacting changes in response to three of the study findings:

- 1 Range of item types
 - ◆ The report distinguished levels of use of multiple-choice items of 50 percent and 75 percent and preferred 50 percent or less. With regard to technology-enhanced items, the report also recommended “using them strategically (read sparingly)” in order to use resources wisely. ACT Aspire will apply its research to make high-quality items of all types and look to expand what is possible in directions that involve technology.
- 2 Content focus
 - ◆ The report recommended “an increase in the number of items focused on the major work of the grade.” We will be working to increase this focus and gathering data to understand the balance in terms of promoting college and career readiness.
- 3 Depth of Knowledge
 - ◆ The study’s review panel recommended “the addition of more items at grade 8 that assess standards at DOK 1.” ACT Aspire currently has a plan in place that will increase the number of DOK 1 questions.

2015–16 Test Program Changes

In an effort to continuously improve ACT Aspire, we have already made adjustments in the following three categories for 2015–16:

- 1 Timing Adjustments – Based on customer feedback and in order to allow all students a better opportunity to show what they know and can do, we will be adjusting the time per test by five to ten minutes (Writing will not change). (See Tables G-1 and G-2 for more information on timing and point adjustments by grade and category.)
- 2 Adjustments to English Test – Adding six multiple-choice items for grades 3, 4, and 5.
- 3 Adjustments to Math Test – Adding six multiple choice items for grades 3, 4, and 5; removing one constructed-response (CR) item from grades 3, 4, and 5.

TABLE G 1

Timing Adjustments

ACT Aspire Summative Testing Time Adjustments (in minutes)								
Grade	English (Current)	English (New)	Math (Current)	Math (New)	Reading (Current)	Reading (New)	Science (Current)	Science (New)
3	30	40	55	65	60	65	55	60
4	30	40	55	65	60	65	55	60
5	30	40	55	65	60	65	55	60
6	35	40	60	70	60	65	55	60
7	35	40	60	70	60	65	55	60
8	35	40	65	75	60	65	55	60
*EHS	40	45	65	75	60	65	55	60

*Early High School

TABLE G 2

English and Mathematics: Number of Points by Reporting Category

GRADE										
	3 (new)	3 (old)	4 (new)	4 (old)	5 (new)	5 (old)	6	7	8	EHS
Reporting Category	# of Points									
English										
Production of Writing	12–14	9–11	8–10	6–8	8–10	6–8	11–13	9–11	9–11	12–14
Knowledge of Language			3–5	2–4	3–5	2–4	2–4	4–6	4–6	6–8
Conventions of Standard English	17–19	14–16	17–19	14–16	17–19	14–16	19–21	19–21	19–21	29–31
Total for English	31	25	31	25	31	25	35	35	35	50
Mathematics										
Number & Operations in Base 10		5–7	5–8	3–5	5–8	3–5	1–3	1–3	1–3	0–2
Number & Operations - Fractions	3–5	2–4	6–8	4–6	6–8	4–6	1–3	1–3	1–3	0–2
The Number System							3–5	3–5	2–4	1–3
Number & Quantity										1–3
Operations & Algebraic Thinking	6–8	3–5	4–6	3–5		3–5	1–3	1–3	0–2	0–2
Expressions & Equations							3–5	3–5	5–7	2–4
Ratios & Proportional Reasoning							3–5	3–5	0–2	1–3
Algebra										2–4
Functions									3–5	3–5
Measurement & Data (measurement)							0–2	0–2	1–3	1–3
Geometry		3–5		3–5	4–6	3–5	5–7	4–6	6–8	5–7
Measurement & Data	5–7	3–5		3–5		3–5				
Measurement & Data (data)							0–2	1–3	1–3	1–3
Statistics & Probability							3–5	3–5	4–6	4–7
Justification & Explanation	12	16	12	16	12	16	16	16	20	20
Total for Mathematics	39	37	39	37	39	37	46	46	53	53

MCAS

Response to Report

Our goal as a Commonwealth is to ensure that every Massachusetts student is prepared to succeed in postsecondary education and compete in the global economy. We have been administering annual assessments in Massachusetts since 1998 as our way of holding ourselves accountable for our progress toward this goal. The Massachusetts Comprehensive Assessment System (MCAS) tests are generally considered the gold standard of state assessments. They hold students to high expectations—in most cases, equivalent to the proficiency standard on the National Assessment of Educational Progress (NAEP)—and use a variety of question formats to ensure that we assess the full range of student abilities. Over the years we have refined the assessments to adapt to changes in the curriculum frameworks, most notably the incorporation of the Common Core State Standards into our 2010 frameworks, and to improve the quality of the assessment over time.

Our students and educators have accomplished incredible things under this system. Massachusetts' NAEP scores have moved from middle of the pack to leading the nation, and our students have scored well on international assessments. We have also made substantial progress toward closing the proficiency gaps between student subgroups, and we have dramatically reduced our dropout rate and increased our cohort graduation rate. That success would not have been possible without a high-quality assessment providing feedback on student, school, district, and state achievement and progress.

The Massachusetts Comprehensive Assessment Systems was a terrific twentieth-century assessment—but it has reached a point of diminishing returns. In 2015, MCAS was administered for the eighteenth year. We have a better understanding now than we did a decade or two ago about learning progression in mathematics, text complexity and the interplay of reading and writing, and the academic expectations of higher education and employers. And we now know that nearly one-third of our public high school students who go on to enroll in Massachusetts public colleges take at least one remedial course in their first semester, suggesting that the curriculum and assessments they have experienced have not adequately prepared them for the world beyond high school. Indeed, MCAS was never designed to be an indicator of college and career readiness. We joined the Partnership for Assessment of Readiness for College and Careers (PARCC) consortium specifically in order to partner with other states in developing an assessment that is more closely aligned to these expectations.

Thus, we were not surprised by this report's conclusion that the MCAS does not always measure well what's most important today. This report also confirms that in many ways, PARCC sets a higher bar than MCAS for student performance. This is particularly true as students move up the grades into middle and high school. This higher bar is not simply about being harder: PARCC provides more opportunities for critical thinking, applying knowledge, research, and making connections between reading and writing. More and more schools have upgraded curriculum and instruction to align with our 2010 frameworks. While we adjusted MCAS to test those frameworks, PARCC was built around them. Classroom instruction is now increasingly focused on the knowledge and skills in the frameworks, rather than how to pass a test.

We are proud of what we have accomplished in Massachusetts in the nearly two decades that we have been administering the MCAS. Now that we have the benefit of that experience and have revised our curriculum frameworks to reflect our upgraded learning expectations, it is time to upgrade our assessments too. Our state Board of Elementary and Secondary Education voted in November 2015 to do exactly that.

2015–16 Test Program Changes

Over the next few years, we will transition to a new statewide assessment system that will take much of what this report identifies as the strengths of PARCC—high-quality content aligned strongly to college and career ready standards—and combine it with elements of MCAS in the context of a Massachusetts-specific governance system that will allow us to set our own policies on test content, administration, and reporting. With this approach, we will continue to benefit from a high-quality, next-generation assessment while ensuring that the test will reflect the Commonwealth's unique needs and concerns. Most importantly, our students will be better prepared for success after high school—our ultimate goal.

PARCC

Response to Report

PARCC would like to thank the study authors and review panelists for a comprehensive, strong study. There are two areas where we would like to present a few additional comments.

ELA/Literacy Content Rating

Panel Recommendation:

The tests could be strengthened by the addition of research tasks that require students to use two or more sources and, as technologies allow, a listening and speaking component.

PARCC Response:

Every PARCC assessment in grades 3–8 requires students to complete a research simulation task where the student reads two or three sources and must integrate or synthesize the ideas in a written essay. All students also read two literary texts and write a literary analysis (literary research) essay.

The PARCC assessment measures many aspects that are key to the Speaking and Listening standards. PARCC uses multimedia texts to measure comprehension for all students taking its tests online (providing students with opportunities to demonstrate strengths and needs in comprehending audio and audiovisual texts). The CCSS build coherence across the ELA strands and identify similar skills built into both the reading comprehension standards (standards RI.7 and RL.7) and the listening standards. PARCC chose to report students' speaking and listening performance in relation to the reading standards.

The PARCC assessment system includes a robust set of Speaking and Listening tools. All schools administering PARCC in 2015–2016 have access to a comprehensive set of formative assessments and instructional tools to support educators, parents, and students in better understanding students' strengths and needs in speaking and listening. Further information about the PARCC Speaking and Listening tools can be found on PARCC's Partnership Resource Center: <https://prc.parcconline.org/library/speaking-and-listening-overview>.

Cognitive Demand

Panel Recommendation:

The program could better meet the Depth criteria by ensuring that the distribution of cognitive demand on the assessments provides sufficient information across the range.

PARCC Response:

It is important to note that students who meet Level 1/Level 2 Depth of Knowledge (DOK) for items situated at higher DOK levels are given partial credit points for demonstrating skills that require lower cognitive complexity. Reviewers did not consider the possibility that scoring, rather than adding more Level 1 or Level 2 items, could allow for the balance of item complexities. For more information on the PARCC scoring rubrics and to view released items, visit: <https://prc.parcconline.org/assessments/parcc-released-items>.

The PARCC assessment uses a cognitive complexity framework that was developed by the PARCC consortium to more accurately reflect the demands of the CCSS. This framework received recognition from AERA (2014 Outstanding Contribution to Practice in Cognition and Assessment award). An article detailing the innovations of this framework and potential next steps in research around cognitive complexity has been published in a new book titled *The Next Generation of Testing: Common Core Standards, Smarter-Balanced*.¹⁰⁶

¹⁰⁶ H. Jiao and R. Lissitz, eds., *The Next Generation of Testing: Common Core Standards, Smarter-Balanced, PARCC, and the Nationwide Testing Movement* (Charlotte, NC: Information Age Publishing, Inc., 2016).

2015–16 Test Program Changes

In May 2015, the chief state school officers from the PARCC states unanimously voted to streamline the assessment. They accomplished this goal while retaining all the key elements of the test—a strong commitment to quality and reliability, measurement of the full range of the standards, and the ability to get results back to teachers and parents quickly, so that they can help meet the needs of students for the coming school year. The following changes to the test design will be instituted in the 2015–16 school year:

- ◆ The two testing windows (the performance-based and end-of-year components) in mathematics and English language arts/literacy (which includes reading and writing) will be consolidated into one. The single testing window will simplify administration of the test for states and schools. Schools will have up to thirty school days to administer the test, and the testing window will extend from roughly the 75 percent mark to the 90 percent mark of instructional time. Most schools will complete testing in one to two weeks during that window.

The testing time for students will be reduced by about ninety minutes overall (sixty minutes in mathematics; thirty minutes in English language arts/literacy). The result will be that the total testing time for ELA/Literacy and mathematics will be approximately 8.5 hours at grades 3–5, 9.2 hours at grades 6–8, and 9.7 hours at grade 11. There will also be more uniformity of test unit times, allowing for easier scheduling in schools.

- ◆ Each PARCC assessment is administered in multiple sections, called units. The number of test units was reduced for all students, and includes three units in English language arts/literacy and three or four units in mathematics.

The testing time was shortened by reducing the number of score points and items in both subject areas. The tables below show a comparison of score points between the previous test design and the redesign.

TABLE G 3

Comparison of Score Points in the Previous ELA/Literacy Design and the Redesign

		Previous Two Administrations	Adopted Single Administration
Grade 3	Reading Points	64	58
	Writing Points	36	36
	Total Points	100	94
	Units	4	3
	Total Testing Time	4.75 hours	4.25 hours*
Grades 4–5	Reading Points	70	62
	Writing Points	36	36
	Total Points	106	98
	Units	4	3
	Total Testing Time	5.0 hours	4.5 hours**
Grades 6–11	Reading Points	94	76
	Writing Points	45	45
	Total Points	139	121
	Units	5	3
	Total Testing Time	5.75 hours	5.2 hours***

* Add 1.5 hours for field test unit

** Add 1.5 hours for field test unit

***Add 1.8 hours for field test unit

TABLE G 4

Comparison of Score Points in the Previous Mathematics Design and the Redesign

		Previous Two Windows	Adopted Single Administration
Grade 3-8	Short Items	56 pts	40 pts
	Reasoning Items	14 pts	14 pts
	Modeling Items	12 pts	12 pts
	Total Points	82 pts	66 pts
	Units	4 @ varies	4 @ 60 min.
	Total Time on Task	5 hours	4 hours
Algebra I, Geometry, Algebra II, and Integrated Math I, II, III	Short Items	65 pts	49 pts
	Reasoning Items	14 pts	14 pts
	Modeling Items	18 pts	18 pts
	Total Points	97 pts	81 pts
	Units	4 @ varies	3 @ 90 min.
	Total Time on Task	5.3 - 5.5 hours	4.5 hours

- ◆ Standalone field testing will be eliminated. As with all similar testing, field test items—items that could be used in future years—are embedded in each student's test. Because the performance tasks in English language arts/literacy are longer, a sampling of students had to take a standalone field test unit for these tasks in spring 2015. To further streamline the testing process for all schools, the PARCC field test will now be wrapped into the testing window. Each year, a small percentage of students will take an additional ELA/Literacy unit. Schools and classrooms selected in one year—per the process determined in their state—will in almost all cases not have to field test again for several years.
- ◆ The test design changes do not result in the loss of any performance tasks in English language arts/literacy (there are still three performance tasks). Additionally, there are now two or three text sets included in the units, depending on the grade level (one text set was removed for grades 6–11).

The test design changes do not result in the loss of any reasoning and modeling mathematics items, with the exception of Algebra II and Integrated Math III at the high school level. Short answer items were removed.

For more information, visit <http://parcconline.org/assessments/test-design/design-changes>.

Smarter Balanced

Response to Report

Tony Alpert, Executive Director

Luci Willits, Deputy Executive Director

December 9, 2015

Thank you to the Thomas B. Fordham Institute and HumRRO for its diligent work to evaluate the quality of the Smarter Balanced summative assessment and its alignment to the Common Core State Standards.

While this report focused on the end-of-year test, Smarter Balanced is more than a summative assessment: it's a system to improve teaching and learning. Our system includes optional and flexible interim assessments available throughout the year to help teachers monitor student progress, as well as a Digital Library with thousands of educator-approved classroom resources. Nearly 5,000 educators from across the country helped build the Smarter Balanced system. Smarter Balanced assessments are designed to be administered online and are customized for every child using built-in accessibility resources.

This report recognizes many of these strengths and gives Smarter Balanced an Excellent or Good Match in all but one subcategory. In addition, the report recognizes that Smarter Balanced is the only assessment that measures students' listening skills. We are proud of these ratings. We also recognize that there is always room for improvement. However, one of the greatest strengths of Smarter Balanced, the computer-adaptive feature of the summative assessment, is not addressed in this report. Because it is an adaptive test that is customized for each student, it is difficult to compare the Smarter Balanced summative assessment on an item-per-item basis to a fixed form test that is static. In addition, this study did not consider some other important features of the Smarter Balanced assessment, including the ability of states to work with the service provider of their choice. Finally, Smarter Balanced is arguably the most accessible large-scale assessment system that includes supports for over 90 percent of the consortium's English language learners' primary languages. Individually, these elements are historic; collectively, they are unprecedented.

It is important to note that due to the timing of this study, reviewers were not able to access all of the interactive features that are available to all students during a live test. For instance, reviewers did not interact with features such as highlighting text in passages and test questions, zooming in and out of test pages, making notes about a test question in the notepad, and using strikethrough for answer options. In addition, the study's version of the system did not provide some of the helpful built-in student tools, such as error messages when students use incorrect keys, the ability to mark items, or move forward in the test without answering all the questions on a page.

2015–16 Test Program Changes

For the 2015–2016 summative tests, Smarter Balanced members have the flexibility to determine whether classroom activities will be given prior to the administration of performance tasks.

The following table describes the differences between the system used in the study versus the one actually used for students:

TABLE G 5

Differences in Student and Study Interface

Students' Actual Version:	Study Version:
Calculators available in grades 6 through high school for items when not measuring computation.	No calculators.
Tutorials to show students how to use the available tools and to interact with all of the different types of items they might see on a test.	No tutorials.
Grade appropriate and item-specific English glossaries are available for mathematics and English language arts.	No glossaries.
For mathematics, grade appropriate item-specific translated glossaries are available in ten different languages, plus dialects.	No glossaries.
Error messages given to students when they try to enter characters that aren't allowed.	No error messages when equation editor and fill-in blank items are incorrectly populated.
Verdana Font (14 pt)	Times New Roman font (12 pt)
Formatted for best results according to student cognitive labs and field testing.	Format not consistent.
All the research-based tools available as appropriate to the content area and item (as shown in the practice tests).	Limited availability of tools. For example, notepad, underline, highlight, etc. were not available.

As part of the development process, Smarter Balanced collaborated with national experts and local teachers to determine how to best measure critical thinking and problem solving skills as part of college and career ready standards. For example, at times, it is most appropriate to ask students to solve engaging items within a real-world scenario; while at other times, presenting students with an equation to solve is a better way to measure student knowledge. This is reflected in the test blueprint.

It is important to note that with the adaptive test, Smarter Balanced can measure more complex skills for low- and high-performing students alike. In mathematics, Smarter Balanced chose to emphasize the more complex skill sets with the understanding that students must have the procedural knowledge to do well on the test. With English language arts, we will discuss the report's findings with our membership and consider changes.

Smarter Balanced is committed to including only high-quality questions on our tests. We were disappointed to see that reviewers found a handful of questions that needed improvement and received a rating of Limited/Uneven Match. Smarter Balanced has an extensive process for question development to ensure each test item

is extensively reviewed prior to being included on a student's test. Educators including national mathematical, English language arts, and accessibility experts write questions and review them for content accuracy as well as for any potential bias or lack of sensitivity. Questions that do not meet a very high standard are revised or are removed. However, we will use this study to improve our item development review processes. Immediately, Smarter Balanced will initiate a detailed review of the existing test questions based on the feedback from this report.

In addition to this positive review of Smarter Balanced, we were pleased to note that the National Network of State Teachers of the Year echoed complimentary feedback in their report as well. That review looked at many of the same questions as this review. The nation's best teachers said Smarter Balanced provides a better picture of student performance, is grade-level appropriate, and supports great teaching and learning throughout the year.

Thank you again for your review and for the opportunity to provide more context into the reviewers' findings.

Sincerely,

A handwritten signature in black ink, appearing to read "Anthony Alpert". The signature is fluid and cursive, with the first name "Anthony" written in a larger, more prominent script than the last name "Alpert".

Anthony Alpert
Executive Director
Smarter Balanced Assessment Consortium