



Technical Report

Spring 2018

Michigan Student Test of Educational Progress

(M-STEP)

TABLE OF CONTENTS

Executive Summary	9
Chapter 1: Background of Spring 2018 M-STEP Assessments	14
1.1 Background of M-STEP	14
1.2 Purpose and Design of ELA and Mathematics M-STEP with Respect to the Smarter Balanced Assessment	14
1.2.1 Background on Smarter Balanced	16
1.2.2 Test Blueprints	17
1.3 Purpose and Design of the Social Studies M-STEP	17
Chapter 2: Uses of Test Scores	18
2.1 Uses of Test Scores	18
2.2 Test-Level Scores	19
2.2.1 Scale Scores	19
2.2.2 Levels of Performance	20
2.2.3 Use of Test-Level Scores	20
2.3 Claim-level Sub-scores for ELA and Mathematics	20
Chapter 3: Test Design and Item Development	23
3.1 Overview	23
3.1.1 A Brief Description of Smarter Balanced Content Structure for ELA and Mathematics	23
3.1.2 Evidence-Centered Design in Constructing Smarter Balanced Assessments	24
3.1.3 A Brief Description of Content Structure for Social Studies	25
3.2 Test Blueprints	26
3.2.1 Test Specifications	26
3.2.2 Item Writer Training	27
3.2.3 Item Development	28
3.2.4 Graphics Creation	29
3.2.5 Item Review	30
3.2.6 Field-Testing	31
3.2.7 Range-Finding	31
3.2.8 Data Review	31

3.3	Operational Test Construction	32
3.3.1	ELA	32
3.3.2	Mathematics	33
3.3.3	Social Studies	34
3.3.4	Science	35
3.3.5	Accommodations	35
3.4	Sources of Items and Metadata	36
3.4.1	ELA and Mathematics	36
3.4.2	Social Studies	36
3.5	Import into DRC INSIGHT Test Engine	37
3.6	Psychometric Review During Assessment Construction	37
3.6.1	ELA and Mathematics	37
3.6.2	Social Studies	38
3.7	Item Types Included	40
3.8	Field-Test Selection and Administration	40
3.8.1	Field Test Item Selection	40
3.8.2	Field Test Administration	41
3.9	Online Form Building and Rendering Process	41
3.9.1	Overview of Rendering Process	41
3.9.2	Form Preparation and Rendering in INSIGHT	44
3.10	Paper/Pencil Form Building and Review Process	44
3.11	Summary	45
Chapter 4:	Test Administration Plan	46
4.1	Universal Tools, Designated Supports, and Accommodations	46
4.2	Online Accommodations	48
4.3	Paper/Pencil Universal Tools, Designated Supports, and Accommodations	52
4.4	Online Test Platform	54
4.5	Test Administration Training	56
4.6	Test Security	57
4.6.1	Prevention	58
4.6.2	Detection	59
4.6.3	Investigation and Remediation	60

4.7	Summary of M-STEP Administration Best Practices	61
4.8	Test Materials	63
4.9	Summary	65
Chapter 5: Test Delivery and Administration		66
5.1	Online Administration Details	66
5.1.1	Online Administration Reports	67
5.1.2	Online User Manuals and Reference Documents.	67
5.2	Paper/Pencil Administration Details	70
5.3	OEAA Secure Site	72
5.4	eDIRECT	73
5.4.1	Michigan Users.	73
5.4.2	Administrative Functions	73
5.4.3	Online Testing Resources	74
5.5	Return Material Processing.	74
5.6	Testing Window and Length of Assessment	78
Chapter 6: Operational CAT.		79
6.1	Entry Point.	79
6.2	Theta Estimates and Standard Error of Measurement	80
6.3	Item Selection	83
6.3.1	Test Blueprint	83
6.3.2	Item Information Function	83
6.3.3	Passage Related Concerns	84
6.4	Test Navigation	85
6.5	Termination	86
6.6	Forced Submission	86
6.7	Summary of Simulation Results Evaluating the CAT Algorithm	87
6.7.1	Adherence to the Test Blueprint.	87
6.7.2	Controlling for Item Exposure.	98
6.8	Summary of Simulation Results for the Student Ability Estimates	101
6.8.1	Ability Estimates at the Extremes.	101
6.8.2	Standard Error of Measurement.	102
6.8.3	Statistical Measures of Bias	105
6.9	Summary	110

Chapter 7: Scoring	111
7.1 Online Scoring	111
7.1.1 Autoscoring	111
7.1.2 Multiple Choice Scoring	112
7.2 Handscoring	112
7.2.1 Security	112
7.2.2 Measurement Incorporated Reader and Team Leader Hiring	113
7.2.3 Preparation of Training Materials for M-STEP	114
7.2.4 Training and Qualifying Readers and Team Leaders	114
7.2.5 Virtual Scoring Center	115
7.2.6 Quality Control and Reliability of Scoring	116
7.2.7 Validity	119
7.2.8 Alerts	120
7.3 Summary	120
Chapter 8: Operational Data Analyses	121
8.1 Operational Analysis of ELA and Mathematics	121
8.1.1 CAT Item Pool Characteristics	121
8.1.2 Item Pool IRT Statistics	125
8.2 Operational CAT ELA and Mathematics Implementation	130
8.2.1 The Scale	130
8.2.2 Lowest and Highest Obtainable Scale Scores (LOSS and HOSS)	130
8.2.3 Item-Pattern Scoring	130
8.2.4 Blueprint Fidelity Summary	131
8.3 Operational Analysis of Social Studies	131
8.3.1 CTT Statistics Social Studies	131
8.3.2 IRT Statistics: Social Studies	133
8.3.3 Item Calibration for Social Studies	134
8.3.4 Anchor Item Evaluation for Social Studies	135
8.3.5 Evidence of Model Fit for Social Studies	138
8.3.6 Test Characteristic Curves (TCCs) and Conversion Tables	138
8.3.7 IRT Statistics	142
8.4 Summary	143

Chapter 9: Test Results	144
9.1 Test Completion	144
9.2 Current Administration Data Scale Score Summaries	144
9.3 Description of Reports	145
9.3.1 Student-Level Data Reports and Data Files	146
9.3.2 Aggregate Data Reports and Data Files	147
9.4 Interpretive Guides	149
9.5 Summary	149
Chapter 10: Performance-Level Setting	170
10.1 Cut Score Validation for English Language Arts and Mathematics	171
10.2 Statistical Articulation for Social Studies	172
10.3 Scale Scores	172
10.4 Cut Scores	173
10.5 Claim Cut Scores	174
10.6 Performance Level Descriptors	174
10.7 Summary	176
Chapter 11: Fairness	177
11.1 Minimizing Bias through Careful Test Development	178
11.1.1 ELA and Mathematics	179
11.1.2 Social Studies	179
11.2 Evaluating Bias through Differential Item Functioning (DIF)	180
11.3 DIF Statistics	181
11.3.1 Flagging Criteria and Results for ELA and Mathematics	184
11.3.2 Flagging Criteria and Results for Social Studies	187
11.4 Summary	189
Chapter 12: Reliability and Evidence of Construct-Related Validity	190
12.1 Reliability	190
12.1.1 Reliability and Standard Error of Measurement	191
12.1.2 Cronbach's Coefficient Alpha	192
12.1.3 Standard Error of Measurement	192
12.1.4 Marginal Reliability for ELA and Mathematics	193
12.1.5 Observed Reliability, SEM, and CSEM for ELA and Mathematics	193
12.1.6 Reliability of Claims for ELA and Mathematics	209

12.1.7 Reliability, SEM, and CSEM for Social Studies	212
12.2 Classification Accuracy and Consistency	215
12.2.1 ELA and Mathematics	215
12.2.2 Social Studies	218
12.3 Assumption of Unidimensionality	220
12.3.1 ELA and Mathematics	220
12.3.2 Social Studies	220
12.4 Validity Evidence	221
12.4.1 Minimization of Construct-Irrelevant Variance and Construct Underrepresentation	222
12.4.2 Evidence Based on Test Content	222
12.4.3 Evidence Based on Response Processes	223
12.4.4 Evidence Based on Internal Test Structure	223
12.4.5 Evidence Based on Relations to Other Variables	223
12.4.6 Correlations among Claims as Evidence of Convergent Validity	224
12.4.7 Divergent (Discriminant) Validity	227
12.4.8 Evaluation of Item Exposure for CAT ELA and Mathematics	227
12.4.9 Evidence Based on Consequences of Test Use	229
12.5 Summary	230
References	231
Appendix A: Test Administration Documents	236
Appendix A.1 Guide to State Assessments	236
Appendix A.2 M-STEP Test Administration Manual	259
Appendix A.3 M-STEP Test Administration Directions – Grade 5 Online	368
Appendix A.4 M-STEP Test Administration Directions – Grade 11 Online	433
Appendix A.5 M-STEP Test Administration Directions – Grade 8 Paper	474
Appendix B: Interpretive Guides	539
Appendix B.1 Interpretive Guide to M-STEP Reports	539
Appendix B.2 Interpretive Guide to MME Reports	580
Appendix B.3 M-STEP Student Data File Format	613
Appendix B.4 M-STEP Aggregate Data File Format	623
Appendix C: Target Score Report	626
Appendix D: M-STEP SGP and AGP Report	638

Appendix E: M-STEP Standards Validation	671
Appendix E-1. Validity Evidence for English Language Arts and Mathematics Cut Scores	671
Appendix E-2 Summary	671
Appendix E-3 Background	671
Appendix E-4 Standards Validation Methodology	672
Appendix E-5 Review of the Recommendations Made at the Standards Validation	674
Appendix E.6 References	685
Appendix F: Test Mode Comparison	686
Appendix G: Michigan Assessment System Participant Groups	703
Appendix G.1 Michigan Educators	703
Appendix G.2 Technical Advisory Committee	703
Appendix G.3 Michigan’s Division of Educator, Student, and School Supports (DESSS) Advisory Committee	704

Executive Summary

In June 2014, the Michigan legislature required the Michigan Department of Education (MDE) to develop a new assessment to administer in the spring of 2015. MDE, in conjunction with its testing vendors, worked to create a new assessment system called the Michigan Student Test of Educational Progress, or M-STEP. M-STEP is designed to effectively measure student mastery and growth in comparison to Michigan state standards. The assessment program is made up of three content areas: English Language Arts (ELA), mathematics, and social studies. ELA and mathematics are assessed in grades 3–8, and social studies is assessed in grades 5, 8, and 11. The designs for the ELA and mathematics assessments are based on assessments provided by the Smarter Balanced Assessment Consortium (Smarter Balanced) with Michigan-specific blueprints. The social studies assessments are designed specifically for Michigan. For spring 2018, the fourth content area of science assessments included only field test items aligned to the new Michigan state science standards. Since science was not administered operationally, the content area will not be addressed in this report.

This technical report addresses all phases of the testing cycle with the intention of providing evidence to support the validity of the M-STEP summative assessment program. All subsequent chapters of this report constitute evidence for the validity argument that M-STEP was developed with rigor, implemented with fidelity, and validated psychometrically.

E.1 ELA and Mathematics

MDE partners with Smarter Balanced, utilizing its ELA and mathematics test items, and Data Recognition Corporation (DRC) for the creation of M-STEP ELA and mathematics assessments in grades 3–8. Smarter Balanced member states retain flexibility regarding how to customize the system so that it may best be used as part of each state's approach to improving their local educational systems. Michigan has taken advantage of this in 2018 by not only customizing the M-STEP blueprints, but also adding passage-based writing (PBW) items to the ELA assessments to assess writing standards. This allowed Michigan to reduce the ELA and mathematics assessments to a legislatively mandated three-hour median testing time (combined) while retaining a writing task in all grades. As a hybrid of the Smarter Balanced assessments and Michigan selected PBW prompts for ELA, the M-STEP assessments are a key part of preparing all Michigan students for success in college and career readiness.

E.2 Science and Social Studies

M-STEP items for social studies are written and reviewed by Michigan educators. Teachers receive training in writing items for standardized assessment and write items testing specific Michigan content standards. Committees of educators review the items for content validity and potential bias issues. These reviews take place both before students see the items on a field test and using student data after they have been field tested. MDE staff and contractor content specialists provide guidance and review throughout this process, ultimately selecting the final items that appear on each test form to cover the full range of Michigan content standards.

E.3 MDE Office of Educational Assessment and Accountability (OEAA)

MDE's Office of Educational Assessment and Accountability (OEAA) has the responsibility of carrying out the requirements in state and federal statutes and rules for statewide assessments. The office oversees the planning, scheduling, and implementation of all major assessment activities and supervises MDE's testing contractors (i.e., DRC and Measurement Incorporated). In addition, the MDE staff from OEAA, in collaboration with outside contractors, conducts quality control activities for every aspect of the development and administration of the assessment program. For additional details for those groups, please refer to Appendix F. OEAA is also active in monitoring the security provisions of the assessment program.

E.4 Michigan Testing Contractors

DRC is MDE's item development contractor. DRC is responsible for providing test development content leads who work in conjunction with OEAA's content leads. DRC works with OEAA to develop test items. DRC is also a liaison between the Smarter Balanced item bank and OEAA test development staff. MDE administers online assessments to 99% of the students in grades 3–8 and 11. M-STEP is delivered through DRC's online test engine. DRC test development staff are responsible for rendering test items according to OEAA's style guide. Each item is reviewed by both DRC and OEAA content leads to ensure each student is presented with properly formatted test items that are clear and engaging and to ensure the content of each item replicates how the item appears in the item bank.

Measurement Incorporated is Michigan's contractor for paper/pencil materials, handscoring, and reporting. Measurement Incorporated is responsible for the development, distribution, and collection of all paper/pencil test materials as well as the monitoring of test security. Measurement Incorporated produces accommodated testing materials based on the test maps OEAA provides and in accordance with industry standards. Measurement Incorporated scores all the PBW prompts using Michigan-developed rubrics. Once testing is complete, Measurement Incorporated is responsible for developing and providing student results.

E.5 Michigan's Assessment System

Michigan's assessment system is a comprehensive, standards-based system. M-STEP is an accountability assessment, which means that it is used to evaluate school and district success in Michigan's accountability system. Other assessments exist for special populations of students, such as students with significant cognitive disabilities or English learners. All students in grades 3–8 and 11 are required to take Michigan's standards-based accountability assessments. Michigan's accountability assessments are listed in Table E-1 and are described in more detail in Section 3.3.

Table E-3. Claims for ELA/Literacy and Mathematics

Test	Content	Grades
M-STEP	Mathematics	3–8
M-STEP	ELA	3–8
M-STEP	Social Studies	5, 8, 11
SAT	Mathematics	11
SAT	ELA	11
MI-Access (alternate assessment)	Mathematics	3–8, 11
MI-Access (alternate assessment)	ELA	3–8, 11
MI-Access (alternate assessment)	Science	4, 7, 11
MI-Access (alternate assessment)	Social Studies	5, 8, 11
WIDA	Listening	1–12
WIDA	Reading	K–12
WIDA	Speaking	K–12
WIDA	Writing	1–12

E.6 Overview of This Report

Subsequent chapters of this technical report document the major activities of the testing cycle. This report provides comprehensive details that confirm that the processes and procedures applied in the M-STEP program adhere to appropriate professional standards and practices of educational assessment. Ultimately, this report serves to document evidence that valid inferences about Michigan student performance can be derived from the M-STEP assessments. Note that part of this report is intended to be utilized in tandem with the *Smarter Balanced 2014–15 Technical Report* (2016) and *Smarter Balanced 2017–18 Technical Report* (2018), while providing additional Michigan-specific validity and reliability information.

Each chapter of this report details the procedures and processes applied in M-STEP as well as the results. Each chapter also highlights the meaning and significance of the procedures, processes, and results in terms of validity and the relationship to the *Standards for Educational and Psychological Testing* (AERA, APA & NCME, 2014). Below is a brief overview of the contents of this report.

Chapter 1, “Background of Spring 2018 M-STEP Assessments,” describes the background and history of M-STEP.

Chapter 2, “Uses of Test Scores,” describes the use of the assessment scores and touches on the validity arguments the technical report intends to address.

Chapter 3, “Test Design and Item Development,” describes the involvement of Michigan educators in the item and assessment development process. As indicated, the assessment development process and the involvement of Michigan educators in that process formed an important part of the validity of M-STEP. The knowledge, expertise, and professional judgment

offered by Michigan educators ultimately ensured that the content of M-STEP formed an adequate and representative sample of appropriate content, and that content formed a legitimate basis upon which to derive valid conclusions about student performance. This part of the technical report thus addresses Standard 4.6 of the *Standards* (AERA, APA, & NCME, 2014, p. 87). It shows that the assessment design process, and the participation of Michigan educators in that process, provides a solid rationale for having confidence in the content and design of M-STEP as a tool from which to derive valid inferences about Michigan student performance. This chapter also addresses AERA, APA, and NCME (2014) *Standards* 3.1, 3.2, 4.0, 4.1, 4.2, 4.12, and 7.2. Chapters 4 and 5, “Test Administration Plan” and “Test Delivery and Administration,” describe the processes, procedures, and policies that guided the administration of M-STEP, including accommodations, security measures, and written procedures provided to assessment administrators and school personnel. These chapters address AERA, APA, and NCME (2014) *Standards* 4.15, 4.16, 6.1, 6.2, 6.3, 6.4, 6.6, 6.7, and 6.10.

Chapter 6, “Operational CAT,” supports Chapter 3 in showing how assessment specification documents, derived from earlier developmental activities, guided the final phases of assessment development and ultimately yielded the assessments administered to students. This chapter thus addresses AERA, APA, and NCME (2014) *Standards* 1.11, 3.1, 3.2, 3.5, 4.0, 4.6, 4.7, 4.8, 4.10, 4.12, 7.2, 8.4, 12.4, and 12.8.

Chapter 7, “Scoring,” explains the procedures used for scoring M-STEP autoscored items and handscored items. Chapter 7 adheres to AERA, APA, & NCME *Standards* 4.18, 4.20, 6.8, and 6.9.

Chapter 8, “Operational Data Analyses,” describes the data used for calibration and scaling. For content areas for which they are appropriate, raw-score results and a classical item analysis were provided and served as a foundation for subsequent analyses. This chapter also describes the calibration and scaling processes, procedures, and results. Some references to introductory and advanced discussions of Item Response Theory (IRT) are provided. This chapter thereby demonstrates adherence to AERA, APA, and NCME (2014) *Standards* 1.8, 5.2, 5.13, and 5.15.

Chapter 9, “Test Results,” presents scale-score results and achievement level information. Scale-score results provide a basic quantitative reference to student performance as derived through the IRT models that were applied. This chapter thus addresses AERA, APA, and NCME (2014) *Standards* 5.1, 6.10, 7.0, and 12.18.

Chapter 10, “Performance-Level Setting,” provides background on the standard-setting activities and functions to address *Standards* 5.21 and 5.22 of the *Standards* (AERA, APA, & NCME, 2014).

Chapter 11, “Fairness,” address validity evidence, specifically with respect to issues of bias. It demonstrates adherence to AERA, APA, and NCME (2014) *Standards* 3.1, 3.2, 3.3, 3.4, 3.5, and 3.6.

The first half of Chapter 12, “Reliability and Evidence of Construct-Related Validity,” demonstrates adherence to the AERA, APA, and NCME (2014) *Standards* through several analyses of the reliability of the 2018 M-STEP. Information on reliability/precision, standard error of measurement (SEM), conditional standard error of measurement (CSEM), and a detailed

examination of classification consistency and accuracy are provided. The first half of Chapter 12 thereby addresses AERA, APA, and NCME (2014) *Standards* 2.0, 2.3, 2.13, and 2.19.

The second half of Chapter 12 addresses validity evidence, including assessment content, response processes, issues of bias, dimensionality analysis, relations to other assessments, and consequences of assessment use. It demonstrates adherence to AERA, APA, and NCME (2014) *Standards* 3.16 and 4.3. This chapter ends with a section addressing the development of validity arguments for M-STEP.

MDE and its testing vendors have maintained an unwavering focus on the gathering of validity evidence in support of M-STEP throughout the development, administration, analysis, and reporting of the 2018 M-STEP administration.

Chapter 1: Background of Spring 2018 M-STEP Assessments

1.1 Background of M-STEP

The Michigan Department of Education (MDE), partnering with Smarter Balanced Assessment Consortium (Smarter Balanced), utilizes the ELA and mathematics test items from Smarter Balanced and the passage-based writing (PBW) prompts from Data Recognition Corporation (DRC) for the creation of M-STEP ELA and mathematics assessments. MDE uses test items written by Michigan educators for the M-STEP social studies assessments. MDE also partners with DRC for all online delivery, item development, and some psychometric work for the program; and with Measurement Incorporated for the paper/pencil and reporting portions of the program.

In the spring 2018 administration of M-STEP, 99% of Michigan students took M-STEP online. Paper/pencil tests were available for accommodated testing for individual students and for MDE-approved schools that were unable to test online.

1.2 Purpose and Design of ELA and Mathematics M-STEP with Respect to the Smarter Balanced Assessment

Summative assessments measure students' progress toward college and career readiness in ELA and mathematics. These assessments are given at the end of the school year as a computer adaptive test (CAT).

Page ix of the *Smarter Balanced 2017–2018 Technical Report* (2018) details the purposes of the Smarter Balanced summative assessments. Represented in part for this report, the “assessments are to provide valid, reliable, and fair information about” the following:

- students' ELA and mathematics achievement with respect to those Common Core State *Standards* (CCSS) measured by the ELA and mathematics summative assessments,
- whether students prior to grade 11 have demonstrated sufficient academic proficiency in ELA and mathematics to be on track for achieving college readiness,
- students' annual progress toward college and career readiness in ELA and mathematics,
- how instruction can be improved at the classroom, school, district and state levels,
- students' ELA and mathematics proficiencies for federal accountability purposes and potentially for state and local accountability systems, and
- students' achievement in ELA and mathematics that is equitable for all students and subgroups of students.

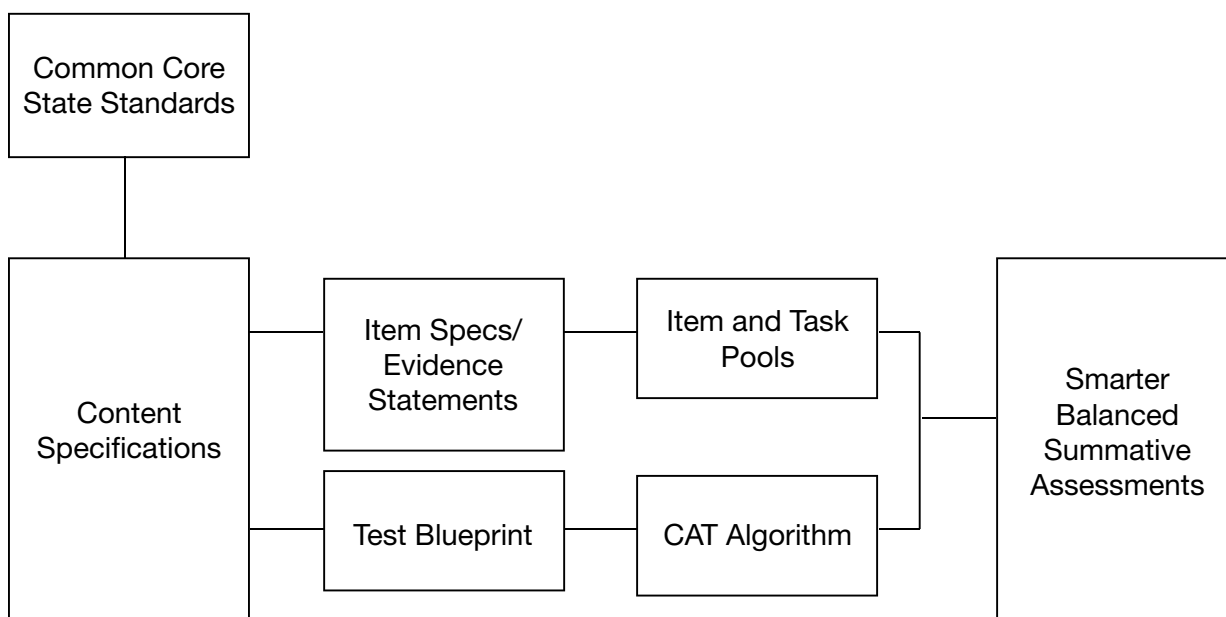
As stated on page 4-4 of the *Smarter Balanced 2017–2018 Technical Report* (2018) summative assessment scores will

- accurately describe both student achievement (i.e., how much students know at the end of the year) and student growth (i.e., how much students have improved since the previous year) to inform program evaluation and school, district, and state accountability systems.
- include writing at every grade and ask students to solve multistep, real-world problems in mathematics.
- capitalize on the strengths of CAT (i.e., efficient and precise measurement with a quick turnaround of results).
- provide valid, reliable, and fair measures of students' progress toward, and attainment of, the knowledge and skills required to be college- and career-ready.
- measure the breadth and depth of the CCSS across the full spectrum of student ability by incorporating a variety of item types (including items and tasks scored by expert raters) that are supported by a comprehensive set of accessibility resources.

The Smarter Balanced assessment system is a valid, fair, and reliable approach to student assessment that provides educators, students, and parents with meaningful results and actionable data to help students succeed.

In developing and maintaining a system of assessments, Smarter Balanced ensures that the assessments' measurement properties reflect industry standards for content, rigor, and performance. A key step in this direction is to ensure that the Smarter Balanced assessments are aligned with the CCSS, which Michigan adopted in 2014. Figure 1-1 (originally from the *Smarter Balanced 2017–2018 Technical Report*, 2018, p. 4-2), shows the components of the assessment.

Figure 1-1. Components of the Smarter Balanced Assessment Design

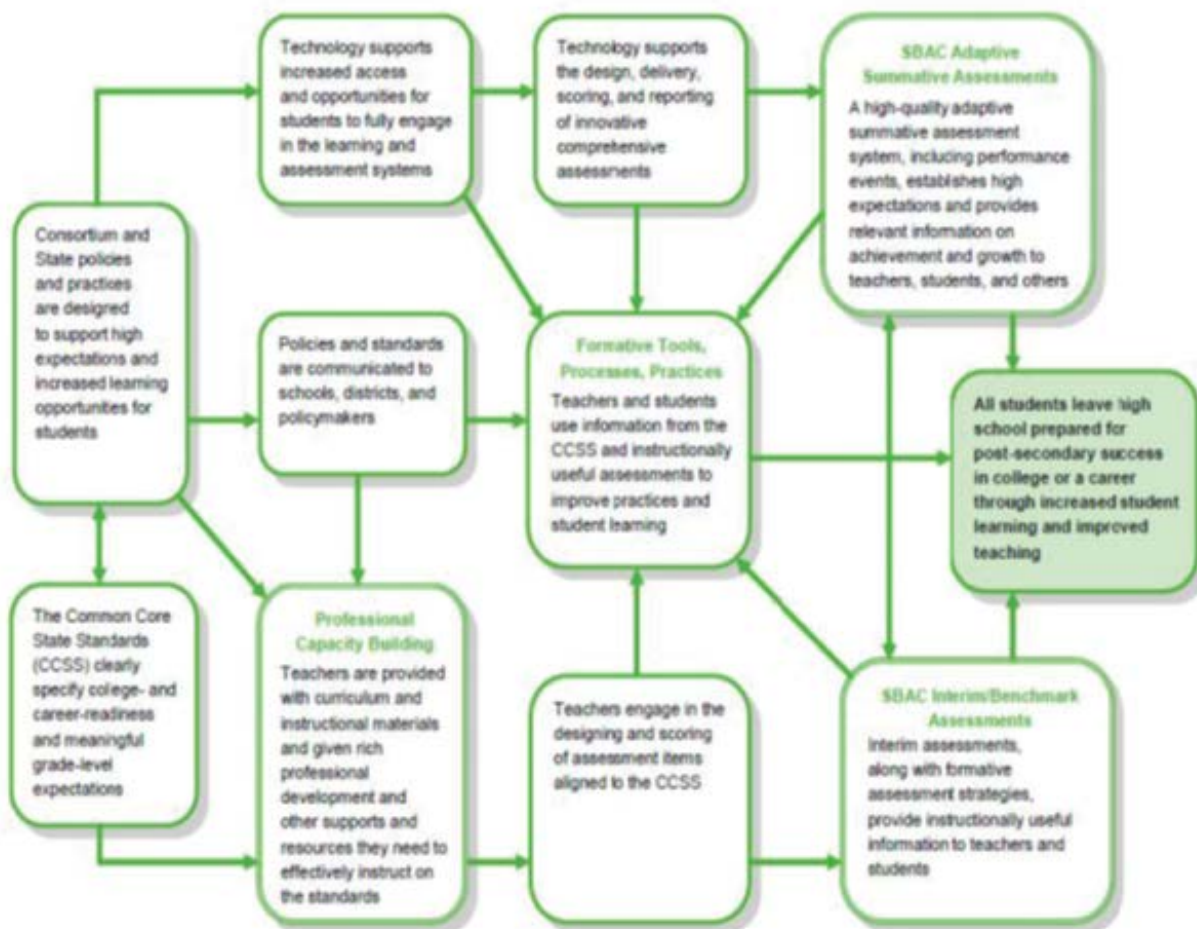


1.2.1 Background on Smarter Balanced

Smarter Balanced supports the development and implementation of learning and assessment systems to reshape education in member states in order to improve student outcomes. Through expanded use of technology and targeted professional development, the Consortium's Theory of Action calls for the integration of learning and assessment systems, leading to more informed decision-making and higher-quality instruction and ultimately increasing the number of students who are well prepared for college and careers.

The ultimate goal of the Smarter Balanced assessment system is to ensure that all students leave high school prepared for postsecondary success in college or a career through increased student learning and improved teaching. This approach suggests that enhanced learning will result from high-quality assessments that support ongoing improvements in instruction and learning. A quality assessment system strategically balances summative, interim, and formative components (Darling-Hammond & Pecheone, 2010). An assessment system must provide valid measurement across the full range of performance on common academic content, including assessment of deep disciplinary understanding and higher-order thinking skills increasingly demanded by a knowledge-based economy. Figure 1-2 presents an overview of the [Smarter Balanced Theory of Action](#) (2011, pg. 7).

Figure 1-2. Overview of Smarter Balanced Theory of Action



1.2.2 Test Blueprints

Part of the innovative aspect of the mathematics and ELA assessments is that the test blueprints sample the content domains using both a CAT engine and a PBW prompt. The test blueprints can be inspected to determine the contribution of the CAT and PBW components in a grade and content area toward the construct intended to be measured. Another aspect of the assessments is the provision of a variety of both autoscored and handscored item types. The contribution of these item types is specified in the Smarter Balanced test blueprints.

In February 2015, the governing members of Smarter Balanced adopted blueprints for the summative assessments of ELA and mathematics for grades 3–8 (Smarter Balanced, 2015a; Smarter Balanced, 2015b). These blueprints were fully implemented in the 2014–15 school year and were in effect in the 2017–18 school year.

For the 2017–18 school year, Michigan slightly modified the Smarter Balanced blueprints for ELA and mathematics. To reduce testing time, the use of Performance Tasks was eliminated for both Mathematics and ELA. In ELA, the PBW prompt was added to assess the writing standards. The net result is that, while the blueprints were modified, all students will continue to receive a writing claim score. In mathematics, Michigan added items to Claims 2 and 4 to address any blueprint gaps caused by the removal of the PT items. Due to the drift from the original Smarter Balanced blueprint, it should be noted that Michigan conducted a standards validation in July 2018 to review the M-STEP cut scores and determine if any changes needed to be made. More information can be found in Chapter 10 and Appendix E.

1.3 Purpose and Design of the Social Studies M-STEP

The summative assessments determine students' progress toward college and career readiness in social studies. These are given at the end of the school year. These assessments are primarily delivered online (99% of Michigan students took the test online) with paper/pencil and accommodated options. The social studies assessments are fixed forms. The summative assessments accurately describe student achievement (i.e., how much students know at the end of the year) to inform program evaluation and school, district, and state accountability systems.

The blueprints for social studies contain no constructed-response items, leading to a quick turnaround of results.

The social studies blueprints are located in Chapter 3, Section 3.3.

Chapter 2: Uses of Test Scores

Validity is an overarching component of M-STEP. The following excerpt is from the *Standards for Educational and Psychological Testing* (the *Standards*) (AERA, APA, & NCME, 2014):

Ultimately, the validity of an intended interpretation of test scores relies on all the available evidence relevant to the technical quality of a testing system. Different components of validity evidence . . . include evidence of careful test construction; adequate score reliability; appropriate test administration and scoring; accurate score scaling, equating, and standard setting; and careful attention to fairness for all test takers, as appropriate to the test interpretation in question. (p. 22)

As stated in the *Standards*, the validity of a testing program hinges on the use of the test scores. Validity evidence that supports the uses of M-STEP scores is provided in this technical report. This chapter examines some possible uses of the test scores.

As the *Standards* note, “validation is the joint responsibility of the test developer and the test user” (AERA, APA, & NCME, 2014, p. 13). For ELA and mathematics, the Smarter Balanced Assessment Consortium (Smarter Balanced) does not control aspects of test administration and test use. The Smarter Balanced members deliver the test, score operational items, and provide reports. Members use Smarter Balanced test scores in their own accountability models. In the *Smarter Balanced 2014–15 Technical Report* (2016)¹ and the *Smarter Balanced 2017–18 Technical Report* (2018), guidelines for administration and use are documented. Please see Chapter 1 of the *Smarter Balanced 2017–18 Technical Report* for the complete validity argument related to ELA and mathematics, member documentation on specific test administration procedures, reporting, and use for the Smarter Balanced assessments.

The following chapters of this technical report provide additional evidence for these uses as well as technical support for some of the interpretations and uses of test scores. The information in Chapters 3 through 12 also provides a firm foundation that M-STEP measures what it is intended to measure. However, this technical report cannot anticipate all possible interpretations and uses of M-STEP scores. It is recommended that policy and program evaluation studies, in accordance with the *Standards*, be conducted to support some of the uses of the test scores.

2.1 Uses of Test Scores

The validity of a test score ultimately rests on how that test score is used. To understand whether a test score is being used properly, the purpose of the test must first be understood. The intended uses of M-STEP scores include

- identifying Michigan students’ strengths and weaknesses;
- communicating expectations for all students;
- evaluating school-, district-, and/or state-level programs; and

¹ <https://portal.smarterbalanced.org/library/en/v2.0/2014-15-technical-report.pdf>

- informing stakeholders (i.e., teachers, school administrators, district administrators, Michigan MDE staff members, parents, and the public) on progress toward meeting state academic performance standards and meeting the requirements of the state's accountability program.

This technical report refers to the use of the test-level scores (i.e., scale scores and performance levels), claim-level scores, and claim performance indicators².

2.2 Test-Level Scores

At the test level, an overall scale score that is based on student performance on the entire test is reported. In addition, an associated performance level is reported. These scores indicate, in varying ways, a student's performance in ELA, mathematics, or social studies. Test-level scores are reported at four reporting levels: the state, school district, school, and student.

Items on the ELA and mathematics test forms were developed by Smarter Balanced. Items on the braille and enlarged print ELA and mathematics forms were also developed by Smarter Balanced. Final pencil/paper and accommodated forms were created using the items developed by Smarter Balanced, but the item selections were finalized by MDE and DRC content development staff. For social studies, all items and test forms were developed by MDE test development staff.

The following sections discuss two types of test-level scores that are reported to indicate a student's performance on M-STEP: (1) the scale score, and (2) its associated level of performance.

2.2.1 Scale Scores

A scale score indicates a student's total performance for each content area on M-STEP. A "scale score" is a statistical conversion of the "raw score" (numbers of questions that are answered correctly and incorrectly) that takes into account differences in questions students might see on different versions of the test across years, forms, or adapted versions of the test. In other words, the scale score represents the student's level of performance, where higher scale scores indicate higher levels of performance on the test and lower scale scores indicate lower levels of performance. Scaling scores permits comparison of assessment results across different test administrations within a particular grade and content area.

Scale scores are not comparable across grade levels or content areas. Scores are scaled within grade levels, so even if the same numbers are used in different grades, it does not mean that the scales form a single "vertical scale." M-STEP is a standards-based test that assesses the standards for each grade, so a very high score on grade 4 standards does not provide a valid estimate of how that student performs on grade 5 standards.

Details of the development of M-STEP scale scores are described in Chapter 10, Section 10.3. The scale score is stable because it allows for students' scores to be reported on the

² Claim scores are only available for ELA and Math.

same scale regardless of which year the students took the assessment and which form of the assessments the student took. Schools can use scale scores to compare the performances of groups of students across years. These comparisons can then be used to assess the impact of changes or differences in instruction or curriculum. The scale scores can be used to determine whether students are demonstrating the same skill and ability across cohorts within a grade and content area.

2.2.2 Levels of Performance

A student's performance on M-STEP is reported in one of the four levels of performance: Not Proficient, Partially Proficient, Proficient, and Advanced. The cut scores for the ELA and mathematics performance levels were established by Smarter Balanced during the standard setting, which occurred in three phases: online panel, in-person workshop, and cross-grade review in October 2014. These cut scores were then evaluated and confirmed in the Michigan standards validation in July 2018 (see Appendix E). The cut scores for the social studies performance levels were established by MDE in August 2015.

M-STEP performance levels reflect the performance standards and abilities intended by the Michigan legislature, Michigan teachers, Michigan citizens, and MDE. Descriptions of each performance level in terms of what a student should know and be able to do are provided by MDE and are referenced in the [M-STEP & MME Performance Level Descriptors](#).³

2.2.3 Use of Test-Level Scores

M-STEP scale scores and performance levels provide summary evidence of student performance. Classroom teachers may use these scores as evidence of student performance in these content areas. At the aggregate level, district and school administrators may use this information for activities such as planning curriculum. The results presented in this technical report provide evidence that the scale scores are valid and reliable indicators of student performance.

2.3 Claim-level Sub-scores for ELA and Mathematics

Claim-level sub-scores are scores on important domain areas within each content area. In most cases, sub-scores correspond to claims, but in mathematics, Claims 2 and 4 are so intertwined that they are reported as a single sub-score. The claims and reporting categories (sub-scores) are primary structural elements in test blueprints and item development. Figures 2.2 through 2.15 from the *Smarter Balanced 2016–17 Technical Report* (2017) provide information on the claims or sub-score reporting categories for ELA and mathematics.

The claim-level performance indicators are reported for ELA and mathematics for each student. A student's performance on each of the ELA and mathematics claims is reported in one of three levels of performance: *Adequate progress*, *Attention may be needed*, and *Most at risk of falling behind*. Performance-level indicator designations are based on the standard error of measurement of the claim-level sub-score and the distance of the claim sub-score from the proficient cut score. If the proficient cut score falls within a 1.5 SEM error band, it is designated

³ https://www.michigan.gov/documents/mde/2015_M-STEP_and_MME_PL_Descriptors_504568_7.pdf

as “Attention may be needed.” If the Level 2/3 cut score is above the error band, the sub-score is designated as “Most at risk of falling behind;” if the cut score is below the error band, the claim level sub-score is “Adequate Progress.”

The purpose of reporting claim-level sub-scores on M-STEP is to show for each student the relationship between the overall performance being measured and the skills in each of the areas delimited by the claims in ELA and mathematics. Teachers may use these sub-scores for individual students as indicators of strengths and weaknesses, but they are best corroborated by other evidence, such as homework, class participation, diagnostic test scores, or observations. Chapter 12 of this technical report provides evidence of content validity and reliability that supports the use of the claim-level sub-scores. Chapter 12 of this technical report also provides evidence of construct validity that further supports the use of these sub-scores (for additional information see *Smarter Balanced 2017–18 Technical Report* (2018), p. 5-18)

Figure 2-1. ELA Claims

Claim #1—Reading

- Students can read closely and analytically to comprehend a range of increasingly complex literary and informational texts.

Claim #2—Writing

- Students can produce effective and well-grounded writing for a range of purposes and audiences.

Claim #3—Speaking and Listening

- Students can employ effective speaking and listening skills for a range of purposes and audiences. At this time, only listening is assessed.

Claim #4—Research

- Students can engage in research/inquiry to investigate topics and to analyze, integrate, and present information.

Figure 2-2. Mathematics Claims

Claim #1—Concepts and Procedures

- Students can explain and apply mathematical concepts and interpret and carry out mathematical procedures with precision and fluency.

Claim #2—Problem Solving/Claim #4-Modeling and Data Analysis

- Students can solve a range of complex well-posed problems in pure and applied mathematics, making productive use of knowledge and problem-solving strategies. Students can analyze complex real-world scenarios and can construct and use mathematical models to interpret and solve problems.
- Students can analyze complex real-world scenarios and can construct and use mathematical models to interpret and solve problems.

Claim #3—Communicating Reasoning

- Students can clearly and precisely construct viable arguments to support their own reasoning and to critique the reasoning of others.

Chapter 3: Test Design and Item Development

3.1 Overview

This chapter is particularly relevant to AERA, APA, & NCME (2014) *Standards* 4.0, 4.1, and 4.7. It also addresses *Standards* 3.1, 3.2, 3.9, 4.12, and 7.4, which will be discussed in pertinent sections of this chapter. *Standards* 4.0, 4.1, and 4.7 are from Chapter 4 of the AERA, APA, & NCME (2014) *Standards*, “Test Design and Development.” AERA, APA, & NCME (2014) Standard 4.0 states the following:

Tests and testing programs should be designed and developed in a way that supports the validity of interpretations of the test scores for their intended uses. Test developers and publishers should document steps taken during the design and development process to provide evidence of fairness, reliability, and validity for intended uses for individuals in the intended examinee population. (p. 85)

The purpose of this chapter is to document the test design and item development process used for M-STEP. This chapter describes steps taken to create M-STEP, from the development of test specifications to the selection of operational items.

3.1.1 A Brief Description of Smarter Balanced Content Structure for ELA and Mathematics

The Common Core State *Standards* (CCSS) are the content standards in ELA and mathematics that many states have adopted. Because the CCSS were not specifically developed for assessment, they contain extensive rationale descriptions and information concerning instruction. Therefore, by adopting previous practices used by many state programs, Smarter Balanced content experts produced content specifications in ELA and mathematics, which distill assessment-focused elements from the CCSS. The Smarter Balanced *Content Specifications for the Summative Assessment of the CCSS for English Language Arts/Literacy* (2015a) and *Content Specifications for the Summative Assessment of the CCSS for Mathematics* (2015b) were expressly created to guide the structure and content of assessment development. Within each of the two content areas in grades 3–8, there are four broad claims. Within each claim, there are a number of assessment targets. The claims in ELA and mathematics are given in Table 3–1 (from the *Smarter Balanced 2017–18 Technical Report* (2018), p. 5-18).

Table 3-1. Claims for ELA and Mathematics

Claim	ELA	Mathematics
1	Reading	Concepts and Procedures
2	Writing	Problem Solving
3	Speaking/Listening	Communicating Reasoning
4	Research	Modeling and Data Analysis

Currently, only the listening part of ELA Claim 3 is assessed. In mathematics, Claims 2 and 4 are reported together as a single sub-score, so there are only three reporting categories for mathematics but four claims.

Because of the breadth in coverage of the individual claims, targets within each claim were needed to define more specific performance expectations. The relationship between targets and CCSS elements is made explicit in the Smarter Balanced content specifications (2015a; 2015b).

The *Item and Task Specifications* (Smarter Balanced, 2015c) for ELA and mathematics provide guidance on how to translate the Smarter Balanced content specifications into assessment items. In addition, guidelines for bias and sensitivity issues, accessibility and accommodations, and style help item developers and reviewers ensure consistency and fairness across the item bank. The specifications and guidelines were reviewed by member states, school districts, higher education representatives, and other stakeholders. The item specifications describe the evidence to be elicited and provide sample task models to guide the development of items that measure student performance relative to the target.

The Smarter Balanced assessment blueprints found in the *Smarter Balanced 2017–18 Technical Report* (2019) describe the content of the ELA and mathematics summative assessments for grades 3–8 administered in the 2017–18 school year and how that content was assessed. The blueprints also describe the composition of the assessment and its scoring. Specific items administered to each student are uniquely determined based on an item-selection algorithm that includes content constraints that correspond to the test blueprint. Developed with broad input from member states, partners, and stakeholders, the summative test blueprints reflect the depth and breadth of the performance expectations of the CCSS. Smarter Balanced governing members adopted the preliminary test blueprints in 2012. The summative test blueprints that were subsequently developed contain refinements and revisions based on the analyses of the pilot and field tests.

3.1.2 Evidence-Centered Design in Constructing Smarter Balanced Assessments

The *Smarter Balanced 2017–18 Technical Report* (2018) discusses the concept of evidence-centered design:

Evidence-centered design (ECD) is an approach to the creation of educational assessments in terms of reasoning about evidence (arguments) concerning the intended constructs. The ECD process begins with identification of claims, or inferences, users want to make concerning student achievement. Evidence needed to support those claims is then specified, and finally, items/tasks capable of eliciting that information are designed (Mislevy, Steinberg, & Almond, 2003). Explicit attention is paid to the potential influence of unintended constructs. The ECD process accomplishes this in two ways. The first is by incorporating an overarching concept of assessment as an argument made from imperfect evidence. This argument makes explicit the claims (i.e., the inferences that one intends to make based on scores) and the nature of the evidence that supports those claims (Hansen & Mislevy, 2008; Mislevy & Haertel, 2006). The second is by distinguishing the activities and structures involved in the assessment enterprise to exemplify an assessment argument in operational processes. By making the underlying evidentiary argument more explicit,

the framework makes operational elements more amenable to examination, sharing, and refinement. Making the argument more explicit also helps designers meet diverse assessment needs caused by changing technological, social, and legal environments (Hansen & Mislevy, 2008; Zhang, Haertel, Javitz, Mislevy, Murray, & Wasson, 2009). The ECD process entails five types of activities. The layers focus in turn on the identification of the substantive domain to be assessed; the assessment argument; the structure of assessment elements such as tasks, rubrics, and psychometric models; the implementation of these elements; and the way they function in an operational assessment, as described below. For Smarter Balanced, a subset of the general ECD elements was used. (p. 4-4)

3.1.3 A Brief Description of Content Structure for Social Studies

M-STEP content in social studies is defined by the knowledge and skills identified in the Michigan state standards. Michigan state standards were approved by the Michigan State Board of Education after consultation and collaboration with educators and the general public, representing consensus of the essential content for Michigan learners. Evidence of validity based on test content includes information about the test specifications, including the test design and test blueprint. Test development involves creating a design framework from the statement of the construct to be measured. The M-STEP social studies test specifications evolve from the tension between the constraints of the assessment program and the benefits sought from the examination of students. These benefits and constraints mix scientific rigor with policy considerations.

The M-STEP test specifications consist of a blueprint and test maps for each grade level and content area. For social studies, the 2018 M-STEP test selection specifications were finalized by MDE and its psychometricians and vendors in 2017.

The key structural aspect is the test blueprint, which specifies the target score points for each discipline in social studies, as shown in Table 3-5. The blueprint represents a compromise among many constraints, including the target weights for each discipline, availability of items from field testing, and results of multiple reviews by content specialists. Test design includes such elements as number and types of items for each of the scores reported. The 2018 M-STEP operational forms matched the test blueprint that was intended for this assessment.

3.2 Test Blueprints

Test specifications and blueprints define the knowledge, skills, and abilities intended to be measured on each student's test event. A blueprint also specifies how skills are sampled from a set of content standards (i.e., the CCSS or Michigan state standards). Other important factors, such as Depth of Knowledge (DOK), are also specified. Specifically, a test blueprint is a formal document that guides the development and assembly of an assessment event/form by explicating the following types of essential information:

- content (i.e., claims/disciplines and assessment targets) that is included for each assessed content area and grade across various levels of the system (i.e., student, classroom, school, district, and state levels)
- the relative emphasis of content standards generally indicated as the number of items or percentage of points per claim and assessment target
- item types used or required, which communicate to item developers how to measure each claim and assessment target and communicate to teachers and students about learning expectations
- DOK, indicating the complexity of item types for each claim and assessment target

The test blueprint is an essential guide for both assessment developers and for curriculum and instruction. For assessment developers, the blueprint and related test-specification documents define how the test will ensure coverage of the full breadth and depth of content and how it will maintain fidelity to the intent of the CCSS and/or Michigan state standards on which the assessments are based. Full content alignment is necessary to ensure that educational stakeholders can make valid, reliable, and unbiased inferences about student, classroom, school, and state performance. At the instructional level, the test blueprint provides a guide to the relative importance of competing content demands and suggests how the content is demonstrated, as indicated by item type and DOK. In summary, an assessment blueprint provides clear development specifications and signals to the broader education community both the full complexity of the standards and how performance on these standards is substantiated.

3.2.1 Test Specifications

AERA, APA, and NCME (2014) Standard 4.1 states the following:

Test specifications should describe the purpose(s) of the test, the definition of the construct or domain measured, the intended examinee population, and interpretations for intended uses. The specifications should include a rationale supporting the interpretations and uses of test results for the intended purpose(s). (p. 85)

The purpose of M-STEP is discussed in Sections 1.2 and 1.3 of Chapter 1. M-STEP tests the knowledge and skills that are identified within Michigan's standards-based accountability system. This framework, in turn, is based on prior consensus among MDE staff, Michigan educators, and experienced content-matter experts that the framework represents content that is important for teachers to teach and students to learn.

The test specifications are discussed in accordance with AERA, APA, and NCME (2014) Standard 4.12, which states the following:

Test developers should document the extent to which the content domain of a test represents the domain defined in the test specifications. (p. 89)

Item and test development are guided by sets of specifications. Details on these specifications for ELA and mathematics can be found in the *Smarter Balanced 2017–2018 Technical Report* (2018), the *Item and Task Specifications* (Smarter Balanced, 2015c), and the Content specifications for the summative assessment of the common core state standards for English language arts and literacy in history/social studies, science, and technical subjects. (2015a). While MDE reviews all Smarter Balanced operational items, MDE utilizes the Smarter Balanced documentation for the technical details of item and test development. The remainder of this section will focus on the details for Michigan-developed assessments and items (operational items and test maps for social studies and field test items for all content areas).

A general description of development activities applying to Michigan-developed assessments (i.e., M-STEP social studies) is provided below. OEAA staff, contractors, and Michigan educators work together to develop these state assessments. Specifically, the development cycle includes the following steps:

- Item writer training
- Item development
- Item review
- Field-testing
- Field-test item review
- Operational test construction

3.2.2 Item Writer Training

Once item specifications are finalized, Michigan's item development contractor uses customized materials approved by OEAA to train item writers to author items specifically for M-STEP. Item writer training can last anywhere from three to five days and is conducted by contractor staff in conjunction with OEAA test development staff. The process of item writing includes cycle(s) of feedback from contractor and OEAA staff and can take between 4 and 8 weeks for an item to move from initial assignment to accepted status. All item writers are Michigan educators who have curriculum and instruction expertise for the grade and content for which they are writing items. In addition, prospective item writers are required to submit three original test items aligned to grade-specific content standards, which OEAA test development staff review and possibly approve for item authoring. Michigan's item writers possess relevant degrees and experience, and many have previous experience in item writing that is M-STEP specific.

3.2.3 Item Development

Item development is discussed in this section in compliance with the AERA, APA, and NCME (2014) *Standards*. Standard 4.7 states the following:

The procedures used to develop, review, and try out items and to select items from the item pool should be documented. (p. 87)

For ELA and mathematics, development of item content for the operational test was completed by Smarter Balanced from 2012 to 2014. Smarter Balanced tested items and refined its approach to item development through three steps: small-scale tryouts in fall 2012, a large pilot test in 2013, and a field test in spring 2014. Items/tasks administered for the 2018 M-STEP operational test complied with Smarter Balanced content specifications and with the item and task specifications that were refined after the pilot test and before the field test. Further details can be found in Chapter 3 in the Item Development section of the *Smarter Balanced 2017–2018 Technical Report*.

For social studies items and Michigan-developed ELA and mathematics items, Michigan item writers drafted test items in accordance with item specifications approved by OEAA test development staff. Contractor staff reviewed items internally and shared with OEAA test development for an additional review. Sections 3.2.6 and 3.3 discuss how the items are selected for field-testing or operational use. The internal review consisted of meeting the following criteria:

Skill:

- Item measures one skill level.
- Item measures skill in manner consistent with specifications.
- Item assesses appropriate (i.e., realistic) level of skill.
- Item makes clear the skill to be employed.

Content:

- Item measures one primary academic standard.
- Item measures academic standard in a manner consistent with specifications.
- Item taps appropriate (i.e., important) aspect of content associated with the academic standard.
- Item makes clear the benchmark or problem to be solved.

Relevance:

- Item is not contrived.
- Item is appropriate for the grade level to be tested.
- Item groups reflect instructional emphasis.

Accuracy:

- Item is factually accurate.
- Multiple-choice (MC) items contain only one correct or best response.
- Multi-select items contain answer choices that are clearly correct or best responses.
- Technology-enhanced (TE) items follow approved style guidelines for each grade and content area.
- If item pertains to disputed content, context for correct answer is clearly defined.
- Item is unambiguously worded.
- Item contains no extraneous material, except as required by the standard.
- Vocabulary is grade-level appropriate and clear.
- Item contains no errors of grammar, spelling, or mechanics.
- Item responses are parallel and related to the stem.
- Item responses are independent.
- Item contains no clues or irrelevant distractors.
- Directions for responding to a PBW prompt are clear.
- PBW prompt and rubric match.
- PBW rubric is clear and easy to apply.
- Item is clearly and conveniently placed on the page.
- Physical arrangement of item is consistent with OEAA style guide.
- Keys for sets of MC items are balanced (e.g., equal numbers of As, Bs, Cs, and Ds).

Bias:

- Item is free of race and gender stereotypes.
- Item contains no material known or suspected to give advantage to any group.
- Item is free of insensitive language.
- Item sets that identify race or gender either directly or indirectly are balanced with reference to race and gender.
- Item content and format are accessible to students with disabilities.
- Item content and format are accessible to students with limited English proficiency.

3.2.4 Graphics Creation

For all Michigan-developed items, MDE has an internal team of media designers that uses the graphic descriptions submitted by the item writers through Michigan's Item Bank System (IBS) to create the pictures, graphs, maps, artwork, etc. that are needed for online test items. MDE and DRC staff review and approve the completed artwork in preparation for the item review.

3.2.5 Item Review

Continuing from Standard 4.7 above, AERA, APA, and NCME (2014) Standard 3.2 is particularly relevant to fairness in item development:

Test developers are responsible for developing tests that measure the intended construct and for minimizing the potential for tests' being affected by construct-irrelevant characteristics, such as linguistic, communicative, cognitive, cultural, physical, or other characteristics. (p. 64)

The Bias and Sensitivity Review Committees (BSC) are comprised of representatives from various backgrounds whose purpose was to screen the items for racial, socioeconomic, gender, and other sensitivity issues. This follows AERA, APA, and NCME (2014) Standard 3.1, which states the following:

Those responsible for test development, revision, and administration should design all steps of the testing process to promote valid score interpretations for intended score uses for the widest possible range of individuals and relevant subgroups in the intended population. (p. 63)

Panels of educators, including those from Michigan, under Smarter Balanced patronage, reviewed all Smarter Balanced items and item stimuli for accessibility, bias/sensitivity, and content. (Item stimuli include the reading passages used on the ELA assessments and the figures and graphics used on the Mathematics assessments.) During the accessibility reviews, panelists identified issues that could negatively affect a student's ability to access stimuli or items or to elicit valid evidence about an assessment target. During the bias and sensitivity review, panelists identified content in stimuli and items that could negatively affect a student's ability to produce a correct response because of their background. The content review focused on developmental appropriateness and alignment of stimuli and items to the content specifications and appropriate depths of knowledge. Panelists in the content review also checked the accuracy of the content, answer keys, and scoring materials. Items flagged for accessibility, bias/sensitivity, and/or content concerns were either revised to address the issues identified by the panelists or removed from the item pool. The final and approved selection by Smarter Balanced educators became the Smarter Balanced computer adaptive item pool and was used for M-STEP ELA and mathematics tests.

For Michigan-developed items, after the internal reviews take place, all M-STEP items are reviewed by Michigan educators through the Content Advisory Committee (CAC) and BSC. Contractor staff trains the CAC and BSC participants using OEAA-approved materials and facilitates the committee meetings under the leadership of OEAA test development staff. All newly written test items are typically reviewed first by the BSC and then by the CAC.

An item rejected by the BSC may or may not get passed on to the CAC for review. Each review is led by experienced contractor staff, with test development staff in attendance, using the following prescribed guidelines to indicate the final status of each item:

- **Accept:** The criteria outlined in the review were met in all areas (i.e., skill, content, relevance, accuracy, and bias), and the item appears suitable for field-testing.

- **Revise:** One or more of the criteria have not been met or the item needs minor changes to make it acceptable. Reviewers provide recommendations on changes to be made to the item that will make the item suitable for field-testing.
- **Reject:** Several category conditions have not been met, are suspect, or need radical changes to make the item acceptable. In such cases, the item may be vague or ambiguous, inappropriate, or not clearly related to the text or the standard. Without extensive revisions, it is unlikely to be salvaged. Reviewers provide comments to explain why the item should be rejected.

Items that have passed bias/sensitivity and content reviews are eligible for field-testing.

3.2.6 Field-Testing

Before an item can be used on an operational test or added to the operational item pool, it must be field-tested. OEAA uses two approaches to administer field-test items: embed field-test items in an operational administration or embed field-test items in a stand-alone field-test administration. Items that have passed bias/sensitivity and content review are eligible for field-testing.

OEAA embeds field-test items in multiple forms of operational fixed-form assessments or randomly assigns field-test items to students across the state during the computer adaptive test (CAT) administrations. Administering field-test items this way ensures that they are randomly distributed, and this allows a large representative sample of responses to be gathered under operational conditions for each item. Enough field-test items are administered annually to replenish and improve the item pools.

When MDE implements testing at new grade levels, for new content areas, or for revised academic standards, it is necessary to conduct a separate stand-alone field test to obtain performance data. In 2018, MDE administered a stand-alone field test in science.

3.2.7 Range-Finding

After the student responses to the field-tested PBW prompts are collected, a range-finding is conducted to determine scoring guidelines and score-point ranges for the different score points for each field-tested writing item. This information is then used in the preparation of materials to guide the handscoring of the PBW item student responses by a trained team of readers, as described in Chapter 7 of this report.

3.2.8 Data Review

After field-testing, MDE psychometric staff analyze results. Contractor staff and test development staff convene data review committee meetings with Michigan educators. Significant effort goes into ensuring that these committee members represent the state demographically with respect to ethnicity, gender, school district size, and geographical region. These committees receive training on interpreting the psychometric data compiled for each field-test item by OEAA psychometric staff. Content experts (usually teachers) and group facilitators apply this training to the data review process. During these data review meetings, participants review the items with field-test statistics. Data provided to the data review committees are separated by BSC and CAC. The data that are reviewed during BSC include

- *N*-count;
- adjusted *p*-value (i.e., adjusted item mean in the range of 0–1 for all items);
- Differential Item Functioning (DIF) flag;
- favored group; and
- percentage of students who choose each option (option-total correlation), omit a response (omit-total correlation), and in paper/pencil tests, submit multiple marks (multiple marks-total correlation).

The data that are reviewed during CAC include

- overall *N*-count;
- adjusted *p*-value;
- difficulty flag;
- item-total correlation;
- item-total flag; and percentage of students who choose each option (option-total correlation), omit a response (omit-total correlation), and in paper/pencil tests, submit multiple marks (multiple marks-total correlation).

As mentioned above, specific directions are provided on the use of the statistical information and how to use Michigan's IBS. BSC members evaluate each test item for fairness issues with respect to culture, ethnicity, gender, geographic location, and economic status, using the data listed above for this group. CAC members evaluate each test item with regard to alignment to the academic content standard, grade-level appropriateness, and level of DOK, using the data information listed above for this group. Both committees then recommend that the item be accepted, revised for additional field-testing, or rejected.

After new items have survived all reviews and field-testing, they are saved in the Michigan IBS as "Ready for Operational," meaning they are now eligible for operational use.

3.3 Operational Test Construction

OEAA test development staff build test maps that meet the test specifications (i.e., blueprint and psychometric specifications) inside Michigan's IBS. All test maps are reviewed for correct answer key, accurate content standards, and appropriate statistic/psychometric information for each item. In addition, comparability of the overall test across forms and across adjacent years is also examined for social studies. Corresponding details for the three content areas are presented below.

3.3.1 ELA

M-STEP ELA is based on Michigan's ELA academic content standards, which were adopted by the State Board of Education in 2010. M-STEP ELA consists of four claims: Reading, Writing, Listening, and Research. The assessment is administered in grades 3–8.

M-STEP ELA is a CAT using Smarter Balanced items, all of which are reviewed and approved by OEAA staff for use in Michigan's CAT. Also, each CAT form distributes one PBW prompt per student. The PBW prompts were developed by DRC and reviewed by OEAA staff. In addition, Michigan embeds five ELA field-test items in each form for grades 3–8.

In the CAT at all grades, Claim 1 (Reading) consists of both informational and literary passages, each with related items. Passages are assessed using MC items and a variety of technology-enhanced items, such as hot text, drop-down menus, and multi-select items. Claim 2 (Writing) includes student writing samples with a set of associated items, some independent items, and one PBW item. PBW prompts cover all Claim 2 content categories. Claim 3 (Speaking/Listening) consists of 3 or 4 listening passages, each with 2 or 3 associated items. Claim 4 (Research) consists of 8 or 9 independent items. The ELA assessment structure is summarized in Table 3-2.

Table 3-2. ELA Structure for Grades 3–8

Claim/Score Reporting Category	Content Category	CAT Stimuli	PBW Stimuli	CAT Items	PBW prompts
1. Reading	Literary	2	0	7–8	0
1. Reading	Informational	2	0	7–8	0
2. Writing	Organization/Purpose and Evidence/Elaboration	0	1	6–8	1
2. Writing	Conventions	0	0	5	0 ¹
3. Speaking/Listening	Listening	3–4	0	8–9	0
4. Research	Research	0	0	8–9	0

3.3.2 Mathematics

M-STEP mathematics is based on Michigan’s mathematics academic content standards, which were adopted by the State Board of Education in 2010. M-STEP mathematics consists of four claims: Concepts and Procedures, Problem Solving, Communicating Reasoning, and Modeling and Data Analysis. The assessment is administered in grades 3–8.

There are non-calculator portions of the mathematics assessment embedded throughout the online test. All items in grades 3–5 are non-calculator items.

M-STEP mathematics is a CAT using Smarter Balanced items, all of which are reviewed and approved by OEAA staff for use in Michigan’s CAT. Michigan embeds five mathematics field-test items in the CAT in each form in grades 3–8.

In the mathematics assessment, the Claim 1 (Concepts and Procedures) section consists of 20 items (MC or TE) in the CAT. Details of the various TE types can be found in Section 3.7. The Claim 2 (Problem Solving) section consists of 4 items. The Claim 3 (Communicating Reasoning) section consists of 8 items. The Claim 4 (Modeling and Data Analysis) section consists of 4 items. Claims 2 and 4 are combined in the blueprint and reporting structure because of content similarity and to provide flexibility for item development. There are still four claims, but only three claim scores are reported with the overall mathematics score. The mathematics assessment structure is summarized in Tables 3-3 and 3-4.

¹ PBW prompts cover all Claim 2 content categories but are listed under only one in Table 3-2 to avoid double-counting.

Table 3-3. Mathematics Overall Structure: Number of Items Claim/Reporting Category

Claim/Score Reporting Category	Grades 3–8
1. Concepts and Procedures	20
2. Problem Solving and 4. Modeling and Data Analysis	8
3. Communicating Reasoning	8

Table 3-4. Mathematics Structure for Grades 3–8

Claim/Score Reporting Category	Content Category	CAT Items
1. Concepts and Procedures	Priority Cluster	15
1. Concepts and Procedures	Supporting Cluster	5
2. Problem Solving and 4. Modeling and Data Analysis	Problem Solving, Modeling and Data Analysis	8
3. Communicating Reasoning	Communicating Reasoning	8

3.3.3 Social Studies

M-STEP social studies is based on Michigan's social studies academic content standards, which were adopted by the State Board of Education in 2007. The assessment is administered in grades 5, 8, and 11. The M-STEP social studies assessment in grade 5 consists of five domains: History, Geography, Civics and Government, Economics, and Public Discourse. There are 45 operational items and 15 embedded field-test items. The M-STEP social studies assessment in grade 8 consists of four domains: History, Geography, Civics and Government, and Economics. There are 44 operational items and 22 embedded field-test items. The M-STEP social studies assessment in grade 11 social studies assessment consists of four domains: U.S. History and Geography, World History and Geography, Civics, and Economics. There are 38 operational items and 16 embedded field-test items. The social studies assessment structure is summarized in Table 3-5.

Table 3-5. Social Studies Structure for Grades 5, 8, and 11

Grade	Domain	# of Operational Items
5	History	19
5	Geography	7
5	Civics and Government	10
5	Economics	7
5	Public Discourse	2
8	History	21
8	Geography	14
8	Civics and Government	4
8	Economics	5
11	U.S. History and Geography	12
11	World History and Geography	12
11	Civics	7
11	Economics	7

3.3.4 Science

M-STEP science is based on Michigan’s science academic content standards, which were adopted by the State Board of Education in 2015. Because science is tested once in a three-year grade band, it would not be appropriate to test under the new standards until schools have had three years of instruction under the new standards, nor would it be appropriate to test under the previous standards. Instead, M-STEP science was a statewide field test for spring 2018.

3.3.5 Accommodations

Michigan is committed to ensuring all students, including English learners and students with disabilities, have access to a wide array of tools across M-STEP. Sections 4.1–4.3 in this report detail the Universal Tools, Designated Supports, and Accommodations Michigan provides. It is important to note that M-STEP is available to students who require Accommodations according to their Individualized Education Program (IEP) or 504 plan. Paper/pencil accommodated tests are available in many different forms to meet the needs of Michigan’s students by providing the tests in contracted and uncontracted braille, enlarged print, and by even having translated forms of the tests. Students may also test online with many different options such as video sign language (American Sign Language and Signed Exact English), stacked Spanish, English text-to-speech, and closed captioning (in the Listening claim). Whether students take a test with or without Universal Tools, Designated Supports, and Accommodations, the M-STEP assessments are administered during the same testing window as regular operational tests.

3.4 Sources of Items and Metadata

3.4.1 ELA and Mathematics

M-STEP ELA and mathematics have two sources for test items:

1. The Smarter Balanced Assessment Consortium
2. The Michigan Item Bank System (IBS)

Smarter Balanced worked with a variety of assessment vendors, state education departments, and educators throughout 2012 to create a pool of ELA and mathematics test items in preparation for pilot testing. In the process of creating the test items, the item writers were provided trainings in evidence-centered design, universal design, DOK, accessibility, and issues of bias and sensitivity. The item writers also received content and item specifications to guide their development. Each test item passed through approval from a content committee, an accessibility committee, and a bias and sensitivity committee before being added to the item pool.

In 2013, Smarter Balanced conducted a pilot test in a small number of schools across states participating in Smarter Balanced, using items from the existing item pool. Smarter Balanced used feedback from this pilot test in preparation for further item development and testing. In 2014, Smarter Balanced administered a field test of the existing item pool to more than 4 million students across states participating in the consortium, including Michigan. Smarter Balanced conducted a subsequent data review using educator committees to evaluate the performance of the test items across the country and to ensure that the items met the quality levels required in terms of content, accessibility, and issues of bias and sensitivity to be included in the operational item pool. The items were then made available to Michigan for inclusion in the CAT item pool. Each year, additional items are field-tested to replenish the general item pool.

The Michigan IBS contains items that have been developed and reviewed by Michigan teachers using processes described earlier in this chapter. The items from both sources (i.e., Smarter Balanced and Michigan IBS) contained a mixture of MC and TE item types, with a DRC-developed PBW prompt added for each student.

3.4.2 Social Studies

The item development process for M-STEP social studies utilizes the Michigan IBS as its main resource. The Michigan IBS is a secure, web-based application that allows users to create contexts and test items. It leads users through all the steps of the item development process, including context review, item review, and data review as described in Section 3.2.

3.5 Import into DRC INSIGHT Test Engine

M-STEP is administered through the DRC INSIGHT test engine. The test items must be imported into INSIGHT from the various sources noted earlier. Once the items are loaded into INSIGHT, they can be rendered for review in the identical formatting structure in which a student would see the item in a test. After the items have been formatted and rendered, they can be assembled into online test forms based on the sequence and information provided in the test maps.

3.6 Psychometric Review During Assessment Construction

Content specialists and psychometricians both from MDE and from Smarter Balanced followed psychometric guidelines and targets for operational forms construction. The foremost guideline was for item content to match the test blueprint for the given content. Both groups used item flagging criteria (discussed below) to guide the assessment construction. Items with flags were avoided when possible.

Details for psychometric reviews are described below by content area groups. Such reviews for ELA and mathematics are done by the Smarter Balanced psychometrician(s), while social studies reviews are carried out by an MDE psychometrician.

3.6.1 ELA and Mathematics

The psychometric review for the items in the M-STEP CAT pool and fixed forms was conducted by Smarter Balanced. Smarter Balanced flagged items based on the following content criteria (Smarter Balanced, 2016, p. 4–22):

- The following items were flagged based on item difficulty and score distribution:
 - items with a low average item score (i.e., less than .10)
 - items with a high average item score (i.e., greater than .95)
 - items with a proportion obtaining any score category less than 0.03
- The following items were flagged based on item discrimination:
 - items with a low item-total correlation (i.e., less than .30)
 - items with a higher mean criterion score for students in a lower score-point category
- The following multiple-choice items were flagged:
 - items where higher ability students (i.e., those in the top 20% on overall score) select a distractor more often than the key
 - items with a higher criterion score mean for students choosing a distractor than the mean for those choosing the key
 - items with a positive correlation between distractor and total score

Items are also classified into three Differential Item Functioning (DIF, for corresponding details please see Chapter 11) categories of A, B, or C. The focus group was indicated by a positive value (e.g., C+), and the reference group was noted with a negative value (e.g., C-). The positive and negative values were reported for items with C DIF. DIF comparison was not done if the

sample size for either group was less than 100 or if the combined sample size for the groups being compared was less than 400 (Smarter Balanced, 2017, p. 3–15.)

DIF was evaluated for eight subgroup comparisons, shown here with the focal groups listed first and the reference groups listed second.

- Gender: Female – Male
- Race/Ethnicity: Asian – White
- Race/Ethnicity: Black – White
- Race/Ethnicity: Hispanic – White
- Race/Ethnicity: Native American – White
- Individualized Education Program: Yes – No
- Limited English Proficiency: Yes – No
- Title 1: Yes – No

Items with C+ or C- DIF were flagged for data review.

Items that were not flagged for content or bias statistical issues were eligible for use in the operational pools. Flagged items became eligible for the operational pools if they were approved by a multidisciplinary panel of experts during data review.

3.6.2 Social Studies

For social studies, the following analyses were carried out for psychometric review (note that the listed analyses are routine annual procedures):

1. Content standard distribution check: This check is to ensure that operational (OP) items on each form have the desired content coverage (i.e., the reporting categories are the same as depicted in the test blueprint; and within each reporting category, the content standards have as much variety as possible.)
2. Item position check: Equating items and common items (i.e., non-equating items that appear on multiple forms) need to appear in the same test positions across forms. Moreover, equating items are checked to make sure they are within +/-2 position change from the previous year's positions.
3. Across year comparability check: For this check, distributions of item difficulty and item discrimination (p-values and adjusted item-total correlations, see Section 8.3.1.2 for details) are checked across adjacent years for unique items to make sure they are comparable. Moreover, when Item Response Theory (IRT) item difficulty and item discrimination (b-parameters and a-parameters) (see Section 6.2 and Equation 6-2 for details) are available for all OP items, test characteristic curves (TCCs), test information function (TIF) curves, and test standard error (TSE) curves are plotted to check the comparability across years
4. Across mode comparability check: Comparability of OP items across modes (paper/pencil, online) is checked using the same approaches listed above in the across year comparability check.

5. Comparability of equating items and other OP items per form: Two analyses are involved in this comparison on each form: (1) content coverage homogeneity test (to make sure that equating items and other OP items have comparable content coverage) and (2) distributions of item difficulty and adjusted item-total correlation comparability check. These analyses are conducted to make sure that the equating items function as a mini-test (i.e., they are representative of the overall test, both statistically and in terms of content).
6. Item key distribution check: This check involves all multiple-choice (MC) items on the test (i.e., OP and field-test items). Here the desired result is for all four key options to appear relatively equally on each test map, with no same key option appearing three times consecutively. Although it is desirable to have unique field-test items on each form, if a field-test item must be repeated on multiple forms, a check is carried out to ensure that it appears in the same test position across forms.
7. Overall OP item set quality check: This check ensures that no OP items have problematic flags. Specifically, DIF results are checked to make sure that no OP items are with “B” or “C” DIF flags. All OP items that appear on the final form have been scrutinized to make sure that there are no bias or sensitivity issues involved. Moreover, adjusted item-total correlations, various item statistics flags (e.g., key option-total correlation being negative, distractor option-total correlation being positive, omit-total correlation being positive, key option percentage not being the highest), and IRT item parameters are also checked to see if items are free of concerns (i.e., adjusted item-total correlation should be ≥ 0.2 , a-parameter should be > 0 , b-parameter should be in the range of $[-3, +3]$, and there should be no item statistics flags).

All identified problems are documented and communicated to the corresponding content leads. Content leads then revise and resubmit test maps for another round of review. This iterative process continues until all issues have been resolved or the problematic item selections are proven to be the best selections given various constraints (e.g., content coverage considerations, and the need to avoid possible clueing).

3.7 Item Types Included

In addition to the traditional MC items, TE items were included in M-STEP. The following is a list of the TE item types used:

- Drag and Drop—Students drag pictures or words into boxes or “drop zones” to indicate an answer.
- Choice Interaction—This is similar to an MC item, but the item can have more than four options, and any number of the options can be correct.
- Hotspot (Count or Selection)—Students answer by selecting graphics, either a particular number of hotspots (Count) or a specific hotspot (Selection).
- Matching Interaction—Students select areas of an interaction grid to match options in rows and columns.
- Matching—Students make line connections between options from two sets.
- Keypad Input—Students use an embedded keyboard with mathematical functions to answer mathematics questions.
- Drop-Down—Students select options from a drop-down list.
- Hot Text Highlight (Line and Paragraph)—Text is selectable and, once selected, will become highlighted for the students. Students select one or more lines of text (Line) or words or sentences from a block of text (Paragraph).
- Order—Students answer by rearranging a list of items or sentences.
- Coordinate Graph Input—Students plot points, lines, and shapes on a coordinate grid.
- Number Line Graph—Students plot points on a number line.
- Text Input—Students enter values in a response box.
- Bar Graph—Students answer by selecting amounts to complete a bar graph.

Not all the TE item types are used in every content area.

3.8 Field-Test Selection and Administration

3.8.1 Field Test Item Selection

The OEAA content leads are tasked with selecting field-test items. The blueprints specify the number of field-test items by grade level and content area. The content leads work within Michigan’s IBS to monitor the number of operational items available for each content standard. Where there are gaps in the numbers available, content leads may decide to field-test items assessing that standard. The content leads also monitor the number of items that may be overexposed and need replacement items as one way to select field-test items.

Responses on field-test items do not contribute to a student’s score on the operational tests. The specific locations of the embedded items in the assessment are not disclosed. These data are free from the effects of differential student motivation that might characterize stand-alone field-test designs since the items are answered by students taking operational tests under standardized test administration procedures.

3.8.2 Field Test Administration

3.8.2.1 Mathematics and ELA

MDE-developed field-test items are embedded within the ELA and mathematics CAT assessments at all grade levels. The items are not designated as field-test items to the students, so the field-test items are not distinguishable from the operational items. This ensures that the students give the same effort to the field-test items as the operational items. All the students taking the CAT receive the same number of field-test items, and the selection and delivery of the field-test items are not affected by a student's performance on the test or the difficulty level of the field-test items. To avoid complications due to position placement, the field-test items are not distributed in the first five or the last five positions on the test.

For mathematics, the field-test items are placed at the same sequence positions throughout the CAT testing experience. For ELA, the field-test items are positioned in similar locations; however, due to the inclusion of passage sets in ELA, the field-test positions are shifted as necessary to accommodate a preceding passage set. Students receive either stand-alone field-test items or a field-test passage set containing associated items.

3.8.2.2 Social Studies

Social studies assessments consist entirely of MDE-developed operational and embedded field-test items for all grade levels.

The operational item set is the same across all online forms in a grade level, appearing in the same test positions. The remaining form positions are used for field-test items, which are unique to each form. The three online forms at each grade level are randomly administered to the student population.

The paper/pencil forms for social studies share all the equating items with the online forms. However, since TE items cannot be presented on paper/pencil forms, items in those positions are replaced by items assessing the same content standards and having similar item statistical profiles that were presentable on paper/pencil forms and in braille format.

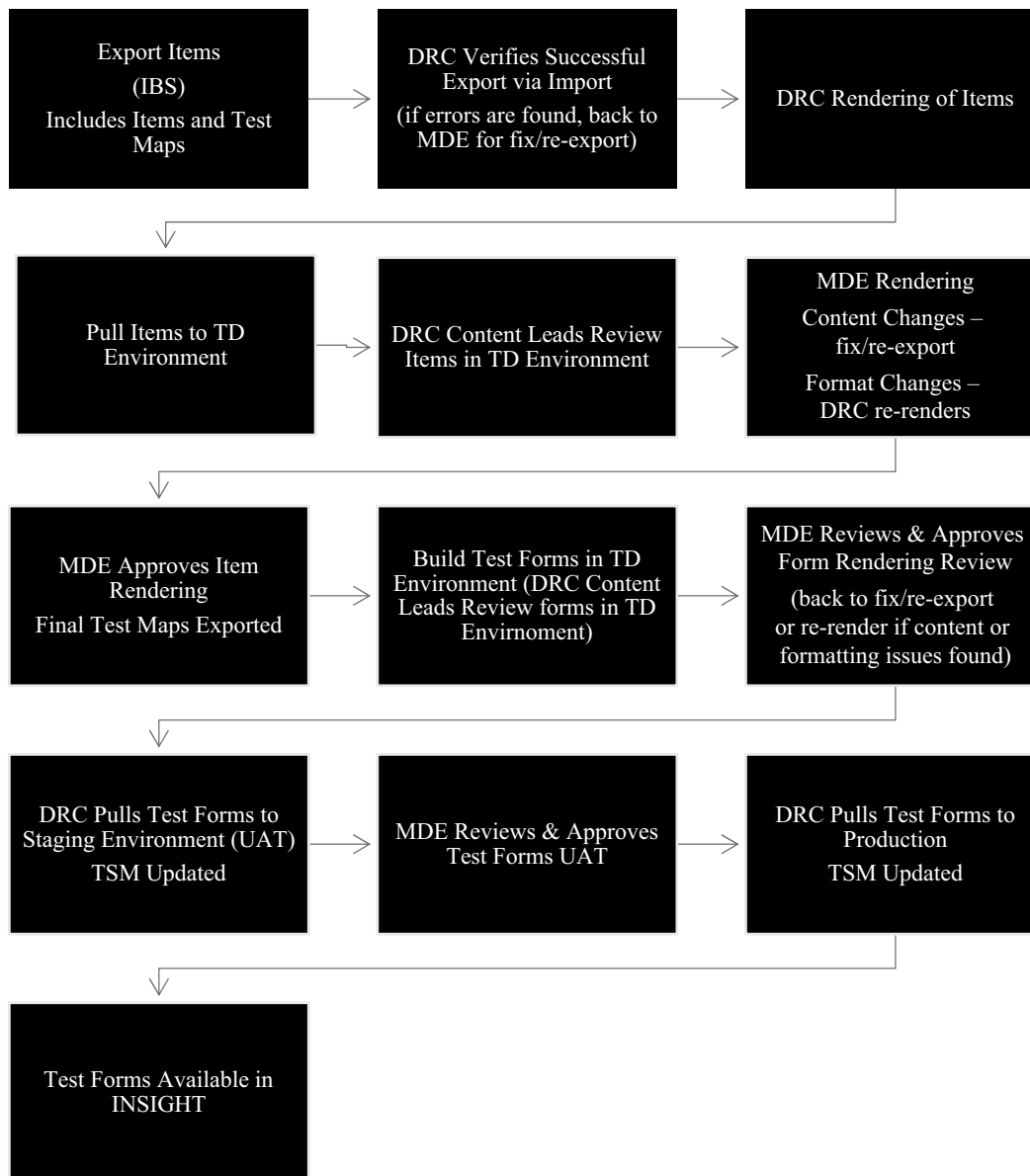
Details on constructing forms and follow in Sections 3.9 and 3.10.

3.9 Online Form Building and Rendering Process

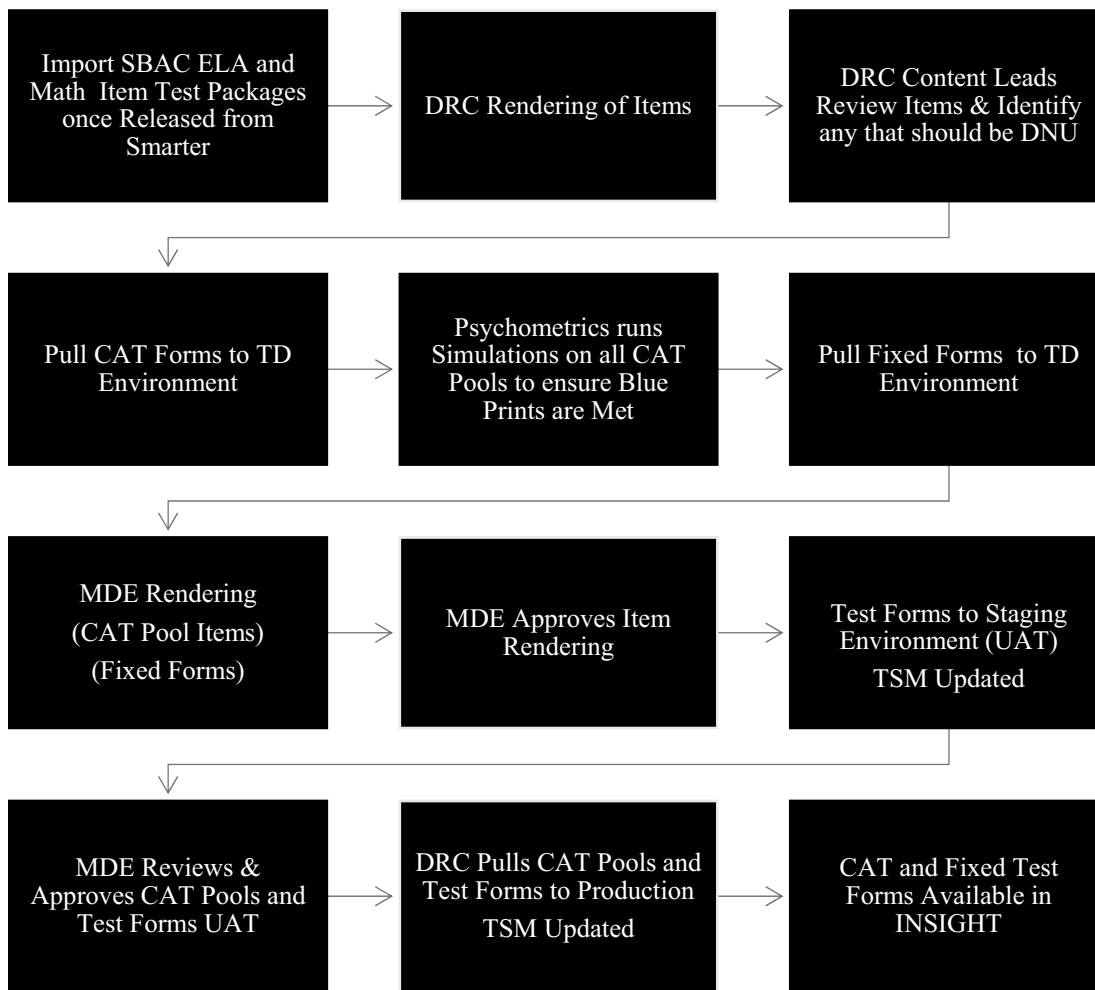
3.9.1 Overview of Rendering Process

DRC and MDE follow a very rigorous rendering process for all items on the 2018 M-STEP. Using the web-based application LeanKit, DRC and MDE monitor the progress of each grade and content batch. The process begins right after the import of items from the Michigan and Smarter Balanced item banks. All parts of the rendering process are completed a month prior to the start of testing to ensure time for User Acceptance Testing (UAT) of all grades and contents. Figure 3-1 below shows the entire process for M-STEP field-test items and social studies items that are imported from the Michigan IBS.

Figure 3-1. Rendering Process of Michigan-Built Items



The rendering process for the Smarter Balanced items is slightly different. Figure 3-2 shows the process followed for all items that are imported from Smarter Balanced to use for M-STEP ELA and mathematics.

Figure 3-2. Rendering Process of Smarter Balanced Items

Requirements are established and reviewed with MDE prior to the imports of the 2018 M-STEP items. The requirements include the QTI 2.2 import specs between the IBS and DRC's IDEAS system as well as specific rules when importing each type of item. Detailed rendering requirements are also documented and reviewed.

3.9.2 Form Preparation and Rendering in INSIGHT

For all fixed forms, after the individual items are formatted and rendered, online test forms are assembled in the INSIGHT test engine based on the sequence and information provided in the test maps created by MDE. The test maps provide test-form data, item form sequence location, and metadata (e.g., content standard, DOK, item position, *p*-value, IRT parameters, answer key, points possible) for each test form for each test type (i.e., program, content, grade). DRC applies the appropriate styles and formatting to the fixed forms based on the previously set style and formatting guidelines.

The assembled fixed forms are then reviewed by content leads at DRC and MDE in a UAT setting to ensure that the forms match the exact design and data displayed in the test maps and that the forms, features, and functionality of INSIGHT appear and operate correctly. The UAT is conducted using the same INSIGHT test delivery system as the students use so the forms appear and function just as the students see them. The forms include features such as the online tools provided for each item, test directions, help files, calculators, and reference materials. Detailed information on student tools can be found in Chapter 4.

3.10 Paper/Pencil Form Building and Review Process

Although approximately 99% of Michigan students test online, there will always be paper/pencil forms available for those students who may not be able to test online and for student groups that require specific Accommodations or tests in other languages. Michigan offers the following Accommodations for students with disabilities and the following accessibility features for English learners delivered through paper/pencil assessments: enlarged print; braille; audio supports, such as reader scripts for teacher read-aloud Accommodations; audio CDs; and DVDs in Arabic and Spanish. The ELA and mathematics paper/pencil tests are provided by Smarter Balanced and align to Michigan's ELA and mathematics blueprints. OEAA's composition unit assembles the test booklets. There are several rounds of reviews conducted by OEAA content leads, OEAA assessment specialists, and OEAA's editor. Once the initial test booklets are approved, they are posted for printing by Measurement Incorporated, and the paper/pencil test maps are provided to Measurement Incorporated for use in creating braille and enlarged print forms using the American Printing House (APH) for the Blind.

The social studies paper/pencil tests are developed by OEAA's content leads using Michigan's IBS. They mirror their online counterparts with modifications, i.e., only TE items are replaced. The content leads review each item in the test map to check for text and/or graphic errors, clueing, correct answer keys, and a balance of answer keys. Once the test map is approved by the content lead, the psychometric lead reviews the test map in a similar way as mentioned above for online forms, but with more focus on comparability of paper/pencil forms to their online counterparts. Once the test maps are approved by both the content lead and the psychometric lead, the composition unit creates one item per page (i.e., "one-per") for review by both the OEAA content lead and the OEAA editor. A one-per is created for each item on the test map, showing how each item will appear in a test booklet. Content leads ensure the one-per matches the item as it is in the IBS, which is the source of truth. The item as it appears on the one-per must also follow OEAA's style guide and be free of errors. After the content lead approves the one-pers, they are reviewed by OEAA's editor. Once the editor approves the one-

pers, test booklets are created. The draft printed test booklets are reviewed first by the editor and then by the content lead. Both the content leads and the editor use OEAA's Proofing Tools Guide and its task checklists to ensure each step is followed. Once the test booklet has final approval, the test maps and approved test booklets are sent to Measurement Incorporated for mass printing and accommodated format production of enlarged print, braille, reader scripts, audio CDs, and DVDs in Spanish and Arabic.

3.11 Summary

In summary, the overall purpose of this chapter is to explicate the procedures used in the development of M-STEP. The efforts by MDE and its vendors address multiple best practices of the test industry, particularly the following AERA, APA, and NCME (2014) *Standards*:

- Standard 3.1—Those responsible for test development, revision, and administration should design all steps of the testing process to promote valid score interpretations for intended score uses for the widest possible range of individuals and relevant subgroups in the intended population.
- Standard 3.2—Test developers are responsible for developing tests that measure the intended construct and for minimizing the potential for tests' being affected by construct-irrelevant characteristics, such as linguistic, communicative, cognitive, cultural, physical, or other characteristics.
- Standard 4.0—Tests and testing programs should be designed and developed in a way that supports the validity of interpretations of the test scores for their intended uses. Test developers and publishers should document steps taken during the design and development process to provide evidence of fairness, reliability, and validity for intended uses for individuals in the intended examinee population.
- Standard 4.1—Test specifications should describe the purpose(s) of the test, the definition of the construct or domain measured, the intended examinee population, and interpretations for intended uses. The specifications should include a rationale supporting the interpretations and uses of test results for the intended purpose(s).
- Standard 4.7—The procedures used to develop, review, and try out items and to select items from the item pool should be documented.
- Standard 4.12—Test developers should document the extent to which the content domain of a test represents the domain defined in the test specifications.

Chapter 4: Test Administration Plan

Chapter 4 reviews the test administration process for both the online and paper/pencil administrations of the M-STEP assessment. Detailed information on supports, accommodations, test materials, and training and test security practices can be found outlined throughout this chapter. According to the AERA, APA, & NCME *Standards* (2014), “[t]he usefulness and interpretability of test scores require that a test be administered and scored according to the developer’s instructions” (p. 111). Chapter 4 examines how test administration procedures implemented for M-STEP strengthen and support the intended score interpretations and reduce construct-irrelevant variance that could threaten the validity of score interpretations.

The online platform components of eDIRECT and INSIGHT, which were necessary for all online test administrations, are discussed in Section 4.4. The web-based application known as eDIRECT was used for all test preparation and test monitoring, while INSIGHT was the online test delivery system used by students when taking online assessments. More information on the online components can be found in Chapter 4.

4.1 Universal Tools, Designated Supports, and Accommodations

To allow all students the ability to fully demonstrate their knowledge and skills on the statewide assessments, a variety of tools are made available across all grades, content areas, and modes of testing. The variety of tools offered attempts to ensure that an equal opportunity for a student to demonstrate what he or she knows on a test is not negatively impacted by the student’s disability or English language proficiency.

MDE categorizes tools into three levels: Universal Tools, Designated Supports, and Accommodations. Universal Tools can be used by students at their own discretion. Use of a Designated Support requires an educator identify that support type for a student because of an instructional need. Tools listed as Accommodations require that a student has an Individualized Education Program (IEP) or 504 plan and that the need to use that support is identified within that document.

Regardless of the level of the tool type, MDE requires educators to make decisions about use on an individual basis. The decision for use should be based on the individual student’s instructional needs for each content area. Some tools may be classified as nonstandard, as described in the Supports and Accommodations documentation, in which case the use of those tools by students may result in invalid test scores. School districts may contact MDE if an IEP or 504 team wants to use an Accommodation that is not on the approved list. MDE will consider allowing that Accommodation for the current administration and in future administrations pending literature and research reviews and discussions with MDE’s assessment content leads.

MDE’s policies related to the use of Accommodations are in compliance with AERA, APA, and NCME (2014) Standard 6.2, which states the following:

When formal procedures have been established for requesting and receiving accommodations, test takers should be informed of these procedures in advance of testing. (p. 115)

Additional information about Michigan’s accommodations framework and a list of which Universal tools, Designated Supports, and Accommodations are considered allowable and valid for students to use can be found in the [Student Supports and Accommodations Table](#).¹

4.1.2.1 Educator Guidelines

Many of the allowable Designated Supports and Accommodations require educators to perform an action for the student or on behalf of the student. For example, a student needing a scribe may have one provided to them as long as the educator is using the guidelines for scribing outlined in MDE’s Scribing Protocol. Additional documents exist to ensure educators are providing these Designated Supports and Accommodations in a consistent and reliable manner. Additional guidelines include *Read-Aloud Guidelines*, *Spanish Read-Aloud Guidelines*, and *Arabic Read-Aloud Guidelines*.

4.1.2.2 Research Base for Supports and Accommodations

Smarter Balanced has published multiple literature reviews that support the use of MDE’s Universal Tools, Designated Supports, and Accommodations. Because MDE uses Smarter Balanced test content, the framework upon which the assessments have been built was based on the development efforts of Smarter Balanced. These [Smarter Balanced Literature Reviews](#) address research related to tools for students with disabilities and English learners.

4.1.2.3 Monitoring the Use of Designated Supports and Accommodations

The 2018 administration included a pilot monitoring of Designated Supports and Accommodations used by students to ensure high reliability and validity of test results. Data audits included verification that students receiving Accommodations on the assessment had an Individualized Education Program or 504 plan. In the event that students received accommodations without an IEP or 504 plan, the school was contacted and asked to verify the use of Accommodations and make a plan to improve their process for future student use of Designated Supports and Accommodations. Interviews were conducted with schools after assessment monitoring to verify the decision-making processes used in providing Designated Supports and Accommodations to students for use on the assessment.

¹ https://www.michigan.gov/documents/mde/M-STEP_Supports_and_Accommodations_Table_477120_7.pdf

4.2 Online Accommodations

Appropriate Universal Tools, Designated Supports, and Accommodations were available for students to use while taking the assessment. Students with an IEP or 504 plan are required to have their assessment needs formalized in those documents prior to using any Universal Tools, Designated Supports, or Accommodations. Some embedded Designated Supports and the assessments with embedded Accommodations were delivered via fixed forms. Embedded refers to supports that were provided within the online delivery platform and non-embedded refers to supports that were provided externally from the online delivery platform. The Designated Supports and Accommodations that were embedded in the online fixed forms and used for the spring 2018 M-STEP were as follows.

- Audio Sign Language (applicable for ELA and Math) was available to students at grades 3–8. For ELA, Audio Sign Language video was available for only Listening passage stimuli and items.
- Stacked Translation (applicable for Spanish Math) was available to students at grades 3–8.
- Closed Captioning (applicable for ELA) was available to students at grades 3–8. ELA Closed Captioning was available for only Listening passages stimuli.

Embedded and non-embedded Designated Supports were also available and could be selected for each student via eDIRECT.

Embedded Designated Supports and embedded Accommodations were available within the CAT assessments and the fixed-form assessments. The available embedded Designated Supports and Accommodations are listed below.

- Text-to-Speech (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11. This Designated Support reads aloud items only. Text-to-Speech with Passages (applicable for ELA only) was available to students at grades 6–8 as an Accommodation. This embedded Accommodation reads aloud both items and passages.
- Masking (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11 (Designated Support).
- Color Choice (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11 (Designated Support).
- Contrasting Color (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11 (Designated Support).

In addition to the Designated Supports and Accommodations delivered by the test engine, there are a number of non-embedded Designated Supports and Accommodations available to students. The use of these non-embedded Designated Supports and Accommodations can be indicated in eDIRECT.

The list of non-embedded Designated Supports and Accommodations that are listed in eDIRECT can be found below. This is not a full list of non-embedded allowable Designated Supports and Accommodations but is only a list of what MDE considers the most frequently used non-embedded Designated Supports and Accommodations.

- Administered Individually/Small Group (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11.
- Noise Buffers (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11.
- Oral Translated Test Directions (applicable for mathematics) was available to students at grades 3–8.
- Read Aloud (Human Reader) (applicable for ELA and mathematics) was available to students at grades 3–8.
- Bilingual Word-to-Word Dictionary (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11.
- Auditory Amplification (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11.
- Visual Aids (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11.
- Scribe (non-writing items) (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11.
- Scribe (PBW prompt) (applicable for ELA) was available to students at grades 3–8.
- OEAA Multiplication Table (applicable for mathematics) was available to students at grades 4–8.
- Abacus (applicable for mathematics and social studies) was available to students at grades 3–8 and 11.
- Non-embedded Calculator (applicable for mathematics and social studies) was available to students at grades 4–8 and 11.
- Administrator Sign Test Directions in ASL (applicable for ELA, mathematics, and social studies) was available to students at grades 3–8 and 11.
- Administrator Sign Test Content in ASL (applicable for social studies) was available to students at grades 5, 8, and 11.
- Alt Communication Devices (ACD) was available to students at grades 5, 8, and 11.



















Table 4-1 below presents more details for DRC INSIGHT student tools. The following tools are available only on some fixed forms or in certain content areas.



Table 4-1. DRC INSIGHT Student Tools by Grade, Content Area, and Test Type

Assessment	Gr	Pointer	Crossoff	Highlighter	Magnifier	Line Guide	Sticky Notes	Ruler	Protractor	Calculator	Graphing Tool	Dictionary/Thesaurus	Periodic Table	Help	Flag for Review	Pause	Writing Tools
ELA – CAT	3	x	x	x	x	x	x					x		x		x	x
ELA – CAT	4	x	x	x	x	x	x					x		x		x	x
ELA – CAT	5	x	x	x	x	x	x					x		x		x	x
ELA – CAT	6	x	x	x	x	x	x					x		x		x	x
ELA – CAT	7	x	x	x	x	x	x					x		x		x	x
ELA – CAT	8	x	x	x	x	x	x					x		x		x	x
Math – CAT	3	x	x	x	x	x	x							x		x	
Math – CAT	4	x	x	x	x	x	x		x					x		x	
Math – CAT	5	x	x	x	x	x	x							x		x	
Math – CAT	6	x	x	x	x	x	x			x				x		x	
Math – CAT	7	x	x	x	x	x	x			x				x		x	
Math – CAT	8	x	x	x	x	x	x			x				x		x	
Science	5	x	x	x	x	x	x							x	x	x	
Science	8	x	x	x	x	x	x						x	x	x	x	
Science	11	x	x	x	x	x	x						x	x	x	x	
Social Studies	5	x	x	x	x	x	x							x	x	x	
Social Studies	8	x	x	x	x	x	x							x	x	x	
Social Studies	11	x	x	x	x	x	x							x	x	x	

Figure 4-1 provides descriptions of system tools that help with navigation, may be Universal Tools, and some that are made available based on the item type and/or when a Designated Support or Accommodation is enabled.

Figure 4-1. DRC INSIGHT Student Tools Descriptions

TOOL	DESCRIPTION/FUNCTION
Navigation Tools	
	Back and Next —Move to the next question or a previous question. (Back is only available in CAT within passage and listening sets.)
	Go To Question —Jump to any item or passage set on the test by choosing the item from a drop-down list (only available in fixed forms).
	Pause —Pause the test for a short period of time (e.g., restroom break) and resume upon return.
	Flag —Mark a question for review at a later point (only available in fixed forms).
	Test Review —Review and change answers by section and indicate whether the test is ready to be scored (only available in fixed forms).
Standard Test-Taking Tools (available at all times)	
	Pointer —Select, change, or unselect an answer option; select other user tools; and navigate through the test. When moved over an answer choice, the pointer converts to a pencil image.
	Cross-Off Tool —Cross out an MC answer selection believed to be incorrect. This tool includes an eraser to remove the cross off if a student changes his or her mind.
	Highlighter —Highlight a portion of text or a graphic and remove highlights.
	Magnifier —Magnify/enlarge a portion of the screen (i.e., object, image, or text) by two times for better viewing.
	Line Guide —Movable, straightedge line used to follow along with each line of text. Student can drag the guide up or down on the screen as an aid in reading an item or passage.
	Help —The Help Library provides information on tool usage, test directions, helpful hints, and other topics. Also includes a “What’s This?” feature that allows a student to access contextual help for a specific tool or button.
	Sticky Note —Creates and places a small note in which a student can type a short message for later reference (multiple notes can be created for each item or passage).
	Calculator —Basic four-function and scientific options are available as required, either individually or together.
	Measurement Tools —Includes a Protractor for measuring angles that can be moved over any object on the screen and rotated.
	Reference Materials —Includes a Periodic Table for grades 7 and 11 science only.
	Graphing Tool —Used to graph one or several functions. Includes zoom and trace features.
	Click to Respond —Allows for placing various types of response areas in a snapshot view that a student expands to respond to the question. For example, a large graphing item can be placed in an item where it might not normally fit.
	Click to Enlarge —Allows for large graphics by using a thumbnail image of the graphic that can be enlarged for viewing. Student can interact with the test item and other tools simultaneously.

TOOL	DESCRIPTION/FUNCTION
Accommodations Tools (determined at the student level)	
	Audio/Video tools —Includes a Text-to-Speech Synthesizer that allows all test-related information (e.g., test directions, questions and answers, formula sheets) to be read aloud to the student. VSL fixed forms provide video for sign language administration .
	Display Options —Can be made available for all students or just those with a specific accommodation, such as Color Overlays , that allows a student to change the background color for text, graphics, and response areas.

4.3 Paper/Pencil Universal Tools, Designated Supports, and Accommodations

As noted earlier, OEAA provides a multitude of opportunities for students to demonstrate their knowledge on the M-STEP assessment with appropriate Universal Tools, Designated Supports, and Accommodations on the paper/pencil forms as well. Below is a list of available paper/pencil Designated Supports and Accommodations that require a specific form of the assessment.

- Stacked Spanish, available for Mathematics in all grades
- Arabic, Spanish, and English DVD, available for Science and Social Studies in all grades, to be used with form 1
- Braille, contracted and uncontracted for all content areas and grades
- Reader Script, available for Science and Social Studies in all grades, to be used with form 1
- English Audio CD, available for Science and Social Studies in all grades, to be used with form 1

Referenced in Table 4-2 is the Designated Support and Accommodation information that is tracked (i.e., bubbled in) on each content area's booklet. This is not a full list of allowable Universal Tools, Designated Supports, and Accommodations but is only a list of what MDE considers the most frequently used Designated Supports and Accommodations.

Table 4-2. Paper/Pencil Accommodations Table

Accommodation	ELA	Math	Social Studies
Directions Read in Native Language	✓	✓	
Oral Translation in Native Language		✓	✓
Spanish Booklet		✓	
Enlarged Print	✓	✓	✓
Multiple-Day Testing	✓	✓	✓
Audio CD			✓
English DVD			✓
Spanish DVD			✓
Arabic DVD			✓
Reader Script			✓
Alternate Response	✓	✓	
American Sign Language (ASL)	✓	✓	
Noise Buffers	✓	✓	
Read-Aloud (see Supports and Accommodations Table for specifics)	✓	✓	
Scribe	✓	✓	
Speech-to-Text	✓	✓	
Abacus		✓	
L1 Glossary		✓	
Other	✓	✓	✓
Nonstandard Accommodation/Support	✓	✓	✓

4.4 Online Test Platform

The secure web-based test engine DRC INSIGHT Online Learning System is downloaded onto computers that students access for all online assessments. Test items and forms can only be accessed using a valid test ticket. Automatic updates are suggested to be turned to “Enable” for the software to be automatically updated as needed. From the INSIGHT landing page, students have access to the test via the “Test Sign In” link and to the sample item sets via the “Online Tools Training” link.

DRC’s client portal, eDIRECT, is used to manage the test setup functions of student assessments and provide the downloads available for installation. The INSIGHT secure browser software is downloaded from eDIRECT and installed on student testing devices. The secure browser can be installed on computers individually, or it can be downloaded to a central location, copied, and simultaneously distributed to multiple computers using common network distribution tools. Everything needed for testing is found within the secure browser, eliminating the need for districts to coordinate updates to third-party software.

Technology coordinators install testing site manager(s) (TSM[s]) to manage the content (test content, responses, and audio files) and regulate traffic between testing sites and DRC’s servers. The System Readiness Check helps troubleshoot any issues that may have occurred during INSIGHT installation or while INSIGHT was running. This application is installed when INSIGHT is installed, and the System Readiness Check performs a series of tests that can be used to diagnose and prevent or correct most errors.

The Load Simulation Tool is also available for sites to use for pre-planning purposes. The software is used by technology coordinators to perform load simulation tests that help estimate the amount of time needed to download tests and upload responses based on the number of students testing at the same time, the current network traffic, the amount of available bandwidth, and other site-specific factors.

The TSM software features Load Balancing, which allow the ability to monitor content caching availability. Load Balancing solutions also allow a district to quickly add or remove TSM servers when required without reconfiguring testing clients or redirecting or reassigning addresses. This tool also allows for an easier method to distribute testers between servers; each testing client is not dependent on a single TSM server having enough capacity.

Prior to an assessment’s operational use, DRC’s quality assurance staff perform full system-level tests in an independent test environment that simulates the production configuration. Tests are run on all supported computer platforms and browsers and include a comprehensive review of system functionality, usability, reliability, security, and overall performance. Test content is also validated during this process.

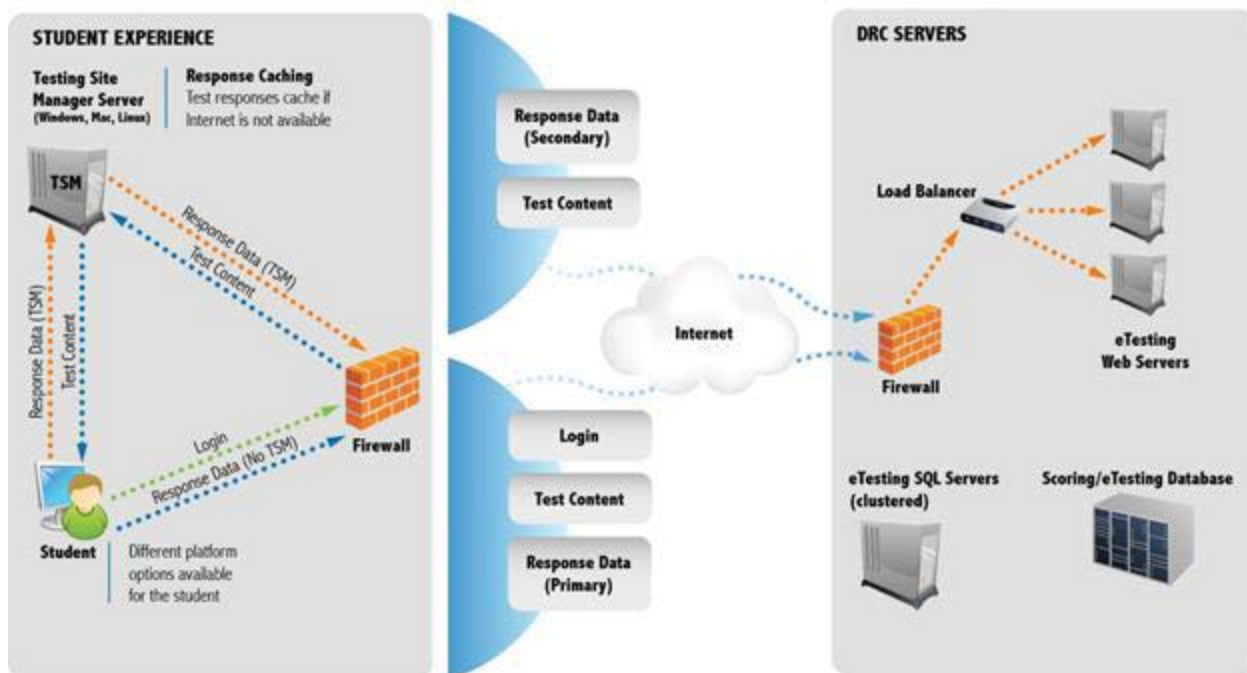
Multiple methods are used to ensure secure data transfer, including encryption technologies and Secure Sockets Layer (SSL) protocol through Hypertext Transfer Protocol Secure (HTTPS). Test content is encrypted at the host server and remains encrypted throughout all network transmissions; content is decrypted only after the student login is validated. Decrypted test content on the student workstation is stored in memory only during each test session. After the session has ended (i.e., the test is completed, or the student logs out), computer memory is

purged to ensure the security of test content.

During testing, responses are sent to a DRC server each time the student navigates away from an item or clicks the “Next” button to submit an answer. Responses are saved automatically every 45 seconds during testing, when the student navigates away from an item, or when the student answers a selected-response item, depending on whichever comes first. If an item takes the student longer than 45 seconds to answer, then the partial, incomplete response is submitted at 45-second intervals until the student completes the item. This autosave helps safeguard against students losing their work on longer items, such as Passage-Based items. When the student returns to the test after a break or interruption, the student is returned to the point at which he or she left off without having to navigate through all previously answered questions.

Figure 4-2 illustrates the secure transfer of online test responses between the student and DRC.

Figure 4-2. Architecture of the Student Testing Experience



4.5 Test Administration Training

All staff involved in the administration of M-STEP are required to receive training based on the role they will serve during the test administration. Districts provide training for Building Assessment Coordinators, and districts or Building Assessment Coordinators provide training for Test Administrators and Proctors. MDE provides test administration training resources for District Assessment Coordinators, Building Assessment Coordinators, and Test Administrators.

DRC, in conjunction with MDE, held a WebEx training presentation on March 6, 2018, with the District and Building Coordinators and Test Administrators. The presentation included pertinent information for all M-STEP online testing. The presentation was recorded and posted to eDIRECT for Michigan users to reference throughout the testing window.

MDE held a New Assessment Coordinator Preconference Workshop for both paper/pencil and online M-STEP administrations at the 2018 Michigan School Testing Conference on February 13, 2018. This presentation provided detailed information for new assessment coordinators administering both the paper/pencil assessment and the online assessment. This training was structured into before-, during-, and after-testing activities and included the following:

- Before Testing
 - Universal Tools, Designated Supports, and Accommodations
 - Pre-identification of students
 - Materials ordering
 - Providing training to test administrators and proctors
 - Scratch paper and calculator policies
 - How to prepare students for testing (M-STEP tutorials, Online Tools Training (OTTs))
 - Off-Site testing requirements and requests
 - eDIRECT training
 - Test security and the Assessment Integrity Guide (AIG)
 - Test materials and handling of secure materials
 - Test schedules and test session setup
 - How to address a testing irregularity
- During Testing
 - Test directions
 - Testing irregularities
 - Active monitoring during testing
 - Materials allowed/not allowed in a test session
- After Testing
 - Materials return
 - Preliminary reports
 - Data files
 - Final reports

MDE also provided three webcasts with accompanying PowerPoint presentations organized

into sections that discuss what administrators should do before, during, and after M-STEP administration. These webcasts followed the format used in the New Assessment Coordinator Workshop. These presentations are available on the MDE [YouTube channel](#).²

Training materials are provided to districts to use for training purposes. These materials include the following:

- M-STEP Test Administration Manual (TAM)
- Secure Site training and resource materials—provide training on pre-identification for testing, materials ordering, student scores and reporting, and using each function in the OEAA Secure Site ([OEAA Secure Site](#))
- Test Directions—offered for each test mode (online and paper/pencil) and for each grade
- M-STEP List of Important Dates
- Supports and Accommodations Guidance Document
- eDIRECT mini-modules—provide training for all functions used in eDIRECT
- eDIRECT User Guide
- INSIGHT Tools poster—displays the tools available for students and describes how to use each tool
- AIG
- Scratch Paper Policy
- Calculator Policy

All these materials were available to schools during the 2017–2018 academic year on the [M-STEP Home Page](#).

Additionally, OEAA publishes a weekly online newsletter called “Spotlight on Assessment and Accountability” throughout the year. The newsletter takes a two-week break in late December/early January, and no issues are published in July. The newsletter provides districts and schools with timely information regarding the M-STEP assessments and test administration, including training opportunities, document availability, and date reminders.

4.6 Test Security

The primary goal of test security is to protect the integrity of the assessment and to ensure that results are accurate and meaningful. OEAA uses four test security goals to maintain the integrity of the State of Michigan Assessment System. These goals are

1. to provide secure assessments that result in valid and reliable scores,
2. to adhere to high professional test administration standards,
3. to maintain consistency across all testing occasions and sites, and
4. to protect the investment of resources, time, and energy.

² https://www.youtube.com/channel/UC7cyZmw_5Q6_5bkfXDilquA

4.6.1 Prevention

Prevention of breaches in test security includes following standards and best practices for test integrity and ensuring security aspects of the design, development, operation, and administration of M-STEP are met to prevent irregularities from occurring. Operational and administrative security policies and procedures apply to both online and paper/pencil test administrations.

Online testing uses DRC's INSIGHT Online Learning System. This is a secure browser that locks the student into the testing environment, preventing access to other applications or websites. The software must be installed on each device used for testing. Test content is held securely in a TSM, which is an encrypted local cache. The TSM also provides backup response storage in the event of network issues. All students are assigned to test sessions and require an individual test ticket for every online test session. Each ticket has a username and a unique password. Access to test tickets is controlled through DRC's eDIRECT site, and eDIRECT access is controlled through locally administered permissions in the OEAA Secure Site.

For the paper/pencil test administration, OEAA and Measurement Incorporated design forms to assist the district and building assessment coordinators with the successful receipt and return of test materials. These forms provide security and accountability during fulfillment and distribution, test administration, and collection processes. Secure packaging and distribution of materials for M-STEP are provided to ensure prompt, accurate, and secure delivery of test materials to districts and schools. All materials that contain test questions or student responses are considered secure materials and must be handled in a way that maintains their security before, during, and after testing. Handling of secure materials for paper/pencil and online testing is discussed at length in Chapter 5. As part of professional test administration practices, OEAA provides test security resources for state, district, and school personnel to use in the prevention of testing irregularities. These include the AIG, TAM, online and paper/pencil administration directions, test security training modules, and incident reporting guides.

All school staff members involved in testing are required to be trained in test administration and security prior to the opening of the assessment window. Training resources are available on a statewide basis. Districts and schools can customize trainings by role and location, using state-provided materials and including local plans. The AIG is intended to be used by districts and schools in the fair and appropriate administration of state assessments. It includes guidelines on the expected professional conduct of educators who administer state assessments to ensure proper test administration and academic integrity. Four assessment security training modules are available as a supplement to the AIG. The modules are intended to be used as an online training program for district and building assessment coordinators, test administrators, and test proctors. These modules explain why test security is important, describe different staff roles in test administration, and detail how to plan for and handle incidents that compromise test security. The M-STEP TAM helps staff administering the assessment understand the administration process, key dates for specific assessment activities, the roles of school personnel in the administration process, and the ways to use available Universal Tools, Designated Supports and Accommodations. Test administrators have online and paper/pencil test directions to follow when administering M-STEP. District assessment coordinators are required to file an incident report in the case of any testing irregularity. The incident reports are filed on the OEAA Secure Site. The test security specialist and other MDE assessment

administrative staff review the incidents and determine what the required remediation will be through the use of internal and independent investigations.

4.6.2 Detection

Detection practices include guidelines for assessment monitoring, testing, and reporting irregularities. Detection resources and practices include the AIG, incident reporting, random/targeted test administration monitoring, administration observation, Universal Tools, Designated Supports and Accommodations monitoring, social media monitoring, and data forensic analysis. Districts are instructed to monitor test sessions for proper test administration and to enforce the policies and guidelines in the AIG to promote fair, approved, and standardized practices. OEAA uses random and targeted assessment monitoring to ensure the security and confidentiality of state assessments and to ensure testing personnel adhere to proper procedures. Targeted assessment monitoring is used when schools have had a previous irregularity or show unusual results from previous state assessment data analyses. Random assessment monitoring uses a sample of schools that are randomly selected for quality and integrity checks. Specific requirements of assessment monitoring are documented in the *Assessment Observation Requirements Document* created with OEAA's vendor Measurement Incorporated. The AIG details the process for monitoring district and school personnel. Internet and media monitoring occurs during testing windows. The goal of this monitoring is to combat breaches and disclosure of secure assessment materials. These monitoring activities include monitoring comments on the internet for test items captured and shared either from testing computer screens or from paper/pencil test booklets. Social media sites are also monitored for posts discussing or exposing test material. Requirements for social media monitoring are documented in the *Social Media Monitoring Requirements Document* created with OEAA's vendor Measurement Incorporated. The AIG details the process for monitoring the social media sites of district and school personnel.

DRC provides MDE with online forensic telemetry data via a secure table data load. The table below references the data that are captured and sent to MDE on a weekly basis during the testing windows.

Table 4-3. INSIGHT Forensic Data

Attribute of Forensic Data	Description
Test Interrupted Stopped Flag	Test was interrupted/stopped
Test Interrupted Stopped Count	Number of times the test was interrupted/stopped
Total Item Time	Total time spent on an item
Item Visit Count	Total number of times the item was visited
Wrong to Right	Item's response was changed from wrong to right (within or across item visits)
Wrong to Right Count	Total number of times the item's response was changed from wrong to right (within or across item visits)
Right to Wrong	Item's response was changed from right to wrong (within or across item visits).
Right to Wrong Count	Total number of times the item's response was changed from right to wrong (within or across item visits)
Wrong to Wrong	Item's response was changed from wrong to wrong (within or across item visits).
Wrong to Wrong Item Count	Total number of times the item's response was changed from wrong to wrong (within or across item visits)
Total Enters Net Total Exits	Records total enters are greater than or less than total exits.

During and after online and paper/pencil test administrations, OEAA conducts multiple analyses on student assessment results. These statistical analyses help in flagging potential testing irregularities. The types of data forensic analyses conducted in spring 2018 included unusual score gains and losses, online right-to-wrong changes, proficiency level gains, occurrence of perfect scores, and response time analysis.

4.6.3 Investigation and Remediation

District assessment coordinators are required to notify OEAA as soon as they are made aware of an alleged or suspected violation or misadministration of M-STEP. Testing irregularities are reported to OEAA via an online incident report form. The M-STEP TAM and AIG provide an incident reporting guide for districts and schools. Each testing irregularity report is reviewed by MDE test administration staff, and corrective action is taken based on policies and procedures for test administration outlined in the M-STEP TAM and AIG.

OEAA also has a phone and online “tip line” to report unethical behavior. Reports can be made anonymously. This provides a means for school staff members to report test integrity issues within their chain of command when they do not feel comfortable reporting the issues to their superior.

All incident reports and supporting documentation are reviewed by MDE, and a determination is made regarding the disposition of each incident. If OEAA determines that the irregularity caused no consequences affecting security, validity, or fraud and that the school took appropriate actions to correct the situation, OEAA may consider the issue resolved and log or close the case. If OEAA determines that questions remain regarding the security, validity, or authenticity of

the test administration, OEAA will request either a school self-investigation or, if the problem is considered potentially severe, an independent investigation.

After investigations have taken place, OEAA will create a summary report of the findings. Determination of the investigation is provided in the report.

Remediation of the incidents reported and investigated differ based on the severity of a confirmed allegation or misadministration. Minor mistakes receive recommendations of best practices. Isolated security incidents or negligence provide good candidates for targeted monitoring the next year. Individual student tests tainted by misadministration are typically invalidated. More serious incidents can lead to invalidating entire classes of tests, retraining staff, or barring staff from participating in statewide testing. When possible, remediation happens within the testing window so that students with invalidated tests can be retested if appropriate. Further information on remediation is available in the AIG.

4.7 Summary of M-STEP Administration Best Practices

The elements discussed in previous sections align with MDE's prevention practices that help maintain the integrity of the assessment and also adhere to the testing practices and AERA, APA, & NCME (2014) *Standards* relevant to test administration. The previous sections also demonstrate how information in the MDE trainings and manuals addresses *Standards* 4.15 and 6.1:

Standard 4.15 The directions for test administration should be presented with sufficient clarity so that it is possible for others to replicate the administration conditions under which the data on reliability, validity, and (where appropriate) norms were obtained. Allowable variations in administration procedures should be clearly described. The process for reviewing requests for additional testing variations should also be documented. (p. 90)

The M-STEP TAM, Test Directions, and AIG provide instructions for before-, during-, and after-testing activities with sufficient detail and clarity to support reliable test administrations by qualified test administrators. To ensure uniform administration conditions throughout the state, instructions in the TAM, Test Directions, and AIG describe the following: general rules of online and paper/pencil testing; pause and break rules; test scheduling; assessment duration, timing, and sequencing information and recommendations; handling of secure materials; and materials that the test administrator and students need for testing.

Standard 6.1 Test administrators should follow carefully the standardized procedures for administration and scoring specified by the test developer and any instructions from the test user. (p. 114)

To ensure the usefulness and interpretability of test scores and to minimize sources of construct-irrelevant variance, it is essential that M-STEP is administered according to the directions provided in the TAM, Test Directions, and AIG.

MDE's protocol, discussed in Section 4.6, stresses incident reporting and adheres to *Standards* 6.3, 6.6, and 6.7.

Standard 6.3 Changes or disruptions to standardized test administration procedures or scoring should be documented and reported to the test user. (p. 115)

Incident reporting by district assessment coordinators is required when there is any type of misadministration or problem with test administration. MDE provides an Incident Reporting Guide within the TAM that details incidence categories and subcategories that are used in the Secure Site Reporting tool and provides sample scenarios for each category or subcategory. MDE staff review the incident reports and respond with corrective action when appropriate.

Standard 6.6 Reasonable efforts should be made to ensure the integrity of test scores by eliminating opportunities for test takers to attain scores by fraudulent or deceptive means. (p. 116)

MDE requires that the testing environment for both online and paper/pencil testing be conducive to a proper test environment. All information regarding the content being measured or test-taking strategies displayed in the testing room must be removed or covered. Students must be seated so there is enough space between them to minimize opportunities to view each other's work. Test administrators and Proctors are encouraged to frequently move through the testing room and monitor the students' work areas during testing. Only staff involved in administering the test and students taking the test can be in the testing room. Students are not permitted to access any electronic devices used for communication, for capturing images of the test or testing room, or for data storage during testing. Testing materials are required to be kept secure at all times before, during, and after the testing sessions. Certain secure materials are required to be returned to Measurement Incorporated or securely destroyed.

Standard 6.7 Test users have the responsibility of protecting the security of test materials at all times. (p. 117)

The AIG and TAM describe the ethical practices that testing staff and students must follow during test administration. Students are reminded at the start of the testing session that in order for their results to be valid, they must not talk to or help other students; look at or copy other students' answers; ask for or accept any help from other students; use their cell phones or any other electronic devices, including an eBook; take pictures or make copies of any test materials; cause a disturbance; remove a test booklet, test ticket, or answer document from the room; or post or chat about any part of the test through social media. All staff who participate in a state assessment or handle secure assessment materials must be fully trained and sign an Assessment Security Compliance Form. By signing the Assessment Security Compliance Form, staff certify that they will follow test administration directions, maintain security and confidentiality of the tests, and report any suspected violations of test security.

Throughout the manuals, test coordinators and examiners are reminded of test security requirements and procedures to maintain test security. Specific actions that are direct violations of test security are so noted. Detailed information about test security procedures are presented in Section 4.6.

4.8 Test Materials

A list of available test materials can be found below in Table 4-4.

Table 4-4. M-STEP Paper Test Materials

Material Description	Product Type
Blank Labels	Ancillary
DVD Information Sheet	Ancillary
FedEx Return Air Bills	Ancillary
Instruction for Materials Return	Ancillary
OEAA Security Compliance Form	Ancillary
Outgoing Box Labels (M-STEP Materials Label)	Ancillary
Packing List Enclosed Label	Ancillary
PreID Labels	Ancillary
Return Kit Cover Sheet	Ancillary
Scorable Labels	Ancillary
Special Handling Envelopes	Ancillary
ELA Answer Document	Answer Document
ELA Emergency Answer Document	Answer Document
Mathematics Answer Document	Answer Document
Mathematics Emergency Answer Document	Answer Document
Social Studies Answer Document	Answer Document
ELA AABB	Braille
ELA Braille—Contracted Test Booklet	Braille
ELA Braille—Uncontracted Test Booklet	Braille
ELA Braille—Uncontracted Print to Braille Correspondence Document	Braille
Mathematics AABB	Braille
Mathematics Braille—Contracted Test Booklet	Braille
Mathematics Braille—Uncontracted Test Booklet	Braille
Mathematics Braille—Uncontracted Print to Braille Correspondence Document	Braille
Social Studies AABB	Braille
Social Studies-Contracted Braille Test Booklet	Braille
Social Studies-Uncontracted Braille Test Booklet	Braille
Social Studies—Uncontracted Print to Braille Correspondence Document	Braille
ELA Listening Audio CD	CD
Social Studies Audio CD	CD
Social Studies Arabic DVD	DVD
Social Studies English DVD	DVD

Chapter 4: Test Administration Plan

Material Description	Product Type
Social Studies Spanish DVD	DVD
ELA Enlarged Print Test Booklet	Enlarged Print
Mathematics Enlarged Print Test Booklet	Enlarged Print
Social Studies Enlarged Print Test Booklet	Enlarged Print
Glossary Reference Sheets	Glossary
Graph Paper	Graph Paper
ELA Listening Script	Listening Script
ELA Listening Script, Emergency	Listening Script
M-STEP Test Administration Manual	Manual
M-STEP Paper/Pencil Test Directions	Manual
M-STEP Online Test Directions	Manual
M-STEP Emergency Test Administration Directions Addendum	Manual
Social Studies Emergency Reader Script (English)	Reader Script
Social Studies Reader Script (English)	Reader Script
ELA Emergency Test Booklet	Test Booklet
ELA Test Booklet	Test Booklet
Mathematics Emergency Test Booklet	Test Booklet
Mathematics Spanish Test Booklet	Test Booklet
Mathematics Test Booklet	Test Booklet
Social Studies Test Booklet	Test Booklet
Social Studies Emergency Test Booklet	Test Booklet
Math Test Booklet	Test Booklet
Science Test Booklet	Test Booklet
Social Studies Test Booklet	Test Booklet
Social Studies Emergency Test Booklet	Test Booklet

4.9 Summary

In summary, the overall purpose of the test administration documentation and training opportunities is to keep districts informed about policies and procedures related to testing in general and the M-STEP program. The information imparted is clearly related to maintaining the integrity of the administration of M-STEP, maintaining the security of the assessment, allowing access to the assessments for special populations by clearly delineating appropriate Universal Tools, Designated Supports or Accommodations, and providing guidance on appropriate interpretations of the test results. These communication and training efforts by MDE and its test vendors are in alignment with multiple best practices of the testing industry, particularly the following standards (AERA, APA, & NCME, 2014):

- Standard 4.15—The directions for test administration should be presented with sufficient clarity so that it is possible for others to replicate the administration conditions under which the data on reliability, validity, and (where appropriate) norms were obtained. Allowable variations in administration procedures should be clearly described. The process for reviewing requests for additional testing variations should also be documented.
- Standard 6.1—Test administrators should follow carefully the standardized procedures for administration and scoring specified by the test developer and any instructions from the test user.
- Standard 6.2—When formal procedures have been established for requesting and receiving accommodations, test takers should be informed of these procedures in advance of testing.
- Standard 6.3—Changes or disruptions to standardized test administration procedures or scoring should be documented and reported to the test user.
- Standard 6.6—Reasonable efforts should be made to ensure the integrity of test scores by eliminating opportunities for test takers to attain scores by fraudulent or deceptive means.
- Standard 6.7—Test users have the responsibility of protecting the security of test materials at all times.

Chapter 5: Test Delivery and Administration

5.1 Online Administration Details

In conjunction with DRC, MDE delivered 99% of M-STEP online via DRC's online testing platform, INSIGHT, in spring 2018 when 854 Michigan school districts administered M-STEP online to 3,155 Michigan schools.

For the third consecutive administration, M-STEP ELA and mathematics were administered as computer adaptive tests (CATs). M-STEP social studies was administered as a fixed form, the same as in spring 2017. Additionally, accommodated forms in all content areas were delivered as fixed-form assessments.

The spring 2018 M-STEP was administered to enrolled students in grades 3–8 and 11, along with grade 12 students who missed testing in 2017. Table 5-1 presents content areas tested by grade.

Table 5-1. Content Areas Tested by Grade

Grade Tested	Content Areas Tested
Grade 3	ELA and Mathematics
Grade 4	ELA and Mathematics
Grade 5	ELA, Mathematics, and Social Studies
Grade 6	ELA and Mathematics
Grade 7	ELA and Mathematics
Grade 8	ELA, Mathematics, and Social Studies
Grade 11 & 12	Social Studies

The number of students tested online for the spring 2018 M-STEP can be found in Table 5-2 below.

Table 5-2. Number of Students Tested Online

Grade	Subject	Online Students Tested
3	ELA	101,461
4	ELA	103,968
5	ELA	107,987
6	ELA	107,846
7	ELA	107,334
8	ELA	109,751
3	Mathematics	101,739
4	Mathematics	104,188
5	Mathematics	108,131

Grade	Subject	Online Students Tested
6	Mathematics	107,950
7	Mathematics	107,390
8	Mathematics	109,767
5	Social Studies	107,951
8	Social Studies	109,537
11 & 12	Social Studies	107,567

5.1.1 Online Administration Reports

DRC and MDE outlined requirements for all online administration reporting prior to administering the 2018 assessments. Administration reports were delivered to MDE daily or weekly based on the established requirements. Table 5-3 shows the types of administration reports that were delivered to MDE during the 2018 M-STEP testing windows.

Table 5-3. Online Administration Reports

Report Name	Delivery Frequency	Description of Report
After-Hours Report	Daily throughout the testing window	Shows online tests that have test login times and/or stop times within the defined after-hours time
Form Distribution Report	Weekly throughout the testing window	Shows fixed-form assignments for monitoring equal distribution of fixed forms per grade and content area
Testing Times Report	Daily throughout the testing window	Daily summary of testing times to allow MDE to monitor how long students take to complete tests
Cumulative Student Status	Daily throughout the testing window	Status of student testing by site; allows MDE to monitor how students are progressing with testing by grade and content area
Excessive Logins Report	Daily throughout the testing window	Shows online tests that have been logged into more than four times

5.1.2 Online User Manuals and Reference Documents

To help assist with the administration of the online M-STEP, numerous manuals and documents were created. Some of these include the test administration manuals, online test directions by grade, and the Technology User Guide, as well as many additional reference documents.

The M-STEP Test Administration Manual (TAM) is available for all test modes, grades, and content areas of M-STEP tests. It provides an overview of the assessments, important testing dates, information on when and how to assign and use Universal Tools, Designated Supports and Accommodations, guidelines on who must test, testing policies and procedures including scratch paper and calculator policies, and resources for assessment coordinators and administrators.

Chapter 5: Test Delivery and Administration

The TAM provides detailed information regarding the roles involved in administering a test and responsibilities for each role: District Assessment Coordinator, Building Assessment Coordinator, and Test Administrator.

Information provided in the M-STEP TAM includes the following:

- Important Dates
 - Paper/Pencil Testing Dates
 - Online Testing Windows
 - Important Pre-testing Activities
 - Materials Return Dates
- Overview
 - M-STEP Assessments
 - Scratch Paper Guidelines
 - Supports and Accommodations
 - Resources for Students to Prepare for Testing
 - M-STEP Roles and Responsibilities
 - Valid, Equitable, and Ethical Assessment
 - OEAA Security Compliance Forms
 - Students to Be Tested
 - Incident Reporting
 - Testing Irregularities
 - Training Tools
 - Security
 - Materials Permitted or Required during Testing
- District Assessment Coordinators
 - Roles and Responsibilities
 - Training Requirements and Resources
- Building Assessment Coordinators
 - Roles and Responsibilities
 - Training Requirements and Resources
- Test Administrators
 - Roles and Responsibilities
 - Test Administrator Resources
 - Managing Test Sessions
 - Test Directions
 - Ending a Test Session
- Supports and Accommodations
 - What Are Supports and Accommodations?
 - Ordering Accommodated Materials
 - Supports and Accommodations Tracking Sheet
 - Where to Find More Information on Supports and Accommodations
 - Nonstandard Accommodations

- Read-Aloud Guidelines
 - Scribing Protocol
- Materials Return Instructions
- Appendices:
 - Calculator Policy
 - Scratch Paper Policy
 - Incident Reporting Guide
 - Important Dates
 - M-STEP Administration Resources

Online Test Directions documents are provided for each grade level. These documents provide information for testing including materials needed during testing, items permitted in testing rooms, test scratch paper and calculator policy information, and test directions for each content area test in the grade level.

Information provided in the Online Test Directions for each grade includes the following:

- Introduction
 - Key
 - Online Tools Training (OTT) and Student Tutorials
- Before Testing
 - Test Materials Needed for M-STEP
 - Before Testing Checklist
- During Testing
 - Permitted Items in Testing Room
 - Procedures for Testing Breaks, Interruptions, or Pauses
 - Test Directions – Introduction
 - Test Sign-In
 - Welcome Screen
 - System Check and Test Security
 - Introduction
 - Answering Questions
 - Test Directions by grade and content area, including directions for accommodated assessments
- After Testing
 - Completing the Test Session
 - Exiting the Test Engine

5.2 Paper/Pencil Administration Details

MDE delivered paper/pencil assessments to meet individual students' needs and for buildings that applied and were approved for a waiver of online testing.

Online testing waivers were available for the following reasons:

- Buildings that were not technologically ready
- Buildings that were under construction or otherwise had a disrupted technological environment
- Locations testing a center-based program
- Locations testing in a juvenile justice facility
- Buildings that had other instructional obstacles

The paper/pencil test was available in Enlarged Print and in both contracted and uncontracted braille versions. A Spanish language paper/pencil test was also available for mathematics in each grade.

There were three forms for each test, including the braille form. These forms are listed in the table below.

Table 5-4. Paper/Pencil Test Forms by Content Area

Content Area	Paper Pencil Forms Available
ELA	Form 1—administered to all students testing paper/pencil
ELA	Form 2—Emergency form
ELA	Form 88—Braille form
Mathematics	Form 1—administered to all students testing paper/pencil
Mathematics	Form 2—Emergency form
Mathematics	Form 88 - Braille form
Social Studies	Form 1—administered to all students testing paper/pencil
Social Studies	Form 2—Emergency form
Social Studies	Form 88—Braille form

The paper/pencil test was provided for the same grades and content areas that had online counterparts (see Table 5-1).

The M-STEP Test Administration Manual (TAM) is common for all test modes, grades, and content area M-STEP tests. It provides an overview of the assessments, important testing dates, information on when and how to assign and use Universal Tools, Designated Supports and Accommodations, guidelines on who must test, testing policies and procedures including scratch paper and calculator policies, and resources for assessment coordinators and administrators. See Section 5.1.2 for information about the content included in the M-STEP

TAM.

Paper/Pencil Test Directions documents are provided for each grade level. These documents provide information for testing including materials needed during testing, items permitted in testing rooms, test scratch paper and calculator policy information, and test directions for each content area test in the grade level.

Information provided in the Paper/Pencil Test Administrations Directions for each grade includes the following:

- Paper/Pencil Test Schedule
- Introduction
 - Test Security
 - Establishing Appropriate Testing Conditions
 - Food, Drink, Snacks
 - Requirements of Test Environment
- Pre-identification Label Directions
- Student Data Grid Information and Administration Directions
 - Directions for Completing the Student Demographic Page
 - Administration Directions for Completing the Student Data Grid
- Test Directions by Content Area for Each Grade
 - Participation of Students with Disabilities and/or English Learners
 - Scratch Paper Policy
 - Calculator Policy for the Grade Tested
 - Testing Times and Schedules
 - General Rules for the Paper/Pencil Assessment
- After Testing
 - Assemble Materials for Return
 - Test Administrator Checklist

The number of students tested using the spring 2018 paper/pencil M-STEP can be found in the table below.

Table 5-5. Number of Students Tested with Paper/Pencil

Grade	Content Area	Number of Students Tested with Paper/Pencil
3	ELA	788
4	ELA	884
5	ELA	870
6	ELA	789
7	ELA	554
8	ELA	592
3	Mathematics	848
4	Mathematics	925
5	Mathematics	926
6	Mathematics	839
7	Mathematics	610
8	Mathematics	616
5	Social Studies	923
8	Social Studies	654
11 & 12	Social Studies	448

5.3 OEAA Secure Site

The OEAA Secure Site is a web-based application used for state assessments and accountability. The primary functions of the Secure Site include pre-identification of students for both paper/pencil and online assessments; ordering paper/pencil materials (see Chapter 4.8), including accommodated versions of the assessments; incident reporting; review of accountable students and test verification; and retrieval of data score files and score reports.

The Secure Site is available to authorized district and school personnel only. The [MDE Secure Site training page](#)¹ includes a complete list of Secure Site functions and how to use them.

The Secure Site takes student information from the Michigan Student Data System (MSDS) and provides a secure, centralized interface for using student information in testing. District-identified permissions for school staff carry through to vendor systems. To prepare for testing, students can be pre-identified for tests, moved between general and alternate assessments, arranged in online testing sessions, and have test materials ordered. During testing, the Secure Site provides tools for correcting issues missed during preparation and reporting incidents. After testing, the Secure Site provides access to student test score reports and data files, as well as providing tools for accountability and test verification.

¹ https://www.michigan.gov/mde/0,4615,7-140-22709_57003---,00.html

5.4 eDIRECT

5.4.1 Michigan Users

DRC uses MDE's Secure Site to pull and load Michigan users to eDIRECT based on Secure Site Test Cycle IDs. For the 2017–18 school year, the M-STEP *Test Cycle ID* was 157. Users were identified by their *Security Role ID* and pulled into eDIRECT according to the established requirements. The mapping of users from the Secure Site to eDIRECT can be found below in Table 5-6.

Table 5-6. Mapping of Building Users from Secure Site to eDIRECT

Security Role ID	eDIRECT Role and Permission Set
17—Public School Administrator	School
20—District Administrator	School
40—Public Online Test Administrator	School
31—Nonpublic School Administrator	School
41—Private School Online Test Administrator	School
42—District Test Administrator	School
45—State	State
38—District Technology Coordinator	District Technology Coordinator
39—School Technology	District Technology Coordinator
43—Public School Technology	District Technology Coordinator
44—Private School Technology	District Technology Coordinator

All users were identified by the site code(s) they had access to within eDIRECT. Users were only able to access student and test information based on their site permissions in the MDE Secure Site.

5.4.2 Administrative Functions

Online administration is managed through the DRC eDIRECT client portal that provides tiered, secure access to all required administrative functions. Within eDIRECT, users manage student information and create test sessions.

Student information for M-STEP is imported into eDIRECT via automatic loading of data. DRC utilizes the MDE Secure Site to pull new and updated student records for import into eDIRECT. Student data is pulled three times a day so that any new student records or updated student records are loaded in a timely manner. Building users are able to view all the demographic information associated with the students from the Secure Site before placing them in test sessions for test tickets.

Once the student data is loaded into the Test Setup application within eDIRECT, users organize students into test sessions. Test sessions can be created by content area, class, grade, or school. Through Test Setup, users can also update student Designated Support/Accommodation information, print test tickets, and monitor student testing status.

The student login ticket contains unique login credentials used by the student to access the testing software. For a selected test session, users can download and print a PDF document containing instructions, a roster of student tickets, and the actual test tickets. Student test tickets are considered secure materials, and test administrators are required to keep printed tickets in a predetermined, locked, secure storage area.

5.4.3 Online Testing Resources

eDIRECT houses an assortment of testing resources available to the district and school users as well as the technology coordinators. The INSIGHT installables and requirements are maintained on eDIRECT, as are all technology guides and information necessary for setting up schools' computers and servers.

Video tutorials containing mini-chapters on how to use eDIRECT applications are available to help users familiarize themselves with the different administrative applications within eDIRECT. An eDIRECT user guide is also available for reference.

For more information on M-STEP-specific online testing resources, visit the [MDE website](#).²

5.5 Return Material Processing

Each box of materials shipped to schools contains a box list, which shows each item in the box. Each order contains a packing list, which shows a complete list of items, quantities, and box locations for the entire order. When an order contains secure materials, a security list is also included that shows a complete list of secure items and the associated shrink-wrapped pack barcodes.

All M-STEP scorable and non-scorable secure testing materials are to be returned via FedEx Express Saver to Measurement Incorporated to be processed.

When boxes of returned materials arrive at Measurement Incorporated, the warehouse team scans the boxes into the Measurement Incorporated tracking system database, where they are checked against the tracking numbers that are assigned to each school. FedEx also scans each of its tracking barcodes to record each box as it was delivered to Measurement Incorporated. This provides immediate information on the number of boxes received and points of origin of the boxes. Once this procedure is completed, the boxes are opened and all materials are sorted.

Scorable and non-scorable materials are securely scanned in using Measurement Incorporated's Security Barcode Check-In Application. This application allows IT Operations to scan the security identifier on individual secure materials or the security identifier located on the outside of an intact pack of shrink-wrapped documents using Measurement Incorporated's automated security scanning process. Scanning the security identifier on the shrink-wrapped pack is equivalent to scanning all the individual security identifiers included in the shrink-wrapped pack and is more efficient than scanning each individual test booklet in the shrink-wrapped pack.

As each security identifier is securely scanned, it is checked against the original list of identifiers

² <http://www.michigan.gov/mde/>

that were entered into the Measurement Incorporated database. Any discrepancies are noted, and a security report is generated for MDE.

For scorable answer documents, the same scanning process that captured the security identifier information also captures information from the student pre-identification label, bubbled demographic information on the answer document cover, bubbled student responses, and images of constructed responses to be sent on to handscoring.

All loose (i.e., individual) test booklets are securely scanned into the Measurement Incorporated database by IT Operations using Measurement Incorporated's automated security scanners.

Warehouse personnel securely scan in all returned accommodated materials using a human-operated computer station equipped with a barcode reader and entered those materials into the ObjectTracker database.

The accommodated materials include CDs, DVDs, braille test booklets, Enlarged Print test booklets, and Reader Scripts. Although they are not accommodated materials, ELA Listening CDs and Reader Scripts for M-STEP are also scanned in.

After all returned secure materials are checked in, Measurement Incorporated's IT team prepares the initial security report data by comparing the security barcodes of checked-in materials with the barcodes of all secure materials.

The initial missing materials and security report data are provided to MDE in a spreadsheet. All schools that were sent materials by Measurement Incorporated are included in the summary, regardless of whether the schools are active or inactive entities.

For public school districts that are missing secure materials, district coordinators are shipped security reports to be further distributed to building coordinators.

For public school academies and nonpublic schools that are missing secure materials, each building coordinator is shipped a security report.

Missing materials reported as destroyed or never received are not included on the security report sent to the district or school. Missing materials reported as lost remain on the security report, and the comment "Reported Lost" is added to the comment section of the security report.

FedEx Ground Package Returns Program labels are provided in case any secure materials need to be returned. Schools that find no additional secure materials are directed to return the summaries of missing secure materials and any additional information.

The Measurement Incorporated IT team updates the security report data using the spreadsheet of issues reported to the Call Center, which includes materials that were lost, destroyed, or never received. This spreadsheet is maintained by the Measurement Incorporated management team. MDE staff forwards to the Measurement Incorporated management team any information collected via phone calls or incident reports regarding materials that were lost, destroyed, or never received.

Chapter 5: Test Delivery and Administration

If a summary of missing secure materials is accompanied by a corresponding explanation letter, the two are stapled together. All summaries of missing secure materials are checked in using the district/building code barcode and are filed in order by assessment, district code, and building code. Any returned secure materials are checked in by security barcode and are stored with the other secure materials.

After the initial response window ends and the returned letters and secure materials are processed, the IT team refreshes the security report data for each assessment, indicating schools that responded with newly returned secure materials and/or letters and schools that did not respond. Follow-up security reports are generated.

A second round of cover letters and security reports is sent to districts and schools that still have outstanding missing materials and have not returned a letter or a security report with comments. This procedure is the same as the ones used for the first round of security reports. Schools that return a letter, materials, or both in the first round are not included in the second round.

Measurement Incorporated checks in and files any returned summaries of missing secure materials, secure materials, and additional information received. When MDE determines that schools have had sufficient time to respond, Measurement Incorporated generates and provides to MDE a final missing materials report.

The final security report spreadsheet sent from Measurement Incorporated to MDE includes all schools and districts that were tested. The Excel filter feature is used to list those that still have outstanding missing materials. The “Returned Letter or Additional Items or Both” column reflects letters and items returned in response to both the initial round and the second round of security reports.

Tables 5-7 through 5-9 show shipped M-STEP material information. The amount of material shipped was and should be expected to be higher than the number of students testing on paper/pencil. Each student needs at least two secure materials for testing, plus some secure accommodated materials for students testing either online or on paper/pencil.

Table 5-7. Count of Secure M-STEP Materials Shipped

Grade	ELA	Mathematics	Social Studies
3	2,949	2,805	N/A
4	3,025	2,804	N/A
5	3,045	2,879	1,880
6	2,827	2,672	N/A
7	2,413	2,255	N/A
8	2,443	2,242	1,560
11	N/A	N/A	1,904

Table 5-8. Count of Secure M-STEP Materials Returned

Grade	ELA	Mathematics	Social Studies
3	2,643	2,529	N/A
4	2,808	2,619	N/A
5	2,789	2,646	1,844
6	2,669	2,500	N/A
7	2,163	2,060	N/A
8	2,280	2,089	1,444
11	N/A	N/A	1,862

Table 5-9. Count of Secure M-STEP Materials Not Returned

Grade	ELA	Mathematics	Social Studies
3	306	276	N/A
4	217	185	N/A
5	256	233	15
6	158	172	N/A
7	250	195	N/A
8	163	153	45
11	N/A	N/A	70

5.6 Testing Window and Length of Assessment

The four-week testing windows for the online 2018 operational M-STEP were as follows:

- Grades 5 and 8 were administered the ELA, mathematics, and social studies assessments from April 9 through May 4, 2018.
- Grade 11 was administered the social studies assessment from April 9 through May 4, 2018.
- Grades 3, 4, 6, and 7 were administered the ELA and mathematics assessments from April 30 through May 25, 2018.

All online accommodated and standard assessments were administered in these time frames; there were no specific make-up windows for online assessments.

Paper/pencil testing dates for grades 5 and 8 were as follows:

- ELA Days 1 and 2: April 10 and 11, 2018
- Math: April 17, 2018
- Social Studies: April 19, 2018
- Makeup days:
 - ELA: April 12, 13, and 16
 - Any content area: April 20–27, 2018

Paper/pencil testing dates for grade 11 were as follows:

- Social Studies: April 5, 2018
- Makeup days: April 13–27, 2018

Paper/pencil testing dates for grades 3, 4, 6, and 7 were as follows:

- ELA Days 1 and 2: May 1 and 2, 2018
- Math: May 8, 2018
- Makeup days:
 - ELA: May 3, 4, and 7, 2018
 - Any content area: May 9–18, 2018

The spring 2018 M-STEP was not timed and was paced by students. Schools scheduled test sessions and determined the appropriate amount of time for students to spend testing in a single test session. Any students needing more time were able to complete the test in a later test session during the four-week grade-level testing windows. Further information on test session timing is provided on pages 4–8 of the [2017–2018 Guide to State Assessments](#).

Chapter 6: Operational CAT

This chapter mainly covers elements of the CAT algorithm, including entry point, ability estimation and standard error of measurement (SEM), passage selection, test navigation, test termination, and forced submission. M-STEP CAT configurations and simulations for ELA and mathematics are reported toward the end of this chapter. Information on the Smarter Balanced Summative CAT configurations can be found in the *Smarter Balanced 2017–2018 Technical Report* (2018), and the 2018 Smarter Balanced Summative CAT simulations can be found [here](#).¹

Before a CAT administration, the configurations and the item pool need to be loaded into the CAT engine. The configurations define the operational test blueprint with different content rules (e.g., Min and Max number of items in one or more content standards and/or item types), field-test blueprint (e.g., number of items in each claim and item position), scoring algorithm (e.g., theta estimation method, scaling constants, highest obtainable scale score [HOSS], and lowest obtainable scale score [LOSS]), and passage-selection criteria (e.g., Min and Max number of items in a passage, passage Min percentage, distinct passage Min for ranking, passage ranking criteria, passage randomization options, and options to fulfill rules). The details of the configurations for each grade and content are presented after the descriptions of the processes. Note that specific information related to the psychometric background can be found in the *Smarter Balanced 2017–2018 Technical Report* (2018).

6.1 Entry Point

The M-STEP CAT algorithm for ELA and mathematics is designed to administer items targeted for an individual student based on his or her performance. However, students' performance is unknown at the beginning of the test. With no prior information about a student, DRC has determined, based on simulation studies prior to operational administrations, that using a starting point one standard deviation (SD) below the average item difficulty of the M-STEP ELA and mathematics CAT pools provides students with a better test-taking experience at the beginning of the test, particularly for those who are at the lower end of the achievement continuum. Table 6-1 lists the initial values used in the 2017–18 M-STEP CAT.

The M-STEP CAT algorithm includes a randomization component when selecting items to control item exposure. That is, one item is selected from among a set of items that is near the targeted item difficulty. This is especially important at the beginning of the test when no prior information is available. Randomization of items and rules defined by the test blueprint ensure that students will not see the same set of items in the same order even when all the students are assumed to perform the same or have the same initial theta at the beginning of the test.

¹ <https://portal.smarterbalanced.org/library/en/2016-17-summative-assessments-simulation-results.pdf>

Table 6-1. Initial Thetas for the CAT

Content	Grade	Initial Theta
ELA	3	-1.634
ELA	4	-1.262
ELA	5	-0.893
ELA	6	-0.421
ELA	7	-0.233
ELA	8	-0.159
Mathematics	3	-1.978
Mathematics	4	-1.195
Mathematics	5	-0.540
Mathematics	6	-0.312
Mathematics	7	0.600
Mathematics	8	0.560

6.2 Theta Estimates and Standard Error of Measurement

After each item response, the theta estimate and SEM are calculated via the maximum likelihood estimation (MLE) for the total test and each claim. Note that only responses to autoscored items are accounted for in the theta estimate used by the CAT algorithm. The items in the item bank are calibrated based on the Generalized Partial Credit Model (GPCM) (Muraki, 1992) (see Equation 6-1).

$$P_{im}(\theta_j) = \frac{\exp[\sum_{k=0}^m Da_i(\theta_j - b_{ik})]}{\sum_{v=0}^{M_i-1} \exp[\sum_{k=0}^v Da_i(\theta_j - b_{ik})]}, \quad (6-1)$$

where $a_i(\theta_j - b_{i0}) \equiv 0$; $P_{im}(\theta_j)$ is the probability of an examinee with ability θ_j getting score m on item i ; M_i is the number of score categories of item i with possible scores as consecutive integers from 0 to $M_i - 1$; D is the scaling constant, 1.7; a_i is the discrimination parameter of item i ; b_{ik} is the location parameter or threshold of category k . The GPCM is equivalent to the 2 Parameter Logistic (2PL) Model (Birnbaum, 1968) (see Equation 6-2) when the item is scored dichotomously.

$$P_i(\theta_j) = \frac{1}{1 + \exp[-Da_i(\theta_j - b_i)]}, \quad (6-2)$$

where $P_i(\theta_j)$ is the probability of an examinee with ability θ_j answering item i correctly; D is the scaling constant, 1.7; a_i and b_i are the discrimination and difficulty parameters of item i .

For a general MLE, the likelihood combines both dichotomously and polytomously scored items as shown below:

$$L(\theta_j | U) = \left(\prod_{i=1}^n P_i(\theta_j)^{u_i} Q_i(\theta_j)^{1-u_i} \right) \cdot \left(\prod_{i=n+1}^N \prod_{m=0}^{M_i-1} P_{im}(\theta_j)^{u_{im}} \right), \quad (6-3)$$

where $Q_i(\theta_j)$ is $1 - P_i(\theta_j)$ and the response matrix U contains the response of dichotomously scored items

$$u_i = \begin{cases} 1, & \text{if correct,} \\ 0, & \text{otherwise,} \end{cases}$$

for $i = 1, \dots, n$, and the responses of polytomously scored items

$$u_{im} = \begin{cases} 1, & \text{if scored } m, \\ 0, & \text{otherwise,} \end{cases}$$

for $i = n + 1, \dots, N$ and $m = 0, 1, \dots, M_i - 1$.

The modified version of the Newton-Raphson equation used by DRC for estimating theta at iteration t is given as below:

$$[\hat{\theta}]_t = [\hat{\theta}]_{t-1} + \frac{L'_1 + L'_2}{ABS(L''_1 + L''_2)}. \quad (6-4)$$

where ABS stands for the absolute value. L'_1 and L'_2 are the first and second derivatives of the likelihood function of polytomously scored items:

$$L'_1 = \sum_{i=1}^m D a_i (u_i - p_i) \quad \text{and} \quad (6-5)$$

$$L''_1 = \sum_{i=1}^m \frac{D^2 a_i^2 (-p_i^2)(1 - p_i)}{p_i}, \quad (6-6)$$

where u_i is the score a student gets from a dichotomously scored item and the possible values are 1 or 0. L'_2 and L''_2 are the first and second derivative of the likelihood function of polytomously scored items:

$$L'_2 = \sum_{i=n+1}^N D a_i \sum_{m=0}^{M_i-1} u_{im} \left(m - \sum_{m=0}^{M_i-1} m P_{im}(\theta_j) \right) \quad \text{and} \quad (6-7)$$

$$L''_2 = - \sum_{i=n+1}^N D^2 a_i^2 \left[\sum_{m=0}^{M_i-1} m^2 P_{im}(\theta_j) - \left(\sum_{m=0}^{M_i-1} m \cdot P_{im}(\theta_j) \right)^2 \right], \quad (6-8)$$

where u_{im} is the value 1 or 0.

During the M-STEP CAT administration process, in the case of scores that are zero (i.e., all items are incorrect) and perfect (i.e., all items are correct), a correction factor is applied before computing the relevant MLEs because the corresponding thetas cannot be estimated. The correction factor can be configured as any fractional value between 0 and 1 (e.g., 0.3). However, for the final scoring, the LOSS and the HOSS are assigned to the “all incorrect” and “all correct” cases according to the scoring specifications.

For each theta estimate, the corresponding SEM is calculated. SEM is the inverse of the square root of the test information function (TIF), which is the sum of the item information functions (IIFs). The IIF for dichotomously and polytomously scored items can be calculated by using the following equations, respectively:

$$IIF_i = D^2 a_i^2 (1 - P_i) P_i \quad (6-9)$$

and

$$IIF_i = D^2 a_i^2 \left[\sum_{m=0}^{M_i-1} m^2 P_{im} - \left(\sum_{m=0}^{M_i-1} m P_{im} \right)^2 \right] \quad (6-10)$$

6.3 Item Selection

After the initial item set is administered, the M-STEP CAT algorithm is designed to administer items targeted at an individual student's current performance, given content coverage boundaries. Specifically, the M-STEP CAT algorithm makes selection decisions each time based on the interim theta estimates while also taking many other factors, including test blueprint, item information function, and/or passage-related factors, into consideration. The details related to these factors are discussed below.

6.3.1 Test Blueprint

The adaptive item selection algorithm is designed to cover a standards-based blueprint, which includes the content standards, DOKs, item types, and score-point constraints. The M-STEP CAT algorithm closely resembles a modified constrained CAT (MCCAT) design (Leung, Chang, & Hau, 2003). The general idea is that the CAT algorithm is configured with upper and lower bounds that specify the minimum and maximum numbers of items that will be administered to students at the total-test, claim, content-category, assessment-target, and/or item-type levels. For the set of items configured, further configurations can be set up so only items at the specified DOK level and/or score point will be selected. The configurations specified in the test blueprint can be prioritized to ensure that the blueprint is met for each administration.

6.3.2 Item Information Function

After a content rule, among the content rules with the same priority level, is selected randomly or by the highest need, the M-STEP CAT algorithm targets the top N-ranked items, which are configurable, with the higher information function at the theta estimate. In general, the most efficient way to run an M-STEP CAT is to select items with the highest information function, which contains the smallest standard error for any given number of items. However, the consequence is that the items with high discriminations tend to be used more frequently. At the beginning of the test, it may not be necessary to select an item with the highest information function because the theta estimate used for calculating the information function contains a large measurement error. To control the item exposure rate for the high-discriminating items in the bank, a randomization process is introduced. Instead of the item with the highest information function being selected, the item to be used next is randomly selected from the top N (e.g., $N = 5$) number of items ranked by the item information (given interim theta estimate) and content-related criteria (see the "Distinct Top-Ranked Passage #" column in Table 6-2). Table 6-3 provides the scaling constants (i.e., slope A and intercept B), LOSS, and HOSS. All are fixed for each grade and content, and all were finalized before the test administration. They are used to convert students' estimated scores to the scale scores. More information about the scale transformations can be found in Chapter 8.

Table 6-2. Passage Selection Criteria

Content Area	Grade	Item Range	Passage Min #	Passage Max #	Passage Min %	Distinct Top-Ranked Passage #	Percentage of Items Used (% of Weight)	Items Delivered (% of Weight)	Max Information (% of Weight)
ELA	3–8	1–4	1	1	100	10	100	0	0
ELA	3–8	5–6	2	3	60	5	0	100	0
ELA	3–8	7	1	1	100	4	100	0	0
ELA	3–8	8–9	3	4	100	5	0	50	50
ELA	3–8	≥ 10	1	4	50	5	0	40	60
Mathematics	3–6	1–3	1	1	100	15	0	100	0
Mathematics	3–6	≥ 4	1	1	66	10	0	0	100
Mathematics	7–8	1–2	1	1	100	15	0	100	0
Mathematics	7–8	≥ 3	1	1	66	10	0	0	100

Table 6-3. Scoring Algorithm

Content	Grade	Slope A	Intercept B	LOSS	HOSS
ELA	3	26.0061	1322.5934	1203	1357
ELA	4	24.6036	1409.5875	1301	1454
ELA	5	25.8718	1501.3628	1409	1560
ELA	6	24.5491	1592.9699	1508	1655
ELA	7	23.8151	1687.3543	1618	1753
ELA	8	24.1951	1782.9264	1721	1857
Mathematics	3	26.3725	1325.7407	1217	1361
Mathematics	4	25.2608	1409.0233	1310	1455
Mathematics	5	23.3374	1495.6493	1409	1550
Mathematics	6	20.4573	1589.9260	1518	1650
Mathematics	7	19.6292	1686.6036	1621	1752
Mathematics	8	18.5194	1782.8881	1725	1850

6.3.3 Passage Related Concerns

Each passage in the ELA test has one or more associated items. The M-STEP CAT algorithm does not require that all items associated with a passage be administered; instead, it evaluates all possible combinations of items within a passage. Item sequencing within a passage is preserved when items are presented to the student. For example, if a six-item passage is selected and items 1 and 4 are not administered, then the items administered in order will be 2, 3, 5, and 6.

The configurable elements of a passage-based M-STEP CAT include the following:

Passage Minimum Percentage—This element defines the minimum percentage of the items associated with a passage to be used.

For example, if the distinct passage minimum percentage is set at 80, then the selection routine will consider passage combinations such as 1 of 1 (100%), 4 of 5 (80%), 5 of 6 (83%), and 6 of 6 (100%). It will not consider combinations such as 1 of 2 (50%), 3 of 4 (75%), 3 of 5 (60%), etc. Near the end of a test, the passage minimum percentage constraint may need to be loosened by a configurable reduction factor to meet content constraints such as the number of items per assessment target.

Passage Minimum and Maximum Number—This element defines the minimum and maximum numbers of items in a passage combination.

In the example above, 6 of 6 (100%) meets the passage minimum percentage (i.e., $\geq 80\%$); however, this passage combination may not be selected if the maximum number of items in a passage is specified as 5.

Passage Evaluation Criteria—Multiple factors are considered when evaluating and ranking each passage combination to determine the best combination to administer to a student. Passage combinations with higher criteria rankings are more likely to be administered. The criteria used in M-STEP CAT were as follows:

- Percentage of items used—the percentage of items associated with the passage selected for consideration
- Items delivered—the total number of items associated with a passage relative to the number of items selected to be delivered per passage
- Max information of passage combination—the higher the item information, the higher the combination is ranked

Different weights may be assigned to each of the factors mentioned above. For example, if 100% of the weight is assigned to the number of items delivered, then the algorithm will select the passages with the highest number of associated items and administer all those items until the maximum number of items is reached. Based on the simulation results, the criteria shown in Table 6-2 provided a better result that balanced the psychometric and test blueprint specifications.

6.4 Test Navigation

Due to a variety of reasons, many versions of CAT engines do not allow students to skip items in a test or return to previously answered items to change answers. Currently, all mathematics tests do not allow students to skip items or return to items to change answers. However, in the ELA tests, students are allowed to skip items within a passage. For example, when presented with a passage and five associated items, a student does not have to answer questions 1–5 in that order. However, if the student tries to navigate to the next passage without answering all items associated with the previous passage, the test engine will prompt the student to answer all items and will not move to the next passage until all are answered.

6.5 Termination

The CAT algorithm allows for both a fixed- and a variable-length test. With a fixed-length test, the test ends when a student has taken a predefined fixed number of items. With a variable-length test, in some cases, the algorithm stops administering items when the threshold of SEM or the maximum number of items is reached. Following the criteria set by Smarter Balanced, which MDE adopted for M-STEP, the algorithm stops administering items when a student has taken a predefined minimum number of items and the test blueprint specifications have been met.

6.6 Forced Submission

Tests are considered “complete” when students respond to the minimum number of operational items specified in the blueprint. Otherwise, the tests are “incomplete.” MLE is used to score the incomplete tests, counting unanswered items as incorrect.

When tests are adaptive, the specific unanswered items are unknown; thus, simulated items are used in place of administered items. Simulated items are generated with the following rules:

- Minimum operational test length is used to determine the test length of the incomplete tests.
- It is assumed that all unanswered operational items are dichotomously scored items. The item parameters of all unanswered operational items are equal to the average values of all the dichotomously scored operational items in the bank for discrimination and difficulty parameters.
- All unanswered operational items are scored as “incorrect.”

Table 6-4 lists the average discrimination and difficulty item parameter estimates and the minimum number of items needed to calculate the scores of the forced submitted students.

Table 6-4. Key Values Used for the 2017–18 Forced Submission

Content	Grade	Mean Discrimination	Mean Difficulty	Minimum Number of Items
ELA	3	0.663	-0.518	44
ELA	4	0.586	0.041	44
ELA	5	0.600	0.412	44
ELA	6	0.545	0.915	44
ELA	7	0.527	1.170	44
ELA	8	0.523	1.234	44
Mathematics	3	0.842	-0.809	36
Mathematics	4	0.825	-0.096	36
Mathematics	5	0.763	0.501	36
Mathematics	6	0.690	1.148	36
Mathematics	7	0.711	1.877	36
Mathematics	8	0.573	2.316	36

6.7 Summary of Simulation Results Evaluating the CAT Algorithm

This section summarizes the CAT simulation results with regard to the evaluation of the operational 2017–18 CAT algorithm. It is described in two subsections: (1) adherence to the test blueprint and (2) control of item exposure. Overall, the results are as expected and meet the acceptable psychometric requirements given the available item pool. For comparisons, the [Smarter Balanced 2017 Simulation Document](https://portal.smarterbalanced.org/library/en/2016-17-summative-assessments-simulation-results.pdf)² can be used.

6.7.1 Adherence to the Test Blueprint

During the M-STEP CAT simulations, blueprint constraints were used to ensure that the test blueprint was adhered to for all grades and content areas (see Figures 6-1 to 6-7). Note that for all the ELA tests, given the available items in the item pool and the fact that all the items are passage based in Claim 1, the number of items in Claim 1 can be met only at the content-category level. The simulation results show that every student received the number of items configured.

Tables 6-5 and 6-7 summarize the minimum and maximum numbers of items and points by claim and total for ELA and mathematics. The minimum and maximum numbers of passages and the number of items per passage per claim and per content category in Claim 1 are also summarized for ELA in Table 6.6. The results indicate that the CAT engine offered students the expected number of items and points, the expected number of passages, and a reasonable number of items delivered per passage.

² <https://portal.smarterbalanced.org/library/en/2016-17-summative-assessments-simulation-results.pdf>

Figure 6-1. Blueprint Target Sampling, Grade 3 through Grade 8 ELA

Claim (Goal1)	Content Category (Goal3)	Assessment Targets (Goal 4)	DOK	# of CAT Items Needed		# of CAT Items Configured		# PBW Items
				Min	Max	Min	Max	
1	LT	2	2,3	1	2	7	8	0
		4	3	1	2			
		1	1,2	3	6			
		3	1,2					
		5	3,4					
		6	2,3					
		7	2,3					
	IT	9	2,3	1	2	7	8	
		11	3	1	2			
		8	1,2	3	6			
		10	1,2					
		12	3,4					
		13	2,3					
		14	2,3					
2	O	1b/3b/6b	3	2	3	2	3	1
	E	1b/3b/6b	3	2	3	2	3	
	E	8	1,2	2	2	2	2	
	C	9	1,2	5	5	5	5	
3	L	4	1,2,3	8	9	8	9	0
4	CR	2	2	8	9	8	9	0
		3	2					
		4	2					

Figure 6-2. Blueprint Target Sampling, Grade 3 Mathematics

Claim (Goal1)	Content Category (Goal3)	Assessment Targets (Goal 4)	DOK	# of CAT Items Needed		# of CAT Items Configured		# of Multi- points Items
				Min	Max	Min	Max	
1	P	B	1	5	6	6	6	0-3
		C	1					
		I	1,2					
		G	1,2					
		D	2	5	6	6	6	
		F	1,2					
		A	1,2					
	S	E	1	3	4	4	4	
		J	1					
		K	1,2					
		H	2,3	1	1	1	1	
2	OA,NBT,NF,MD,G	A	2,3	2	2	2	2	
		B	1,2,3	2	2	2	2	
		C	1,2,3					
		D	1,2,3					
4	OA,NBT,NF,MD,G	A	2,3	2	2	2	2	
		D	2,3					
		B	2,3,4	1	1	1	1	
		E	2,3,4					
		C	1,2,3	1	1	1	1	
		F	1,2,3					
		G	3,4	0	0	0	0	
3	OA,NBT,NF,MD,G	A	2,3	3	3	3	3	
		D	2,3					
		B	2,3,4	3	3	3	3	
		E	2,3,4					
		C	2,3	2	2	2	2	
		F	2,3					

Figure 6-3. Blueprint Target Sampling, Grade 4 Mathematics

Claim (Goal1)	Content Category (Goal3)	Assessment Targets (Goal 4)	DOK	# of CAT Items Needed		# of CAT Items Configured		# of Multi- points Items
				Min	Max	Min	Max	
1	P	A	1,2	8	9	9	9	0-3
		E	1,2					
		F	1,2					
		G	1,2	2	3	3	3	
		D	1,2	1	2	2	2	
		H	1,2	1	1	1	1	
	S	I	1,2	2	3	3	3	
		K	1,2					
		B	1,2	1	1	1	1	
		C	2,3					
		J	1,2					
		L	1,2	1	1	1	1	
2	OA,NBT,NF,MD,G	A	2,3	2	2	2	2	
		B	1,2,3	2	2	2	2	
		C	1,2,3					
		D	1,2,3					
4	OA,NBT,NF,MD,G	A	2,3	2	2	2	2	
		D	2,3					
		B	2,3,4	1	1	1	1	
		E	2,3,4					
		C	1,2,3	1	1	1	1	
		F	1,2,3					
		G	3,4	0	0	0	0	
3	OA,NBT,NF,MD,G	A	2,3	3	3	3	3	
		D	2,3					
		B	2,3,4	3	3	3	3	
		E	2,3,4					
		C	2,3	2	2	2	2	
		F	2,3					

Figure 6-4. Blueprint Target Sampling, Grade 5 Mathematics

Claim (Goal1)	Content Category (Goal3)	Assessment Targets (Goal 4)	DOK	# of CAT Items Needed		# of CAT Items Configured		# of Multi- points Items
				Min	Max	Min	Max	
1	P	E	1,2	5	6	6	6	0-3
		I	1,2					
		F	1,2	4	5	5	5	
		D	1,2	3	4	4	4	
		C	1,2					
	S	J	1	2	3	3	3	
		K	2					
		A	1	2	2	2	2	
		B	2					
		G	1					
		H	1,2					
2	OA,NBT,NF,MD,G	A	2,3	2	2	2	2	
		B	1,2,3	2	2	2	2	
		C	1,2,3					
		D	1,2,3					
4	OA,NBT,NF,MD,G	A	2,3	2	2	2	2	
		D	2,3					
		B	2,3,4	1	1	1	1	
		E	2,3,4					
		C	1,2,3	1	1	1	1	
		F	1,2,3					
		G	3,4	0	0	0	0	
3	OA,NBT,NF,MD,G	A	2,3	3	3	3	3	
		D	2,3					
		B	2,3,4	3	3	3	3	
		E	2,3,4					
		C	2,3	2	2	2	2	
		F	2,3					

Figure 6-5. Blueprint Target Sampling, Grade 6 Mathematics

Claim (Goal1)	Content Category (Goal3)	Assessment Targets (Goal 4)	DOK	# of CAT Items Needed		# of CAT Items Configured		# of Multi- points Items Needed
				Min	Max	Min	Max	
1	P	E	1	5	6	6	6	0-3
		F	1,2					
		A	1,2	3	4	4	4	
		G	2	3	3	3	3	
		B	1,2					
		D	1,2	2	2	2	2	
	S	C	1,2	4	5	5	5	
		H	1,2					
		I	2					
		J	1,2					
2	RP,NS,EE,G,SP	A	2,3	2	2	2	2	
		B	1,2,3	2	2	2	2	
		C	1,2,3					
		D	1,2,3					
4	RP,NS,EE,G,SP	A	2,3	2	2	2	2	
		D	2,3					
		B	2,3,4	1	1	1	1	
		E	2,3,4					
		C	1,2,3	1	1	1	1	
		F	1,2,3					
		G	3,4	0	0	0	0	
3	RP,NS,EE,G,SP	A	2,3	3	3	3	3	
		D	2,3					
		B	2,3,4	3	3	3	3	
		E	2,3,4					
		C	2,3	2	2	2	2	
		F	2,3					
		G	2,3					

Figure 6-6. Blueprint Target Sampling, Grade 7 Mathematics

Claim (Goal1)	Content Category (Goal3)	Assessment Targets (Goal 4)	DOK	# of CAT Items Needed		# of CAT Items Configured		# of Multi- points Items
				Min	Max	Min	Max	
1	P	A	2	8	9	9	9	0-2
		D	1,2					
		B	1,2	5	6	6	6	
		C	1,2					
	S	E	1,2	2	3	3	3	
		F	1,2					
		G	1,2	1	2	2	2	
		H	2					
		I	1,2					
2	RP,NS,EE,G,SP	A	2,3	2	2	2	2	
		B	1,2,3	2	2	2	2	
		C	1,2,3					
		D	1,2,3					
4	RP,NS,EE,G,SP	A	2,3	2	2	2	2	
		D	2,3					
		B	2,3,4	1	1	1	1	
		E	2,3,4					
		C	1,2,3	1	1	1	1	
		F	1,2,3					
		G	3,4	0	0	0	0	
3	RP,NS,EE,G,SP	A	2,3	3	3	3	3	
		D	2,3					
		B	2,3,4	3	3	3	3	
		E	2,3,4					
		C	2,3	2	2	2	2	
		F	2,3					
		G	2,3					

Figure 6-7. Blueprint Target Sampling, Grade 8 Mathematics

Claim (Goal1)	Content Category (Goal3)	Assessment Targets (Goal 4)	DOK	# of CAT Items Needed		# of CAT Items Configured		# of Multi- points Items
				Min	Max	Min	Max	
1	P	C	1,2	5	6	6	6	0-3
		D	1,2					
		B	1,2	5	6	6	6	
		E	1,2					
		G	1,2					
		F	1,2	2	3	3	3	
		H	1,2					
	S	A	1,2	4	5	5	5	
		I	1,2					
		J	1,2					
2	NS,EE,F,G,SP	A	2,3	2	2	2	2	
		B	1,2,3	2	2	2	2	
		C	1,2,3					
		D	1,2,3					
4	NS,EE,F,G,SP	A	2,3	2	2	2	2	
		D	2,3					
		B	2,3,4	1	1	1	1	
		E	2,3,4					
		C	1,2,3	1	1	1	1	
		F	1,2,3					
		G	3,4	0	0	0	0	
3	NS,EE,F,G,SP	A	2,3	3	3	3	3	
		D	2,3					
		B	2,3,4	3	3	3	3	
		E	2,3,4					
		C	2,3	2	2	2	2	
		F	2,3					
		G	2,3					

Table 6-5. Summary of Items and Points for ELA

Grade	Level	Min # of Items	Max # of Items	Min # of Points	Max # of Points
3	Total	45	46	48	49
3	Claim 1	16	16	16	16
3	Claim 2	13	13	16	16
3	Claim 3	8	9	8	9
3	Claim 4	8	8	8	8
3	Claim 1_LT	8	8	8	8
3	Claim 1_IT	8	8	8	8
4	Total	45	46	48	49
4	Claim 1	16	16	16	16
4	Claim 2	13	13	16	16
4	Claim 3	8	9	8	9
4	Claim 4	8	8	8	8
4	Claim 1_LT	8	8	8	8
4	Claim 1_IT	8	8	8	8
5	Total	44	46	47	49
5	Claim 1	15	16	15	16
5	Claim 2	13	13	16	16
5	Claim 3	8	9	8	9
5	Claim 4	8	8	8	8
5	Claim 1_LT	8	8	8	8
5	Claim 1_IT	7	8	7	8
6	Total	44	46	47	49
6	Claim 1	15	16	15	16
6	Claim 2	13	13	16	16
6	Claim 3	8	9	8	9
6	Claim 4	8	8	8	8
6	Claim 1_LT	7	8	7	8
6	Claim 1_IT	8	8	8	8
7	Total	45	46	48	49
7	Claim 1	16	16	16	16
7	Claim 2	13	13	16	16
7	Claim 3	8	9	8	9
7	Claim 4	8	8	8	8
7	Claim 1_LT	8	8	8	8
7	Claim 1_IT	8	8	8	8

Chapter 6: Operational Computer Adaptive Test (CAT)

Grade	Level	Min # of Items	Max # of Items	Min # of Points	Max # of Points
8	Total	45	46	48	49
8	Claim 1	16	16	16	16
8	Claim 2	13	13	16	16
8	Claim 3	8	9	8	9
8	Claim 4	8	8	8	8
8	Claim 1_LT	8	8	8	8
8	Claim 1_IT	8	8	8	8

Table 6-6. Summary of Passages and Items per Passage for ELA

Grade	Level	Min # of Passages	Max # of Passages	Min # of Items per Passage	Max # of Items per Passage
3	Total	7	8	2	4
3	Claim 1	4	4	4	4
3	Claim 3	3	4	2	3
3	Claim 1_LT	2	2	4	4
3	Claim 1_IT	2	2	4	4
4	Total	7	8	2	4
4	Claim 1	4	4	4	4
4	Claim 3	3	4	2	3
4	Claim 1_LT	2	2	4	4
4	Claim 1_IT	2	2	4	4
5	Total	7	8	2	4
5	Claim 1	4	4	3	4
5	Claim 3	3	4	2	3
5	Claim 1_LT	2	2	4	4
5	Claim 1_IT	2	2	3	4
6	Total	7	8	2	4
6	Claim 1	4	4	3	4
6	Claim 3	3	4	2	3
6	Claim 1_LT	2	2	3	4
6	Claim 1_IT	2	2	4	4
7	Total	7	8	2	4
7	Claim 1	4	4	4	4
7	Claim 3	3	4	2	3
7	Claim 1_LT	2	2	4	4
7	Claim 1_IT	2	2	4	4

Grade	Level	Min # of Passages	Max # of Passages	Min # of Items per Passage	Max # of Items per Passage
8	Total	7	8	2	4
8	Claim 1	4	4	4	4
8	Claim 3	3	4	2	4
8	Claim 1_LT	2	2	4	4
8	Claim 1_IT	2	2	4	4

Table 6-7. Summary of Items and Points for Mathematics

Grade	Level	Min # of Items	Max # of Items	Min # of Points	Max # of Points
3	Total	36	36	36	40
3	Claim 1	20	20	20	22
3	Claim 2	4	4	4	5
3	Claim 3	8	8	8	11
3	Claim 4	4	4	4	8
3	Claim 2 & 4	8	8	8	12
4	Total	36	36	36	40
4	Claim 1	20	20	20	23
4	Claim 2	4	4	4	4
4	Claim 3	8	8	8	10
4	Claim 4	4	4	4	6
4	Claim 2 & 4	8	8	8	10
5	Total	36	36	36	40
5	Claim 1	20	20	20	20
5	Claim 2	4	4	4	5
5	Claim 3	8	8	8	11
5	Claim 4	4	4	4	7
5	Claim 2 & 4	8	8	8	11
6	Total	36	36	36	40
6	Claim 1	20	20	20	22
6	Claim 2	4	4	4	6
6	Claim 3	8	8	8	11
6	Claim 4	4	4	4	7
6	Claim 2 & 4	8	8	8	11
7	Total	36	36	36	39
7	Claim 1	20	20	20	20
7	Claim 2	4	4	4	5
7	Claim 3	8	8	8	11

Grade	Level	Min # of Items	Max # of Items	Min # of Points	Max # of Points
7	Claim 4	4	4	4	6
7	Claim 2 & 4	8	8	8	10
8	Total	36	36	36	40
8	Claim 1	20	20	20	23
8	Claim 2	4	4	4	6
8	Claim 3	8	8	8	12
8	Claim 4	4	4	4	6
8	Claim 2 & 4	8	8	8	10

6.7.2 Controlling for Item Exposure

A common concern when implementing a CAT is the exposure rate of the items. It is important to control the item exposure rate while balancing the other constraints of the CAT. Tables 6-8 and 6-9 show the item exposure rates for ELA and mathematics, respectively. Each table provides the number and proportion of items for each of the six exposure rate categories, including no exposure. For example, an exposure rate of (0.0, 0.1] means that between 0% (excluding 0, as it forms its own category) and 10% of the students took that item. For grade 3 ELA, 598 items, or 69% of the items in the pool, had an exposure rate between 0% and 10%. For both ELA and mathematics, most items had a low exposure rate and were categorized with an exposure rate between 0% and 10%.

Table 6-8. Summary of Item Exposure Rate by Grade and Level for ELA

Grade	Level	Number Items	Proportion of Items
3	0	113	0.13
3	(0.0, 0.1]	598	0.69
3	(0.1, 0.2]	96	0.11
3	(0.2, 0.3]	30	0.03
3	(0.3, 0.4]	17	0.02
3	> 0.4	9	0.01
4	0	88	0.11
4	(0.0, 0.1]	593	0.72
4	(0.1, 0.2]	70	0.08
4	(0.2, 0.3]	51	0.06
4	(0.3, 0.4]	19	0.02
4	> 0.4	5	0.01
5	0	99	0.13
5	(0.0, 0.1]	534	0.68
5	(0.1, 0.2]	92	0.12
5	(0.2, 0.3]	36	0.05
5	(0.3, 0.4]	19	0.02
5	> 0.4	6	0.01
6	0	84	0.11
6	(0.0, 0.1]	525	0.7
6	(0.1, 0.2]	57	0.08
6	(0.2, 0.3]	45	0.06
6	(0.3, 0.4]	24	0.03
6	> 0.4	12	0.02
7	0	91	0.14
7	(0.0, 0.1]	428	0.66
7	(0.1, 0.2]	57	0.09
7	(0.2, 0.3]	34	0.05
7	(0.3, 0.4]	28	0.04
7	> 0.4	15	0.02
8	0	89	0.12
8	(0.0, 0.1]	500	0.68
8	(0.1, 0.2]	68	0.09
8	(0.2, 0.3]	38	0.05
8	(0.3, 0.4]	28	0.04
8	> 0.4	9	0.01

Table 6-9. Summary of Item Exposure Rate by Grade and Level for Mathematics

Grade	Level	Number Items	Proportion of Items
3	0	197	0.16
3	(0.0, 0.1]	926	0.74
3	(0.1, 0.2]	118	0.09
3	(0.2, 0.3]	2	0
3	(0.3, 0.4]	0	0
3	> 0.4	0	0
4	0	145	0.11
4	(0.0, 0.1]	1026	0.8
4	(0.1, 0.2]	101	0.08
4	(0.2, 0.3]	5	0
4	(0.3, 0.4]	0	0
4	> 0.4	0	0
5	0	146	0.12
5	(0.0, 0.1]	932	0.77
5	(0.1, 0.2]	127	0.11
5	(0.2, 0.3]	1	0
5	(0.3, 0.4]	0	0
5	> 0.4	0	0
6	0	147	0.13
6	(0.0, 0.1]	867	0.77
6	(0.1, 0.2]	111	0.1
6	(0.2, 0.3]	8	0.01
6	(0.3, 0.4]	0	0
6	> 0.4	0	0
7	0	74	0.07
7	(0.0, 0.1]	807	0.8
7	(0.1, 0.2]	110	0.11
7	(0.2, 0.3]	22	0.02
7	(0.3, 0.4]	0	0
7	> 0.4	0	0
8	0	148	0.16
8	(0.0, 0.1]	614	0.68
8	(0.1, 0.2]	105	0.12
8	(0.2, 0.3]	35	0.04
8	(0.3, 0.4]	0	0
8	> 0.4	0	0

6.8 Summary of Simulation Results for the Student Ability Estimates

For Smarter Balanced tests with an adaptive component, test reliability is estimated through simulations conducted using the operational summative item pool. For fixed-form tests, reliability and SEM are calculated using the items on the forms and their psychometric properties relative to the population. DRC conducted simulation studies for the 2017–18 tests using the 2016–17 M-STEP ability estimates, which had the means and SDs shown in Table 6-10.

Table 6-10. Mean and Standard Deviation of the Sample Used in the Simulation Study

Content	Grade	Mean	SD
ELA	3	-1.07	1.01
ELA	4	-0.62	1.05
ELA	5	-0.07	1.02
ELA	6	0.06	1.09
ELA	7	0.30	1.12
ELA	8	0.57	1.08
Mathematics	3	-1.10	0.98
Mathematics	4	-0.63	1.02
Mathematics	5	-0.28	1.08
Mathematics	6	-0.10	1.23
Mathematics	7	0.08	1.35
Mathematics	8	0.26	1.39

6.8.1 Ability Estimates at the Extremes

The examinee ability in the simulation study was estimated using MLE. To provide a limit to the score range for extreme values, the test scoring algorithm used the HOSS and LOSS that were derived during the Smarter Balanced 2014 achievement level setting. Scores above HOSS or below LOSS are assigned HOSS and LOSS values respectively. Table 6-11 presents the LOSS and HOSS values that were used in the simulation and the percentage of the affected scores at those values.

Table 6-11. HOSS/LOSS and Percentages of Affected Scores from Simulation Results

Content	Grade	LOSS	HOSS	Percentage of Scores at LOSS	Percentage of Scores at HOSS
ELA	3	-4.59	1.34	0.07	0.70
ELA	4	-4.40	1.80	0.07	1.00
ELA	5	-3.58	2.25	0.00	0.83
ELA	6	-3.48	2.51	0.10	0.83
ELA	7	-2.91	2.75	0.30	0.97
ELA	8	-2.57	3.04	0.40	0.70
Math	3	-4.11	1.33	0.43	0.60
Math	4	-3.92	1.82	0.33	0.67
Math	5	-3.73	2.33	0.10	0.30
Math	6	-3.53	2.95	0.87	0.40
Math	7	-3.34	3.32	1.67	0.50
Math	8	-3.15	3.63	1.43	0.73

6.8.2 Standard Error of Measurement

The SEM, in the theta metric, is calculated for each reportable claim score and the total score. Note that for mathematics, the combined score for Claims 2 and 4 is reported, so the SEM for the combined score is calculated. Tables 6-12 and 6-13 provide statistical summaries (including the minimum, maximum, mean, median, and SD values) of the SEMs for claim scores and total scores. For all the tests, the average SEMs for claim scores are larger than the SEMs for the total scores. This is expected because the number of items in each claim is smaller than the number of items in the total test. As the grade increases, the average SEM increases. This is possibly due to the mismatch between the item difficulty distributions and student ability distributions in higher grades. The 3,000 simulated students' abilities or scores were randomly selected from the previous year's operational results on M-STEP. It was found that the items in the higher grades were relatively harder for the students. The SEMs are reasonable given the length of the total test and the claim level of the test.

Table 6-12. Summary of Standard Error of Measurement by Grade and Level for ELA

Grade	Level	Mean	SD	Min	Max	Median
3	Total	0.24	0.04	0.20	0.95	0.22
3	Claim 1	0.41	0.15	0.31	2.44	0.37
3	Claim 2	0.46	0.11	0.36	2.20	0.44
3	Claim 3	0.81	0.30	0.52	2.79	0.72
3	Claim 4	0.54	0.21	0.39	3.80	0.48
4	Total	0.26	0.03	0.23	0.71	0.25
4	Claim 1	0.46	0.17	0.36	3.19	0.42
4	Claim 2	0.50	0.15	0.36	3.17	0.49
4	Claim 3	0.81	0.30	0.55	3.43	0.72
4	Claim 4	0.62	0.20	0.45	2.17	0.56
5	Total	0.26	0.02	0.22	0.62	0.25
5	Claim 1	0.45	0.10	0.36	2.50	0.44
5	Claim 2	0.53	0.11	0.38	2.00	0.52
5	Claim 3	0.81	0.24	0.61	3.29	0.75
5	Claim 4	0.53	0.15	0.42	3.43	0.49
6	Total	0.28	0.04	0.24	1.01	0.26
6	Claim 1	0.50	0.16	0.40	2.94	0.46
6	Claim 2	0.56	0.17	0.43	2.84	0.52
6	Claim 3	0.82	0.34	0.54	4.60	0.72
6	Claim 4	0.64	0.24	0.46	4.11	0.56
7	Total	0.30	0.05	0.25	1.02	0.29
7	Claim 1	0.53	0.19	0.37	3.40	0.49
7	Claim 2	0.63	0.15	0.51	2.57	0.59
7	Claim 3	0.85	0.31	0.56	3.17	0.76
7	Claim 4	0.72	0.27	0.50	2.57	0.65
8	Total	0.29	0.04	0.26	0.81	0.28
8	Claim 1	0.54	0.18	0.44	3.71	0.50
8	Claim 2	0.61	0.14	0.50	2.44	0.58
8	Claim 3	0.89	0.30	0.60	3.51	0.80
8	Claim 4	0.61	0.20	0.48	2.53	0.54

Table 6-13. Summary of Standard Error of Measurement by Grade and Level for Mathematics

Grade	Level	Mean	SD	Min	Max	Median
3	Total	0.21	0.05	0.18	1.17	0.20
3	Claim 1	0.27	0.06	0.23	2.29	0.26
3	Claim 3	0.58	0.28	0.33	3.26	0.49
3	Claim 2 & 4	0.56	0.30	0.33	3.55	0.45
4	Total	0.21	0.06	0.16	1.21	0.20
4	Claim 1	0.27	0.06	0.22	1.18	0.25
4	Claim 3	0.57	0.3	0.33	2.99	0.46
4	Claim 2 & 4	0.57	0.26	0.33	2.85	0.49
5	Total	0.25	0.09	0.17	1.45	0.21
5	Claim 1	0.32	0.12	0.21	2.19	0.28
5	Claim 3	0.62	0.31	0.36	4.54	0.51
5	Claim 2 & 4	0.71	0.44	0.33	3.28	0.54
6	Total	0.28	0.11	0.19	2.54	0.25
6	Claim 1	0.34	0.12	0.25	2.48	0.31
6	Claim 3	0.84	0.49	0.40	6.74	0.67
6	Claim 2 & 4	0.85	0.62	0.36	6.06	0.63
7	Total	0.33	0.21	0.19	4.35	0.27
7	Claim 1	0.40	0.27	0.24	3.88	0.33
7	Claim 3	1.03	0.68	0.42	5.87	0.78
7	Claim 2 & 4	1.01	0.70	0.35	6.11	0.74
8	Total	0.36	0.13	0.22	2.08	0.35
8	Claim 1	0.43	0.17	0.27	5.22	0.42
8	Claim 3	1.23	0.64	0.47	4.98	1.03
8	Claim 2 & 4	1.20	0.76	0.44	4.90	0.89

6.8.3 Statistical Measures of Bias

This section presents the statistics calculated for the annual Michigan simulation investigation. Note that these statistics are the same as those reported in the *Smarter Balanced 2017–2018 Technical Report* (2018). Therefore, a direct quote from this Smarter Balanced report is used here for describing these statistics.

- Bias: [T]he statistical bias of the estimated theta parameter. This is a test of the assumption that error is randomly distributed around true ability. It is a measure of whether scores systematically underestimate or overestimate ability.
- Mean squared error (MSE): This is a measure of the magnitude of difference between true and estimated theta.
- Significance of bias [“Bias Sig” in Tables 6-14 and 6-15]: [A]n indicator of the statistical significance of bias.
- Average standard error of the estimated theta: This is the average of the simulated standard error of measurement [SEM] over all examinees. It is the marginal reliability for the simulated population.
- Standard error of estimates of theta at the 5th, 25th, 75th, and 95th percentiles
- Percentage of students’ estimated theta falling outside the 95% and 99% confidence intervals [Miss Rate]. (p. 2–3)

For detailed mathematical formulas in computing these statistics, please refer to pages 2–4 of the *Smarter Balanced 2017–2018 Technical Report* (2018).

Tables 6-14 and 6-15 present the bias of the estimated abilities for ELA and mathematics, respectively. As was found in the Smarter Balanced simulation study (Smarter Balanced, 2016), the bias in the overall scores is both small and insignificant. It should also be noted that claim scores do have some systematic bias. This is likely caused by the application of HOSS and LOSS values.

Table 6-14. Bias of the Estimated Theta from Simulation Results: ELA

Level	Grade	Mean Bias	SE of Mean Bias	Bias Sig	MSE	95% CI Miss Rate	99% CI Miss Rate
Overall	3	-0.01	0.02	0.66	0.06	5.50	1.37
Overall	4	0.00	0.02	0.96	0.07	4.47	0.83
Overall	5	-0.01	0.02	0.79	0.07	5.20	0.87
Overall	6	0.00	0.02	0.90	0.08	4.03	0.87
Overall	7	0.01	0.02	0.79	0.10	5.43	1.10
Overall	8	-0.01	0.02	0.64	0.09	5.17	0.67
Claim 1	3	-0.02	0.02	0.18	0.22	4.57	0.90
Claim 1	4	-0.03	0.02	0.08	0.23	3.40	0.53
Claim 1	5	-0.02	0.02	0.27	0.23	3.83	0.50
Claim 1	6	-0.02	0.02	0.31	0.26	3.07	0.27
Claim 1	7	0.02	0.02	0.38	0.32	3.87	0.57
Claim 1	8	-0.01	0.02	0.52	0.35	3.83	0.73
Claim 2	3	-0.01	0.02	0.48	0.25	3.53	0.23
Claim 2	4	0.02	0.02	0.27	0.28	3.00	0.23
Claim 2	5	0.03	0.02	0.10	0.30	3.33	0.40
Claim 2	6	-0.01	0.02	0.57	0.36	3.07	0.23
Claim 2	7	-0.02	0.02	0.41	0.43	3.33	0.27
Claim 2	8	-0.03	0.02	0.12	0.41	3.00	0.30
Claim 3	3	-0.06	0.02	0.00	0.70	2.07	0.17
Claim 3	4	-0.02	0.02	0.36	0.72	2.20	0.10
Claim 3	5	0.01	0.02	0.71	0.71	2.07	0.17
Claim 3	6	-0.01	0.02	0.50	0.80	1.77	0.10
Claim 3	7	-0.01	0.02	0.66	0.83	2.13	0.17
Claim 3	8	-0.04	0.02	0.07	0.87	1.73	0.17
Claim 4	3	-0.04	0.02	0.04	0.33	2.20	0.13
Claim 4	4	-0.04	0.02	0.05	0.46	2.03	0.20
Claim 4	5	-0.02	0.02	0.27	0.32	2.23	0.13
Claim 4	6	-0.04	0.02	0.02	0.46	2.10	0.27
Claim 4	7	-0.07	0.02	0.00	0.56	2.47	0.27
Claim 4	8	-0.03	0.02	0.12	0.40	2.37	0.10

Table 6-15. Bias of the Estimated Theta from Simulation Results: Mathematics

Level	Grade	Mean Bias	SE of Mean Bias	Bias Sig	MSE	95% CI Miss Rate	99% CI Miss Rate
Overall	3	-0.01	0.02	0.61	0.05	5.23	0.93
Overall	4	0.00	0.02	0.79	0.05	4.77	0.83
Overall	5	-0.02	0.02	0.40	0.07	5.27	0.87
Overall	6	-0.03	0.02	0.26	0.10	4.73	0.77
Overall	7	-0.05	0.02	0.04	0.19	5.10	1.20
Overall	8	-0.05	0.03	0.07	0.18	5.17	1.07
Claim 1	3	0.00	0.02	0.87	0.08	4.77	0.73
Claim 1	4	0.00	0.02	0.83	0.08	4.97	0.77
Claim 1	5	-0.03	0.02	0.14	0.12	4.80	0.80
Claim 1	6	-0.03	0.02	0.24	0.15	4.77	0.77
Claim 1	7	-0.07	0.02	0.00	0.27	4.63	0.93
Claim 1	8	-0.05	0.03	0.04	0.26	4.27	1.30
Claim 3	3	-0.08	0.02	0.00	0.39	2.57	0.30
Claim 3	4	-0.07	0.02	0.00	0.35	2.27	0.57
Claim 3	5	-0.09	0.02	0.00	0.43	2.60	0.47
Claim 3	6	-0.14	0.02	0.00	0.70	2.57	0.40
Claim 3	7	-0.19	0.02	0.00	1.05	2.83	0.53
Claim 3	8	-0.25	0.03	0.00	1.47	2.33	0.63
Claims 2 & 4	3	-0.07	0.02	0.00	0.33	3.00	0.37
Claims 2 & 4	4	-0.05	0.02	0.01	0.34	2.67	0.37
Claims 2 & 4	5	-0.09	0.02	0.00	0.48	3.10	0.60
Claims 2 & 4	6	-0.12	0.02	0.00	0.56	2.40	0.40
Claims 2 & 4	7	-0.15	0.02	0.00	0.74	2.60	0.40
Claims 2 & 4	8	-0.21	0.03	0.00	1.12	2.93	0.73

Tables 6-16 and 6-17 below present marginal reliability coefficients and precisions for the overall tests and for reported claims. As expected, estimated reliability coefficients for the overall tests are high and are in the acceptable range for a large-scale, high-stakes test. Reliability estimates at the claim level are lower, and corresponding errors are higher. Claims with smaller numbers of items and fewer points from the adaptive section of the test exhibit the lowest reliability. This shows the importance of incorporating error data in claim-level reports.

Table 6-16. Overall Score and Claim Score Precision/Reliability of Simulation Results: ELA

Level	Grade	Mean # Items	Mean SEM	Reliability	RMSE	SD theta
Overall	3	45.23	0.24	0.95	0.25	1.06
Overall	4	45.03	0.26	0.94	0.26	1.09
Overall	5	45.40	0.26	0.94	0.26	1.06
Overall	6	45.15	0.28	0.94	0.28	1.13
Overall	7	45.06	0.30	0.93	0.31	1.18
Overall	8	45.31	0.29	0.93	0.30	1.14
Claim 1	3	16.00	0.41	0.86	0.47	1.17
Claim 1	4	16.00	0.45	0.84	0.48	1.19
Claim 1	5	15.96	0.46	0.83	0.48	1.16
Claim 1	6	15.98	0.50	0.83	0.51	1.25
Claim 1	7	16.00	0.53	0.83	0.57	1.33
Claim 1	8	16.00	0.54	0.80	0.59	1.29
Claim 2	3	13.00	0.47	0.83	0.50	1.17
Claim 2	4	13.00	0.50	0.82	0.53	1.22
Claim 2	5	13.00	0.53	0.80	0.55	1.19
Claim 2	6	13.00	0.56	0.79	0.60	1.27
Claim 2	7	13.00	0.63	0.76	0.65	1.32
Claim 2	8	13.00	0.61	0.77	0.64	1.30
Claim 3	3	8.23	0.80	0.61	0.84	1.37
Claim 3	4	8.03	0.81	0.64	0.85	1.43
Claim 3	5	8.45	0.80	0.62	0.84	1.35
Claim 3	6	8.17	0.82	0.64	0.89	1.49
Claim 3	7	8.06	0.86	0.65	0.91	1.56
Claim 3	8	8.31	0.89	0.61	0.93	1.50
Claim 4	3	8.00	0.54	0.77	0.57	1.22
Claim 4	4	8.00	0.63	0.74	0.68	1.31
Claim 4	5	8.00	0.54	0.79	0.57	1.21
Claim 4	6	8.00	0.64	0.74	0.68	1.34
Claim 4	7	8.00	0.71	0.70	0.75	1.38
Claim 4	8	8.00	0.61	0.76	0.63	1.31

Table 6-17. Overall Score and Claim Score Precision/Reliability of Simulation Results: Mathematics

Level	Grade	Mean # Items	Mean SEM	Reliability	RMSE	SD theta
Overall	3	36	0.21	0.95	0.23	1.03
Overall	4	36	0.21	0.96	0.23	1.06
Overall	5	36	0.25	0.95	0.27	1.13
Overall	6	36	0.28	0.95	0.32	1.30
Overall	7	36	0.33	0.93	0.44	1.49
Overall	8	36	0.37	0.93	0.43	1.50
Claim 1	3	20	0.27	0.93	0.29	1.04
Claim 1	4	20	0.27	0.93	0.29	1.08
Claim 1	5	20	0.32	0.91	0.35	1.16
Claim 1	6	20	0.34	0.92	0.39	1.33
Claim 1	7	20	0.40	0.90	0.52	1.53
Claim 1	8	20	0.44	0.90	0.51	1.53
Claim 3	3	8	0.58	0.71	0.62	1.22
Claim 3	4	8	0.57	0.72	0.59	1.24
Claim 3	5	8	0.62	0.73	0.66	1.34
Claim 3	6	8	0.84	0.62	0.84	1.56
Claim 3	7	8	1.03	0.52	1.02	1.79
Claim 3	8	8	1.25	0.44	1.21	1.90
Claim 2 & 4	3	8	0.56	0.72	0.57	1.19
Claim 2 & 4	4	8	0.57	0.73	0.58	1.22
Claim 2 & 4	5	8	0.70	0.63	0.70	1.34
Claim 2 & 4	6	8	0.84	0.52	0.75	1.48
Claim 2 & 4	7	8	1.01	0.43	0.86	1.62
Claim 2 & 4	8	8	1.23	0.34	1.06	1.81

One of the advantages of adaptive tests is that SEM can be controlled for all ability levels. Ideally, the SEM should be similar throughout the ability distribution. Table 6-18 presents average error by decile of the true thetas, which were generated based on the Michigan population. For both ELA and mathematics, the results show that the error at the lower end of the test tends to be the highest (except for Grade 5 ELA), indicating that there is more error associated with the ability estimation at the lower end of the ability distribution, which is caused by the relative difficulty of the item pools.

Table 6-18. Average Standard Errors by Grade and by Deciles of True Proficiency Scores of Simulation Results

Subject	Grade	Decile 1	Decile 2	Decile 3	Decile 4	Decile 5	Decile 6	Decile 7	Decile 8	Decile 9	Decile 10
ELA	4	7.53	6.31	6.06	6.01	6.02	6.01	6.00	6.01	6.06	6.85
ELA	5	7.28	6.15	6.02	6.07	6.11	6.26	6.55	6.76	6.96	7.47
ELA	6	8.94	7.29	6.75	6.37	6.17	6.07	6.05	6.07	6.28	7.02
ELA	7	9.40	7.63	7.09	6.84	6.73	6.70	6.70	6.82	6.91	7.40
ELA	8	9.09	7.31	7.01	6.96	6.94	6.94	6.90	6.97	7.00	7.41
Math	3	8.02	5.97	5.50	5.25	5.11	5.05	5.02	5.01	5.02	5.49
Math	4	8.44	6.04	5.46	5.12	5.04	4.95	4.76	4.53	4.52	5.10
Math	5	9.75	7.72	6.74	5.99	5.29	4.91	4.44	4.12	4.23	4.76
Math	6	10.12	6.88	6.19	5.75	5.34	5.07	4.90	4.44	4.15	4.22
Math	7	14.21	8.77	7.15	6.28	5.79	5.18	4.79	4.24	4.02	4.07
Math	8	11.67	8.39	7.53	7.17	6.65	6.28	5.76	5.24	4.91	4.62
Math	8	13.91	9.23	8.20	7.69	7.18	6.82	6.36	5.78	5.24	5.20

6.9 Summary

In summary, Chapter 6 of this report demonstrates M-STEP's adherence to AERA, APA, & NCME (2014) *Standards* regarding construct-related validity and reliability. The analyses described above are related to the following standards:

- Standard 2.0— Appropriate evidence of reliability/precision should be provided for the interpretation for each intended score use.
- Standard 2.1 —The range of replications over which reliability/precision is being evaluated should be clearly stated, along with a rationale for the choice of this definition, given the testing situation.
- Standard 4.3—Test developers should document the rationale and supporting evidence for the administration, scoring, and reporting rules used in computer-adaptive, multistage-adaptive, or other tests delivered using computer algorithms to select items. This documentation should include procedures used in selecting items or sets of items for administration, in determining the starting point and termination conditions for the test, in scoring the test, and in controlling item exposure.

Chapter 7: Scoring

Chapter 7 shows how M-STEP scoring adhered to the AERA, APA, & NCME *Standards*. Standard 4.18 provides some general guidance for Chapter 7:

Procedures for scoring and, if relevant, scoring criteria, should be presented by the test developer with sufficient detail and clarity to maximize the accuracy of scoring. Instructions for using rating scales or for deriving scores obtained by coding, scaling, or classifying constructed responses should be clear. This is especially critical for extended-response items such as performance tasks, portfolios, and essays. (p. 91)

Chapter 7 explains the procedures used for autoscoring and handscoring, with the latter applicable to the Passage-Based Writing items. There was no AI scoring of student writing in 2018. The scoring criteria used for each item are not presented in this chapter to preserve the integrity of the items for future use.

7.1 Online Scoring

7.1.1 Autoscoring

All content areas of M-STEP contain items that required autoscoring. Autoscoring was used for Technology Enhanced items which could involve combining many components to form a single correct answer. Scoring rules for each item were set up prior to the start of testing. These rules listed all the different correct components per item. DRC ensured that all rubrics and scoring rules were verified for accuracy before scoring began. Quality checks were run against all autoscored items using the autoscoring simulator tool to ensure the item was scored as designed. The autoscoring simulator tool allowed specialists to respond to items with expected responses. In some cases, the simulator generated all possible responses and expected values. Once student responses were entered, the scoring engine was run against the student response and the expected value. The tool alerted the specialist of any mismatches, which could then be updated. All of this quality check occurred before the start of the testing window. During the testing window, the autoscoring process ran daily, and all available, completed items were scored. After testing was complete, a secondary check was completed by the psychometrics team. Any items that did not perform as expected were communicated back to the autoscoring specialists, who reran the simulations to assure the autoscoring was set up as requested. If an autoscoring setup issue was found at this point, the items are updated and rescored. This occurred before reporting.

DRC provided MDE with complete item frequency reports, which includes the following information for each response pattern/combination: (1) the number/percentage of students gave that response pattern/combination, and (2) the score provided by the scoring system.

7.1.2 Multiple Choice Scoring

The online scoring process includes the scoring of multiple-choice items, in which students chose only one correct answer from choices A–D. The items were scored against a scoring key that was prepared and validated before the start of each testing window. Responses to multiple-choice items were captured during the online test administration, and items were scored as “right,” “wrong,” or “blank” (i.e., not answered). Additional answer key checks were conducted during the testing windows to ensure that the items were scored based on the provided key.

7.2 Handscoring

Measurement Incorporated performed all required scoring of paper/pencil and online items needing handscoring. For M-STEP ELA, these were Passage-Based Writing (PBW) items for grades 3–8. For M-STEP mathematics, these included short-text and short-text fill-in table items for grades 3–8.

M-STEP items were scored by readers working in Taylor, Michigan; Grand Rapids, Michigan; and at other scoring centers (i.e., Durham, Greensboro, Wilmington, and Charlotte, North Carolina; Nashville, Tennessee; Tampa, Florida). Readers also scored remotely through Virtual Scoring Center (VSC Score) (i.e., distributive scoring).

AERA, APA, & NCME (2014) Standard 4.20 specifies the following:

The process for selecting, training, qualifying, and monitoring scorers should be specified by the test developer. The training materials, such as the scoring rubrics and examples of test takers’ responses that illustrate the levels on the rubric score scale, and the procedures for training scorers should result in a degree of accuracy and agreement among scorers that allows the scores to be interpreted as originally intended by the test developer. Specifications should also describe processes for assessing scorer consistency and potential drift over time in raters’ scoring. (p. 92)

Sections 7.2.1 through 7.2.5 explain how scorers are selected and trained for the M-STEP handscoring process. Sections 7.2.6 and 7.2.7 describes how the scorers are monitored throughout the M-STEP handscoring process.

7.2.1 Security

All Measurement Incorporated scoring rooms are designated secure areas with stringent security regulations that are vigorously enforced. Measurement Incorporated routinely implements a number of measures to help safeguard the security of student responses while they are in Measurement Incorporated’s possession and to maintain the confidentiality of student identity.

In the scoring rooms, the use of cellphones, tablets, MP3 players, laptops, or recording or photographic equipment is prohibited. The copying of materials for anything other than the training purposes that are expressly permitted by MDE is prohibited.

All buildings that house student responses, including Measurement Incorporated headquarters, scoring centers, and warehouses, utilize an electronic security system during nonbusiness hours.

All readers scoring remotely are required to work from a private, password-protected environment. No free or public Wi-Fi can be used. Readers can access a project website only from a secure, password-protected network. Readers cannot access any project website from a public computer or a public network, such as a wireless network at a hotel or restaurant. While in VSC Score, readers are unable to take screenshots or to access e-mail or other applications. Maintaining a secure workstation is a condition for employment for all remote employees.

Before receiving any training materials, all scoring project staff are required to sign a confidentiality and proprietary agreement, which indicates that no participant in training and/or scoring may reveal any specific information about the test or about the criteria and methods for scoring to any person as part of his or her contractual obligation to score student responses.

At scoring centers, all training materials remain on the premises during a project and are collected at the end of each workday to be secured. All materials are collected and accounted for at the end of the scoring project.

Readers who score remotely access training materials from an online resource library. The software does not allow readers to print or download data.

No identifying student information is provided on the images sent to readers via VSC Score software.

Readers do not have the ability to access training materials or student responses unless they and their team leader are logged on to the system.

Violation of any portion of the Measurement Incorporated security policy results in termination.

7.2.2 Measurement Incorporated Reader and Team Leader Hiring

Measurement Incorporated recruits, interviews, and hires a pool of readers to ensure sufficient staff for scoring projects.

All readers must have a minimum of a bachelor's degree. MDE has the right to review the names, demographics, educational backgrounds, and experience (including scoring experience) of all readers. Reader degrees are verified before the applicants are interviewed. Applicants must provide either an official transcript with a seal (no copies accepted), an official letter from a registrar's office (which would be mailed to the Site Manager), or access to a third-party company such as Parchment or Student Clearing House. Reader applicants can also bring their original diploma with a seal when they come for an interview.

Team leaders are selected and recruited from experienced reader staff. Each team leader supervises a group of 10–12 readers during live scoring.

7.2.3 Preparation of Training Materials for M-STEP

Three types of sets of student responses were used in training readers and team leaders:

- Anchor sets consisted of typical student responses at each score point, with examples of what would barely earn that point, a median answer for that point, and a high response within that point without quite reaching the next point. These sets were used to show readers and team leaders how the rubric was applied to each response.
- Training sets consisted of atypical student responses and were used to further demonstrate application of the rubric to actual student responses.
- Qualifying sets consisted of student responses similar to those in the anchor and training sets. These sets were used for readers to demonstrate their understanding of the application of the rubric to student responses.

Measurement Incorporated scoring directors used MDE-approved training materials. Anchor sets consisted of three responses at each score point. Each response was annotated to explain how the rubric criteria were applied. Training sets contained 5–10 papers. There was a training set for each trait for analytic scoring and a training set that combined the traits. The responses in each of these sets were arranged in random score-point order, and all score points were represented.

7.2.4 Training and Qualifying Readers and Team Leaders

AERA, APA, & NCME (2014) Standard 6.9 specifies the following:

Those responsible for test scoring should establish and document quality control processes and criteria. Adequate training should be provided. The quality of scoring should be monitored and documented. Any systematic source of scoring errors should be documented and corrected. (p. 118)

Readers and team leaders were trained by the scoring director on the scoring criteria approved by MDE and were required to achieve qualifying standards set by MDE.

Readers were divided into teams consisting of one team leader and 10–15 readers.

For brief write, research, and mathematics scoring, the scoring director presented the item and anchor set and then discussed each score point as readers and team leaders took notes.

Following the presentation of these anchor sets, readers and team leaders scored a training set and then one or two qualifying sets.

For full write scoring, the scoring director introduced the readers and team leaders to the three analytic traits (i.e., Organization/Purpose, Evidence/Elaboration, and Conventions) using a unique anchor set for each trait so that readers and team leaders were fluent in the individual traits before they scored the traits simultaneously.

Following the presentation of each trait anchor set, a training set was scored for the trait and discussed; readers and team leaders then prepared to score all traits concurrently. Readers and team leaders took two qualifying sets, in which scores were assigned for all three analytic traits

on each student response.

Readers and team leaders were provided a copy of anchor sets, training sets, and qualifying sets. Readers and team leaders were required to refer to the anchor sets and their notes when taking training sets and qualifying sets.

Readers and team leaders scored the qualifying set and submitted their scores. The percentage of correct scores was recorded. After the set was completed, the scoring director discussed the set with the group.

If a particular response or type of response generated numerous questions across teams, the scoring director discussed the problem with the group or posted a note to chat to ensure that everyone heard the same explanation.

Once the group had finished discussing the first qualifying set, the readers and team leaders scored the next set. Training continued until all training sets and qualifying sets were scored and discussed.

Readers were required to demonstrate their ability to score accurately by attaining the qualifying agreement percentage approved by MDE before they gained access to actual student responses.

Any reader or team leader unable to meet the qualifying standards set by MDE was released.

Reference Tables 7-1 and 7-2 for additional information.

Table 7-1. Qualifying Sets

Content	Number of Qualifying Sets per Item
Math	1 or 2
Research	1
Brief Write	1
Full Write	2 for each trait

Table 7-2. Qualifying Standards

Score Point Range	Qualifying Standard (Exact Agreement)
0–1	90%; no nonadjacent scores
0–2	80%; no nonadjacent scores
0–3	70%; no nonadjacent scores

7.2.5 Virtual Scoring Center

Measurement Incorporated used its VSC Score system for the image-based scoring of paper/pencil responses and for the scoring of online responses transferred to Measurement Incorporated from DRC.

Readers and team leaders accessed the VSC Score system through a secure web-based interface with the use of a unique user ID and password. Each team leader and reader was assigned a unique number for easy identification of his or her scoring work throughout the scoring session. VSC Score enabled readers and team leaders to score only those items that they were trained and qualified to score.

Each PBW response was randomly assigned to be read by one reader. A random sample of all student responses (i.e., 10% of responses) was then randomly assigned to a second reader. VSC Score managed readers' individual workloads and allowed readers to review and submit their scores.

Readers were trained on how to use the VSC Score performance assessment scoring system—how to assign scores, how to adjust the image for legibility, how to “flag” responses that were atypical from the anchor sets, training sets, and qualifying sets for review by the team lead and scoring director, etc.

Readers logged in and “checked out” a scoring set of student responses. This scoring set was generated by randomly selecting student responses from the pool of unscored student responses. The reader evaluated the first response, entered the score by clicking the appropriate value on the scoring toolbar, and clicked the “Submit” button. The next response in the scoring set then appeared for the reader to score and submit. This process continued until all responses in the set had been scored. After scoring all responses in a set, the reader had the option to review any of the responses and modify the scores before submitting them to the system.

Once the scores had been submitted, the set was “checked in” and responses were routed to other qualified readers as necessary. The requirements for subsequent readings were defined in the system during setup, and student responses were not marked as complete until the requisite number of independent readers had scored the response.

When a reader had a question about a response, he or she could transfer the image (along with the question and/or comments) from the current scoring set to a review set, which was assigned to a team leader. The team leader could forward the question to the scoring director, submit the appropriate score, or return the response to the reader with comments. This procedure was used whenever a reader had scoring concerns or encountered apparent non-scorable responses. Readers could mark completely blank responses as non-scorable, but otherwise only scoring directors or the project director could assign a non-scorable condition code to a student response.

7.2.6 Quality Control and Reliability of Scoring

AERA, APA, & NCME (2014) Standard 6.8 states the following:

Those responsible for test scoring should establish scoring protocols. Test scoring that involves human judgment should include rubrics, procedures, and criteria for scoring. When scoring of complex responses is done by computer, the accuracy of the algorithm and processes should be documented. (p. 118)

Section 7.2.6 explains the monitoring procedures that Measurement Incorporated uses to ensure that handscoring evaluators follow established scoring criteria while items are being scored. Detailed scoring rubrics are available for all PBW prompts, which specify the criteria for scoring those PBW prompts. These rubrics will not be presented in this report in order to preserve the integrity of the items for use in future MAP forms.

MDE reader production and reliability statistics, including reader training results, were available to MDE via a suite of VSC reports, which could be accessed online using secure credentials supplied to MDE staff.

Detailed Reader Status Reports were generated for each scoring project, utilizing a comprehensive system for collecting and analyzing score data. Daily analyses of the Reader Status Reports alerted management personnel to individual or group retraining needs.

After the readers' scores were submitted in the VSC Score system, the data was uploaded into the primary Scoring Resource Center servers. The scores were then validated and processed.

Updated real-time reports that showed both daily and cumulative data (i.e., project-to-date data) were available 24 hours a day via a secure website. Reports included data on the number of responses scored by each reader, the percentage of responses scored that day in exact agreement or adjacent agreement with a second reader, and the total number of responses scored at each score point.

For M-STEP performance assessment scoring, a random sample of 10% of all student responses are scored a second time to generate agreement data.

Readers were required to consistently demonstrate the ability to assign scores according to the rubric and anchor papers that were introduced during training. Their scoring accuracy is under scrutiny using validity responses that were included daily with the actual student responses (for details, see Section 7.2.7).

If questionable reader reliability indications were found, the affected responses were scored again.

The monitoring and retraining process was sustained throughout the project to promote strict adherence to MDE-approved scoring criteria and consistency throughout the scoring effort.

Scoring directors and team leaders provided consistent monitoring of the scoring patterns of each reader throughout the project, responded to questions, spot-checked (i.e., read behind) reader scoring, provided feedback, and counseled readers who were having difficulty with the criteria.

Scoring directors continued to look for atypical types of responses that were not covered in the initial training and presented further instruction about handling these types of responses when necessary.

The inter-rater reliability information for the handscored ELA and mathematics items is presented in Table 7-3.

Table 7-3. Human-to-Human Inter-rater Reliability

Content	Grade	Item ID	Maturity	% Perfect Plus	N Perfect	% Perfect	N Adjacent	% Adjacent	N Nonadjacent	% Nonadjacent
ELA	3	945967	OP	99.4	1611	87	229	12.4	11	0.6
ELA	3	945969	OP	98.1	1541	78.7	381	19.4	37	1.9
ELA	3	945975	OP	99	1497	77.6	414	21.5	19	1
ELA	3	945976	OP	98.9	1309	78.4	341	20.4	19	1.1
ELA	4	945962	OP	99.4	2061	78	565	21.4	15	0.6
ELA	4	945979	OP	98.5	968	73.8	324	24.7	20	1.5
ELA	4	945980	OP	99.1	2296	74.4	760	24.6	28	0.9
ELA	4	945981	OP	98.9	2059	78.1	549	20.8	28	1.1
ELA	5	945968	OP	99.3	3154	74	1076	25.3	30	0.7
ELA	5	945974	OP	98.8	2,943	75	937	23.9	46	1.2
ELA	5	945982	OP	98.8	2882	77.3	802	21.5	46	1.2
ELA	5	945983	OP	99	3,385	76.4	1001	22.6	46	1
ELA	6	945964	OP	99.1	3,192	79.3	799	19.8	35	0.9
ELA	6	945965	OP	98.6	2839	75.5	868	23.1	54	1.4
ELA	6	945972	OP	98.4	3072	71.5	1157	26.9	67	1.6
ELA	6	945984	OP	98.5	2840	72.1	1037	26.3	60	1.5
ELA	7	945966	OP	99.3	3114	74.1	1057	25.2	30	0.7
ELA	7	945970	OP	98.7	2941	69.5	1234	29.2	57	1.3
ELA	7	945971	OP	99.3	3000	74.9	979	24.4	28	0.7
ELA	7	945985	OP	99.2	3,287	75.6	1025	23.6	36	0.8
ELA	8	945963	OP	99.4	4306	78.5	1142	20.8	35	0.6
ELA	8	945973	OP	98.8	3468	78.2	914	20.6	55	1.2
ELA	8	945977	OP	99.2	3,494	77.6	974	21.6	36	0.8
ELA	8	945978	OP	99.2	3,615	77.1	1039	22.1	37	0.8
Math	3	1906	OP	100	82	94.3	5	5.7		
Math	3	1928	OP	100	87	98.9	1	1.1		
Math	3	76505	OP	100	82	94.3	5	5.7		
Math	3	76529	OP	100	87	98.9	1	1.1		
Math	4	2620	OP	100	96	100				
Math	4	2666	OP	100	96	100				
Math	4	78588	OP	100	94	100				
Math	5	2228	OP	100	94	100				
Math	5	5545	OP	100	88	95.7	4	4.3		
Math	5	78500	OP	100	88	95.7	4	4.3		

Content	Grade	Item ID	Maturity	% Perfect Plus	N Perfect	% Perfect	N Adjacent	% Adjacent	N Nonadjacent	% Nonadjacent
Math	6	10839	OP	100	80	96.4	3	3.6		
Math	6	684	OP	100	83	98.8	1	1.2		
Math	6	77674	OP	100	82	100				
Math	7	12069	OP	100	61	100				
Math	7	25989	OP	100	60	98.4	1	1.6		
Math	7	5577	OP	100	60	100				
Math	7	7180	OP	100	60	98.4	1	1.6		
Math	7	9499	OP	100	61	100				
Math	8	21111	OP	100	56	100				

7.2.7 Validity

Measurement Incorporated used validity responses, similar to the student responses found in the qualifying sets, during live scoring to monitor readers' accuracy in scoring. Preselected validity responses were approved by MDE. Scoring directors also had the ability to select live responses as validity responses, which were also subject to MDE approval. The true scores for these responses were entered into a validity database.

Validity responses were randomly incorporated into readers' sets each day of the project. Team leaders reviewed the validity results and provided feedback to the readers.

A validity report was generated that included the response identification number, the scores assigned by the readers, and the "true" scores. Measurement Incorporated provided MDE with daily and project-to-date summaries of what percentages of papers scored by readers matched the validity checks or were high or low at each score point. Five percent of the responses that a reader scored were validity papers. These responses appeared to the reader daily throughout the entire scoring project. The validity standards can be found in Table 7-4.

Table 7-4. Validity Standards

Score Point Range	Validity Standard (Exact Agreement)
0–1	90%
0–2	80%
0–3	80%
0–4	70%

7.2.8 Alerts

Measurement Incorporated implemented a formal process for notifying MDE when student responses reflected a possibly dangerous situation for the student, which may include responses indicating endangerment, abuse, or psychological and/or emotional difficulties.

Measurement Incorporated also alerted MDE if there appeared to be possible instances of teacher or proctor interference or student collusion with other students.

Measurement Incorporated always takes immediate action following a scoring alert.

7.3 Summary

The information presented in this chapter summarizes the scoring procedures for different types of items and the steps taken by DRC and Measurement Incorporated to ensure accuracy in the technology-enhanced item scoring and handscoring processes. The reliability statistics presented in Sections 7.2.6 and 7.2.7 demonstrate that the items are scored reliably. These efforts follow multiple best practices of the testing industry, particularly AERA, APA, & NCME (2014) *Standards* 4.18 4.20, 6.8, and 6.9:

- Standard 4.18—Procedures for scoring and, if relevant, scoring criteria, should be presented by the test developer with sufficient detail and clarity to maximize the accuracy of scoring. Instructions for using rating scales or for deriving scores obtained by coding, scaling, or classifying constructed responses should be clear. This is especially critical for extended-response items such as performance tasks, portfolios, and essays.
- Standard 4.20—The process for selecting, training, qualifying, and monitoring scorers should be specified by the test developer. The training materials, such as the scoring rubrics and examples of test takers' responses that illustrate the levels on the rubric score scale, and the procedures for training scorers should result in a degree of accuracy and agreement among scorers that allows the scores to be interpreted as originally intended by the test developer. Specifications should also describe processes for assessing scorer consistency and potential drift over time in raters' scoring.
- Standard 6.8—Those responsible for test scoring should establish scoring protocols. Test scoring that involves human judgment should include rubrics, procedures, and criteria for scoring. When scoring of complex responses is done by computer, the accuracy of the algorithm and processes should be documented.
- Standard 6.9—Those responsible for test scoring should establish and document quality control processes and criteria. Adequate training should be provided. The quality of scoring should be monitored and documented. Any systematic source of scoring errors should be documented and corrected.

Chapter 8: Operational Data Analyses

This chapter describes the analyses conducted with the operational (OP) data. Item/test analyses from both the Classical Test Theory (CTT) and the item response theory (IRT) frameworks are used (when appropriate) and reported here.

This chapter demonstrates adherence of M-STEP to AERA, APA, & NCME (2014) *Standards* 1.8, 5.2, 5.13, and 5.15. Each standard will be explicated within the appropriate section of this chapter. Standard 7.2 provides general guidance that is relevant to this chapter:

The population for whom a test is intended and specifications for the test should be documented. (p. 126)

Chapter 3 presents the test specifications. Information regarding reported data is discussed in detail in Chapter 9.

8.1 Operational Analysis of ELA and Mathematics

The *Smarter Balanced 2016–2017 Technical Report* (2017) states that part of the Smarter Balanced Theory of Action is to leverage appropriate technology and innovation. The use of CAT methodologies helps ensure that students across the range of proficiency levels have an assessment experience with items well targeted to their skill level. Adaptive testing allows average-, low-, and high-performing students to stay engaged in the assessment because they respond to items specifically targeted to their skill level. CAT tests are also efficient because they provide a higher level of score precision than fixed-form tests with the same number of items. For the CAT component, there are both content constraints (e.g., a long reading passage in ELA must be administered) and psychometric criteria that must be optimized for each student.

8.1.1 CAT Item Pool Characteristics

8.1.1.1 CAT Item Types

This section presents different item types used by Smarter Balanced (*Smarter Balanced 2017–2018 Technical Report*, 2018, p. 4–28) to compose the summative item pools. These different item types are listed in Table 8-1.

Table 8-1. Item Types Found in the Summative Item Pools

Item Types	ELA	Mathematics
Multiple Choice (MC)	X	X
Multi-Select (MS)	X	X
Evidence-Based Selected Response (EBSR)	X	
Match Interaction (MI)	X	X
Hot Text (HTQ)	X	
Passage Based Writing (PBW)	X	
Equation Response (EQ)		X
Grid-Item Response (GI)		X
Table Interaction (TI)		X
Table Interaction (TI)		X
Constructed Response (CR)		X

The Smarter Balanced item/task type characteristics are defined as sufficient to ensure that the content measured the intent of the Common Core State *Standards* (CCSS) and that there was consistency across item/task writers and editors. This included item types such as selected-response, constructed-response, and technology-enhanced.

As shown in Table 8-1, the common item types for both ELA and mathematics are MC, MS, and MI. In addition, ELA also included the following item types: EBSR, HTQ, and PBW. Mathematics also included the following item types: EQ, GI, and TI.

For ELA, PBW prompts are included in the pool. For more information on PBW prompts, please see Chapter 3. Additionally, it should be noted that the following sections provide information about the ELA and mathematics item pools that were administered in Michigan.

8.1.1.2 CAT Item Pool Specification

“An item pool refers to a collection of test questions (known as items) that support the test blueprint for a particular content area and grade. The Consortium takes multiple steps to ensure the quality of the items in the Smarter Balanced item pool. Building on the continuing process of developing item/task specifications and test blueprints, the Consortium uses an iterative process for creating and revising each item as well as the collection of items” (*Smarter Balanced 2017–2018 Technical Report*, 2018, p. 4–17).

8.1.1.3 CAT Distribution of Item Types

The M-STEP distribution of item types is shown in Tables 8-2 and 8-3.

Table 8-2. Distribution of ELA Item Types by Grade and Claim

Grade	Claim	MC	MS	EBSR	MI	HTQ	PBW	Total
3	1	167	48	47		54		316
3	2	119	57			56	4	236
3	3	77	39	47	20			183
3	4	64	40		4	22		130
3	Total	427	184	94	24	132	4	865
4	1	108	49	48		52		257
4	2	130	48			57	4	239
4	3	89	38	47	21			195
4	4	64	49		2	21		136
4	Total	391	184	95	23	130	4	827
5	1	118	65	56		48		287
5	2	116	62			45	4	227
5	3	65	27	37	16			145
5	4	57	46			27		130
5	Total	356	200	93	16	120	4	789
6	1	78	50	38		59		225
6	2	91	68			57	4	220
6	3	78	24	41	19			162
6	4	68	58			18		144
6	Total	315	200	79	19	134	4	751
7	1	85	45	36		45		211
7	2	86	62			51	4	203
7	3	68	30	42	14			154
7	4	26	27		3	31		87
7	Total	265	164	78	17	127	4	655
8	1	84	50	42		52		228
8	2	94	77			48	4	223
8	3	117	36	24	5			182
8	4	41	32		2	32		107
8	Total	336	195	66	7	132	4	740

Table 8-3. Distribution of Mathematics Item Types by Grade and Claim

Grade	Claim	MC	MS	MI	EQ	GI	TI	Total
3	1	124	4	77	500	71	33	809
3	2	14	5	6	54	22		101
3	3	87	32	22	12	60		213
3	4	32	11	8	48	20		119
3	Total	257	52	113	614	173	33	1242
4	1	111		191	453	58	14	827
4	2	30	4	7	58	6	2	107
4	3	66	33	19	22	77		217
4	4	60	10	4	32	13	6	125
4	Total	267	47	221	565	154	22	1276
5	1	232	1	88	438	45		804
5	2	8	2	2	64	12	2	90
5	3	87	24	19	15	58	2	205
5	4	27	6	6	34	30	4	107
5	Total	354	33	115	551	145	8	1206
6	1	70	136	105	367	70	14	762
6	2	7	12	3	51	12	2	87
6	3	47	41	31	18	51		188
6	4	10	12	3	52	13	4	94
6	Total	134	201	142	488	146	20	1131
7	1	58	128	74	390	32		682
7	2	10	13	6	54	8	3	94
7	3	34	33	17	16	43		143
7	4	14	11	2	49	16	1	93
7	Total	116	185	99	509	99	4	1012
8	1	164	80	82	255	39	16	636
8	2	5	3	4	33	9	3	57
8	3	31	33	21	8	47		140
8	4	11	8	5	18	17	4	63
8	Total	211	124	112	314	112	23	896

8.1.1.4 Item Pool Calibration and Model Fit Evaluation

Item parameters contained in ELA and mathematics tests were estimated using a marginal maximum-likelihood procedure with either the 2-parameter logistic (2PL) model for MC items or the generalized partial credit model (Muraki, 1992) for technology-enhanced (TE) items administered after the 2013–14 Smarter Balanced field-test administration. Additionally, for model fit, the evaluation of goodness-of-fit used the likelihood ratio test in PARSCALE (Muraki & Bock, 2003).

For details on item calibration and model fit for ELA and mathematics, please refer to Chapter 9 of the Smarter Balanced 2013–2014 Technical Report (2015), which was published on the [Smarter Balanced website](#).¹

With the exception of the PBW prompts, Smarter Balanced ELA and mathematics operational item parameters were used to score Michigan students who took ELA and mathematics assessments.

To place the PBW prompts on the same scale as the Smarter Balanced item pool that is used for M-STEP ELA, DRC used a common-item, non-equivalent groups design to link the PBW prompts to the established scale. After the initial IRT item calibration, item parameters were linked to the existing M-STEP scale using the Stocking & Lord (1983) equating procedure, where all other items in the pool were used as anchor items.

8.1.2 Item Pool IRT Statistics

The distributions of item parameters by grade and claim are shown in Tables 8-4 and 8-5. Item difficulty is represented by the b-parameter, and discrimination is represented by the a-parameter. Note that there is a wide range of difficulty in each category.

¹ <http://www.smarterbalanced.org/wp-content/uploads/2015/08/Chapter-9-Field-Test-IRT.pdf>

Table 8-4. Distribution of Item Difficulty (b-parameter) and Discrimination (a-parameter) for ELA

Grade	Claim	N Items	Difficulty Mean	Difficulty Min	Difficulty Max	Discrimination Mean
3	1	317	-0.531	-2.596	4.693	0.704
3	2	237	-0.796	-2.896	4.115	0.685
3	3	183	-0.163	-2.920	3.815	0.544
3	4	130	-0.380	-2.216	1.699	0.683
3	Total	867	-0.503	-2.920	4.693	0.662
4	1	260	0.304	-2.529	6.233	0.619
4	2	239	-0.383	-3.252	2.935	0.590
4	3	195	0.027	-2.822	4.254	0.552
4	4	136	0.394	-1.996	3.727	0.562
4	Total	830	0.056	-3.252	6.233	0.585
5	1	288	0.673	-2.360	5.651	0.611
5	2	228	0.013	-2.278	3.294	0.600
5	3	145	0.526	-2.403	3.481	0.520
5	4	130	0.453	-1.494	3.832	0.665
5	Total	791	0.420	-2.403	5.651	0.600
6	1	225	1.001	-1.636	4.779	0.578
6	2	220	0.860	-2.719	5.542	0.536
6	3	163	0.870	-1.497	7.385	0.496
6	4	144	0.941	-1.305	3.609	0.560
6	Total	752	0.920	-2.719	7.385	0.545
7	1	212	1.233	-1.836	6.630	0.551
7	2	203	1.131	-2.019	5.305	0.514
7	3	154	0.883	-1.706	5.885	0.502
7	4	88	1.692	-0.815	5.525	0.533
7	Total	657	1.181	-2.019	6.630	0.526
8	1	228	1.495	-1.170	6.421	0.557
8	2	224	1.122	-2.192	4.558	0.502
8	3	185	0.885	-2.119	3.871	0.484
8	4	107	1.472	-1.788	5.188	0.566
8	Total	744	1.228	-2.192	6.421	0.523

Table 8-5. Distribution of Item Difficulty (b-parameter) and Discrimination (a-parameter) for Mathematics

Grade	Claim	N Items	Difficulty Mean	Difficulty Min	Difficulty Max	Discrimination Mean
3	1	810	-1.134	-4.338	4.163	0.839
3	2	101	-0.300	-2.537	1.967	0.991
3	3	213	-0.115	-2.424	5.116	0.724
3	4	119	-0.076	-2.677	3.201	0.790
3	Total	1243	-0.790	-4.338	5.116	0.827
4	1	827	-0.296	-3.260	4.483	0.848
4	2	107	0.214	-2.248	2.574	0.887
4	3	218	0.248	-2.083	5.184	0.740
4	4	126	0.255	-2.148	3.284	0.691
4	Total	1278	-0.106	-3.260	5.184	0.817
5	1	804	0.305	-2.791	3.606	0.774
5	2	90	1.082	-1.267	3.939	0.957
5	3	205	0.764	-1.903	5.278	0.658
5	4	108	1.156	-1.232	4.634	0.691
5	Total	1207	0.517	-2.791	5.278	0.761
6	1	767	0.853	-2.846	9.158	0.691
6	2	87	1.567	-2.978	5.498	0.769
6	3	188	1.859	-2.161	8.754	0.581
6	4	95	1.849	-0.715	6.440	0.782
6	Total	1137	1.157	-2.978	9.158	0.686
7	1	694	1.724	-1.792	7.801	0.723
7	2	94	2.146	-1.085	5.596	0.800
7	3	144	2.264	-1.645	6.594	0.563
7	4	93	2.178	-0.789	4.777	0.708
7	Total	1025	1.880	-1.792	7.801	0.706
8	1	637	2.066	-1.868	7.752	0.578
8	2	59	2.862	0.046	5.751	0.716
8	3	142	3.190	-0.993	9.022	0.438
8	4	64	2.425	-1.364	6.476	0.623
8	Total	902	2.320	-1.868	9.022	0.568

It is also beneficial to examine the distribution of item difficulty compared to the distribution of abilities across the student population. This can be used to ensure that the item pool is deep enough to measure the abilities of the student population without item exposure rates being too high. Figures 8-1 and 8-2 show the comparison of item difficulty, student scores, and cut scores for ELA and mathematics, respectively. For most grades, the item pool has good alignment with the student ability distribution. However, in grades 6 to 8 for mathematics, the item pool appears to be more difficult when compared to the corresponding student ability distribution.

Figure 8-1. ELA Item Pool Difficulty in Comparison to the Student Ability Distribution

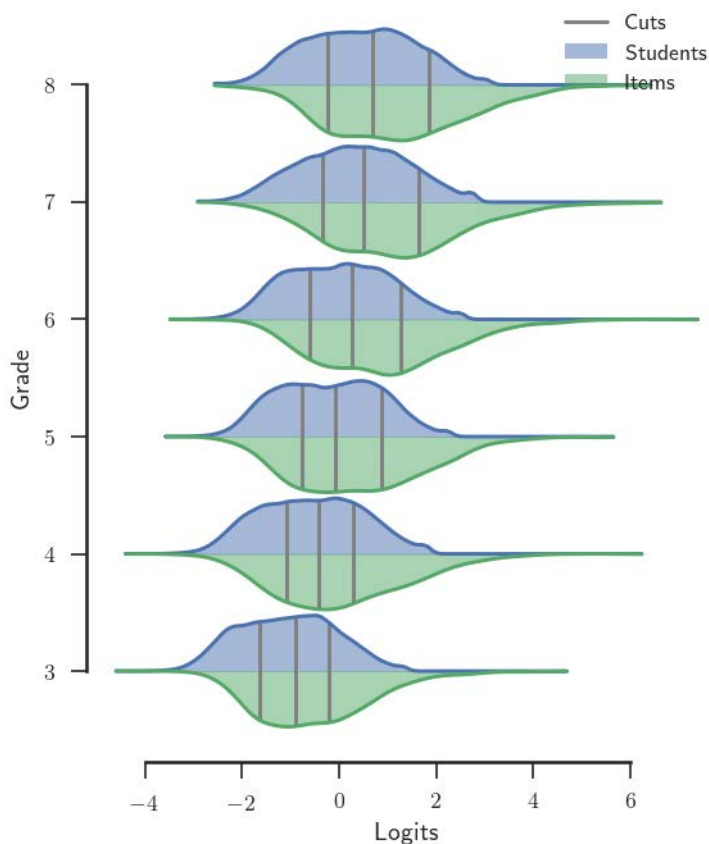
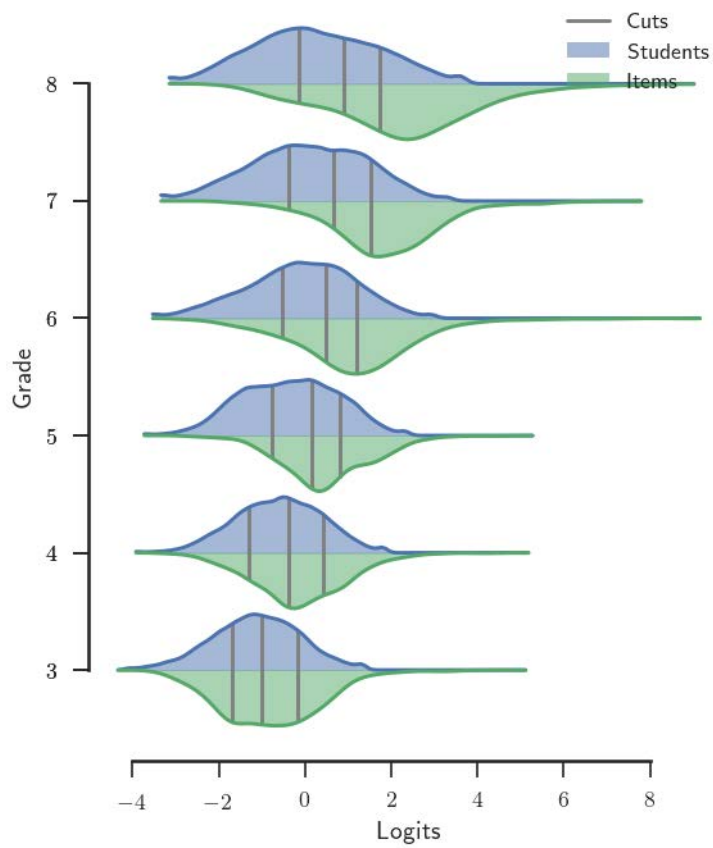


Figure 8-2. Mathematics Item Pool Difficulty in Comparison to the Student Ability Distribution



8.2 Operational CAT ELA and Mathematics Implementation

8.2.1 The Scale

The scales on which M-STEP ELA and mathematics scale scores are reported were established by Smarter Balanced after the 2014 field test. The underlying scales are not unique to Michigan but have been adapted by several states that were members of the Smarter Balanced Consortium. Michigan has used the underlying scale to create state-specific M-STEP scales used solely by Michigan.

The Smarter Balanced ELA and mathematics scores are reported on vertical scales, sometimes referred to as growth scales, showing student progress from grade to grade. For details on ELA and mathematics vertical scale development, refer to Chapter 9 of the *Smarter Balanced 2013–2014 Technical Report* (2015), which is posted on the [Smarter Balanced website](http://www.smarterbalanced.org).² However, the scale scores reported by Michigan should not be considered vertical scale scores.

Additional information regarding M-STEP scale scores can be found in Chapter 10.

8.2.2 Lowest and Highest Obtainable Scale Scores (LOSS and HOSS)

A maximum likelihood procedure cannot produce scale-score estimates for students with perfect scores or scores below the level expected by guessing. In addition, although maximum likelihood estimates are available for students with extreme scores other than zero or perfect, occasionally these estimates have standard errors of measurement that are very large and differences between these extreme values have little meaning. Therefore, scores are established for these students based on a rational but necessarily non-maximum likelihood procedure. These values, which are set separately by grade, are called the lowest obtainable theta (LOT) and the highest obtainable theta (HOT). For reporting purposes, the LOT and HOT are transformed, using a linear transformation, to the LOSS and the HOSS. For more information on the LOSS and HOSS, see Chapter 10 and Table 10-1.

8.2.3 Item-Pattern Scoring

M-STEP scale scores are derived using item-pattern scoring; thus, these scale scores are based on the student's responses to all items on a given test and account for the characteristics of the test items, such as item difficulty. A scale score can be interpreted as a highly probable estimate of a student's ability in a given content area.

Using item-pattern scoring, a student's scale score is based on the student's response to each item (i.e., his or her item-response vector). Each item uses optimal item weights in terms of item information, meaning that items do not contribute equally to the overall scale score. Students with the same raw score may be assigned to different scale scores, depending on which items they answered correctly.

² <http://www.smarterbalanced.org/wp-content/uploads/2015/08/Chapter-9-Field-Test-IRT.pdf>

8.2.4 Blueprint Fidelity Summary

The *Smarter Balanced 2017–2018 Technical Report* notes that “A key design document of the summative assessments is the test blueprint, which specifies the number and nature of items to be administered” (p. 6–9). Chapter 6 of that same report states that “The analyses showed that the operational tests delivered in the 2016–2017 administration fulfilled the blueprint requirements very well” (p. 6–9).

For M-STEP ELA and mathematics implementation, a review of the blueprint fulfillment was completed for both the simulation (see Chapter 6 of this report) and the OP tests.

8.3 Operational Analysis of Social Studies

This section describes analyses conducted for social studies and reports corresponding results. As mentioned above, item/test analyses from the CTT and IRT frameworks have been carried out. They are reported below separately. If the IRT models fit the empirical item-response data for the population for which generalizations are made (i.e., Michigan students), then it is likely that the scores are valid indicators of an underlying ability.

8.3.1 CTT Statistics Social Studies

This section presents test-level summary statistics for each form and grade of social studies. This is followed by item-level statistics for each form and grade of social studies. These statistics were produced using census data.

8.3.1.1 Test-Level Analysis

This section presents the test-level summary statistics for social studies. In addition to the number of students taking the form (N), Table 8-6 provides the following raw score descriptive statistics for a given grade and form: mean, standard deviation (SD), minimum (Min), and maximum (Max).

Table 8-6. Test-Level Descriptive Statistics by Form: Social Studies Raw Score Distribution

Grade	N OP Items	Form	N	Mean	SD	Min	Max
5	45	1	36009	22.45	8.00	3	45
5	45	2	36069	22.60	7.99	3	45
5	45	3	36062	22.46	7.97	3	45
5	45	4	936	18.57	7.21	5	43
8	44	1	36570	22.19	8.57	0	44
8	44	2	36599	22.15	8.58	2	44
8	44	3	36557	22.06	8.58	2	44
8	44	4	653	17.58	7.40	4	40
11	38	1	35012	21.16	8.07	0	38
11	38	2	35062	21.18	8.05	0	38
11	38	3	35171	21.08	8.05	0	38
11	38	4	728	19.43	7.99	3	38

8.3.1.2 Item-Level Analysis

This section presents various item-level statistics for all OP items³ on the spring 2018 M-STEP social studies tests. Specifically, item difficulty and adjusted item-total correlations defined by the CTT are reported here.

Since all items on the spring 2018 M-STEP social studies tests are dichotomously scored, the p -value is computed as an indicator for item difficulty. The p -value equals the proportion of students who answer an item correctly. A high p -value means that an item is easy, and a low p -value means that an item is difficult.

The adjusted item-total correlation is an index of the association between students' performance on an item and their performance on the test as a whole; however, the item of interest is excluded from the total raw score. A high adjusted item-total correlation is desired, as high correlations indicate that students with high scores on all other test items (i.e., students with high ability) tend to get a correct answer on the item under consideration and that the students with low scores on all other test items (i.e., students with low ability) tend to get this specific item incorrect. Since all items are dichotomously scored, the adjusted point biserial correlation is computed.

The item-level descriptive statistics (by grade and form) for all OP items are presented in Tables 8-7 and 8-8 for social studies. For each grade, forms 1 through 3 are online forms and form 4 is the paper/pencil form.⁴ All online forms for social studies have the same set of OP

³ All statistics for field-test items are excluded from this report.

⁴ One emergency form and one braille form per grade were also created. However, responses from these forms are excluded from any analysis due to negligible occurrences (for braille forms, which are exactly the same as the corresponding paper/pencil forms) and a different calibration approach (for emergency forms, banked values from the item pool would be activated for scale-score computation). In 2018, no emergency forms were used at all.

items; therefore, forms 1 through 3 are reported together per grade in Tables 8-7 and 8-8. Each paper/pencil form for social studies per grade has a few different OP items from its online counterpart because the TE items could not be presented on paper/pencil forms. As shown in Tables 8-7 and 8-8, for both item difficulty (p -value) and adjusted item-total correlation (adjusted point biserial), the following descriptive statistics are reported: number of OP items (N OP Items), mean, SD, minimum (Min), and maximum (Max), for a given grade by form.

Table 8-7. Item-Level Descriptive Statistics by Form: Social Studies P -Value

Grade	N OP Items	Form	Mean	SD	Min	Max
5	45	1–3	0.50	0.12	0.28	0.86
5	45	4	0.48	0.12	0.28	0.86
8	44	1–3	0.50	0.10	0.32	0.78
8	44	4	0.48	0.11	0.27	0.69
11	38	1–3	0.56	0.12	0.21	0.76
11	38	4	0.56	0.12	0.33	0.76

Table 8-8. Item-Level Descriptive Statistics by Form: Social Studies Adjusted Point Biserial

Grade	N OP Items	Form	Mean	SD	Min	Max
5	45	1–3	0.31	0.06	0.19	0.42
5	45	4	0.31	0.06	0.19	0.42
8	44	1–3	0.35	0.09	0.22	0.53
8	44	4	0.33	0.09	0.15	0.53
11	38	1–3	0.39	0.08	0.17	0.53
11	38	4	0.39	0.08	0.24	0.53

8.3.2 IRT Statistics: Social Studies

The unidimensional 2PL model is used for M-STEP social studies at each grade level, as all items are dichotomously scored. For this model, the probability that person j answers item i correctly is defined as follows (adapted from Embretson & Reise, 2000, p. 70):

$$P(X_{ij} = 1 | \theta_j, b_i, a_i) = \frac{\exp(a_i(\theta_j - b_i))}{1 + \exp(a_i(\theta_j - b_i))}, \quad (8.1)$$

where θ_j , b_i , and a_i represent person j 's ability, item i 's difficulty, and item i 's discrimination, respectively. Note that $P(X_{ij} = 1 | \theta_j, b_i, a_i)$ is referred to as for simplicity $P_i(\theta)$ hereafter.

8.3.3 Item Calibration for Social Studies

The common-item nonequivalent groups design (Kolen & Brennan, 2004) and the fixed item parameter calibration approach were used to put all items onto the base scale. The IRT software used was flexMIRT (Cai, 2017). An outline of the annual calibration, equating, and scaling procedures for M-STEP social studies is presented below:

- A free run was carried out using flexMIRT (i.e., a 2PL model) with all online OP items for each grade.
- After each free run, the obtained item parameters for anchor items were compared with their banked values (all banked values have already been put onto the base scale). In addition to checking the scatterplots of item difficulty parameters and item-discrimination parameters, a simple linear regression was carried out using free run results as outcomes (i.e., one regression for item difficulty parameters and one regression for item-discrimination parameters) and the corresponding banked values as the predictor. Model diagnostics analyses focusing on finding unusual points were then carried out, which included checking leverage points, outliers, and influential observations. The OEAA psychometrician then made the anchor item inclusion/exclusion decisions and shared with Assessment and Evaluation Services (AES), which functions as an independent party validating the psychometric work done by the OEAA psychometrician on M-STEP social studies.
- AES psychometricians conducted their independent anchor item stability check (with their own methods) and compared their conclusion with what the OEAA psychometrician obtained.
- After the OEAA psychometrician and AES agreed on the anchor item inclusion/exclusion decisions, the OEAA psychometrician carried out a mean/mean method to transform the parameters from free run to the base scale for all anchor items. The constants A and B are computed as follows (adapted from formulas presented in Kolen & Brennan, 2004, p. 163):

$$A = \frac{\bar{a}_{free}}{\bar{a}_{base}},$$

$$B = \bar{b}_{base} - A * \bar{b}_{free} \quad (8.2)$$

where \bar{a}_{free} = mean of anchor set item discrimination parameters from the free run,

\bar{a}_{base} = mean of anchor set item discrimination parameters from the bank,

\bar{b}_{free} = mean of anchor set item difficulty parameters from the free run,

\bar{b}_{base} = mean of anchor set item difficulty parameters from the bank.

After obtaining the constants A and B as mentioned above, the following formulas are used to transform all anchor item parameters onto the base scale (adapted from formulas presented in Kolen & Brennan, 2004, p. 162):

$$a_{equated} = \frac{a_{free}}{A},$$

$$b_{equated} = Ab_{free} + B, (8.3)$$

where \bar{a}_{free} = item discrimination parameter from the free run for an anchor item,

$a_{equated}$ = transformed item discrimination parameter for that anchor item,

\bar{b}_{free} = item difficulty parameter from the free run for an anchor item,

$b_{equated}$ = transformed item difficulty parameter for that anchor item.

- A validation check is then carried out by AES to confirm the transformed item difficulty and item-discrimination parameters. After the anchor item values are verified per grade, a fixed item parameter calibration was used to put all OP items onto the base scale.
- Summed score to Expected *A Posteriori* (EAP) conversion tables were then created using flexMIRT. For social studies, one conversion table was used for all online forms as all OP items are the same across all forms at each grade level. Note that in each year, when conversion tables are made, paper/pencil data are not available for equating and calibration; thus, the online form conversion tables are applied to the paper/pencil forms.
- When the final data (both online and paper/pencil data) are available later in the year, a fixed item parameter calibration (where all online OP items are fixed at the values found during the conversion table creation stage) is carried out once again using the final data for all OP items. Then the obtained OP items' parameters are used in a fixed item parameter calibration to put all field-test items onto the same scale.

8.3.4 Anchor Item Evaluation for Social Studies

There are various methods for evaluating anchor item stability. As mentioned above, model diagnostic analyses were used by the OEAA psychometrician in checking the stability of anchor items at the conversion table creation stage. In this section, an ad hoc approach is reported, which evaluates the anchor quality using anchor item response patterns. This method uses all possible information about student performance that is shared between the 2017 and 2018 online⁵ administrations of M-STEP social studies tests. The annually used evaluation steps are as follows:

- Obtain the item response patterns in both the 2017 and the 2018 online administrations for the anchor items used in 2018. Note that only the same response patterns appearing in both years are used for this evaluation.
- Aggregate these item response patterns to obtain the number of unique item response patterns per grade as well as the mean scale score for each specific item response pattern in 2017 and 2018.
- Plot the mean scale score in 2018 against the mean scale score in 2017 by grade for each anchor item response pattern.
- Plot a 45-degree line on that scatterplot. The observations plotted should cluster relatively tightly and be randomly distributed around the 45-degree line.

⁵ This analysis limits its scale to online responses only because the equating procedures are carried out with online data only.

- Plot the “proficient” cut score on both the vertical and horizontal axes to divide the graph into four quadrants (i.e., item response patterns that are scored proficient in both years, those that are scored proficient in 2017 but not 2018 and vice versa, and those that are scored not proficient in both years).

The final steps in the analysis are to evaluate the degree to which the scatterplot for each grade deviates from expectations for good equating (i.e., deviation from tight clustering and random distribution around the 45-degree line) and to evaluate the distribution of item score patterns in the four quadrants by grade.

Table 8-9 presents the anchor points (same as the number of anchor items) per grade on each form. The results of the anchor quality evaluation are presented in Figures 8-3 to 8-5 and in Table 8-10. As shown in Figures 8-3 to 8-5, the points plotted on the scatterplot for each grade tend to lie along the 45-degree line, indicating that the majority of students who shared the same item response patterns on the anchor set also obtained similar mean scale scores (per item response pattern) across the two years. Therefore, these anchor items are considered to be stable across these two years.

Table 8-9. Number of Anchor Items by Content and Grade for Each Form

Content Area	Grade	Total Points	Anchor Points	Percentage of Anchor Points
Social Studies	5	45	12	26.67
Social Studies	8	44	11	25.00
Social Studies	11	38	10	26.32

Figure 8-3. Social Studies Grade 5

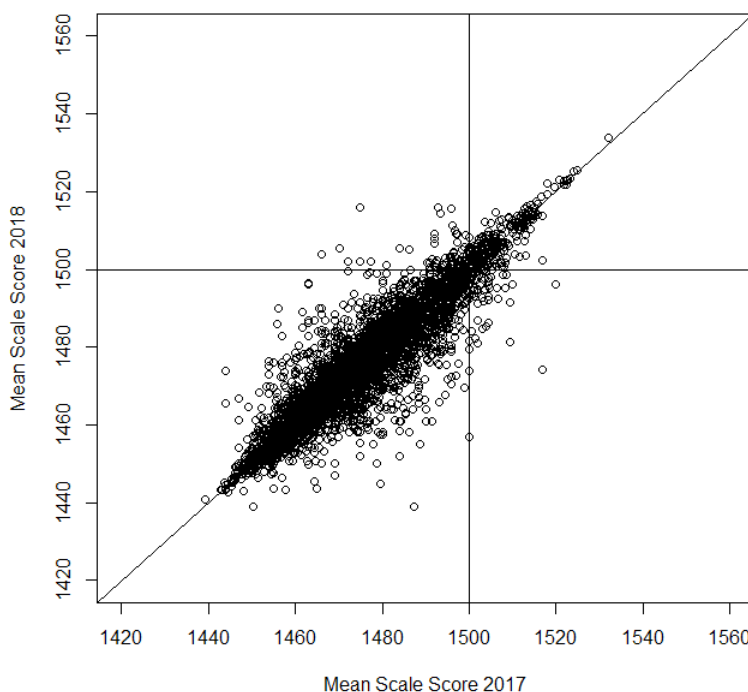
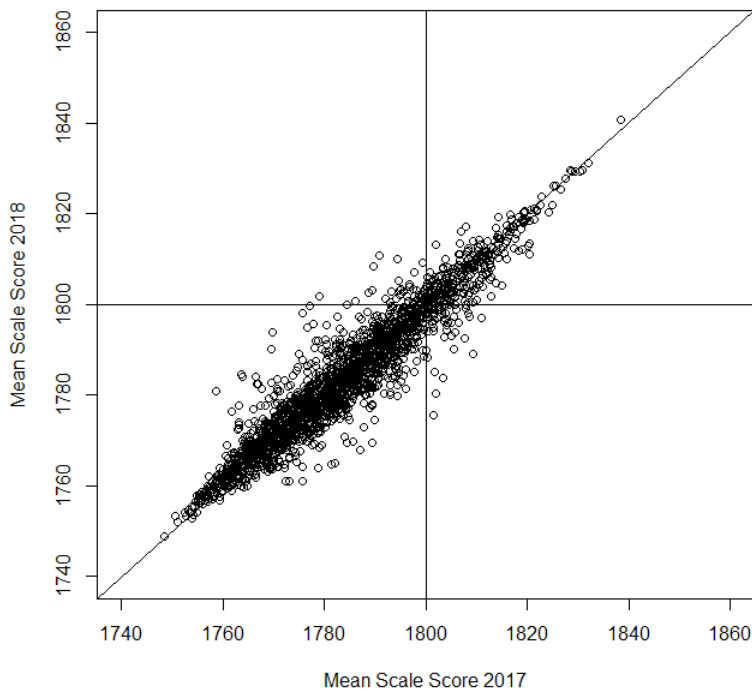
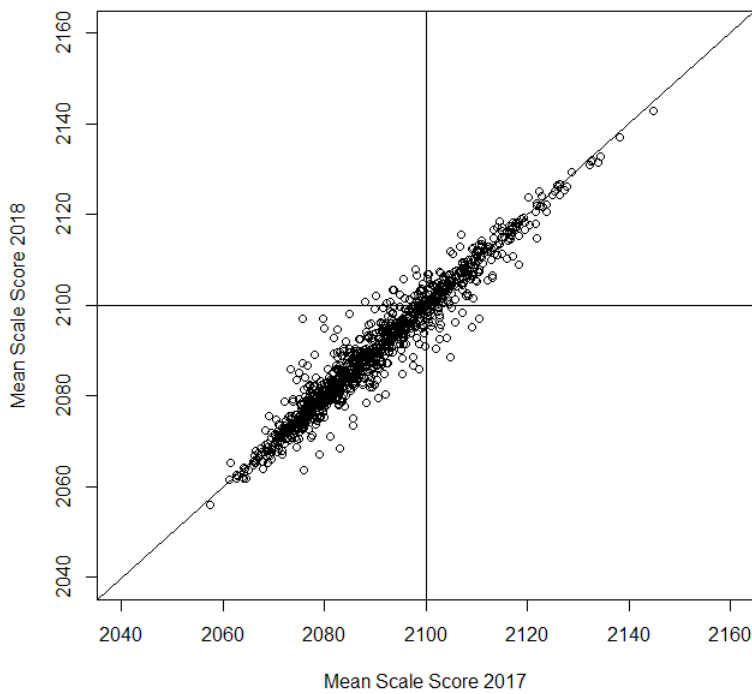


Figure 8-4. Social Studies Grade 8**Figure 8-5. Social Studies Grade 11**

The number and percentage of these anchor item response patterns that fall into each of the four quadrants by grade are summarized in Table 8-10. The percentage of response patterns that are associated with consistent performance categorization (based on the mean scale score for each item response pattern) across the two administrations ranged from 94.04% to 96.63%. According to this table, grade 5 had the highest consistency rate (96.63%), while grade 11 had the lowest consistency rate of about 94.04%.

Table 8-10. Evaluation of Equating Quality Using Linking Item Response Patterns

Content Area	Grade	Item Response Pattern	Proficient in Both Years	Not Proficient in Both Years	Proficient in 2017 Only	Proficient in 2018 Only	Consistent Classification	Inconsistent Classification
Social Studies	5	Count	230	3635	62	73	3865	135
Social Studies	5	Percentage	5.75	90.88	1.55	1.83	96.63	3.38
Social Studies	8	Count	288	1661	55	41	1949	96
Social Studies	8	Percentage	14.08	81.22	2.69	2.00	95.31	4.69
Social Studies	11	Count	235	727	30	31	962	61
Social Studies	11	Percentage	22.97	71.07	2.93	3.03	94.04	5.96

Note. Some rows have percentages that sum to more than 100 due to rounding.

8.3.5 Evidence of Model Fit for Social Studies

Due to sparse contingency tables, the limited-information fit statistics M_2 (Cai & Hansen, 2013) of the fitted model were requested for each fixed item parameter calibration run in flexMIRT. Due to the large sample size (>35,000 per online form), the model selection index tends to prefer more complex models (Cudeck & Henly, 1991). Taking model parsimony into considerations, the RMSEA values are considered rather than the M_2 statistics. The RMSEA values are 0.01 for social studies at grade 5, and 0.02 for social studies at grades 8 and 11. The fact that the RMSEA values are small in magnitude (i.e., close to 0) is evidence to support the use of the 2PL fixed item parameter calibration.

8.3.6 Test Characteristic Curves (TCCs) and Conversion Tables

The TCC is the graphical representation of the test characteristic function (TCF), and TCF is the expected raw total score given θ . Since all items are dichotomously scored, the expression of TCF is as follows (adapted from Yen & Fitzpatrick, 2006, p. 125):

$$E(X_i|\theta) = \sum_{i=1}^n E(X_i|\theta) = \sum_{i=1}^n P_i(\theta) \quad (8.4)$$

Figures 8-6 to 8-8 display the TCCs for the spring 2018 M-STEP social studies tests by grade. These graphs were made using the item parameter estimates obtained from the post-administration calibration in 2018 (based on unidimensional 2PL models). Two TCCs are shown for social studies at each grade level (one for online forms 1–3 and one for paper/pencil form 4) (see Figures 8-6 to 8-8). Note that these curves were created for OP items per form based on the item parameter estimates obtained from the last step mentioned in Section 8.3.3. Due to item differences between online forms and paper/pencil forms (i.e., TE items that appear on online forms cannot appear on paper/pencil forms), slight differences in TCCs can be seen. In general, for each grade, the TCCs across all forms are very close to each other.

Table 8-11 presents the summed scores to EAP conversion tables by grade for social studies. These are used for operational reporting. Note that when conversion tables were made, no paper/pencil data were available for calibration; thus, an operational decision was made to apply the conversion tables from the online form to the corresponding paper/pencil forms. Whether such decision is reasonable is examined in the mode comparison study, which can be found in Appendix F. Note that these tables present very similar results as those shown in the corresponding TCC graphs.

Figure 8-6. TCC for Social Studies Grade 5 by Form

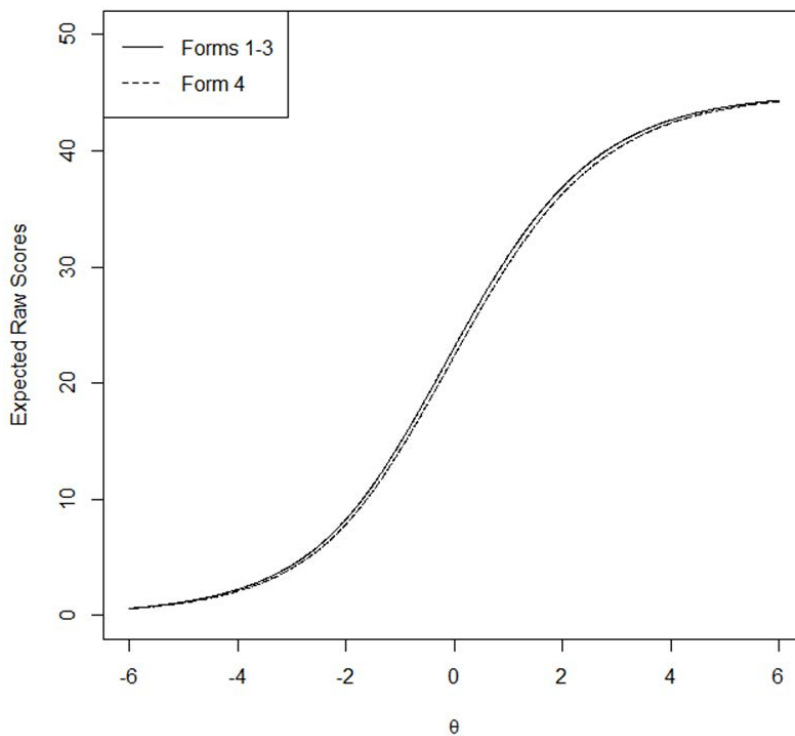


Figure 8-7. TCC for Social Studies Grade 8 by Form

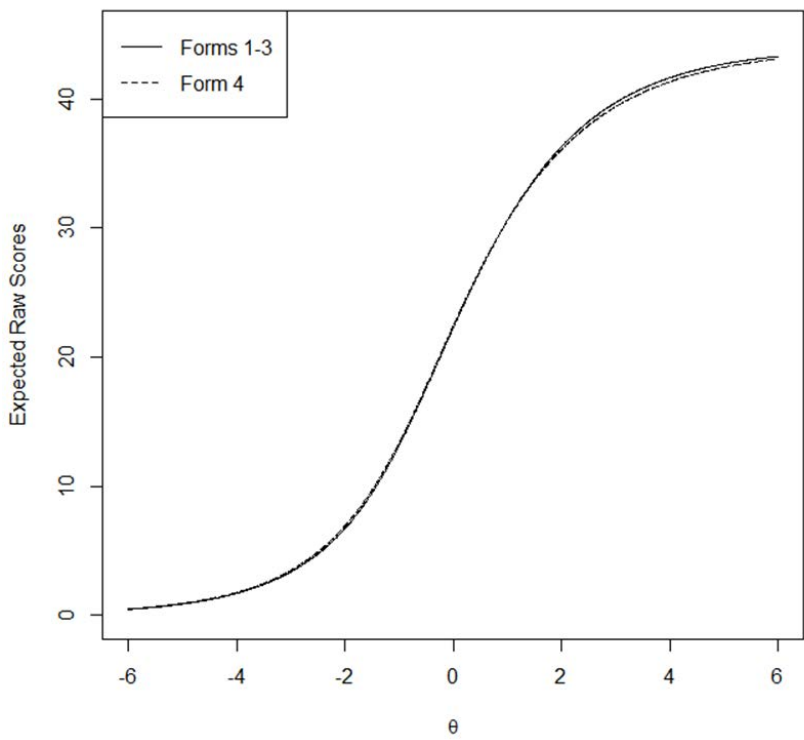


Figure 8-8. TCC for Social Studies Grade 11 by Form

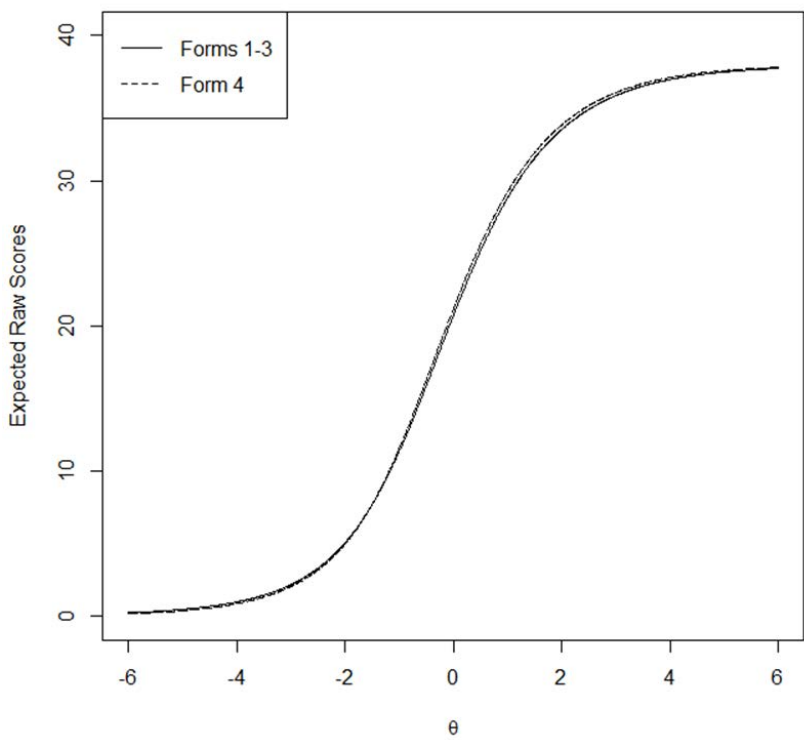


Table 8-11. Social Studies Summed Score to EAP Conversion Tables by Grade

Raw Score	Grade 5 Theta	Grade 5 SE	Grade 8 Theta	Grade 8 SE	Grade 11 Theta	Grade 11 SE
0	-3.0770	0.4840	-2.8990	0.5110	-2.7460	0.5090
1	-2.8990	0.4960	-2.6990	0.5110	-2.5180	0.4940
2	-2.7140	0.4950	-2.5010	0.4990	-2.2990	0.4690
3	-2.5300	0.4850	-2.3100	0.4800	-2.0980	0.4410
4	-2.3520	0.4700	-2.1320	0.4590	-1.9140	0.4150
5	-2.1830	0.4540	-1.9650	0.4390	-1.7470	0.3930
6	-2.0230	0.4390	-1.8090	0.4210	-1.5940	0.3740
7	-1.8720	0.4260	-1.6630	0.4040	-1.4520	0.3580
8	-1.7270	0.4140	-1.5250	0.3900	-1.3200	0.3440
9	-1.5890	0.4040	-1.3950	0.3780	-1.1960	0.3330
10	-1.4570	0.3950	-1.2700	0.3670	-1.0780	0.3240
11	-1.3300	0.3880	-1.1510	0.3580	-0.9650	0.3160
12	-1.2070	0.3820	-1.0370	0.3500	-0.8570	0.3100
13	-1.0880	0.3760	-0.9260	0.3430	-0.7530	0.3050
14	-0.9730	0.3720	-0.8190	0.3380	-0.6510	0.3020
15	-0.8590	0.3680	-0.7140	0.3340	-0.5510	0.2990
16	-0.7490	0.3650	-0.6120	0.3310	-0.4520	0.2970
17	-0.6400	0.3630	-0.5110	0.3290	-0.3540	0.2960
18	-0.5320	0.3610	-0.4110	0.3280	-0.2570	0.2960
19	-0.4260	0.3600	-0.3120	0.3280	-0.1600	0.2970
20	-0.3200	0.3590	-0.2140	0.3280	-0.0620	0.2990
21	-0.2150	0.3590	-0.1160	0.3290	0.0370	0.3020
22	-0.1100	0.3590	-0.0170	0.3310	0.1380	0.3050
23	-0.0050	0.3600	0.0820	0.3330	0.2400	0.3100
24	0.1010	0.3620	0.1820	0.3370	0.3460	0.3150
25	0.2070	0.3640	0.2830	0.3400	0.4540	0.3220
26	0.3140	0.3660	0.3860	0.3450	0.5670	0.3290
27	0.4220	0.3690	0.4910	0.3500	0.6840	0.3390
28	0.5320	0.3720	0.5980	0.3560	0.8060	0.3490
29	0.6430	0.3770	0.7080	0.3630	0.9360	0.3610
30	0.7570	0.3810	0.8210	0.3700	1.0730	0.3750
31	0.8730	0.3870	0.9380	0.3780	1.2180	0.3910
32	0.9930	0.3930	1.0590	0.3870	1.3750	0.4090
33	1.1150	0.4000	1.1840	0.3970	1.5430	0.4300

Raw Score	Grade 5 Theta	Grade 5 SE	Grade 8 Theta	Grade 8 SE	Grade 11 Theta	Grade 11 SE
34	1.2420	0.4080	1.3140	0.4080	1.7250	0.4530
35	1.3730	0.4170	1.4500	0.4200	1.9240	0.4790
36	1.5100	0.4270	1.5920	0.4340	2.1400	0.5060
37	1.6520	0.4380	1.7410	0.4480	2.3750	0.5310
38	1.8010	0.4500	1.8980	0.4640	2.6220	0.5470
39	1.9580	0.4630	2.0640	0.4800		
40	2.1230	0.4780	2.2390	0.4970		
41	2.2980	0.4930	2.4230	0.5130		
42	2.4810	0.5070	2.6140	0.5230		
43	2.6710	0.5160	2.8060	0.5250		
44	2.8620	0.5150	2.9920	0.5150		
45	3.0470	0.5020				

Note. The possible maximum total raw score is 45 for grade 5, 44 for grade 8, and 38 for grade 11.

8.3.7 IRT Statistics

As discussed above, the 2PL model was used to calibrate the spring 2018 social studies items at each grade level. A summary (i.e., minimum, maximum, and mean values) of the item difficulty (b -parameter) and item discrimination (a -parameter) estimates for all OP items per form for each grade is presented in Table 8-12.

Table 8-12. Item Difficulty (b -Parameter) and Item Discrimination (a -Parameter) for Social Studies by Grade and Form

Grade	Form	Difficulty Minimum	Difficulty Maximum	Difficulty Mean	Discrimination Minimum	Discrimination Maximum	Discrimination Mean
5	1–3	-1.9336	1.4443	-0.0280	0.4934	1.9271	0.8307
5	4	-1.9336	1.4443	0.0773	0.4934	1.9271	0.8217
8	1–3	-0.9856	1.2302	0.0767	0.4981	2.0747	0.9267
8	4	-0.9856	1.8781	0.0955	0.3954	1.8605	0.8912
11	1–3	-0.9879	1.6519	-0.1204	0.3894	1.6423	1.1117
11	4	-0.9879	1.0106	-0.1559	0.5708	2.3057	1.1383

8.4 Summary

In summary, the overall purpose of the OP data analysis is to ensure that the test items, as well as the overall test, are functioning appropriately. It also helps maintain the test scale across years so that test results may be appropriately compared across years. The data analyses undertaken by Smarter Balanced, DRC, and the Michigan Department of Education are in alignment with multiple best practices of the assessment industry, particularly the following *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014):

- Standard 5.2—The procedures for constructing scales used for reporting scores and the rationale for these procedures should be described clearly.
- Standard 5.13—When claims of form-to-form score equivalence are based on equating procedures, detailed technical information should be provided on the method by which equating functions were established and on the accuracy of the equating functions.

Chapter 9: Test Results

This chapter of the Technical Report contains information on the results of the spring 2018 administration of M-STEP along with descriptions of the score reports, data structure, and interpretive guide. The AERA, APA, and NCME (2014) *Standards* addressed in Chapter 9 include 5.1, 6.10, and 7.0. Each standard will be presented in the pertinent section of this chapter.

9.1 Test Completion

The spring 2018 M-STEP was administered to Michigan students in three content areas: ELA, mathematics, and social studies. For the purposes of this technical report, “percent valid” is the percentage of students who received a valid scale score given the total number of students who took the online or paper/pencil test. These test completion rates are summarized in Tables 9-1a through 9-3g.

Test completion information is reported for all students and the following demographic subgroups:

- Gender: Female and Male
- Race/Ethnicity: American Indian or Alaska Native, Asian, Black or African American, Hispanic or Latino, Native Hawaiian or Other Pacific Islander, Two or More Races, and White
- Economically Disadvantaged: Yes, No
- English Language Learners: Yes, No
- Students with Disabilities: Yes, No
- Students Used Standard Accommodations: Yes, No

9.2 Current Administration Data Scale Score Summaries

Summaries of the scale-score (SS) descriptive statistics for the spring 2018 administration of the ELA, mathematics, and social studies assessments are reported in Tables 9-4a through 9-6g, by grade, content area, and demographic subgroup.

Additionally, Tables 9-7a through 9-9b present the scale-score descriptive statistics and the performance level percentages by grade for the 2018 M-STEP ELA, mathematics, and social studies tests. These tables provide the scale-score descriptive statistics (i.e., Mean, SD, Min, Max values) and the percentages of students in each performance level: Not Proficient, Partially Proficient, Proficient, and Advanced.

9.3 Description of Reports

Score reports are the primary means of communicating test scores to relevant district and school administrators, teachers, and parents. AERA, APA, and NCME (2014) Standard 6.10 states the following:

When test score information is released, those responsible for testing programs should provide interpretations appropriate to the audience. The interpretations should describe in simple language what the test covers, what scores represent, the precision/reliability of the scores, and how scores are intended to be used. (p. 119)

Standard 5.1 is also addressed:

Test users should be provided with clear explanations of the characteristics, meaning, and intended interpretation of scale scores, as well as their limitations. (p. 102)

Interpretations related to the test scores are provided in the *Spring 2018 Interpretive Guide to M-STEP Reports* for grades 3 through 8 and in the *Spring 2018 Interpretive Guide to MME Reports* for grade 11. The interpretive guides are provided in Appendix B.1 and Appendix B.2 respectively.

In addition to providing interpretation, it is important that the information is understandable by the target audience. Standard 7.0 of the AERA, APA, & NCME (2014) *Standards* states the following:

Information relating to tests should be clearly documented so that those who use tests can make informed decisions regarding which test to use for a specific purpose, how to administer the chosen test, and how to interpret test scores. (p. 125)

In support of Standard 7.0, the *Spring 2018 Interpretive Guide to M-STEP Reports* (presented in Appendices B.1 and B.2) are accessible to parents, teachers, and laypeople alike.

M-STEP score reports comprise student-level reports and data files and aggregate reports and data files. Depending on the audience, reports and data files are available in several systems.

1. The OEAA Secure Site allows authorized school and district users to download student-level and state, district, and building level aggregate data files. It also provides access to the online Dynamic Score Reporting Site.
2. The Dynamic Score Reporting Site (DSRS) provides authorized school and district users access to interactive online student-level and aggregate reports.
3. The Michigan Linked Educational Assessment Reporting Network (MiLearn) provides teachers, parents, and students direct access to student level data through the district's Student Information System. The reports presented are a subset of those available in the DSRS.
4. [MiSchoolData](#) is Michigan's public portal to education-related data. It contains assessment data for public consumption.

This technical report will address the data files and reports available through the OEAA Secure Site and the Dynamic Score Reporting Site.

Brief descriptions of the reports and data files are provided below. More extensive descriptions with samples are included in the *Spring 2018 Interpretive Guide to M-STEP Reports* (M-STEP IGTR) for grades 3 through 8 and in the *Spring 2018 Interpretive Guide to MME Reports* (MME IGTR) for grade 11. The interpretive guides also include information on how to use the data, limitations of the data, and the online functionality associated with each report.

9.3.1 Student-Level Data Reports and Data Files

- The Student Record Labels provide a summary of student performance levels for individual students. The labels include district and school information, student demographic information, M-STEP administration cycle information, and overall student performance level for tested content areas.

Student Record Labels are provided for inclusion in the student's Cumulative Student Record (CA60) folder. They are printed and shipped to the school in which the student tested in late summer and are available through the Secure Site if the school needs to print additional copies.

Additional information can be found starting on page 18 of the M-STEP IGTR and on page 14 of the MME IGTR.

- The Individual Student Report (ISR) provides information about student performance by content area. Each student will have a separate ISR for each content area assessed. The report is divided into three main sections:
 - Student demographic information
 - Overall content performance and detailed claim data for ELA and mathematics
 - Discipline and content expectation data for social studies
 - Additional information can be found starting on page 19 of the M-STEP IGTR and on page 15 of the MME IGTR. Individual Student Reports are also available in MI-Learn to educators, parents, and students where they are referred to as the M-STEP Student Detail Report.
- Parent Reports are printed and shipped to schools for distribution to parents. The parent report provides information for parents about student performance in tested content areas. This report includes five main sections:
 - Superintendent letter
 - Overall performance level and scale score
 - Detailed claim data for ELA and mathematics and discipline data for social studies
 - Definitions for parents
 - Performance-level descriptors

Additional information can be found starting on page 21 of the M-STEP IGTR and on page 16 of the MME IGTR. Parent Reports are also available in the Dynamic Score Reporting Site for schools to access and print copies.

- Student Roster allows users to view student scale scores and claim performance data for ELA and mathematics or discipline data for social studies by content area and grade. The report is divided into five main sections:
 - Overall proficiency summary of the rostered students in graphic format
 - An alphabetical listing of the selected students
 - Overall content performance in a table format
 - Overall content performance in a graphical format
 - Claim data for ELA and mathematics or discipline data for social studies

Additional information can be found starting on page 23 of the M-STEP IGTR and on page 19 of the MME IGTR. The Student Roster report is also available to educators in MiLearn.

- The Student Overview provides summary information about student performance in all tested content areas in the selected grade. For each selected student, the following data are displayed for each tested content area in both graphical and table format:
 - scale score
 - margin of error
 - performance level
 - claim or discipline performance

Additional information can be found starting on page 27 of the M-STEP IGTR and on page 21 of the MME IGTR. The Student Overview is also available to educators in MiLearn.

- The Student Data File (SDF) contains student level demographics, test scores, and performance data for each tested content area. The Downloadable SDF contains detailed individual student data in an Excel file. This data includes school information, student demographic data, test administration data, and student performance data. The SDF is provided for schools to use as a data resource for school- or district-level data reviews. Schools or districts can use the SDF to manipulate and evaluate data in ways that support school improvement goals or for other data-based decision-making purposes.

Additional information can be found starting on page 37 of the M-STEP IGTR. The SDF can be downloaded from the OEAA Secure Site. The file layout is in Appendix B.3.

9.3.2 Aggregate Data Reports and Data Files

The Target Analysis Report is available at the school, district, Intermediate School District (ISD), and state levels for ELA and mathematics. The report is intended to provide an overview of relative strengths and weaknesses in ELA and mathematics by assessment target as compared to student performance on the test as a whole.

Additional information can be found starting on page 28 of the M-STEP IGTR.

- The Expectation Analysis Report provides the percentage of points earned by grade, the content area expectations in each social studies discipline, and the number of students scoring in each of the four quartiles. The report is intended to provide an overview of performance by content expectation. The report displays the number of students assessed in each expectation (not all students were assessed on every expectation), the average percentage of points earned, and the number of students scoring in one of four bands or quartiles: 0%–25%, 26%–50%, 51%–75%, and 76%–100% of all possible points.

Additional information can be found starting on page 30 of the M-STEP IGTR and page 22 of the MME IGTR. The Expectation Analysis Report is also available in MiLearn for educators.

- The Demographic Report provides a comparison of students by grade and content area, aggregated across selected demographic groups, showing the percentages proficient at each level (i.e., advanced, proficient, partially proficient, and not proficient). The demographic report is available at the school, district, ISD, and state levels. Users can select different populations of students to be displayed. The following student populations may be selected:
 - All Students—this is the default
 - All Except Students with Disabilities—students who are not marked Special Education in the Michigan Student Data System (MSDS) at the time of testing
 - Students with Disabilities—students who are marked Special Education in MSDS at the time of testing

Additional information can be found starting on page 31 of the M-STEP IGTR and on page 23 of the MME IGTR.

- The Comprehensive Report provides a comparison of students by grade and content area, aggregated across schools and districts, showing the percentages proficient at each level (i.e., advanced, proficient, partially proficient, and not proficient). The Comprehensive Report is available at the ISD and district levels. After the user selects a grade to view, all tested content areas for that grade are displayed sequentially in alphabetical order.

Additional information can be found starting on page 33 of the M-STEP IGTR and on page 25 of the MME IGTR.

- The Aggregate Data File contains student performance data aggregated for all students and by select demographic subgroups across buildings, districts, and the state. This data includes school information, student population, demographic group, and student performance data.

The Aggregate Data File is provided for schools and districts to use as a data resource for school- or district-level data reviews. Schools or districts can use the Aggregate Data Files to evaluate data in ways that support school improvement goals or other data-based decision-making purposes.

Additional information can be found starting on page 37 of the M-STEP IGTR. The Aggregated Data Files can be downloaded from the OEAA Secure Site. The file layout is in Appendix B.4.

9.4 Interpretive Guides

For the spring 2018 M-STEP, MDE produced individual and aggregate reports for students, schools, districts, and the state. The information provided in these reports can be interpreted and used in a variety of ways. In addition to providing interpretation, it is important that the information is understandable by the target audience. Standard 7.0 of the AERA, APA, and NCME (2014) *Standards* states the following:

Information relating to tests should be clearly documented so that those who use tests can make informed decisions regarding which test to use for a specific purpose, how to administer the chosen test, and how to interpret test scores. (p. 125)

To aid in interpretation, MDE prepared the *Spring 2018 Interpretive Guide to M-STEP Reports* for grades 3 through 8 and the *Spring 2018 Interpretive Guide to MME Reports* for grade 11.

The spring 2018 edition of the *Spring 2018 Interpretive Guide to M-STEP Reports* can be found in Appendix B.1 of this technical report. The MME interpretive guide can be found in Appendix B.2.

9.5 Summary

In summary, the overall purpose of reporting test results is to communicate information on student performance to stakeholders. These results are presented in the context of score reports that aid the user in understanding the meaning of the test scores. The reports and ancillary information developed by MDE are in alignment with multiple best practices of the testing industry, particularly the following *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014):

- Standard 5.1—Test users should be provided with clear explanations of the characteristics, meaning, and intended interpretation of scale scores, as well as their limitations.
- Standard 6.10—When test score information is released, those responsible for testing programs should provide interpretations appropriate to the audience. The interpretations should describe in simple language what the test covers, what scores represent, the precision/reliability of the scores, and how scores are intended to be used.
- Standard 7.0—Information relating to tests should be clearly documented so that those who use tests can make informed decisions regarding which test to use for a specific purpose, how to administer the chosen test, and how to interpret test scores.

Table 9-1a. M-STEP Test Completion Rates by Grade: English Language Arts—All Students

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
All Students	Total Tested	102,276	104,875	108,885	108,672	107,926	110,417
All Students	Number Valid	102,249	104,852	108,857	108,635	107,888	110,343
All Students	Percent Valid	99.97%	99.98%	99.97%	99.97%	99.96%	99.93%

Table 9-1b. M-STEP Test Completion Rates by Grade: English Language Arts—Gender

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Female	Total Tested	50,176	51,433	53,335	53,683	52,969	53,948
Female	Number Valid	50,168	51,428	53,323	53,670	52,951	53,931
Female	Percent Valid	99.98%	99.99%	99.98%	99.98%	99.97%	99.97%
Male	Total Tested	52,100	53,442	55,550	54,989	54,957	56,469
Male	Number Valid	52,081	53,424	55,534	54,965	54,937	56,412
Male	Percent Valid	99.96%	99.97%	99.97%	99.96%	99.96%	99.90%

**Table 9-1c. M-STEP Test Completion Rates by Grade: English Language Arts—
Race/Ethnicity**

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
American Indian or Alaska Native	Total Tested	586	626	708	667	751	736
American Indian or Alaska Native	Number Valid	585	624	707	667	751	736
American Indian or Alaska Native	Percent Valid	99.83%	99.68%	99.86%	100.00%	100.00%	100.00%
Asian	Total Tested	3,456	3,515	3,630	3,611	3,646	3,798
Asian	Number Valid	3,456	3,515	3,630	3,609	3,644	3,798
Asian	Percent Valid	100.00%	100.00%	100.00%	99.94%	99.95%	100.00%
Black or African American	Total Tested	19,097	18,820	19,219	19,112	18,445	18,313
Black or African American	Number Valid	19,092	18,808	19,205	19,099	18,430	18,285
Black or African American	Percent Valid	99.97%	99.94%	99.93%	99.93%	99.92%	99.85%
Hispanic or Latino	Total Tested	8,263	8,612	8,960	8,694	8,852	8,344
Hispanic or Latino	Number Valid	8,262	8,612	8,959	8,693	8,848	8,338
Hispanic or Latino	Percent Valid	99.99%	100.00%	99.99%	99.99%	99.95%	99.93%
Native Hawaiian or Other Pacific Islander	Total Tested	83	78	91	91	68	92
Native Hawaiian or Other Pacific Islander	Number Valid	83	78	91	91	68	92
Native Hawaiian or Other Pacific Islander	Percent Valid	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%
Two or More Races	Total Tested	4,666	4,723	4,529	4,428	3,971	3,918
Two or More Races	Number Valid	4,664	4,723	4,527	4,428	3,970	3,914
Two or More Races	Percent Valid	99.96%	100.00%	99.96%	100.00%	99.97%	99.90%
White	Total Tested	66,125	68,501	71,748	72,069	72,193	75,216
White	Number Valid	66,107	68,492	71,738	72,048	72,177	75,180
White	Percent Valid	99.97%	99.99%	99.99%	99.97%	99.98%	99.95%

Table 9-1d. M-STEP Test Completion Rates by Grade: English Language Arts—Economically Disadvantaged

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	Total Tested	57,726	58,062	59,221	57,812	55,474	54,684
Yes	Number Valid	57,709	58,044	59,198	57,782	55,442	54,618
Yes	Percent Valid	99.97%	99.97%	99.96%	99.95%	99.94%	99.88%
No	Total Tested	44,550	46,813	49,664	50,860	52,452	55,733
No	Number Valid	44,540	46,808	49,659	50,853	52,446	55,725
No	Percent Valid	99.98%	99.99%	99.99%	99.99%	99.99%	99.99%

Table 9-1e M-STEP Test Completion Rates by Grade: English Language Arts—English Language Learners

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	Total Tested	10,120	9,838	7,586	6,789	6,366	6,268
Yes	Number Valid	10,118	9,838	7,585	6,786	6,362	6,261
Yes	Percent Valid	99.98%	100.00%	99.99%	99.96%	99.94%	99.89%
No	Total Tested	92,156	95,037	101,299	101,883	101,560	104,149
No	Number Valid	92,131	95,014	101,272	101,849	101,526	104,082
No	Percent Valid	99.97%	99.98%	99.97%	99.97%	99.97%	99.94%

Table 9-1f. M-STEP Test Completion Rates by Grade: English Language Arts—Students with Disabilities

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	Total Tested	11,491	12,271	12,445	11,928	11,822	11,798
Yes	Number Valid	11,484	12,265	12,435	11,912	11,816	11,772
Yes	Percent Valid	99.94%	99.95%	99.92%	99.87%	99.95%	99.78%
No	Total Tested	90,785	92,604	96,440	96,744	96,104	98,619
No	Number Valid	90,765	92,587	96,422	96,723	96,072	98,571
No	Percent Valid	99.98%	99.98%	99.98%	99.98%	99.97%	99.95%

Table 9-1g. M-STEP Test Completion Rates by Grade: English Language Arts—Students Used Standard Accommodations

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	Total Tested	335	337	358	6,047	5,952	6,092
Yes	Number Valid	333	337	357	6,043	5,950	6,082
Yes	Percent Valid	99.40%	100.00%	99.72%	99.93%	99.97%	99.84%
No	Total Tested	101,941	104,538	108,527	102,625	101,974	104,325
No	Number Valid	101,916	104,515	108,500	102,592	101,938	104,261
No	Percent Valid	99.98%	99.98%	99.98%	99.97%	99.96%	99.94%

Table 9-2a. M-STEP Test Completion Rates by Grade: Mathematics—All Students

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
All Students	Total Tested	102,610	105,135	109,083	108,821	108,040	110,443
All Students	Number Valid	102,587	105,113	109,057	108,789	108,000	110,383
All Students	Percent Valid	99.98%	99.98%	99.98%	99.97%	99.96%	99.95%

Table 9-2b. M-STEP Test Completion Rates by Grade: Mathematics—Gender

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Female	Total Tested	50,324	51,548	53,405	53,730	53,016	53,952
Female	Number Valid	50,317	51,539	53,397	53,721	53,000	53,941
Female	Percent Valid	99.99%	99.98%	99.99%	99.98%	99.97%	99.98%
Male	Total Tested	52,286	53,587	55,678	55,091	55,024	56,491
Male	Number Valid	52,270	53,574	55,660	55,068	55,000	56,442
Male	Percent Valid	99.97%	99.98%	99.97%	99.96%	99.96%	99.91%

Table 9-2c. M-STEP Test Completion Rates by Grade: Mathematics—Race/Ethnicity

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
American Indian or Alaska Native	Total Tested	585	625	709	668	747	736
American Indian or Alaska Native	Number Valid	585	624	708	668	747	736
American Indian or Alaska Native	Percent Valid	100.00%	99.84%	99.86%	100.00%	100.00%	100.00%
Asian	Total Tested	3,582	3,608	3,707	3,682	3,698	3,832
Asian	Number Valid	3,582	3,608	3,706	3,680	3,696	3,832
Asian	Percent Valid	100.00%	100.00%	99.97%	99.95%	99.95%	100.00%
Black or African American	Total Tested	19,094	18,834	19,210	19,110	18,416	18,285
Black or African American	Number Valid	19,089	18,825	19,201	19,100	18,400	18,255
Black or African American	Percent Valid	99.97%	99.95%	99.95%	99.95%	99.91%	99.84%
Hispanic or Latino	Total Tested	8,325	8,679	9,006	8,743	8,913	8,372
Hispanic or Latino	Number Valid	8,322	8,676	9,005	8,740	8,907	8,368
Hispanic or Latino	Percent Valid	99.96%	99.97%	99.99%	99.97%	99.93%	99.95%
Native Hawaiian or Other Pacific Islander	Total Tested	83	79	92	91	70	93
Native Hawaiian or Other Pacific Islander	Number Valid	83	79	92	91	70	93
Native Hawaiian or Other Pacific Islander	Percent Valid	100.00%	100.00%	100.00%	100.00%	100.00%	100.00%
Two or More Races	Total Tested	4,672	4,722	4,531	4,426	3,966	3,910
Two or More Races	Number Valid	4,670	4,722	4,528	4,426	3,961	3,908
Two or More Races	Percent Valid	99.96%	100.00%	99.93%	100.00%	99.87%	99.95%
White	Total Tested	66,269	68,588	71,828	72,101	72,230	75,215
White	Number Valid	66,256	68,579	71,817	72,084	72,219	75,191
White	Percent Valid	99.98%	99.99%	99.98%	99.98%	99.98%	99.97%

Table 9-2d. M-STEP Test Completion Rates by Grade: Mathematics—Economically Disadvantaged

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	Total Tested	57,889	58,203	59,330	57,888	55,534	54,669
Yes	Number Valid	57,872	58,188	59,310	57,865	55,499	54,618
Yes	Percent Valid	99.97%	99.97%	99.97%	99.96%	99.94%	99.91%
No	Total Tested	44,721	46,932	49,753	50,933	52,506	55,774
No	Number Valid	44,715	46,925	49,747	50,924	52,501	55,765
No	Percent Valid	99.99%	99.99%	99.99%	99.98%	99.99%	99.98%

Table 9-2e. M-STEP Test Completion Rates by Grade: Mathematics—English Language Learners

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	Total Tested	10,400	10,074	7,796	6,971	6,545	6,406
Yes	Number Valid	10,397	10,071	7,795	6,965	6,542	6,401
Yes	Percent Valid	99.97%	99.97%	99.99%	99.91%	99.95%	99.92%
No	Total Tested	92,210	95,061	101,287	101,850	101,495	104,037
No	Number Valid	92,190	95,042	101,262	101,824	101,458	103,982
No	Percent Valid	99.98%	99.98%	99.98%	99.97%	99.96%	99.95%

Table 9- 2f. M-STEP Test Completion Rates by Grade: Mathematics—Students with Disabilities

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	Total Tested	11,574	12,329	12,466	11,914	11,809	11,746
Yes	Number Valid	11,571	12,325	12,457	11,902	11,801	11,734
Yes	Percent Valid	99.97%	99.97%	99.93%	99.90%	99.93%	99.90%
No	Total Tested	91,036	92,806	96,617	96,907	96,231	98,697
No	Number Valid	91,016	92,788	96,600	96,887	96,199	98,649
No	Percent Valid	99.98%	99.98%	99.98%	99.98%	99.97%	99.95%

Table 9-2g. M-STEP Test Completion Rates by Grade: Mathematics—Students Used Standard Accommodations

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	Total Tested	109	1,859	2,630	2,454	2,393	2,454
Yes	Number Valid	109	1,859	2,629	2,450	2,390	2,453
Yes	Percent Valid	100.00%	100.00%	99.96%	99.84%	99.87%	99.96%
No	Total Tested	102,501	103,276	106,453	106,367	105,647	107,989
No	Number Valid	102,478	103,254	106,428	106,339	105,610	107,930
No	Percent Valid	99.98%	99.98%	99.98%	99.97%	99.96%	99.95%

Table 9-3a. M-STEP Test Completion Rates by Grade: Social Studies—All Students

		Grade 5	Grade 8	Grade 11
All Students	Total Tested	109,086	110,475	105,098
All Students	Number Valid	108,874	110,191	104,888
All Students	Percent Valid	99.81%	99.74%	99.80%

Table 9-3b. M-STEP Test Completion Rates by Grade: Social Studies—Gender

		Grade 5	Grade 8	Grade 11
Female	Total Tested	53,422	53,978	52,348
Female	Number Valid	53,338	53,852	52,254
Female	Percent Valid	99.84%	99.77%	99.82%
Male	Total Tested	55,664	56,497	52,750
Male	Number Valid	55,536	56,339	52,634
Male	Percent Valid	99.77%	99.72%	99.78%

Table 9-3c. M-STEP Test Completion Rates by Grade: Social Studies—Race/Ethnicity

		Grade 5	Grade 8	Grade 11
American Indian or Alaska Native	Total Tested	708	735	636
American Indian or Alaska Native	Number Valid	706	732	633
American Indian or Alaska Native	Percent Valid	99.72%	99.59%	99.53%
Asian	Total Tested	3,706	3,830	3,821
Asian	Number Valid	3,687	3,825	3,818
Asian	Percent Valid	99.49%	99.87%	99.92%
Black or African American	Total Tested	19,204	18,271	15,961
Black or African American	Number Valid	19,118	18,144	15,840
Black or African American	Percent Valid	99.55%	99.30%	99.24%
Hispanic or Latino	Total Tested	8,997	8,369	7,147
Hispanic or Latino	Number Valid	8,985	8,349	7,135
Hispanic or Latino	Percent Valid	99.87%	99.76%	99.83%
Native Hawaiian or Other Pacific Islander	Total Tested	92	94	89
Native Hawaiian or Other Pacific Islander	Number Valid	92	93	89
Native Hawaiian or Other Pacific Islander	Percent Valid	100.00%	98.94%	100.00%
Two or More Races	Total Tested	4,531	3,916	2,963
Two or More Races	Number Valid	4,522	3,903	2,954
Two or More Races	Percent Valid	99.80%	99.67%	99.70%
White	Total Tested	71,848	75,260	74,481
White	Number Valid	71,764	75,145	74,419
White	Percent Valid	99.88%	99.85%	99.92%

Table 9-3d. M-STEP Test Completion Rates by Grade: Social Studies—Economically Disadvantaged

		Grade 5	Grade 8	Grade 11
Yes	Total Tested	59,320	54,690	44,807
Yes	Number Valid	59,170	54,475	44,639
Yes	Percent Valid	99.75%	99.61%	99.63%
No	Total Tested	49,766	55,785	60,291
No	Number Valid	49,704	55,716	60,249
No	Percent Valid	99.88%	99.88%	99.93%

Table 9-3e. M-STEP Test Completion Rates by Grade: Social Studies—English Language Learners

		Grade 5	Grade 8	Grade 11
Yes	Total Tested	7,787	6,404	4,427
Yes	Number Valid	7,772	6,381	4,415
Yes	Percent Valid	99.81%	99.64%	99.73%
No	Total Tested	101,299	104,071	100,671
No	Number Valid	101,102	103,810	100,473
No	Percent Valid	99.81%	99.75%	99.80%

Table 9-3f. M-STEP Test Completion Rates by Grade: Social Studies—Students with Disabilities

		Grade 5	Grade 8	Grade 11
Yes	Total Tested	12,505	11,817	9,498
Yes	Number Valid	12,442	11,743	9,475
Yes	Percent Valid	99.50%	99.37%	99.76%
No	Total Tested	96,581	98,658	95,600
No	Number Valid	96,432	98,448	95,413
No	Percent Valid	99.85%	99.79%	99.80%

Table 9-3g. M-STEP Test Completion Rates by Grade: Social Studies—Students Used Standard Accommodations

		Grade 5	Grade 8	Grade 11
Yes	Total Tested	57	62	34
Yes	Number Valid	55	61	34
Yes	Percent Valid	96.49%	98.39%	100.00%
No	Total Tested	109,029	110,413	105,064
No	Number Valid	108,819	110,130	104,854
No	Percent Valid	99.81%	99.74%	99.80%

Table 9-4a. Scale-Score Descriptive Statistics by Grade: English Language Arts—All Students

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
All Students	N	102,249	104,852	108,857	108,635	107,888	110,343
All Students	Mean SS	1,294.7	1,395.4	1,496.1	1,592.8	1,694.3	1,793.3
All Students	SD SS	26.0	26.2	27.2	26.2	26.5	27.3

Table 9-4b. Scale-Score Descriptive Statistics by Grade: English Language Arts—Gender

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Female	N	50,168	51,428	53,323	53,670	52,951	53,931
Female	Mean SS	1,296.8	1,397.5	1,499.1	1,596.1	1,698.2	1,797.7
Female	SD SS	26.0	25.9	26.8	25.8	25.5	26.6
Male	N	52,081	53,424	55,534	54,965	54,937	56,412
Male	Mean SS	1,292.7	1,393.3	1,493.3	1,589.7	1,690.5	1,789.1
Male	SD SS	25.9	26.3	27.2	26.2	26.9	27.3

Table 9-4c. Scale-Score Descriptive Statistics by Grade: English Language Arts—Race/Ethnicity

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
American Indian or Alaska Native	N	585	624	707	667	751	736
American Indian or Alaska Native	Mean SS	1,290.4	1,389.7	1,489.8	1,587.7	1,689.0	1,786.5
American Indian or Alaska Native	SD SS	23.7	24.6	24.4	24.2	25.4	25.3
Asian	N	3,456	3,515	3,630	3,609	3,644	3,798
Asian	Mean SS	1,307.1	1,409.2	1,512.6	1,608.5	1,712.1	1,810.6
Asian	SD SS	25.6	26.5	26.9	26.3	26.1	26.7
Black or African American	N	19,092	18,808	19,205	19,099	18,430	18,285
Black or African American	Mean SS	1,278.4	1,378.9	1,479.1	1,577.0	1,678.7	1,777.9
Black or African American	SD SS	22.9	23.6	23.7	22.9	23.6	24.3
Hispanic or Latino	N	8,262	8,612	8,959	8,693	8,848	8,338
Hispanic or Latino	Mean SS	1,288.7	1,389.1	1,489.9	1,586.8	1,687.9	1,786.8
Hispanic or Latino	SD SS	23.8	23.3	25.0	23.9	24.6	25.2
Native Hawaiian or Other Pacific Islander	N	83	78	91	91	68	92
Native Hawaiian or Other Pacific Islander	Mean SS	1,293.3	1,395.0	1,501.6	1,595.8	1,692.7	1,797.2
Native Hawaiian or Other Pacific Islander	SD SS	22.9	28.4	27.5	26.7	25.9	27.1
Two or More Races	N	4,664	4,723	4,527	4,428	3,970	3,914
Two or More Races	Mean SS	1,293.6	1,394.0	1,494.4	1,591.6	1,692.9	1,791.6
Two or More Races	SD SS	25.6	25.7	26.7	25.9	26.5	27.7
White	N	66,107	68,492	71,738	72,048	72,177	75,180
White	Mean SS	1,299.7	1,400.1	1,500.8	1,597.1	1,698.3	1,797.0
White	SD SS	25.0	25.1	26.2	25.4	25.5	26.5

Table 9-4d. Scale-Score Descriptive Statistics by Grade: English Language Arts—Economically Disadvantaged

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	N	57,709	58,044	59,198	57,782	55,442	54,618
Yes	Mean SS	1,286.2	1,386.4	1,486.8	1,583.7	1,684.8	1,783.4
Yes	SD SS	24.2	24.3	25.2	24.2	24.7	25.2
No	N	44,540	46,808	49,659	50,853	52,446	55,725
No	Mean SS	1,305.8	1,406.5	1,507.3	1,603.2	1,704.3	1,803.0
No	SD SS	24.1	24.1	25.1	24.5	24.5	25.7

Table 9-4e. Scale-Score Descriptive Statistics by Grade: English Language Arts—English Language Learners

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	N	10,118	9,838	7,585	6,786	6,362	6,261
Yes	Mean SS	1,288.4	1,387.0	1,480.7	1,575.5	1,675.8	1,774.5
Yes	SD SS	24.0	23.6	22.3	20.8	21.5	21.7
No	N	92,131	95,014	101,272	101,849	101,526	104,082
No	Mean SS	1,295.4	1,396.2	1,497.3	1,594.0	1,695.4	1,794.4
No	SD SS	26.2	26.3	27.1	26.1	26.3	27.2

Table 9-4f. Scale-Score Descriptive Statistics by Grade: English Language Arts—Students with Disabilities

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	N	11,484	12,265	12,435	11,912	11,816	11,772
Yes	Mean SS	1,279.3	1,377.7	1,474.9	1,571.2	1,671.3	1,768.9
Yes	SD SS	23.0	23.1	22.6	21.2	22.0	21.4
No	N	90,765	92,587	96,422	96,723	96,072	98,571
No	Mean SS	1,296.7	1,397.7	1,498.9	1,595.5	1,697.1	1,796.2
No	SD SS	25.7	25.7	26.5	25.5	25.6	26.4

Table 9-4g. Scale-Score Descriptive Statistics by Grade: English Language Arts—Students Used Standard Accommodations

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	N	333	337	357	6,043	5,950	6,082
Yes	Mean SS	1,279.8	1,376.8	1,474.4	1,569.5	1,669.3	1,767.0
Yes	SD SS	22.3	25.0	23.6	19.0	19.4	19.3
No	N	101,916	104,515	108,500	102,592	101,938	104,261
No	Mean SS	1,294.8	1,395.4	1,496.2	1,594.2	1,695.7	1,794.8
No	SD SS	26.0	26.2	27.1	26.0	26.1	26.9

Table 9-5a. Scale-Score Descriptive Statistics by Grade: Mathematics—All Students

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
All Students	N	102,587	105,113	109,057	108,789	108,000	110,383
All Students	Mean SS	1,296.1	1,393.5	1,487.1	1,587.8	1,688.1	1,787.5
All Students	SD SS	27.0	25.7	26.4	25.8	26.6	27.0

Table 9-5b. Scale-Score Descriptive Statistics by Grade: Mathematics—Gender

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Female	N	50,317	51,539	53,397	53,721	53,000	53,941
Female	Mean SS	1,295.0	1,392.3	1,486.2	1,587.8	1,688.4	1,789.1
Female	SD SS	26.2	24.6	25.3	24.8	25.5	26.0
Male	N	52,270	53,574	55,660	55,068	55,000	56,442
Female	Mean SS	1,297.1	1,394.8	1,488.0	1,587.7	1,687.9	1,786.0
Female	SD SS	27.7	26.6	27.4	26.8	27.6	27.9

Table 9-5c. Scale-Score Descriptive Statistics by Grade: Mathematics—Race/Ethnicity

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
American Indian or Alaska Native	N	585	624	708	668	747	736
American Indian or Alaska Native	Mean SS	1,290.7	1,388.8	1,482.1	1,583.6	1,683.2	1,781.2
American Indian or Alaska Native	SD SS	24.5	23.1	22.2	23.8	23.9	23.7
Asian	N	3,582	3,608	3,706	3,680	3,696	3,832
Asian	Mean SS	1,314.7	1,413.0	1,508.5	1,608.3	1,710.8	1,811.2
Asian	SD SS	26.1	25.8	25.7	26.6	27.6	28.1
Black or African American	N	19,089	18,825	19,201	19,100	18,400	18,255
Black or African American	Mean SS	1,278.3	1,375.5	1,467.8	1,569.0	1,669.4	1,769.7
Black or African American	SD SS	24.4	23.2	22.6	23.0	22.8	22.8
Hispanic or Latino	N	8,322	8,676	9,005	8,740	8,907	8,368
Hispanic or Latino	Mean SS	1,289.1	1,386.6	1,480.0	1,580.7	1,680.3	1,779.6
Hispanic or Latino	SD SS	24.4	23.4	23.4	23.7	24.1	24.1
Native Hawaiian or Other Pacific Islander	N	83	79	92	91	70	93
	Mean SS	1,292.3	1,388.9	1,488.1	1,588.0	1,685.3	1,792.4
	SD SS	25.3	25.9	25.1	27.0	25.2	27.6
Two or More Races	N	4,670	4,722	4,528	4,426	3,961	3,908
Two or More Races	Mean SS	1,293.7	1,391.3	1,484.3	1,585.3	1,685.3	1,784.5
Two or More Races	SD SS	26.5	25.6	26.3	25.5	26.8	27.4
White	N	66,256	68,579	71,817	72,084	72,219	75,191
White	Mean SS	1,301.3	1,398.5	1,492.3	1,592.7	1,692.9	1,791.7
White	SD SS	25.4	23.9	24.8	23.9	24.9	25.8

Table 9-5d. Scale-Score Descriptive Statistics by Grade: Mathematics—Economically Disadvantaged

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	N	57,872	58,188	59,310	57,865	55,499	54,618
Yes	Mean SS	1,286.9	1,384.3	1,477.4	1,577.9	1,677.8	1,776.8
Yes	SD SS	25.3	23.9	24.3	24.2	24.3	24.1
No	N	44,715	46,925	49,747	50,924	52,501	55,765
No	Mean SS	1,307.9	1,405.0	1,498.7	1,599.0	1,699.0	1,798.0
No	SD SS	24.4	23.1	24.0	22.9	24.5	25.6

Table 9-5e. Scale-Score Descriptive Statistics by Grade: Mathematics—English Language Learners

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	N	10,397	10,071	7,795	6,965	6,542	6,401
Yes	Mean SS	1,292.8	1,388.1	1,475.9	1,573.6	1,672.5	1,772.4
Yes	SD SS	25.7	24.8	23.2	23.5	23.8	23.1
No	N	92,190	95,042	101,262	101,824	101,458	103,982
No	Mean SS	1,296.4	1,394.1	1,488.0	1,588.7	1,689.2	1,788.4
No	SD SS	27.1	25.7	26.4	25.7	26.4	27.0

Table 9-5f. Scale-Score Descriptive Statistics by Grade: Mathematics—Students with Disabilities

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	N	11,571	12,325	12,457	11,902	11,801	11,734
Yes	Mean SS	1,277.1	1,374.1	1,465.2	1,563.3	1,662.8	1,762.4
Yes	SD SS	28.0	25.7	24.5	24.2	23.1	21.1
No	N	91,016	92,788	96,600	96,887	96,199	98,649
No	Mean SS	1,298.5	1,396.1	1,489.9	1,590.8	1,691.3	1,790.5
No	SD SS	25.9	24.6	25.3	24.4	25.3	26.1

Table 9-5g. Scale-Score Descriptive Statistics by Grade: Mathematics—Students Used Standard Accommodations

		Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
Yes	N	109	1,859	2,629	2,450	2,390	2,453
Yes	Mean SS	1,275.6	1,365.7	1,458.5	1,555.9	1,657.4	1,757.9
Yes	SD SS	25.5	20.5	19.8	19.5	19.8	17.7
No	N	102,478	103,254	106,428	106,339	105,610	107,930
No	Mean SS	1,296.1	1,394.0	1,487.8	1,588.5	1,688.8	1,788.2
No	SD SS	27.0	25.5	26.1	25.5	26.3	26.8

Table 9-6a. Scale-Score Descriptive Statistics by Grade: Social Studies—All Students

		Grade 5	Grade 8	Grade 11
All Students	N	108,874	110,191	104,888
All Students	Mean SS	1,477.1	1,786.6	2,098.9
All Students	SD SS	24.7	25.3	24.7

Table 9-6b. Scale-Score Descriptive Statistics by Grade: Social Studies—Gender

		Grade 5	Grade 8	Grade 11
Female	N	53,338	53,852	52,254
Female	Mean SS	1,475.9	1,785.4	2,096.9
Female	SD SS	23.7	23.9	22.6
Male	N	55,536	56,339	52,634
Male	Mean SS	1,478.3	1,787.7	2,100.9
Male	SD SS	25.5	26.5	26.5

Table 9-6c. Scale-Score Descriptive Statistics by Grade: Social Studies—Race/Ethnicity

		Grade 5	Grade 8	Grade 11
American Indian or Alaska Native	N	706	732	633
American Indian or Alaska Native	Mean SS	1,474.3	1,781.8	2,095.5
American Indian or Alaska Native	SD SS	22.8	23.3	22.9
Asian	N	3,687	3,825	3,818
Asian	Mean SS	1,489.7	1,800.3	2,109.2
Asian	SD SS	26.7	25.5	26.3
Black or African American	N	19,118	18,144	15,840
Black or African American	Mean SS	1,461.7	1,770.0	2,083.3
Black or African American	SD SS	19.5	19.3	19.8
Hispanic or Latino	N	8,985	8,349	7,135
Hispanic or Latino	Mean SS	1,470.1	1,779.5	2,092.1
Hispanic or Latino	SD SS	21.4	22.4	22.7
Native Hawaiian or Other Pacific Islander	N	92	93	89
Native Hawaiian or Other Pacific Islander	Mean SS	1,479.4	1,788.5	2,098.3
Native Hawaiian or Other Pacific Islander	SD SS	23.5	24.2	23.9
Two or More Races	N	4,522	3,903	2,954
Two or More Races	Mean SS	1,474.6	1,784.9	2,098.1
Two or More Races	SD SS	23.9	24.9	24.0
White	N	71,764	75,145	74,419
White	Mean SS	1,481.7	1,790.8	2,102.4
White	SD SS	24.2	24.9	24.3

Table 9-6d. Scale-Score Descriptive Statistics by Grade: Social Studies—Economically Disadvantaged

		Grade 5	Grade 8	Grade 11
Yes	N	59,170	54,475	44,639
Yes	Mean SS	1,468.6	1,776.9	2,089.9
Yes	SD SS	21.6	22.0	22.3
No	N	49,704	55,716	60,249
No	Mean SS	1,487.3	1,796.0	2,105.6
No	SD SS	24.3	24.7	24.2

Table 9-6e. Scale-Score Descriptive Statistics by Grade: Social Studies—English Language Learners

		Grade 5	Grade 8	Grade 11
Yes	N	7,772	6,381	4,415
Yes	Mean SS	1,463.3	1,769.8	2,077.1
Yes	SD SS	18.7	17.9	16.9
No	N	101,102	103,810	100,473
No	Mean SS	1,478.2	1,787.6	2,099.9
No	SD SS	24.8	25.3	24.5

Table 9-6f. Scale-Score Descriptive Statistics by Grade: Social Studies—Students with Disabilities

		Grade 5	Grade 8	Grade 11
Yes	N	12,442	11,743	9,475
Yes	Mean SS	1,461.1	1,767.4	2,079.8
Yes	SD SS	20.4	19.4	20.0
No	N	96,432	98,448	95,413
No	Mean SS	1,479.2	1,788.8	2,100.8
No	SD SS	24.4	24.9	24.3

Table 9-6g. Scale-Score Descriptive Statistics by Grade: Social Studies—Students Used Standard Accommodations

		Grade 5	Grade 8	Grade 11
Yes	N	55	61	34
Yes	Mean SS	1,462.7	1,776.5	2,078.6
Yes	SD SS	20.2	21.7	21.2
No	N	108,819	110,130	104,854
No	Mean SS	1,477.1	1,786.6	2,098.9
No	SD SS	24.7	25.3	24.7

Table 9-7a. Scale-Score Descriptive Statistics: English Language Arts

Grade	<i>N</i>	Mean	SD	Min	Max
3	102,249	1,294.7	26.0	1,203	1,357
4	104,852	1,395.4	26.2	1,301	1,454
5	108,857	1,496.1	27.2	1,409	1,560
6	108,635	1,592.8	26.2	1,508	1,655
7	107,888	1,694.3	26.5	1,618	1,753
8	110,343	1,793.3	27.3	1,721	1,857

Table 9-7b. Performance-Level Percentages: English Language Arts

Grade	<i>N</i>	Not Proficient	Partially Proficient	Proficient	Advanced
3	102,249	31.01%	24.61%	22.37%	22.01%
4	104,852	33.89%	21.07%	21.53%	23.51%
5	108,857	32.21%	21.26%	28.69%	17.85%
6	108,635	31.42%	27.21%	28.16%	13.21%
7	107,888	29.26%	27.33%	30.63%	12.78%
8	110,343	29.96%	27.26%	31.06%	11.72%

Table 9-8a. Scale-Score Descriptive Statistics: Mathematics

Grade	<i>N</i>	Mean	SD	Min	Max
3	102,587	1,296.1	27.0	1,217	1,361
4	105,113	1,393.5	25.7	1,310	1,455
5	109,057	1,487.1	26.4	1,409	1,550
6	108,789	1,587.8	25.8	1,518	1,650
7	108,000	1,688.1	26.6	1,621	1,752
8	110,383	1,787.5	27.0	1,725	1,850

Table 9-8b. Performance-Level Percentages: Mathematics

Grade	<i>N</i>	Not Proficient	Partially Proficient	Proficient	Advanced
3	102,587	27.94%	26.36%	27.10%	18.60%
4	105,113	24.77%	33.27%	25.64%	16.32%
5	109,057	37.22%	28.46%	17.75%	16.57%
6	108,789	34.54%	30.91%	18.57%	15.98%
7	108,000	36.34%	27.96%	19.41%	16.30%
8	110,383	40.35%	26.04%	15.78%	17.84%

Table 9-9a. Scale-Score Descriptive Statistics: Social Studies

Grade	<i>N</i>	Mean	SD	Min	Max
5	108,874	1,477.1	24.7	1,410	1,561
8	110,191	1,786.6	25.3	1,708	1,867
11	104,888	2,098.9	24.7	2,022	2,166

Table 9-9b. Performance-Level Percentages: Social Studies

Grade	<i>N</i>	Not Proficient	Partially Proficient	Proficient	Advanced
5	108,874	22.27%	59.40%	15.65%	2.67%
8	110,191	30.98%	39.76%	24.08%	5.19%
11	104,888	10.99%	40.52%	36.87%	11.63%

Chapter 10: Performance-Level Setting

This chapter briefly describes the M-STEP performance level setting and presents the cut scores established and the performance level descriptors derived from the performance level setting.

M-STEP is administered to assess Michigan students' mastery of the Michigan state standards. The assessments began as an implementation of the Smarter Balanced's ELA and mathematics tests. The current cut scores for the tests are taken from the SBAC tests. A brief overview of the Smarter Balanced standard-setting procedures during which the cut scores for ELA and mathematics were derived can be found in the report about performance level setting, *Smarter Balanced Assessment Consortium: Achievement Level Setting Final Report* (2015d), which is posted on the [Smarter Balanced library web page](https://portal.smarterbalanced.org/library/en/achievement-level-setting-final-report-with-appendix.pdf).¹

Over the course of several years, important changes have been made to the assessments to make them more meaningful to Michigan educators. These include the alignment of the test items to the Michigan Academic Standards, the implementation of a Michigan-specific test blueprint, and a reduction in the number of performance tasks used in ELA to reduce overall test time. These changes were made cautiously and deliberately with the active involvement of Michigan educators and stakeholders.

In school year 2017–18, the tests in grades 3–8 were shortened to a legislatively mandated median time of three hours to reduce the time burden on students and schools. To do so, all performance tasks in ELA were replaced with PBW items, a form of constructed response (CR) item. The ELA test blueprints were adjusted to accommodate the new item type and the reduction in test length. In grades 3–8 mathematics, the test was also shortened to reduce overall testing time, but this change did not involve adding new test items or significantly altering the test blueprint.

As a result of these changes, MDE partnered with DRC and Michigan educators to evaluate the validity of the cut scores derived by Smarter Balanced for ELA and mathematics with a cut score validation meeting in July 2018.

For social studies, a statistical articulation was used to establish cut scores.

The AERA, APA, & NCME (2014) *Standards* addressed in this chapter are 5.21 and 5.22, which will be presented in the pertinent sections of the chapter.

The AERA, APA, and NCME Standard 5.21 states that:

When proposed score interpretations involve one or more cut scores, the rationale and procedures used for establishing cut scores should be documented clearly. (p. 107)

¹ <https://portal.smarterbalanced.org/library/en/achievement-level-setting-final-report-with-appendix.pdf>

To evaluate the validity of M-STEP score interpretations, it is essential to understand that descriptors and cut scores are established in a collaborative and participatory process. The descriptors clearly establish, in plain language, the proper frame of reference for understanding how to interpret test scores, particularly cut scores.

10.1 Cut Score Validation for English Language Arts and Mathematics

The purpose of the standards validation was to determine whether the current M-STEP cut scores for grades 3–8 ELA and mathematics were still valid for continued use, given the 2018 updates to the tests.

A total of 54 Michigan educators engaged in a modification of the Bookmark Standard Setting Procedure (Lewis, Mitzel, & Green, 1996; Lewis, Mitzel, Mercado, & Schulz, 2012) to validate the cut scores. This method has been used on large-scale assessments in Michigan and across the nation, including by Smarter Balanced.

Participants studied the existing Michigan performance level descriptors (PLDs) and Michigan state standard to review the knowledge, skills, and abilities expected of students in each performance level. The four performance levels on M-STEP are Not Proficient, Partially Proficient, Proficient, and Advanced. Each performance level is associated with a level of mastery of the Michigan Learning *Standards*. Participants then discussed the content-based expectations for students at the threshold of each performance level (e.g., a student who is “just” Proficient). To support their discussions of these threshold students, participants were provided with the Smarter Balanced achievement level descriptors (ALDs). These Smarter Balanced ALDs were used at the original standard setting where the cut scores were established.

Participants studied collections of test items that were ordered in terms of difficulty. The existing cut scores were presented as benchmarks for participants’ consideration: participants were asked to consider the knowledge and skills that students would need to demonstrate on the updated ELA and mathematics tests, as based on the benchmarked (existing) cut scores. Then, participants compared these expectations against the content-based expectations for students at the thresholds of each performance level. Participants were instructed to recommend retaining the existing cut scores if there was good correspondence between the benchmarks and these content-based expectations or to recommend alternative cut scores that reflect better correspondence. Participants engaged in two rounds of individual judgments and group discussion. (The grade 5 mathematics committee engaged in three rounds of judgments to accommodate additional discussion.) The committees’ median judgments were taken as their final recommendations.

The available validity evidence suggests that there were no significant differences between the updated ELA and mathematics assessments and the content assessed by the prior assessments and that the differences between the judgments made at the 2018 standards validation workshop and the existing cut scores were not statistically different. That is, the recommendations made by Michigan educators during the standards validation were consistent with the existing cut scores, and the validity evidence collected during this process supports the continued use of the cut scores. More information can be found in the full *M-STEP Standards Validation 2018 Technical Report* (2019) in Appendix E.

10.2 Statistical Articulation for Social Studies

MDE partnered with DRC to conduct a standard-setting workshop for M-STEP social studies in grades 5, 8, and 11 in June 2015. During the workshop, participants considered the test items, performance level descriptors, and test data. Following the workshop, MDE considered participants’ recommendations and discussed with the state superintendent. MDE found that the participants’ recommended proficiency cuts were much lower than in the past and thus determined that such recommendations were not consistent with the high expectations of career and college readiness. As a result, in consultation with members of Michigan Technical Advisory Committee, MDE used statistically articulated cut scores and considered such approach to be more appropriate.

10.3 Scale Scores

This section describes the slopes and intercepts for transforming thetas to scale scores, as well as the LOSS and the HOSS for various M-STEP content areas. The values for ELA and mathematics were derived by MDE and DRC using the work done by Smarter Balanced.

For a detailed description of the methods used in calibration and scaling *Smarter Balanced 2017–2018 Technical Report* (Smarter Balanced, 2017). After calibration, results were in the theta metric. MDE transformed the theta metric results onto a four-digit scale, which is more meaningful for stakeholders. The equation for this linear transformation is

$$\text{Scale score} = (\text{theta} * \text{slope}) + \text{intercept}$$

Table 10-1 presents the information of slopes and intercepts for all four content areas, along with the LOSS and HOSS values which give the effective range of M-STEP scales for each grade and content area.

Table 10-1. Scale Transformation Slopes and Intercepts for M-STEP Summative Assessments with LOSS and HOSS Values

Content Area	Grade	Slope A	Intercept B	LOSS	HOSS
ELA	3	26.0061	1322.5934	1203	1357
ELA	4	24.6036	1409.5875	1301	1454
ELA	5	25.8718	1501.3628	1409	1560
ELA	6	24.5491	1592.9699	1508	1655
ELA	7	23.8151	1687.3543	1618	1753
ELA	8	24.1951	1782.9264	1721	1857
Math	3	26.3725	1325.7407	1217	1361
Math	4	25.2608	1409.0233	1310	1455
Math	5	23.3374	1495.6493	1409	1550
Math	6	20.4573	1589.9260	1518	1650
Math	7	19.6292	1686.6036	1621	1752
Math	8	18.5194	1782.8881	1725	1850
Social Studies	5	27.2005	1478.3212	1395	1568
Social Studies	8	26.9339	1785.9405	1703	1868
Social Studies	11	26.8528	2095.9989	2016	2166

10.4 Cut Scores

This section presents the cut scores for each grade/content area of M-STEP. Table 10-2 shows the cut scores for ELA and mathematics in grades 3–8 and for social studies in grades 5, 8, and 11. It should be noted that for ELA and mathematics, the Smarter Balanced established cut scores on the theta metric were transformed to the (Michigan specific) M-STEP scales using a linear transformation.

Table 10-2. Cut Scores for M-STEP Summative Assessments

Content Area	Grade	SS Cut between Levels 1 and 2	SS Cut between Levels 2 and 3	SS Cut between Levels 3 and 4
ELA	3	1280	1299.5	1317
ELA	4	1383	1399.5	1417
ELA	5	1481	1499.5	1524
ELA	6	1578	1599.5	1624
ELA	7	1679	1699.5	1726
ELA	8	1777	1799.5	1828
Math	3	1281	1299.5	1321
Math	4	1376	1399.5	1420
Math	5	1478	1499.5	1515
Math	6	1579	1599.5	1614
Math	7	1679	1699.5	1716
Math	8	1780	1799.5	1815
Social Studies	5	1458	1499.5	1530
Social Studies	8	1771	1799.5	1831
Social Studies	11	2069	2099.5	2131

10.5 Claim Cut Scores

As stated in Section 2.3, student performance on ELA and mathematics claims was classified into one of the three performance levels: *Adequate progress*, *Attention may be needed*, and *Most at risk of falling behind*. Detailed rules for calculating performance levels for ELA and mathematics claims can be found in the Smarter Balanced Scoring Specifications, 2014–2015 *Administration Summative and Interim Assessments: ELA/Literacy Grades 3–8, 11* and *Mathematics Grades 3–8, 11, V.7* in Appendix D (AIR, 2016).

10.6 Performance Level Descriptors

The performance level descriptors that were adopted by MDE for reporting purposes can be found in Tables 10-3 and 10-4.

Table 10-3. Performance-Level Descriptors for M-STEP, Grades 3–8

Performance Level	Descriptor
Advanced—PL 4	The student's performance exceeds grade-level content standards and indicates substantial understanding and application of key concepts defined for Michigan students. The student needs support to continue to excel.
Proficient—PL 3	The student's performance indicates understanding and application of key grade-level content standards defined for Michigan students. The student needs continued support to maintain and improve proficiency.
Partially Proficient—PL 2	The student needs assistance to improve achievement. The student's performance is not yet proficient, indicating a partial understanding and application of the grade-level content standards defined for Michigan students.
Not Proficient—PL 1	The student needs intensive intervention and support to improve achievement. The student's performance is not yet proficient and indicates minimal understanding and application of the grade level content standards defined for Michigan students.

Table 10-4. Performance-Level Descriptors for M-STEP, Grade 11

Performance Level	Descriptor
Advanced—PL 4	The student's performance exceeds the high school content standards and indicates substantial understanding and application of key concepts defined for Michigan students. The student needs support to continue to excel and to be college- and career-ready.
Proficient—PL 3	The student's performance indicates understanding and application of key high school content standards defined for Michigan students. The student needs continued support to maintain and improve proficiency and to be college- and career-ready.
Partially Proficient—PL 2	The student needs assistance to improve achievement and to become career and college ready. The student's performance is not yet proficient, indicating a partial understanding and application of the high school content standards defined for Michigan students.
Not Proficient—PL 1	The student needs intensive intervention and support to improve achievement and to become career and college ready. The student's performance is not yet proficient and indicates minimal understanding and application of the high school content standards defined for Michigan students.

10.7 Summary

This chapter presented a brief overview of the process for performance level setting used by Smarter Balanced for derivation of the ELA and mathematics cut scores. It also presents an overview of the procedure used for social studies. These procedures are addressed in more detail in Sections 10.1 and 10.2.

The standard settings undertaken by Smarter Balanced support the following *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014):

- Standard 5.21—When proposed score interpretations involve one or more cut scores, the rationale and procedures used for establishing cut scores should be documented clearly.²
- Standard 5.22—When cut scores defining pass-fail or proficiency levels are based on direct judgments about the adequacy of item or test performances, the judgmental process should be designed so that the participants providing the judgments can bring their knowledge and experience to bear in a reasonable way.

² For ELA and mathematics

Chapter 11: Fairness

As noted in the *Standards* (AERA, APA, & NCME, 2014), there are varying definitions of fairness. This chapter examines fairness as it relates to minimizing bias on a test and looks at test performance among varying subgroups assessed by M-STEP. It should be noted that differences in test performance among subgroups do not mean that a test is unfair—they simply mean that groups perform differently on the test. Even when a test is carefully and properly constructed, differences may exist among subgroups as a result of differences in curriculum or learning by the students in the subgroup.

This chapter is particularly relevant to AERA, APA, & NCME (2014) *Standards* 3.1 through 3.6. These standards are from Chapter 3 of the AERA, APA, & NCME (2014) *Standards*, “Fairness in Testing.” Each of these standards will be presented below. Standard 3.6 states the following:

Standard 3.6 Where credible evidence indicates that test scores may differ in meaning for relevant subgroups in the intended examinee population, test developers and/or users are responsible for examining the evidence for validity of score interpretations for intended uses for individuals from those subgroups. What constitutes a significant difference in subgroup scores and what actions are taken in response to such differences may be defined by applicable laws. (p. 65)

There is no specific research on M-STEP showing that the test scores of examinee subgroups differ in meaning; however, this is an ongoing concern in any large-scale testing program. To lessen the possibility of differences in test score meaning, DRC, MDE, and Smarter Balanced follow several steps in the item development and selection processes as explained in Section 11.1 of this chapter. In addition, DRC, MDE, and Smarter Balanced have conducted content and bias reviews on items, as explained in Chapter 3. These practices adhere to Standard 3.3:

Standard 3.3 Those responsible for test development should include relevant subgroups in validity, reliability/precision, and other preliminary studies used when constructing the test. (p. 64)

DRC and MDE have conducted differential item functioning (DIF) studies following the operational administration of M-STEP. Typically, items are evaluated for possible DIF in the field-test phase of test development, and items flagged for DIF are further examined for possible bias. During test development, Smarter Balanced follows procedures to minimize the inclusion of items that may potentially favor one demographic group over another. DRC and MDE staff do the same for social studies. Section 11.2 of this chapter explains the steps taken to evaluate M-STEP items through the use of DIF to adhere with this standard.

In addition, standardized test administration and training of test administrators for M-STEP comply with *Standards* 3.4 and 3.5:

Standard 3.4 Test takers should receive comparable treatment during the test administration and scoring process. (p. 65)

Standard 3.5 Test developers should specify and document provisions that have been made to test administration and scoring procedures to remove construct-irrelevant barriers for all relevant subgroups in the test-taker population. (p. 65)

Section 11.1 of this chapter is also directly relevant to *Standards* 3.1 and 3.2:

Standard 3.1 Those responsible for test development, revision, and administration should design all steps of the testing process to promote valid score interpretations for intended score uses for the widest possible range of individuals and relevant subgroups in the intended population. (p. 63)

Standard 3.2 Test developers are responsible for developing tests that measure the intended construct and for minimizing the potential for tests' being affected by construct-irrelevant characteristics, such as linguistic, communicative, cognitive, cultural, physical, or other characteristics. (p. 64)

Section 11.1 explains the steps taken by DRC and MDE to minimize words, phrases, and content that may be regarded as offensive by members of particular demographic subgroups. Chapter 3 discusses item content considerations during item development and item bias reviews for items included in M-STEP. These reviews are also critical in fulfilling *Standards* 3.1 and 3.2.

11.1 Minimizing Bias through Careful Test Development

The development of a test that is fair for all examinees begins in the early stages of planning and development. The item and test development processes that are used to minimize bias are summarized below.

First, careful attention is paid to content validity during the item development and item selection processes. Bias can occur only if the test is measuring different things for different groups. By eliminating irrelevant skills or knowledge from the items, the possibility of bias is reduced. Second, item writers and test developers follow several published guidelines for reducing or eliminating bias.

11.1.1 ELA and Mathematics

Smarter Balanced developed *Bias and Sensitivity Guidelines* (ETS, 2012) to help ensure that the assessments are fair for all groups of test takers, despite differences in characteristics that include, but are not limited to, disability status, ethnic group, gender, regional background, native language, race, religion, sexual orientation, and socioeconomic status. Unnecessary barriers can be reduced by following some fundamental rules:

- Measuring only knowledge or skills that are relevant to the intended construct
- Not angering, offending, upsetting, or otherwise distracting test takers
- Treating all groups of people with appropriate respect in test materials

These rules help ensure that the test content is fair for test takers as well as acceptable to the many stakeholders and constituent groups within Smarter Balanced member organizations. The more typical view is that bias and sensitivity guidelines apply primarily to the review of test items. However, fairness must be considered in all phases of test development and use. Smarter Balanced strongly rely on the *Bias and Sensitivity Guidelines* (ETS, 2012) in the development of the Smarter Balanced assessments, particularly in item writing and review. Items must comply with the Bias and Sensitivity Guidelines in order to be included in the Smarter Balanced assessments.

Smarter Balanced assessments are developed using the principles of evidence-centered design (ECD). ECD requires a chain of evidence-based reasoning that links test performance to the claims made about test takers. Fair assessments are essential to the implementation of ECD. If test items are not fair, then the evidence they provide means different things for different groups of students. Under those circumstances, the claims cannot be equally supported for all test takers, which is a threat to validity. As part of the validation process, all items are reviewed for issues of bias and sensitivity using the *Bias and Sensitivity Guidelines* (ETS, 2012) prior to being presented to students. This helps ensure that item responses reflect only knowledge of the intended content domain, are free of offensive or distracting material, and portray all groups in a respectful manner. When the guidelines are followed, item responses provide evidence that supports assessment claims.

11.1.2 Social Studies

DRC and MDE item writers and test developers follow documented bias and sensitivity guidelines to ensure that the items are fair for all groups of test takers, despite differences in characteristics that include, but are not limited to, disability status, ethnic group, gender, regional background, native language, race, religion, sexual orientation, and socioeconomic status. Test developers review all items included in M-STEP and other testing materials with these guidelines in mind.

Careful attention is given to item statistics (if available) throughout the test development process. As part of the test assembly process, attempts are made to avoid using or reusing items with poor statistics. Additional steps to reduce bias, including the use of content and bias committees comprised of Michigan educators, are described in more detail in Chapter 3 of this report.

The goal of fairness in assessment is to ensure that test materials are as free as possible from unnecessary barriers to the success of diverse groups of students.

11.2 Evaluating Bias through Differential Item Functioning (DIF)

An empirical approach known as DIF is used to examine items after they have been administered. The DIF statistics indicate the degree to which members of a particular subgroup perform better or worse than expected on each item as compared to the members of the reference group. Therefore, DIF flags do not necessarily indicate that an item is biased; rather, DIF flags indicate that the item functions differently for equally able members of different groups (Camilli & Shepard, 1994). The DIF procedures and results are described in this section. Note that items are not necessarily suppressed from operational scoring if they are flagged for DIF.

The position of DRC concerning test bias is based on two general propositions. First, students may differ in their background knowledge, cognitive and academic skills, language, attitudes, and values. To the degree that these differences are large, no one curriculum and no one set of instructional materials will be equally suitable for all. Therefore, no one test will be equally appropriate for all. Furthermore, it is difficult to specify what amount of difference can be called large and to determine how these differences will affect the outcome of a particular test. Second, schools have been assigned the tasks of developing certain basic cognitive skills and supporting the development of these skills equitably among all students. Therefore, there is a need for tests that measure the common skills and bodies of knowledge that are common to all learners. The test developers' task is to create assessments that measure these key cognitive skills without introducing extraneous or construct-irrelevant elements into the performances on which the measurement is based. If these tests require that students have culturally specific knowledge and skills not taught in school, differences in performance among students can occur because of differences in student background and out-of-school learning. Such tests are measuring different things for different groups and can be called biased (Camilli & Shepard, 1994; Green, 1975).

To lessen such biases, DRC and MDE strive to minimize the role of extraneous elements, thereby increasing the number of students for whom the test is appropriate. As discussed above and in Chapter 3 of this report, careful attention is given during the test development and form construction processes to lessen the influence of these elements for large numbers of students (including the use of content and bias review committees). Unfortunately, in some cases, extraneous elements may continue to play a substantial role. To assess the extent to which items may be performing differently for various subgroups of interest, DIF analyses are conducted during field testing and after each operational test administration. DIF statistics are used to quantify differences in item performance between two groups after controlling for examinees' overall achievement level. For M-STEP, DIF is conducted for ELA, mathematics, and social studies using very similar procedures. Details in Sections 11.3.1 and 11.3.2 provide DIF results for the following subgroups:

- **Gender:** The focal group is female; the reference group is male.
- **Race/Ethnicity:** The focal groups are students whose race/ethnicity is reported as African American or Black, Hispanic or Latino, or Asian; the reference group is students whose race/ethnicity is reported as White.

- **Disability status:** The focal group is students who are identified as students with disabilities (SWD); the reference group is all others.
- **English Proficiency status:** The focal group is students who are identified as Limited English Proficiency (LEP); the reference group is all others.
- **Socio-economic status:** The focal group is students who are identified as economically disadvantaged (EconDis); the reference group is all others.

11.3 DIF Statistics

Two commonly used DIF statistics were applied to M-STEP items and are described here. They are (1) the Mantel-Haenszel (MH) statistic (Mantel & Haenszel, 1959) for dichotomously scored items and an extension of the $MH\chi^2$ (Mantel, 1963) for polytomously scored items, and (2) the standardized mean difference (SMD) effect size (ES) for polytomously scored items (Dorans & Schmitt, 1991).

For dichotomously scored items (e.g., MC items), the MH statistic is computed as follows (Camilli & Shepard, 1994):

$$MH\chi^2 = \frac{\left\{ \sum_{j=1}^S [A_j - E(A_j)] - 1/2 \right\}^2}{\sum_{j=1}^S VAR(A_j)} \quad (11.1)$$

where $VAR(A_j) = \frac{n_{Rj}n_{Fj}m_{1j}m_{0j}}{T_j^2(T_j-1)}$ and $E(A_j) = \frac{n_{Rj}m_{1j}}{T_j}$.

In Equation 11.1, $A_j - E(A_j)$ represents the difference between the observed number and the expected number of correct responses on the item by the reference group members who have the j th score on the matching variable;¹ n_{Rj} and n_{Fj} represent the number of examinees in the reference and focal groups, respectively, for the j th score on the matching variable; m_{1j} represents the total number of examinees (both reference and focal) with the j th score on the matching variable and with a correct response on the current item; m_{0j} represents the total number of examinees with the j th score on the matching variable and with an incorrect response on the current item. The $MH\chi^2$ is evaluated against the standard χ^2 critical with one degree of freedom.

The $MH\chi^2$ does not indicate the strength of association of the relationship between item performance and group membership. The MH odds ratio can be computed to estimate the strength of this association. The resulting estimate represents the relative likelihood of success on a particular item for members of two different groups of examinees (Camilli, 2006). This odds ratio thus provides an estimate of effect size (ES) with a value of 1.0, indicating no DIF. A value greater than 1.0 indicates that, on average, the reference group members performed better than comparable focal group members did. A value less than 1.0 indicates that, on average, the reference group members performed worse than comparable focal group members did.

¹ Total observed score is used as the matching variable for DIF analysis here.

The odds of a correct response (i.e., proportion passing divided by proportion failing) is P/Q (i.e., $P/[1-P]$). The MH odds ratio is simply the odds of a correct response of the reference group divided by the odds of a correct response of the focal group. The formula for its estimation is as follows (Camilli & Shepard, 1994, p. 116):

$$\hat{\alpha}_{MH} = \frac{\sum_{j=1}^S A_j D_j / T_j}{\sum_{j=1}^S B_j C_j / T_j} \quad (11.2)$$

where $S = K - 1$ and represents the actual number of 2×2 contingency tables (assuming the tables have at least 1 person in each cell); K represents the number of items on the test; j signifies the j th score on the matching variable and runs from 0 to K .² For j th score category, A_j represents the number of reference group members with a correct response, B_j represents the number of reference group members with an incorrect response, C_j represents the number of focal group members with a correct response, and D_j represents the number of focal group members with an incorrect response. T_j represents the total number of examinees who have the j th score on the matching variable.

The corresponding null hypothesis is that the odds of getting the item correct are equal for the two groups (i.e., the odds ratio is equal to 1):

$$H_0: \alpha_{MH} = 1 \quad (11.3)$$

To make the odds ratio symmetrical around zero with its range located in the interval $-\infty$ to $+\infty$, the odds ratio is transformed into a log-odds ratio as follows (Camilli & Shepard, 1994, p.116):

$$\hat{\lambda}_{MH} = \log(\alpha_{MH}) \quad (11.4)$$

The natural logarithm transformation of this odds ratio is symmetrical around zero (where 0 indicates no DIF). This DIF measure is a signed index, where a positive value represents DIF in favor of the reference group and a negative value indicates DIF in favor of the focal group.

The variance of the log-odds ratio estimate (V_{λ}) is computed as follows (Camilli & Shepard, 1994, p. 121):

$$V_{\lambda} = \frac{\sum_{j=1}^S T_j^{-2} (A_j D_j + \alpha_{MH} B_j C_j) [A_j + D_j + \alpha_{MH} (B_j + C_j)]}{2(\sum_{j=1}^S A_j D_j / T_j)^2} \quad (11.5)$$

The terms included in Equation 11.5 correspond to those presented for Equation 11.2. In practice, a standardized MH log-odds ratio is computed by dividing the estimate $\hat{\lambda}_{MH}$ by the estimated standard error. According to Penfield (2007, p.16), “A value greater than 2.0 or less than -2.0 may be considered evidence of the presence of DIF.”

² Although the value of the matching variable runs from 0 to K , the all correct (i.e., K) and all incorrect (i.e., 0) score categories are not included in the DIF analysis in order to avoid having a denominator equal to 0.

In addition, once $\hat{\lambda}_{MH}$ is obtained using Equation 11.4, the delta statistic (MH D-DIF, used by SBAC in flagging criteria) can be computed as

$$\text{MH D-DIF} = -2.35 \times \hat{\lambda}_{MH} \quad (11.6)$$

For polytomously scored items, an extension of the MH χ^2 procedure was computed (Mantel, 1963). The statistic is computed as follows (Zwick, Donaghue, & Grima, 1993, p. 239):

$$\text{Mantel } \chi^2 = \frac{(\sum_k F_k - \sum_k E(F_k))^2}{\sum_k \text{VAR}(F_k)} \quad (11.7)$$

where F_k is the sum of scores for the focal group at the k th level of the matching variable and is defined as

$$F_k = \sum_t y_t n_{Ftk}, \quad (11.8)$$

and the expectation of F_k under the hypothesis of no association is

$$E(F_k) = \frac{n_{F+k}}{n_{++k}} \sum_t y_t n_{+tk}, \quad (11.9)$$

and the variance of F_k under the assumption of no association is

$$\text{Var}(F_k) = \frac{n_{R+k} n_{F+k}}{n_{++k}^2 (n_{++k} - 1)} \left\{ \left(n_{++k} \sum_t y_t^2 n_{+tk} \right) - \left(\sum_t y_t n_{+tk} \right)^2 \right\}. \quad (11.10)$$

Using the Mantel approach for ordered categories, the data are organized into a $2 \times T \times K$ contingency table, where T is the number of response categories and K is the number of levels of the matching variable. y_1, y_2, \dots, y_T represent the T scores that can be obtained on the item; n_{Rik} and n_{Fik} represent the number of examinees in the reference and focal groups, respectively, who are at the k th level of the matching variable and received an item score of y_i . The “+” denotes summation over a particular index (e.g., n_{R+k} denotes the total number of reference group members at the k th level of the matching variable). Under the null hypothesis of no association, the Mantel statistic has a chi-square distribution with one degree of freedom. For dichotomous items, the Mantel statistic reduces to the MH statistic (without the continuity correction).

In addition to the MH statistic, an ES was calculated by dividing the SMD statistics by the overall (i.e., focal and reference groups combined) standard deviation (SD) of the item scores: $\text{ES} = \text{SMD} / \text{SD}$. The SMD compares the mean of the reference and focal groups, adjusting for the distribution of reference and focal group members on the matching variable (Zwick et al., 1993), which for these analyses is the M-STEP raw score. SMD is computed as follows (Zwick et al., 1993):

$$\text{SMD} = \sum_k p_{Fk} (m_{Fk} - m_{Rk}) \quad (11.11)$$

where p_{Fk} is the proportion of the focal group members at the k th level of the matching variable m_{Fk} and m_{Rk} indicate mean item score for the focal group and the reference group at the k th level of the matching variable, respectively.

A negative SMD value implies that the focal group has a lower mean item score than the reference group, whereas a positive value implies that the focal group has a higher mean item score than the reference group, conditioned on the matching test score.

11.3.1 Flagging Criteria and Results for ELA and Mathematics

For ELA and mathematics, according to Smarter Balanced (for more information, see the *Smarter Balanced 2017–2018 Technical Report* [2018]), the minimum case count for each of the two groups (i.e., the focal group and the reference group) was set at 100 and the minimum case count for the combined group was set to 400.

The following flagging criteria were used for dichotomously scored items (e.g., MC items):

- Moderate DIF: significant MH chi-square statistic ($p < 0.05$) and $1.0 \leq |\text{MH D-DIF}| < 1.5$
- Large DIF: significant MH chi-square statistic ($p < 0.05$) and $|\text{MH D-DIF}| \geq 1.5$

The following flagging criteria were used for polytomously scored items:

- Moderate DIF: if the extension of the MH statistic is significant ($p < .05$) and $|\text{ES}|$ is > 0.17 and ≤ 0.25 .
- Large DIF: if the extension of the MH statistic is significant ($p < .05$) and $|\text{ES}| > 0.25$.

A positive MH D-DIF or ES value indicates that the item favors the focal group, while a negative value indicates that the item favors the reference group instead.

Table 11-1 shows the item counts for ELA and mathematics DIF analyses based on the 2018 M-STEP administration. Tables 11-2 and 11-3 summarize the number of items having moderate or large DIF flags (i.e., B or C) by grade for each focal/reference group that included at least 100 students for ELA and mathematics, respectively. For example, consider grade 3 ELA. There were 13 items (or 3.8% of all eligible items) flagged for moderate DIF. Specifically, there were 6 items favoring males and 7 items favoring females.

Again, any items included on the M-STEP ELA and mathematics assessments (including those items flagged for DIF) have been thoroughly reviewed by MDE staff, DRC test development staff, and Smarter Balanced staff.

Table 11-1. Item Counts Used in Differential Item Functioning Analyses: ELA and Mathematics

Content Area	Grade	<i>N</i> Items	Female/ Male	Asian/ White	Black or African American/ White	Hispanic or Latino/ White	SWD/ Non-SWD	LEP/ Non-LEP	EconDis/ Non-EconDis
ELA	3	778	707	672	695	669	694	694	706
ELA	4	780	716	673	700	683	701	697	714
ELA	5	733	654	624	635	625	638	633	649
ELA	6	694	652	619	624	619	627	624	639
ELA	7	591	554	526	539	528	540	538	549
ELA	8	680	629	592	617	592	619	613	627
Math	3	1102	988	648	784	678	817	812	979
Math	4	1190	1019	696	736	702	746	735	1026
Math	5	1127	990	672	828	743	975	967	992
Math	6	1044	910	770	887	810	897	888	903
Math	7	973	867	484	610	504	647	642	860
Math	8	803	685	649	663	652	666	663	675

Table 11-2. Number of Differential Item Functioning Flagged Items: ELA

Grade	DIF Category	Female/ Male	Asian/White	Black or African American/ White	Hispanic or Latino/ White	Disabilities/ Without Disabilities	LEP/ Non-LEP	EconDis/ Non-EconDis
3	b-	6	24	5	4	5	10	1
3	b+	7	20	18	9	37	28	12
3	c-	0	3	0	0	1	0	0
3	c+	0	5	3	3	5	4	0
4	b-	5	21	8	5	12	15	2
4	b+	17	19	25	19	41	33	7
4	c-	2	4	1	0	1	1	1
4	c+	1	1	0	1	7	4	0
5	b-	11	19	9	9	5	10	0
5	b+	17	19	13	10	42	36	7
5	c-	1	1	0	0	1	0	0
5	c+	5	2	1	2	6	9	0
6	b-	10	15	1	5	8	9	4
6	b+	13	22	16	17	41	33	8
6	c-	3	2	0	1	1	3	1
6	c+	5	2	1	0	5	7	0
7	b-	9	9	7	5	10	6	0
7	b+	30	11	23	6	33	21	5
7	c-	4	2	0	0	0	1	0
7	c+	3	1	1	0	2	5	0
8	b-	8	10	4	10	4	8	2
8	b+	24	29	31	24	30	37	4
8	c-	4	5	1	0	1	1	0
8	c+	5	1	2	0	3	12	0

Table 11-3. Number of Differential Item Functioning Flagged Items: Mathematics

Grade	DIF Category	Female/ Male	Asian/White	Black or African American/ White	Hispanic or Latino/ White	Disabilities/ Without Disabilities	LEP/ Non-LEP	EconDis/ Non-EconDis
3	b-	36	34	37	13	5	19	10
3	b+	31	16	49	11	22	19	15
3	c-	1	14	2	4	3	1	0
3	c+	4	6	8	0	3	1	2
4	b-	29	39	13	4	8	12	6
4	b+	45	15	45	7	16	17	23
4	c-	4	10	1	0	0	2	0
4	c+	8	2	7	1	2	2	2
5	b-	41	49	19	7	12	19	9
5	b+	33	15	39	7	21	24	14
5	c-	5	15	5	0	0	4	0
5	c+	7	7	3	2	4	12	0
6	b-	59	44	22	5	11	28	7
6	b+	44	18	47	8	31	27	21
6	c-	3	28	5	1	3	3	0
6	c+	9	7	8	3	3	10	0
7	b-	29	47	19	7	7	24	13
7	b+	24	16	23	6	9	18	15
7	c-	1	32	2	0	1	6	2
7	c+	4	5	14	0	2	10	3
8	b-	5	49	20	3	9	11	5
8	b+	28	11	19	12	12	16	21
8	c-	0	31	3	4	2	7	0
8	c+	2	1	6	0	0	6	1

11.3.2 Flagging Criteria and Results for Social Studies

For science and social studies, the minimum case count was 30 for each of the two groups (i.e., the reference group and the focal group). The following flagging criteria, adapted from Penfield (2007), were used:

- Negligible DIF (a): if either MH common log-odds ratio ($\hat{\lambda}_{MH}$) is not significantly different from zero or $|\hat{\lambda}_{MH}| < 0.426$
- Moderate DIF (b): if $\hat{\lambda}_{MH}$ is significantly different from zero and $|\hat{\lambda}_{MH}| \geq 0.426$ and either (a) $|\hat{\lambda}_{MH}| \leq 0.638$, or (b) $|\hat{\lambda}_{MH}|$ is not significantly greater than 0.426
- Large DIF (C): if $|\hat{\lambda}_{MH}|$ is significantly greater than 0.426 and $|\hat{\lambda}_{MH}| > 0.638$.

Chapter 11: Fairness

Table 11-4 shows the item counts for social studies DIF analyses. Tables 11-5 and 11-6 summarize the number of items having moderate and large DIF flags (i.e., b or c). For example, consider grade 8 social studies. There was 1 item (or 1.8% of all items) flagged for large Asian vs. White DIF, favoring White.

Again, any items included on the M-STEP social studies assessments (including those items flagged for DIF) have been thoroughly reviewed by MDE staff, DRC test development staff, and Michigan content/bias committee members.

Table 11-4. Item Counts used in Differential Item Functioning Analyses: Social Studies

Content Area	Grade	<i>N</i> Items	Female/ Male	Asian/ White	Black or African American/ White	Hispanic or Latino/ White	SWD/ Non-SWD	LEP/ Non-LEP	EconDis/ Non-EconDis
Social Studies	5	52	52	52	52	52	52	52	52
Social Studies	8	55	55	55	55	55	55	55	55
Social Studies	11	44	44	44	44	44	44	44	44

Table 11-5. Number of Differential Item Functioning Flagged Items: Social Studies

Grade	DIF Category	Female/ Male	Asian/White	Black or African American/ White	Hispanic or Latino/ White	Disabilities/ Without Disabilities	LEP/ Non-LEP	EconDis/ Non-EconDis
5	b-	0	1	0	0	1	1	0
5	b+	0	1	0	0	0	1	0
5	c-	0	0	0	0	0	0	0
5	c+	0	0	0	0	0	0	0
8	b-	1	0	0	0	0	0	0
8	b+	0	0	1	1	0	1	1
8	c-	0	1	0	0	0	0	0
8	c+	0	0	0	0	0	0	0
11	b-	0	1	0	0	0	1	0
11	b+	1	1	0	0	0	0	0
11	c-	0	0	0	0	0	1	0
11	c+	0	0	1	0	0	0	0

11.4 Summary

In summary, the overall purpose of this chapter is to address fairness concerns that are relevant to the administration of M-STEP. The information in this chapter supports multiple best practices of the testing industry and particularly the following AERA, APA, & NCME (2014) standards:

- Standard 3.1—Those responsible for test development, revision, and administration should design all steps of the testing process to promote valid score interpretations for intended score uses for the widest possible range of individuals and relevant subgroups in the intended population.
- Standard 3.2—Test developers are responsible for developing tests that measure the intended construct and for minimizing the potential for tests' being affected by construct-irrelevant characteristics, such as linguistic, communicative, cognitive, cultural, physical, or other characteristics.
- Standard 3.3—Those responsible for test development should include relevant subgroups in validity, reliability/precision, and other preliminary studies used when constructing the test.
- Standard 3.4—Test takers should receive comparable treatment during the test administration and scoring process.
- Standard 3.5—Test developers should specify and document provisions that have been made to test administration and scoring procedures to remove construct-irrelevant barriers for all relevant subgroups in the test-taker population.
- Standard 3.6—Where credible evidence indicates that test scores may differ in meaning for relevant subgroups in the intended examinee population, test developers and/or users are responsible for examining the evidence for validity of score interpretations for intended uses for individuals from those subgroups. What constitutes a significant difference in subgroup scores and what actions are taken in response to such differences may be defined by applicable laws.

Chapter 12: Reliability and Evidence of Construct-Related Validity

This chapter presents evidence supporting construct-related validity. Part of the test validity argument is that scores must be consistent and precise enough to be useful for the intended purposes. The concepts of reliability and precision are examined through analysis of measurement error in simulated and operational (OP) conditions.

This chapter demonstrates M-STEP's adherence to AERA, APA, & NCME (2014) *Standards* 2.0, 2.1, 2.3, 2.13, 2.14, 2.16, 2.19, and 4.3. Each standard will be discussed in the pertinent section of this chapter.

12.1 Reliability

Reliability refers to the consistency of the students' test scores on parallel forms of a test. A reliable test is one that produces scores that are expected to be relatively stable if the test is administered repeatedly under similar conditions. Often, however, it is impractical to administer multiple forms of the test, and reliability is estimated on a single administration of the test. This type of reliability, known as internal consistency, provides an estimate of how consistently examinees perform across items within a test during a single test administration (Crocker & Algina, 1986). Reliability is a necessary but insufficient condition of validity.

The AERA, APA, & NCME (2014) *Standards* indicates the following:

The term reliability has been used in two ways in the measurement literature. First, the term has been used to refer to the reliability coefficients of classical test theory, defined as the correlation between scores on two equivalent forms of the test, presuming that taking one form has no effect on performance on the second form. Second, the term has been used in a more general sense, to refer to the consistency of scores across replications of a testing procedure, regardless of how this consistency is estimated or reported (e.g., in terms of standard errors, reliability coefficients per se, generalizability coefficients, error/tolerance ratios, item response theory [IRT] information functions, or various indices of classification consistency). (p. 33)

In accordance with the AERA, APA, & NCME (2014) *Standards* and in developing and maintaining tests of the highest quality, the reliability of each

M-STEP test has been calculated.

There are several specific AERA, APA, & NCME (2014) standards that this chapter addresses. These include *Standards* 2.0, 2.3, 2.13, and 2.19. Each standard is articulated below.

Standard 2.0 Appropriate evidence of reliability/precision should be provided for the interpretation for each intended score use. (p. 42)

Standard 2.3 For each total score, subscore, or combination of scores that is to be interpreted, estimates of relevant indices of reliability/precision should be reported. (p. 43)

The total score reliabilities are reported below in Sections 12.1.5 through 12.1.7 of this chapter. The overall standard errors of measurement (SEMs) and conditional standard errors of measurement (CSEMs) by decile are presented in Section 12.1.5.

Standard 2.13 The standard error of measurement, both overall and conditional (if reported), should be provided in units of each reported score. (p. 45)

The SEM based on scale scores and the CSEM based on scale scores are discussed below in Section 12.1.5.

Standard 2.19 Each method of quantifying the reliability/precision of scores should be described clearly and expressed in terms of statistics appropriate to the method. The sampling procedures used to select test takers for reliability/precision analyses and the descriptive statistics on these samples, subject to privacy obligations where applicable, should be reported. (p. 47)

12.1.1 Reliability and Standard Error of Measurement

According to the classical true score theory, which is a fundamental component of the CTT, an observed score is a sum of two parts—a random component of true score (T) and a random component of error score (E), or mathematically, $X = T + E$ (McDonald, 1999). This model has the following properties: (1) the expected error score is zero, (2) the correlation between the true score and the error score is zero, and (3) the correlation between the error scores on different but parallel forms is zero (Lord & Novick, 1968).

Based on this model, a student's observed test score is an imprecise estimate of his or her actual ability because a portion of that score is attributable to random error. A fundamental theoretical quantity in test theory, the *reliability coefficient* of observed scores is defined as the ratio of the variance of true scores to the variance of observed scores. Tests are therefore most reliable when the proportion of observed score variance that may be attributed to error variance is minimized. According to McDonald (1999), test-retest methods, parallel or alternate-form methods, and internal analysis are the three recognized methods for estimating the reliability coefficient.

Due to practical difficulties in applying the first two above-mentioned methods, only the internal consistency reliability approach is described here. Estimates of internal consistency reliability involve “dividing the test into two or more constituent parts and in some way estimating reliability from the consistency of performance across these part-tests” (Haertel, 2006, p. 71).

12.1.2 Cronbach's Coefficient Alpha

Historically, various internal consistency reliability estimates have been proposed, but the most widely used, for fixed forms, is Cronbach's (1951) coefficient alpha (Haertel, 2006). Using sample statistics, it is computed as follows (adapted from Haertel, 2006):

$$\alpha = \frac{l}{l-1} \left(1 - \frac{\sum_{i=1}^l S_i^2}{S_X^2} \right) \quad (12.1)$$

where l represents the number of items on the test, S_i^2 represents the sample variance of item i , and S_X^2 represents the sample variance of the total raw score.

The use of coefficient alpha has several theoretical advantages (Haertel, 2006). First, since it equals the mean of all possible split-half reliability coefficients, which is another estimate of internal consistency reliability that involves the division of the total test into two "parallel" sub-tests, the use of coefficient alpha avoids the arbitrary choice of a split or division. Second, it is mathematically equivalent to one of the lower bounds of the theoretical reliability coefficient. The implication of this is that the theoretical reliability coefficient is higher than the observed coefficient alpha.

12.1.3 Standard Error of Measurement

SEM is related to reliability and is calculated with sample statistics as follows (Hays, 1994, p. 617):

$$SEM(X) = S_X \sqrt{1 - r_{XX'}} \quad (12.2)$$

where $SEM(X)$ represents the estimated SEM of the observed test score X , S_X denotes the estimated standard deviation (i.e., sample standard deviation) of the observed score, and $r_{XX'}$ represents the estimated reliability coefficient of a test. In this report, the observed coefficient alpha is used as the estimated reliability coefficient for social studies.

According to Equation 12.2, the SEM is inversely related to the reliability of a test: For any standard deviation of the observed score, the SEM decreases when the reliability coefficient increases. Thus, when an SEM is small, one has more confidence in the accuracy, or precision, of the observed test scores.

12.1.4 Marginal Reliability for ELA and Mathematics

In a CAT administration, each student receives a different test form; therefore, the calculation of coefficient alpha is not applicable. An observed reliability can be derived from SEMs, which are computed from the test form each student took. The method of standard error calculation for both total and score reporting category scores, as described in the Smarter Balanced Scoring Specifications for 2014–2015 (AIR, 2014), is displayed below:

The standard error (SE) for student i is

$$SE(\theta_i) = \frac{1}{\sqrt{I(\theta_i)}} \quad (12.3)$$

where $I(\theta_i)$ is the test information function for student i , calculated as

$$I(\theta_i) = \sum_{j=1}^J D^2 a_j^2 \left(\frac{\sum_{l=1}^{m_j} l^2 \exp(\sum_{k=1}^l D a_j (\theta_i - b_{jk}))}{1 + \sum_{l=1}^{m_j} \exp(\sum_{k=1}^l D a_j (\theta_i - b_{jk}))} - \left(\frac{\sum_{l=1}^{m_j} l \exp(\sum_{k=1}^l D a_j (\theta_i - b_{jk}))}{1 + \sum_{l=1}^{m_j} \exp(\sum_{k=1}^l D a_j (\theta_i - b_{jk}))} \right)^2 \right) \quad (12.4)$$

where m_j is the maximum possible score point (starting from 0) for the j th item, D is the scale factor, 1.7. Values of a_j and b_{jk} are item parameters for item j and score level k .

SE is calculated based only on the answered items. The upper bound of the SE is set to 2.5 on the theta metric. Any value larger than 2.5 is truncated at 2.5 on the theta metric.

12.1.5 Observed Reliability, SEM, and CSEM for ELA and Mathematics

The marginal reliability for ELA and mathematics was calculated using the 2018 Michigan administration data. The results are presented in Table 12-1.

Table 12-1. ELA and Mathematics Summative Scale-Score Marginal Reliability Estimates

Content Area	Grade	<i>N</i>	Mean # Items	SD(SS)	Mean SEM	Marginal Reliability
ELA	3	101,516	45.22	25.79	6.10	0.94
ELA	4	104,078	45.03	26.04	6.26	0.94
ELA	5	108,040	45.42	27.12	6.57	0.94
ELA	6	107,916	45.13	26.20	6.67	0.93
ELA	7	107,359	45.05	26.54	7.17	0.93
ELA	8	109,775	45.31	27.40	7.25	0.93
Mathematics	3	101,643	36.00	27.04	5.55	0.96
Mathematics	4	104,125	36.00	25.72	5.36	0.95
Mathematics	5	108,030	36.00	26.43	5.91	0.94
Mathematics	6	107,883	36.00	25.85	5.65	0.95
Mathematics	7	107,318	36.00	26.58	6.12	0.94
Mathematics	8	96,169	36.00	27.97	6.72	0.94

SD(SS) = standard deviation of scale score

Table 12-2 shows that the marginal reliability varies by overall score levels. All students take a similar number of items, but the information delivered by the items differs. The most information occurs where the pool item difficulty and students' ability match the best with abundant items for selection. As shown in Figures 8-1 and 8-2, Smarter Balanced pools, used by Michigan, are difficult relative to the student ability levels of the state population. Consistently, as shown in Table 12-2, students with lower scores (e.g., deciles 1 and 2) have lower reliability than those with higher scores (e.g., deciles 8 and 9), except for students in ELA at grade 5.

Table 12-2. Marginal Reliability Overall and by Decile for ELA and Mathematics

Content Area	Grade	N	Var	Overall	Decile 1	Decile 2	Decile 3	Decile 4	Decile 5	Decile 6	Decile 7	Decile 8	Decile 9	Decile 10
ELA	3	101,516	25.79	0.94	0.91	0.94	0.95	0.95	0.95	0.95	0.95	0.95	0.95	0.94
ELA	4	104,078	26.04	0.94	0.92	0.94	0.95	0.95	0.95	0.95	0.95	0.95	0.95	0.93
ELA	5	108,040	27.12	0.94	0.93	0.95	0.95	0.95	0.95	0.95	0.94	0.94	0.93	0.93
ELA	6	107,916	26.20	0.93	0.89	0.93	0.93	0.94	0.94	0.95	0.95	0.95	0.94	0.93
ELA	7	107,359	26.54	0.93	0.89	0.92	0.93	0.93	0.94	0.94	0.94	0.93	0.93	0.92
ELA	8	109,775	27.40	0.93	0.90	0.93	0.93	0.94	0.94	0.94	0.94	0.94	0.94	0.93
Mathematics	3	101,643	27.04	0.96	0.92	0.95	0.96	0.96	0.96	0.97	0.97	0.97	0.97	0.96
Mathematics	4	104,125	25.72	0.95	0.90	0.95	0.96	0.96	0.96	0.96	0.97	0.97	0.97	0.96
Mathematics	5	108,030	26.43	0.94	0.84	0.91	0.93	0.95	0.96	0.97	0.97	0.98	0.97	0.97
Mathematics	6	107,883	25.85	0.95	0.87	0.93	0.94	0.95	0.96	0.96	0.96	0.97	0.97	0.97
Mathematics	7	107,318	26.58	0.94	0.79	0.91	0.93	0.95	0.95	0.96	0.97	0.98	0.98	0.98
Mathematics	8	96,169	27.97	0.94	0.84	0.91	0.93	0.94	0.94	0.95	0.96	0.97	0.97	0.97

Because of the CSEM differences by score level, demographic groups with lower average scores tend to have lower reliability than the population as a whole. Due to the small sample sizes of some of the subgroups (e.g., American Indian or Alaska Native and Native Hawaiian or Other Pacific Islander), corresponding results should be interpreted with caution. Tables 12-3 to 12-6 show marginal reliability by demographic group and the MSE.

Table 12-3. Marginal Reliability of Total Summative Scores by Ethnic Group—ELA

Grade	Group	<i>N</i>	Var	MSE	Marginal Reliability
3	All	101,516	25.79	6.10	0.94
3	American Indian or Alaska Native	595	23.41	6.10	0.93
3	Asian	3,450	25.32	6.05	0.94
3	Black or African American	18,812	22.58	6.38	0.92
3	Hispanic or Latino	8,191	23.64	6.11	0.93
3	Native Hawaiian or Other Pacific Islander	83	22.90	5.98	0.93
3	Two or More Races	4,647	25.39	6.09	0.94
3	White	65,738	24.78	6.02	0.94
4	All	104,078	26.04	6.26	0.94
4	American Indian or Alaska Native	631	24.60	6.25	0.93
4	Asian	3,502	26.37	6.31	0.94
4	Black or African American	18,506	23.29	6.41	0.92
4	Hispanic or Latino	8,546	23.21	6.21	0.93
4	Native Hawaiian or Other Pacific Islander	102	27.43	6.28	0.95
4	Two or More Races	4,705	25.59	6.26	0.94
4	White	68,086	24.91	6.22	0.94
5	All	108,040	27.12	6.57	0.94
5	American Indian or Alaska Native	718	24.61	6.50	0.93
5	Asian	3,609	26.88	6.78	0.94
5	Black or African American	18,939	23.53	6.54	0.92
5	Hispanic or Latino	8,873	24.96	6.48	0.93
5	Native Hawaiian or Other Pacific Islander	109	27.15	6.61	0.94
5	Two or More Races	4,510	26.70	6.54	0.94
5	White	71,282	26.13	6.58	0.94
6	All	107,916	26.20	6.67	0.93
6	American Indian or Alaska Native	666	24.20	6.71	0.92
6	Asian	3,592	26.17	6.56	0.94
6	Black or African American	18,791	22.68	7.04	0.90
6	Hispanic or Latino	8,628	23.90	6.72	0.92
6	Native Hawaiian or Other Pacific Islander	109	25.75	6.65	0.93
6	Two or More Races	4,411	25.79	6.68	0.93
6	White	71,719	25.43	6.58	0.93
7	All	107,359	26.54	7.17	0.93
7	American Indian or Alaska Native	747	25.52	7.20	0.92

Chapter 12: Reliability and Evidence of Construct-Related Validity

Grade	Group	<i>N</i>	Var	MSE	Marginal Reliability
7	Asian	3,635	26.15	7.12	0.93
7	Black or African American	18,168	23.47	7.44	0.90
7	Hispanic or Latino	8,789	24.69	7.21	0.91
7	Native Hawaiian or Other Pacific Islander	78	25.98	7.21	0.92
7	Two or More Races	3,956	26.57	7.19	0.93
7	White	71,986	25.61	7.09	0.92
8	All	109,775	27.40	7.25	0.93
8	American Indian or Alaska Native	734	25.38	7.28	0.92
8	Asian	3,781	26.95	7.18	0.93
8	Black or African American	17,996	24.09	7.48	0.90
8	Hispanic or Latino	8,283	25.33	7.30	0.92
8	Native Hawaiian or Other Pacific Islander	97	27.38	7.14	0.93
8	Two or More Races	3,910	27.78	7.27	0.93
8	All	74,974	26.68	7.19	0.93

Table 12-4. Marginal Reliability of Total Summative Scores by Ethnic Group—Mathematics

Grade	Group	<i>N</i>	Var	MSE	Marginal Reliability
3	All	101,643	27.04	5.55	0.96
3	American Indian or Alaska Native	595	24.49	5.58	0.95
3	Asian	3,557	26.04	5.38	0.96
3	Black or African American	18,830	24.43	6.02	0.94
3	Hispanic or Latino	8,059	24.52	5.62	0.95
3	Native Hawaiian or Other Pacific Islander	83	25.43	5.60	0.95
3	Two or More Races	4,657	26.51	5.58	0.95
3	White	65,862	25.45	5.41	0.95
4	All	104,125	25.72	5.36	0.95
4	American Indian or Alaska Native	631	23.24	5.44	0.94
4	Asian	3,578	25.79	5.14	0.96
4	Black or African American	18,539	23.21	6.02	0.93
4	Hispanic or Latino	8,414	23.35	5.47	0.94
4	Native Hawaiian or Other Pacific Islander	103	24.81	5.37	0.95
4	Two or More Races	4,700	25.64	5.42	0.95
4	White	68,160	23.84	5.17	0.95
5	All	108,030	26.43	5.91	0.94
5	American Indian or Alaska Native	718	22.39	6.06	0.92
5	Asian	3,670	25.71	5.09	0.96
5	Black or African American	18,941	22.57	7.34	0.88
5	Hispanic or Latino	8,750	23.49	6.24	0.92
5	Native Hawaiian or Other Pacific Islander	111	23.94	5.66	0.94
5	Two or More Races	4,508	26.32	6.10	0.94
5	White	71,332	24.77	5.52	0.95
6	All	107,883	25.85	5.65	0.95
6	American Indian or Alaska Native	668	23.75	5.79	0.94
6	Asian	3,654	26.59	4.97	0.96
6	Black or African American	18,815	22.98	6.76	0.91
6	Hispanic or Latino	8,486	23.71	5.95	0.93
6	Native Hawaiian or Other	109	26.52	5.53	0.95
6	Two or More Races	4,413	25.57	5.76	0.95
6	White	71,738	23.90	5.35	0.95

Chapter 12: Reliability and Evidence of Construct-Related Validity

Grade	Group	<i>N</i>	Var	MSE	Marginal Reliability
7	All	107,318	26.58	6.12	0.94
7	American Indian or Alaska Native	743	23.98	6.39	0.92
7	Asian	3,685	27.54	4.96	0.96
7	Black or African American	18,150	22.77	7.78	0.87
7	Hispanic or Latino	8,693	24.10	6.63	0.91
7	Native Hawaiian or Other Pacific Islander	78	23.38	5.83	0.93
7	Two or More Races	3,947	26.82	6.37	0.93
7	White	72,022	24.91	5.68	0.94
8	All	96,169	27.97	6.72	0.94
8	American Indian or Alaska Native	613	25.01	7.07	0.91
8	Asian	3,537	27.60	5.61	0.96
8	Black or African American	15,440	23.54	8.05	0.87
8	Hispanic or Latino	7,079	25.10	7.21	0.91
8	Native Hawaiian or Other Pacific Islander	91	27.74	6.43	0.94
8	Two or More Races	3,390	28.57	6.96	0.93
8	White	66,019	26.56	6.40	0.94

Table 12-5. Marginal Reliability of Total Summative Scores by Group—ELA

Grade	Group	<i>N</i>	Var	MSE	Marginal Reliability
3	Economically Disadvantaged	56,982	24.00	6.19	0.93
3	LEP	9,903	23.74	6.12	0.93
3	Disabilities	11,273	22.85	6.34	0.92
3	All	101,516	25.79	6.10	0.94
4	Economically Disadvantaged	57,301	24.12	6.29	0.93
4	LEP	9,635	23.35	6.25	0.93
4	Disabilities	11,999	22.99	6.41	0.92
4	All	104,078	26.04	6.26	0.94
5	Economically Disadvantaged	58,437	25.18	6.51	0.93
5	LEP	7,392	22.26	6.46	0.91
5	Disabilities	12,160	22.59	6.58	0.91
5	All	108,040	27.12	6.57	0.94
6	Economically Disadvantaged	57,107	24.16	6.83	0.92
6	LEP	6,638	20.71	7.04	0.88
6	Disabilities	11,689	21.16	7.26	0.88
6	All	107,916	26.20	6.67	0.93
7	Economically Disadvantaged	54,927	24.72	7.29	0.91
7	LEP	6,240	21.41	7.49	0.87
7	Disabilities	11,597	22.02	7.69	0.87
7	All	107,359	26.54	7.17	0.93
8	Economically Disadvantaged	54,089	25.26	7.37	0.91
8	LEP	6,139	21.52	7.54	0.87
8	Disabilities	11,543	21.30	7.73	0.87
8	All	109,775	27.40	7.25	0.93

Table 12-6. Marginal Reliability of Total Summative Scores by Group—Mathematics

Grade	Group	N	Var	MSE	Marginal Reliability
3	Economically Disadvantaged	56,998	25.36	5.73	0.95
3	LEP	9,963	25.78	5.57	0.95
3	Disabilities	11,357	28.15	6.23	0.95
3	All	101,643	27.04	5.55	0.96
4	Economically Disadvantaged	57,295	23.98	5.61	0.94
4	LEP	9,660	24.80	5.49	0.95
4	Disabilities	12,076	25.77	6.22	0.94
4	All	104,125	25.72	5.36	0.95
5	Economically Disadvantaged	58,393	24.42	6.51	0.92
5	LEP	7,416	23.38	6.60	0.91
5	Disabilities	12,199	24.64	7.75	0.89
5	All	108,030	26.43	5.91	0.94
6	Economically Disadvantaged	57,034	24.23	6.17	0.93
6	LEP	6,604	23.56	6.46	0.92
6	Disabilities	11,676	24.29	7.28	0.90
6	All	107,883	25.85	5.65	0.95
7	Economically Disadvantaged	54,868	24.33	6.92	0.91
7	LEP	6,221	23.69	7.47	0.89
7	Disabilities	11,598	23.12	8.68	0.84
7	All	107,318	26.58	6.12	0.94
8	Economically Disadvantaged	46,350	25.14	7.47	0.90
8	LEP	5,276	24.07	7.84	0.89
8	Disabilities	9,839	21.98	8.78	0.83
8	All	96,169	27.97	6.72	0.94

In addition to the SEM, the CSEM expresses the degree of measurement error in scale-score units and are conditioned on the ability of the student. The CSEM is reported in support of AERA, APA, & NCME (2014) Standard 2.14, which states the following:

When possible and appropriate, conditional standard errors of measurement should be reported at several score levels unless there is evidence that the standard error is constant across score levels. Where cut scores are specified for selection or classification, the standard errors of measurement should be reported in the vicinity of each cut score. (p. 46)

The CSEMs are defined as the reciprocal of the square root of the test information function and can be estimated across all points of the ability continuum (Hambleton & Swaminathan, 1985); therefore, Equations 12.3 and 12.4 are used to calculate the CSEM.

In further compliance with Standard 2.14, Table 12-7 shows the median CSEM near the achievement level cut scores for ELA and mathematics.

Table 12-7. . Conditional Standard Errors of Measurement near (± 10 Points) Achievement Level Cut Scores, ELA and Mathematics

Content Area	Level—Cut Score	Grade	<i>N</i>	Median	Standard Deviation
ELA	1—2	3	13,260	5.98	0.20
ELA	2—3	3	13,293	5.53	0.50
ELA	3—4	3	12,943	5.85	0.36
ELA	1—2	4	13,866	6.04	0.21
ELA	2—3	4	12,966	6.01	0.10
ELA	3—4	4	13,475	6.01	0.08
ELA	1—2	5	13,762	6.07	0.26
ELA	2—3	5	12,121	6.19	0.40
ELA	3—4	5	12,833	6.93	0.27
ELA	1—2	6	14,059	6.64	0.49
ELA	2—3	6	14,019	6.09	0.28
ELA	3—4	6	10,316	6.19	0.39
ELA	1—2	7	13,690	6.94	0.36
ELA	2—3	7	14,341	6.73	0.45
ELA	3—4	7	9,855	6.98	0.14
ELA	1—2	8	13,999	7.02	0.19
ELA	2—3	8	13,200	6.91	0.29
ELA	3—4	8	9,461	7.00	0.09
Mathematics	1—2	3	13,351	5.53	0.50
Mathematics	2—3	3	14,470	5.05	0.23
Mathematics	3—4	3	11,738	5.01	0.10
Mathematics	1—2	4	13,793	5.40	0.49
Mathematics	2—3	4	15,434	4.97	0.19
Mathematics	3—4	4	11,723	4.60	0.49
Mathematics	1—2	5	14,229	5.83	0.44
Mathematics	2—3	5	14,342	4.52	0.50
Mathematics	3—4	5	12,032	4.17	0.38
Mathematics	1—2	6	16,142	5.82	0.41
Mathematics	2—3	6	15,742	4.91	0.29
Mathematics	3—4	6	12,404	4.21	0.41
Mathematics	1—2	7	15,797	6.13	0.37
Mathematics	2—3	7	13,448	4.84	0.38

Content Area	Level—Cut Score	Grade	<i>N</i>	Median	Standard Deviation
Mathematics	3—4	7	12,396	4.02	0.15
Mathematics	1—2	8	13,715	7.03	0.35
Mathematics	2—3	8	10,682	5.70	0.47
Mathematics	3—4	8	10,947	5.00	0.14

When using a CAT, the CSEM will vary for the same scale score; therefore, it is necessary to report averages. Table 12-8 presents the overall average CSEM and the average CSEM by scale-score decile for ELA and mathematics.

Table 12-8. Overall Average CSEM and Average CSEM by Decile, ELA and Mathematics

Content Area	Grade	Overall SEM	Decile 1	Decile 2	Decile 3	Decile 4	Decile 5	Decile 6	Decile 7	Decile 8	Decile 9	Decile 10
ELA	3	6.10	7.66	6.34	6.02	5.94	5.70	5.52	5.55	5.78	5.98	6.49
ELA	4	6.26	7.29	6.16	6.05	6.04	6.02	6.01	6.00	6.01	6.03	6.85
ELA	5	6.57	7.33	6.28	6.08	6.07	6.12	6.22	6.48	6.82	6.97	7.28
ELA	6	6.67	8.47	7.14	6.88	6.51	6.25	6.11	6.05	6.06	6.16	7.00
ELA	7	7.17	8.92	7.50	7.04	6.85	6.75	6.73	6.71	6.79	6.96	7.38
ELA	8	7.25	8.82	7.42	7.05	7.00	6.97	6.91	6.92	6.96	6.99	7.35
Mathematics	3	5.55	7.67	6.02	5.69	5.29	5.11	5.04	5.01	5.01	5.01	5.57
Mathematics	4	5.36	7.92	5.98	5.35	5.07	5.02	4.99	4.80	4.56	4.63	5.15
Mathematics	5	5.91	10.29	7.81	6.82	6.00	5.19	4.96	4.47	4.14	4.21	4.78
Mathematics	6	5.65	9.24	6.93	6.15	5.84	5.19	5.01	4.91	4.48	4.16	4.26
Mathematics	7	6.12	11.79	8.11	6.92	6.15	5.74	5.07	4.80	4.13	4.02	4.05
Mathematics	8	6.72	11.01	8.35	7.53	7.11	6.81	6.13	5.47	5.03	4.85	4.59

Figures 12-1 through 12-12 display the CSEM curves by grade and content area. The dashed vertical lines represent the cut scores. The CSEM tends to be higher at the ends of the scale-score range. The measurement error increases when there are few items at a particular ability level. The figures show that the CSEM tends to be at its minimum around cut scores between Levels 2 and 3 and Levels 3 and 4.

Figure 12-1. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 3 English Language Arts

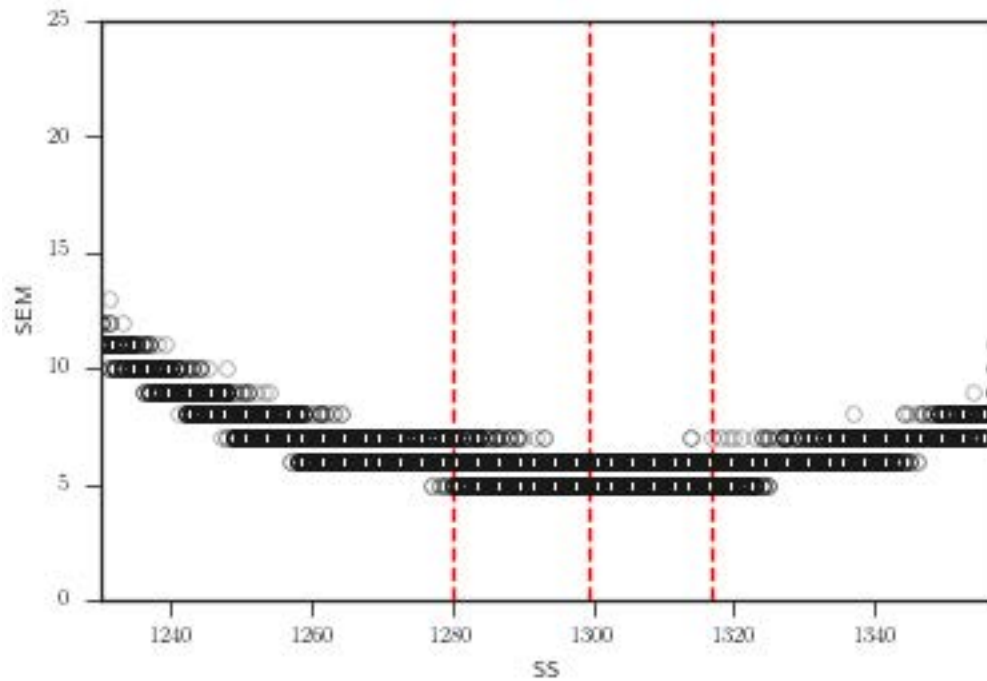


Figure 12-2. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 4 English Language Arts

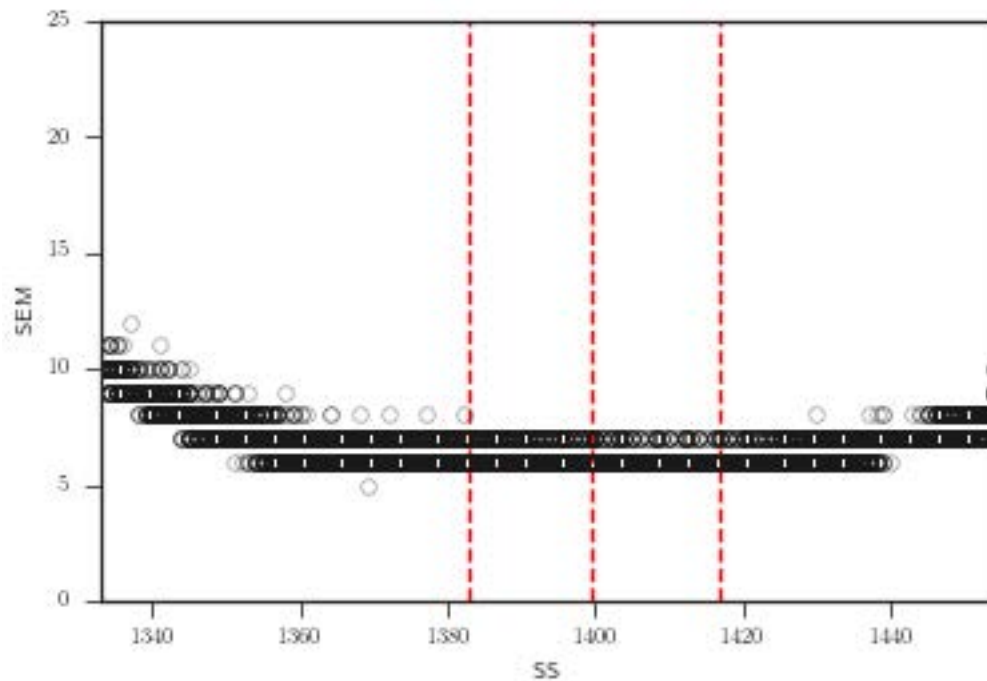


Figure 12-3. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 5 English Language Arts

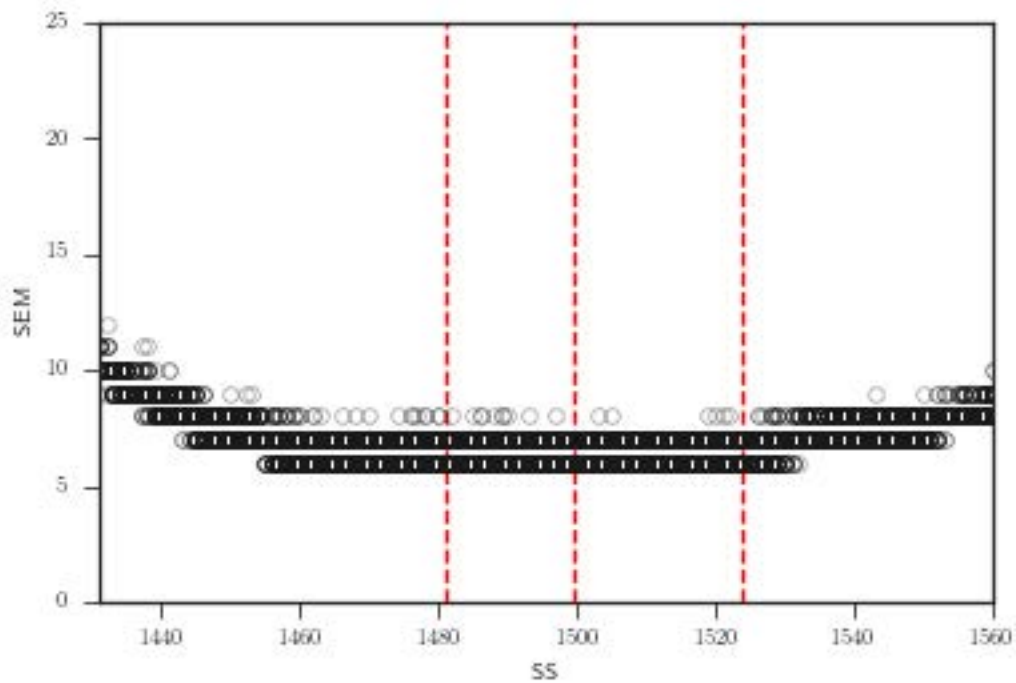


Figure 12-4. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 6 English Language Arts

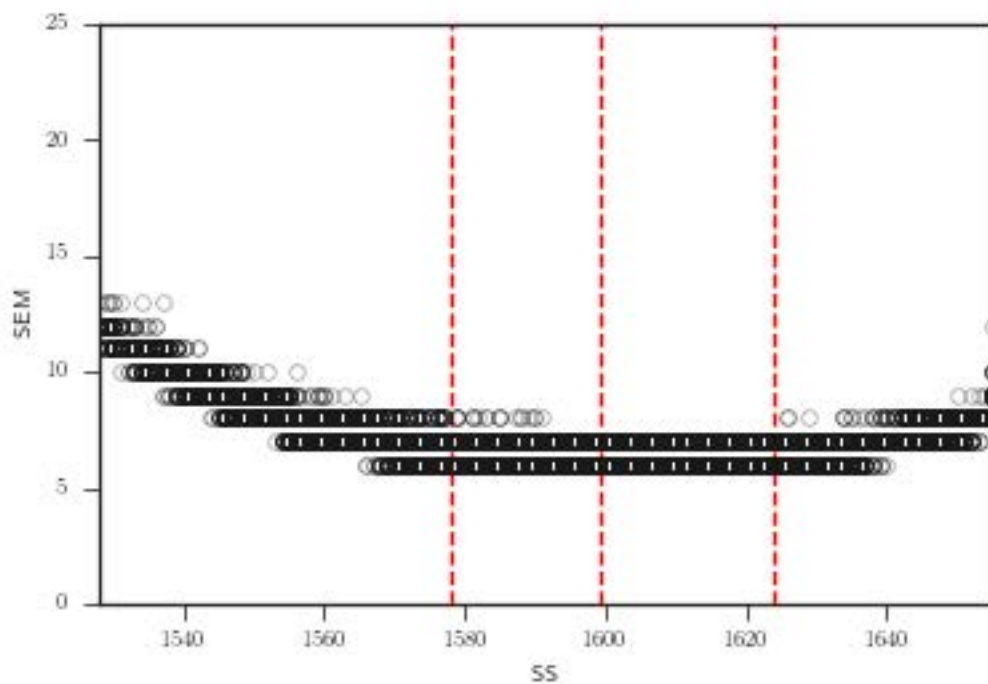


Figure 12-5. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 7 English Language Arts

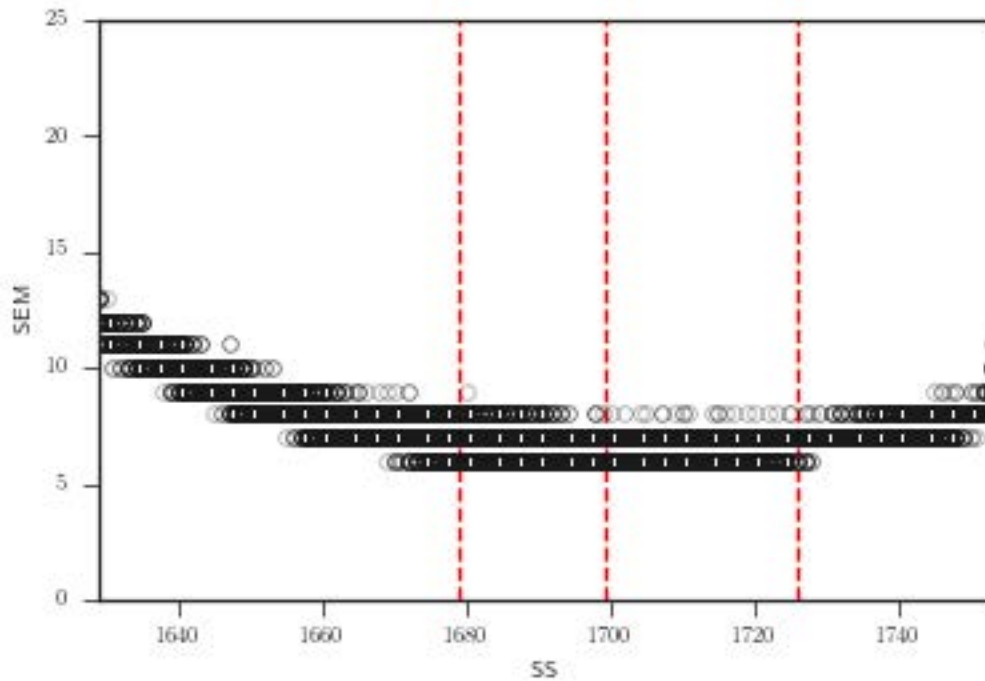


Figure 12-6. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 8 English Language Arts

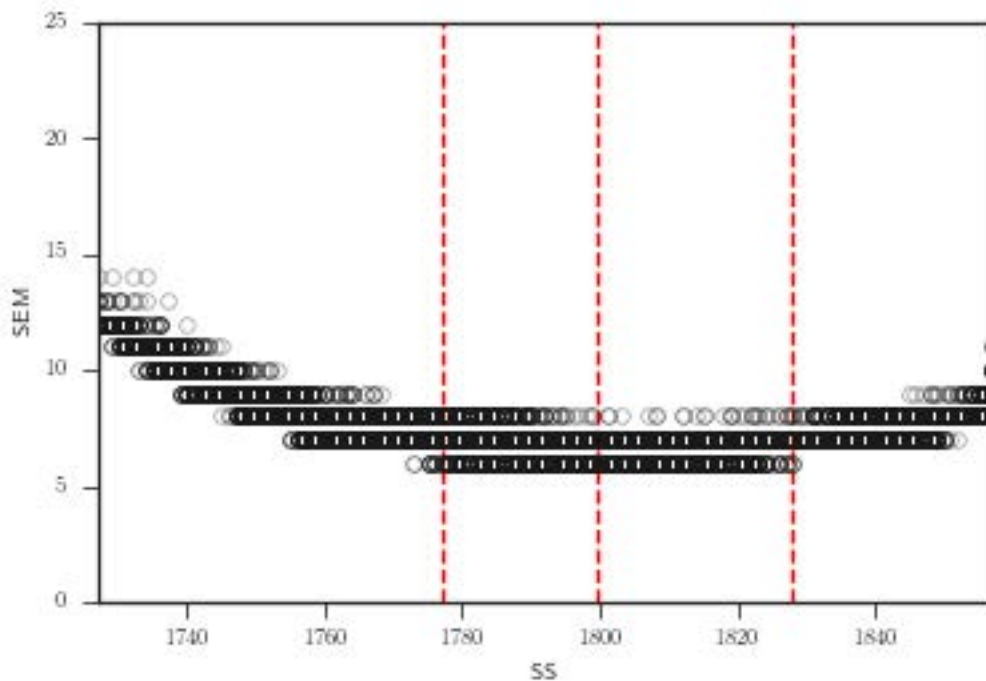


Figure 12-7. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 3 Mathematics

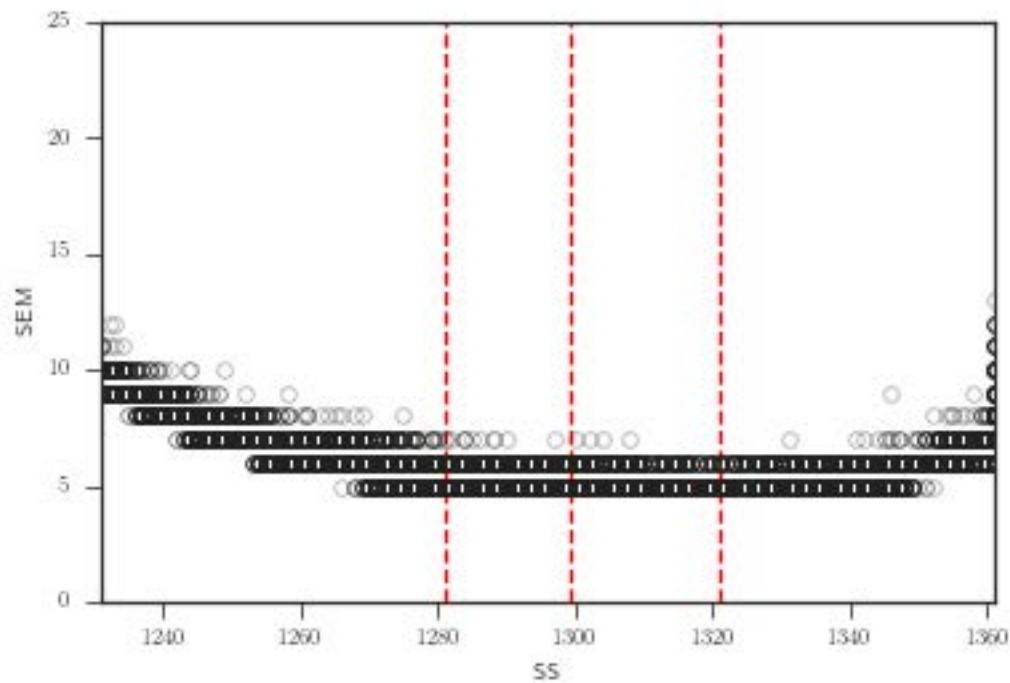


Figure 12-8. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 4 Mathematics

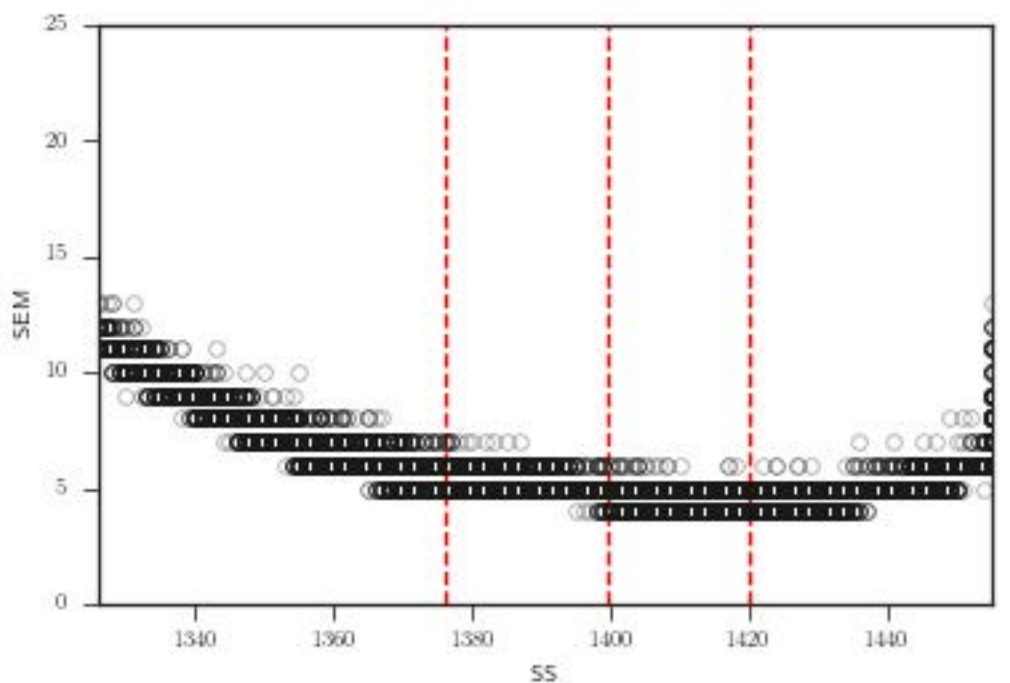


Figure 12-9. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 5 Mathematics

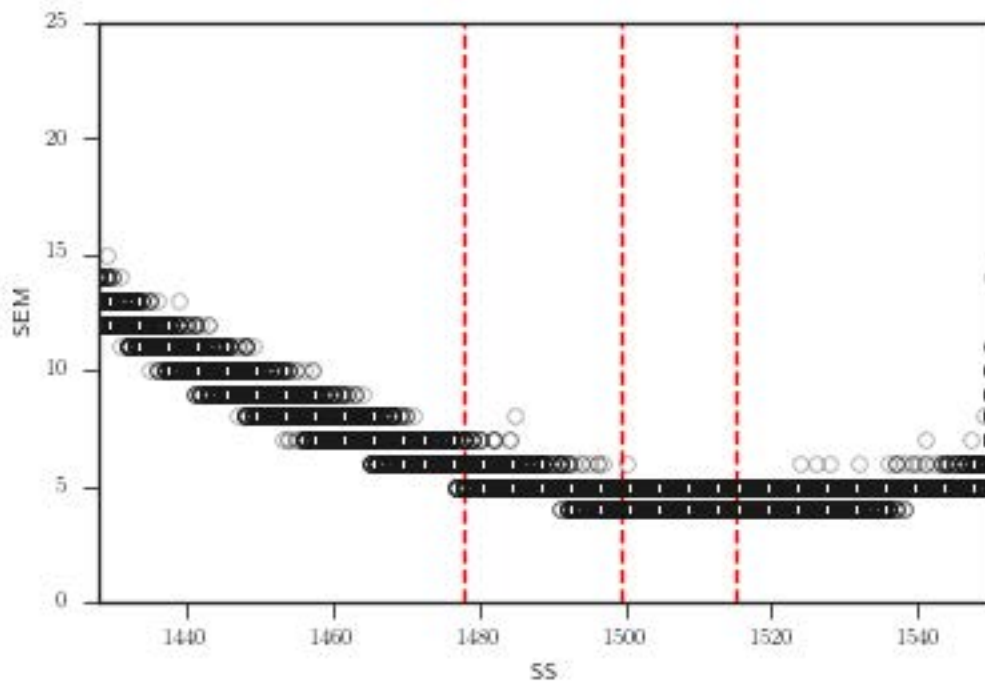


Figure 12-10. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 6 Mathematics

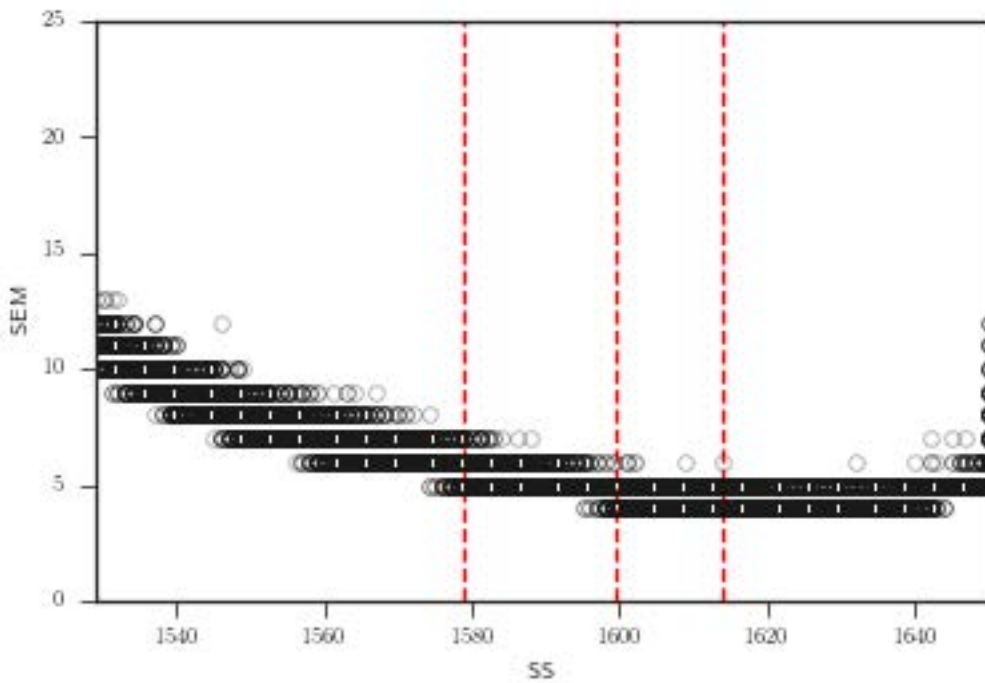


Figure 12-11. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 7 Mathematics

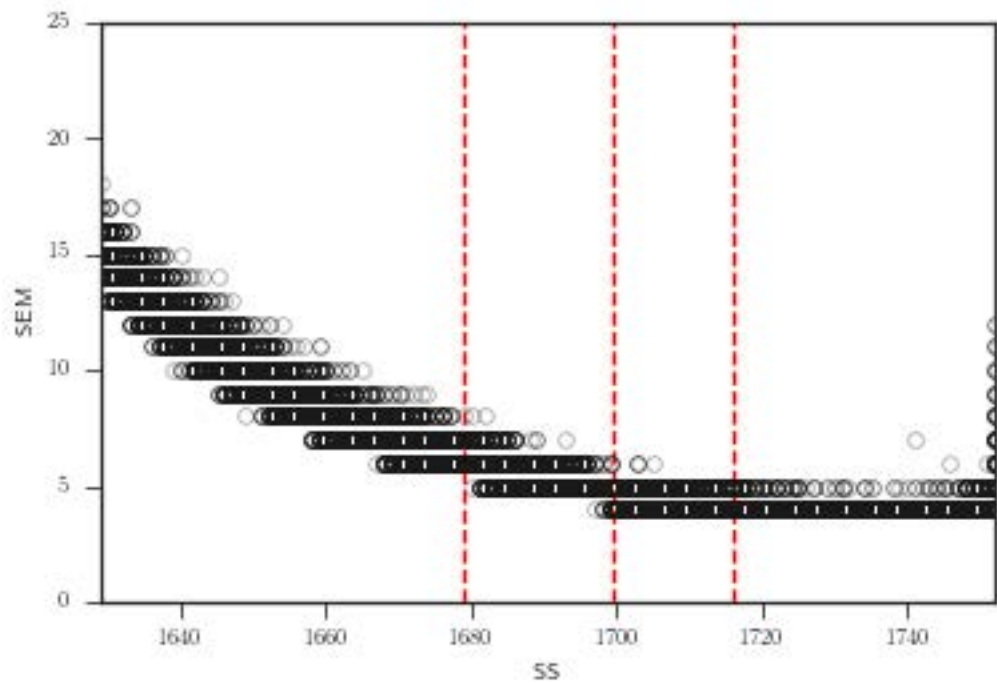
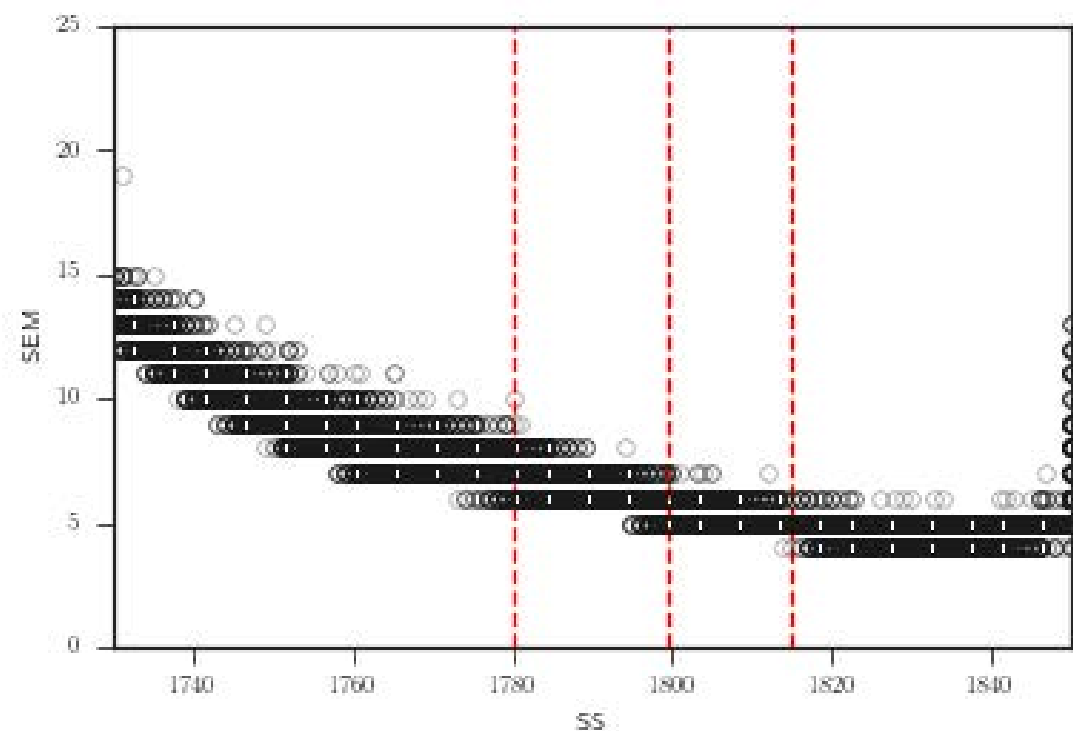


Figure 12-12. Conditional Standard Errors of Measurement for Overall Scale Scores, Grade 8 Mathematics



Smarter Balanced supports using fixed-form paper/pencil tests in schools that lack computer capacity or administering them to address potential religious concerns associated with using technology for assessments. Since the paper/pencil tests for ELA and mathematics consist of Smarter Balanced items, and there are few students who take the paper/pencil forms, DRC has chosen to be consistent with ELA and mathematics online counterparts by calculating the marginal reliability for those forms using Equation 12.6.

Table 12-9 shows the marginal reliability for the paper/pencil forms. As expected, overall estimated reliability coefficients are high and in the acceptable range for a large-scale, high-stakes test.

Table 12-9. Fixed-Form Marginal Reliability: ELA and Mathematics

Content Area	Grade	Number of Items	Reliability	SEM
ELA	3	50	0.89	8.38
ELA	4	49	0.88	8.79
ELA	5	52	0.88	9.31
ELA	6	55	0.88	9.09
ELA	7	50	0.86	8.86
ELA	8	56	0.89	9.51
Mathematics	3	66	0.88	8.84
Mathematics	4	65	0.90	7.62
Mathematics	5	64	0.85	8.83
Mathematics	6	65	0.84	9.45
Mathematics	7	65	0.87	9.65
Mathematics	8	62	0.85	9.43

12.1.6 Reliability of Claims for ELA and Mathematics

Scale-score summary statistics (i.e., mean and standard deviation), marginal reliability coefficients, and mean CSEM were computed for each of the claims by grade and content area using M-STEP data. These statistics are presented in Tables 12-10 and 12-11 for ELA and mathematics, respectively. Reliability indices are a function of the number of test items. As expected, reliability coefficients are lower for a claim assessed by a small number of items compared to a claim assessed by a larger number of items. Consequently, the reliability for claims with larger CSEMs is lower than those with smaller CSEMs. These CSEMs are reported in the scale-score metric.

Table 12-10. Reliability, Mean, Standard Deviation, and Conditional Standard Error of Measurement (CSEM) of ELA Claims

Grade	Claim No.	Claim	Student N Count	Number of Items	Mean	Std. Dev.	Reliability	Mean CSEM
3	1	Reading	101,514	15–16	1295.23	27.19	0.84	10.83
3	2	Writing	101,514	13	1290.77	28.99	0.83	12.11
3	3	Listening	101,514	8–9	1297.49	36.24	0.61	22.59
3	4	Research	101,514	8	1291.25	33.79	0.75	16.74
4	1	Reading	104,078	15–16	1394.53	28.12	0.83	11.76
4	2	Writing	104,078	13	1392.78	28.52	0.83	11.64
4	3	Listening	104,078	8–9	1399.47	33.62	0.62	20.85
4	4	Research	104,078	8	1392.01	34.48	0.74	17.45
5	1	Reading	105,578	15–16	1497.06	28.97	0.82	12.23
5	2	Writing	105,578	13	1493.60	30.67	0.83	12.69
5	3	Listening	105,578	8–9	1500.39	34.98	0.63	21.18
5	4	Research	105,578	8	1494.23	32.32	0.78	15.28
6	1	Reading	105,714	15–16	1591.83	29.39	0.81	12.75
6	2	Writing	105,714	13	1589.60	29.66	0.80	13.38
6	3	Listening	105,714	8–9	1596.31	33.82	0.62	20.84
6	4	Research	105,714	8	1590.61	32.54	0.71	17.54
7	1	Reading	107,359	15–16	1695.54	27.83	0.80	12.40
7	2	Writing	107,359	13	1688.88	31.29	0.79	14.26
7	3	Listening	107,359	8–9	1697.42	33.77	0.64	20.40
7	4	Research	107,359	8	1692.83	32.97	0.70	17.96
8	1	Reading	109,774	15–16	1793.55	28.47	0.78	13.47
8	2	Writing	109,774	13	1790.66	31.41	0.81	13.85
8	3	Listening	109,774	8–9	1797.74	35.56	0.63	21.72
8	4	Research	109,774	8	1792.46	32.72	0.77	15.62

Table 12-11. Reliability, Mean, Standard Deviation, and Conditional Standard Error of Measurement (CSEM) of Mathematics Claims

Grade	Claim No.	Claim	Student N Count	Number of Items	Mean	Std. Dev.	Reliability	CSEM
3	1	Concepts and Procedures	101,643	20	1296.34	28.85	0.93	7.47
3	3	Communicating Reasoning	101,643	8	1294.79	32.55	0.71	17.51
3	2 & 4	Problem Solving and Modeling and Data Analysis	101,643	8	1292.23	32.30	0.64	19.42
4	1	Concepts and Procedures	104,125	20	1393.53	26.98	0.93	6.89
4	3	Communicating Reasoning	104,125	8	1391.35	32.06	0.69	17.71
4	2 & 4	Problem Solving and Modeling and Data Analysis	104,125	8	1390.56	31.82	0.69	17.59
5	1	Concepts and Procedures	108,030	20	1487.12	28.13	0.91	8.43
5	3	Communicating Reasoning	108,030	8	1484.36	32.18	0.70	17.54
5	2 & 4	Problem Solving and Modeling and Data Analysis	108,030	8	1481.44	34.90	0.53	23.89
6	1	Concepts and Procedures	107,883	20	1588.24	26.82	0.93	7.23
6	3	Communicating Reasoning	107,883	8	1584.28	31.23	0.61	19.54
6	2 & 4	Problem Solving and Modeling and Data Analysis	107,883	8	1580.40	35.62	0.47	26.00
7	1	Concepts and Procedures	107,318	20	1688.35	27.89	0.91	8.19
7	3	Communicating Reasoning	107,318	8	1683.37	32.33	0.51	22.66
7	2 & 4	Problem Solving and Modeling and Data Analysis	107,318	8	1680.06	36.20	0.45	26.78
8	1	Concepts and Procedures	96,169	20	1788.52	28.94	0.92	8.26
8	3	Communicating Reasoning	96,169	8	1786.56	33.55	0.51	23.54
8	2 & 4	Problem Solving and Modeling and Data Analysis	96,169	8	1783.82	36.00	0.51	25.22

12.1.7 Reliability, SEM, and CSEM for Social Studies

Table 12-12 provides information on reliability (coefficient alpha) (see Equation 12.1) and SEM (see Equation 12.2) from the classical true score theory for social studies by grade and form, despite the fact that all OP items are the same for forms 1 through 3 per grade. This choice was made for two reasons: (1) conceptually, it makes more sense to report test-level results by form because each form represents one test, and (2) it shows variations of statistics across forms (if there are any) to inform related decisions (i.e., whether to combine online forms per grade for social studies) when computing classification accuracy and consistency.

As shown in Table 12-12, the values of coefficient alpha across forms and grades for social studies range from 0.82 to 0.89. Therefore, in general, based on coefficient alpha, M-STEP social studies tests have an acceptable degree of internal consistency. Moreover, very similar statistics across the three online forms per grade for social studies are observed. This supports the later decision to combine all three online forms per grade when examining classification accuracy and consistency for social studies.

Table 12-12. Test-Level Descriptive Statistics by Form: Social Studies Reliability and Standard Error of Measurement

Grade	N OP Items	Form	N	Reliability	SEM
5	45	1	36,009	0.85	3.06
5	45	2	36,069	0.85	3.06
5	45	3	36,062	0.85	3.07
5	45	4	936	0.82	3.07
8	44	1	36,570	0.88	3.01
8	44	2	36,599	0.88	3.01
8	44	3	36,557	0.88	3.01
8	44	4	653	0.83	3.05
11	38	1	35,012	0.89	2.70
11	38	2	35,062	0.89	2.70
11	38	3	35,171	0.89	2.72
11	38	4	728	0.88	2.76

Additionally, the CSEM was calculated for social studies. Related numerical information can be found in corresponding conversion tables reported in Chapter 8 (i.e., Table 8-11). Graphical representations can be found in Figures 12-13 to 12-15. According to these graphs, the CSEMs are not the lowest at the proficient cut scores (i.e., the vertical line, which indicates the cut between Level 2 and Level 3). However, the ability ranges from -2 to 2 in all graphs appear to have low SE. Note, however, that these graphs are made using the post-administration estimated item parameters.

Figure 12-13. Test (Conditional) Standard Error for Social Studies Grade 5 by Form

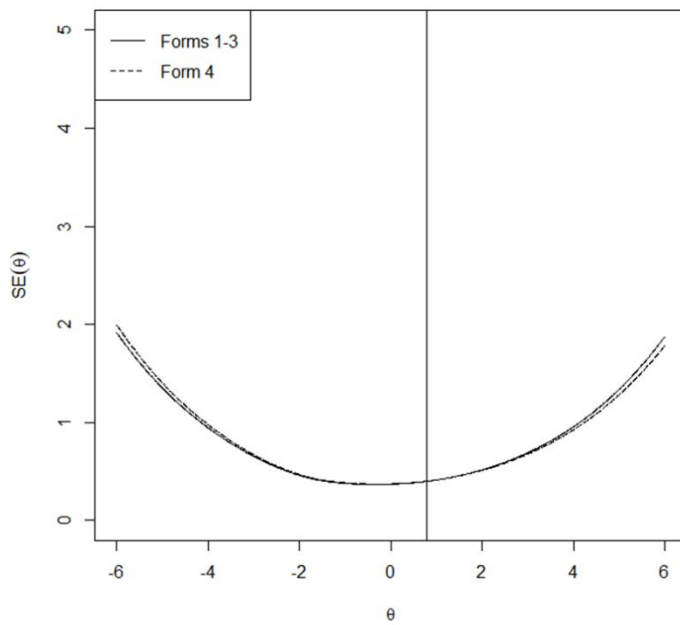


Figure 12-14. Test (Conditional) Standard Error for Social Studies Grade 8 by Form

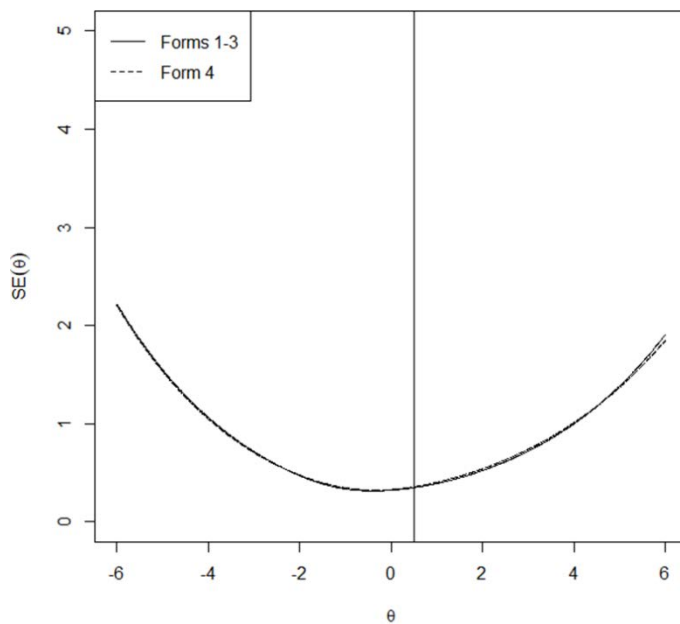
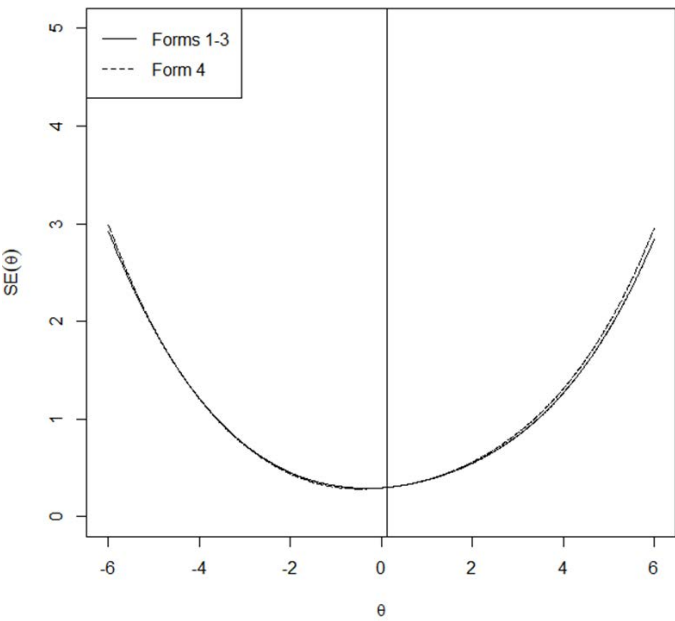


Figure 12-15. Test (Conditional) Standard Error for Social Studies Grade 11 by Form



12.2 Classification Accuracy and Consistency

Based on M-STEP scale scores, student performance in corresponding content areas is classified into one of the four performance levels (i.e., Advanced, Proficient, Partially Proficient, and Not Proficient). Among these, the most important classification is between the Proficient and Partially Proficient categories (i.e., the proficient or not cut). While it is always important to know the reliability of student scores in any examination, it is also important to assess the quality of the decisions, especially regarding the proficient or “not cut” categories. Such evaluation was performed through estimation of the probabilities of accurate and consistent classification of student performance.

Classification accuracy is defined as the extent to which the actual classifications of examinees agree with classifications that would be made on the basis of their true scores (Livingston & Lewis, 1995). It is common to estimate classification accuracy by utilizing a psychometric model to find true scores corresponding to observed scores. The magnitude of classification accuracy measures is influenced by key features of the test design, including the number of items, the number of cut scores, reliability, and associated SEM or CSEM.

12.2.1 ELA and Mathematics

To calculate classification accuracy for each student in ELA and mathematics, the calculations used by Smarter Balanced (see the Smarter Balanced 2017-18 Technical Report, 2018). For each student, the likelihood of scoring in each performance level is calculated. The student likelihoods are used to calculate the accuracy by level and the overall accuracy.

Tables 12-13 through 12-15 provide the classification accuracy for ELA and mathematics. The overall classification accuracy ranges from 0.83 to 0.86, and the accuracy by performance level ranges from 0.70 to 0.92. These results suggest that accurate performance level classifications for ELA and mathematics are being made for students in Michigan based on M-STEP. Note that any inconsistencies between the expected values and accuracy by level or overall accuracy are due to computation rounding error.

Table 12-13. Overall Classification Accuracy: ELA and Mathematics

Content Area	Grade	<i>N</i>	Overall Accuracy
ELA	3	102,490	0.83
ELA	4	105,087	0.83
ELA	5	109,063	0.83
ELA	6	108,850	0.83
ELA	7	108,100	0.82
ELA	8	110,536	0.83
Mathematics	3	102,827	0.85
Mathematics	4	105,347	0.85
Mathematics	5	109,263	0.86
Mathematics	6	109,005	0.86
Mathematics	7	108,214	0.85
Mathematics	8	110,576	0.83

Table 12-14. Classification Accuracy: ELA

Grade	Assigned Level	<i>N</i>	Observed Proportion	Expected Level 1	Expected Level 2	Expected Level 3	Expected Level 4	Accuracy by Level
3	1	31,766	0.31	0.28	0.03	0.00	0.00	0.92
3	2	25,213	0.25	0.03	0.19	0.03	0.00	0.76
3	3	22,929	0.22	0.00	0.03	0.17	0.03	0.76
3	4	22,582	0.22	0.00	0.00	0.03	0.19	0.88
4	1	35,591	0.34	0.31	0.03	0.00	0.00	0.92
4	2	22,131	0.21	0.03	0.15	0.03	0.00	0.70
4	3	22,629	0.22	0.00	0.03	0.16	0.03	0.74
4	4	24,736	0.24	0.00	0.00	0.03	0.21	0.87
5	1	35,117	0.32	0.30	0.03	0.00	0.00	0.92
5	2	23,174	0.21	0.03	0.16	0.03	0.00	0.73
5	3	31,299	0.29	0.00	0.03	0.23	0.03	0.80
5	4	19,473	0.18	0.00	0.00	0.03	0.15	0.84
6	1	34,180	0.31	0.28	0.03	0.00	0.00	0.91
6	2	29,624	0.27	0.03	0.21	0.03	0.00	0.76
6	3	30,659	0.28	0.00	0.03	0.23	0.02	0.81
6	4	14,387	0.13	0.00	0.00	0.02	0.11	0.85
7	1	31,594	0.29	0.26	0.03	0.00	0.00	0.90
7	2	29,543	0.27	0.04	0.20	0.04	0.00	0.74
7	3	33,134	0.31	0.00	0.03	0.25	0.02	0.81

Grade	Assigned Level	N	Observed Proportion	Expected Level 1	Expected Level 2	Expected Level 3	Expected Level 4	Accuracy by Level
7	4	13,829	0.13	0.00	0.00	0.02	0.11	0.83
8	1	33,081	0.30	0.27	0.03	0.00	0.00	0.90
8	2	30,125	0.27	0.04	0.20	0.03	0.00	0.75
8	3	34,356	0.31	0.00	0.03	0.26	0.02	0.82
8	4	12,974	0.12	0.00	0.00	0.02	0.10	0.82

Table 12-15. Classification Accuracy: Mathematics

Grade	Assigned Level	N	Observed Proportion	Expected Level 1	Expected Level 2	Expected Level 3	Expected Level 4	Accuracy by Level
3	1	28,703	0.28	0.26	0.02	0.00	0.00	0.92
3	2	27,112	0.26	0.03	0.20	0.03	0.00	0.77
3	3	27,875	0.27	0.00	0.03	0.22	0.02	0.82
3	4	19,137	0.19	0.00	0.00	0.02	0.16	0.88
4	1	26,083	0.25	0.22	0.02	0.00	0.00	0.91
4	2	35,039	0.33	0.03	0.27	0.03	0.00	0.82
4	3	27,026	0.26	0.00	0.03	0.21	0.02	0.82
4	4	17,199	0.16	0.00	0.00	0.02	0.14	0.88
5	1	40,652	0.37	0.34	0.03	0.00	0.00	0.92
5	2	31,105	0.28	0.03	0.23	0.03	0.00	0.81
5	3	19,401	0.18	0.00	0.02	0.14	0.02	0.79
5	4	18,105	0.17	0.00	0.00	0.02	0.15	0.89
6	1	37,634	0.35	0.32	0.03	0.00	0.00	0.92
6	2	33,684	0.31	0.03	0.25	0.03	0.00	0.80
6	3	20,265	0.19	0.00	0.03	0.14	0.02	0.76
6	4	17,422	0.16	0.00	0.00	0.02	0.14	0.89
7	1	39,285	0.36	0.33	0.03	0.00	0.00	0.91
7	2	30,250	0.28	0.03	0.22	0.03	0.00	0.78
7	3	21,017	0.19	0.00	0.02	0.16	0.02	0.81
7	4	17,662	0.16	0.00	0.00	0.02	0.15	0.89
8	1	44,580	0.40	0.37	0.04	0.00	0.00	0.91
8	2	28,794	0.26	0.04	0.19	0.03	0.00	0.73
8	3	17,449	0.16	0.00	0.02	0.12	0.02	0.74
8	4	19,753	0.18	0.00	0.00	0.02	0.16	0.89

12.2.2 Social Studies

For social studies, each test under consideration consists of equally weighted and dichotomously scored items only, and procedures from Hanson and Brennan (1990) were applied to derive classification accuracy and classification consistency measures. Moreover, the definitions for accuracy and consistency of decisions presented in Young and Yoon (1998) were adopted. Specifically, the *accuracy* of decisions is the extent to which decisions would agree with those made if each student could somehow be tested with all possible forms of an examination; and the *consistency* of decisions is the extent to which decisions would agree with those made if each student had taken a parallel form of the examination, equal in difficulty and covering the same content as the form the student actually took (Young & Yoon, 1998). These ideas are shown schematically in Figures 12-16 and 12-17 using M-STEP social studies as an example. In both figures, “Achieves Proficient Status” refers to the proficient and above category on the total raw score, and “Does Not Achieve Proficient Status” refers to all categories below the proficient cut.

Figure 12-16. Classification Accuracy

		Decision made on a form actually taken	Decision made on a form actually taken
		Does Not Achieve Proficient Status	Achieves Proficient Status
“True status” based on all-forms average	Does Not Achieve Proficient Status	Correct Classification	Misclassification
	Achieves Proficient Status	Misclassification	Correct Classification

Note. Adapted from Young and Yoon (1998).

Figure 12-17. Classification Consistency

		Decision made on the 2nd form taken	Decision made on the 2nd form taken
		Does Not Achieve Proficient Status	Achieves Proficient Status
Decision made on the 1st form taken	Does Not Achieve Proficient Status	Consistent Classification	Inconsistent Classification
	Achieves Proficient Status	Inconsistent Classification	Consistent Classification

Note. Adapted from Young and Yoon (1998).

In Figure 12-16, accurate classification occurs when the decision made on the basis of the form actually taken agrees with the decision made on the basis of the theoretical “all-forms” average. Misclassification occurs, for example, when a student who “Does Not Achieve Proficient Status” based on his or her “all-forms” average is classified incorrectly as “Achieves Proficient Status.”

Consistent classification occurs (see Figure 12-17) when two possible alternate forms agree on the classification of a student as either “Achieves Proficient Status” or “Does Not Achieve Proficient Status,” whereas inconsistent classification occurs when the decisions made by the forms differ.

The analyses made use of the techniques outlined and implemented by Hanson and Brennan (1990) and Brennan (2004). Specifically, a four-parameter beta distribution was used to model the true score, and Lord’s (1965) two-term approximation to the compound binomial distribution was used to model the conditional error. The BB-CLASS software (Version 1.1) was used to complete these analyses (Brennan, 2004).

Table 12-16 presents the analysis results of decision accuracy and consistency for classifying students at each grade level per test form as “Achieves Proficient Status” or “Does Not Achieve Proficient Status” based on their M-STEP social studies total raw scores. As mentioned above, the three online forms for social studies were combined (see Table 12-16) due to the fact that all OP items are exactly the same across these forms and the raw score statistics are very similar across forms (see Table 8-6).

In addition to classification accuracy and consistency, Table 12-16 provides information on the proportion of false positives and false negatives (i.e., the two types of misclassification). The false positive is the type of misclassification in which students should be classified in the “Does Not Achieve Proficient Status” category based on their “all-forms” average but end up in the “Achieves Proficient Status” category based on the actual form. The false negative is just the opposite: students who should be in the “Achieves Proficient Status” category based on their “all-forms” average end up in the “Does Not Achieve Proficient Status” category based on the actual form. The sum of the proportion values for accuracy, false positives, and false negatives should be equal to 1.00. Due to rounding, however, the sum of these values in the table may not be equal to 1.00.

As shown in Table 12-16, the proportion of false positives (i.e., labeling a student as proficient when he or she should be categorized as not proficient) ranged from 0.03 to 0.06 for social studies. Moreover, the proportion of false negatives (i.e., labeling a student as not proficient when he or she should be categorized as proficient) ranged from 0.01 to 0.04 for social studies.

The last column in Table 12-16 reports the proportion of students predicted by the model that would be assigned to the same category (i.e., either proficient or not proficient) if an alternate form of M-STEP social studies assessments (with similar content coverage and item difficulty as the actual form) had been administered. These values range from 0.87 to 0.95.

Table 12-16. Decision Accuracy and Consistency on M-STEP Social Studies Total Raw Score by Grade and Form

Grade	Form	Accuracy	False Positive	False Negative	Consistency
5	1–3	0.93	0.05	0.02	0.90
5	4	0.97	0.03	0.01	0.95
8	1–3	0.92	0.05	0.03	0.89
8	4	0.95	0.03	0.01	0.94
11	1–3	0.91	0.06	0.04	0.87
11	4	0.91	0.06	0.03	0.88

12.3 Assumption of Unidimensionality

Another measure of construct validity is unidimensionality. One of the underlying assumptions of the IRT models used to scale M-STEP content area tests is that the items being calibrated are unidimensional; that is, items comprising M-STEP in each grade/content area measure a single construct. For example, mathematics items should measure mathematics ability and not reading skills. Standard 1.13 of the AERA, APA, & NCME (2014) *Standards* states the following:

If the rationale for a test score interpretation for a given use depends on premises about the relationships among test items or among parts of the test, evidence concerning the internal structure of the test should be provided. (p. 26–27)

12.3.1 ELA and Mathematics

Smarter Balanced examined the unidimensionality for the Smarter Balanced/M-STEP ELA and mathematics assessments. Based on the findings of the dimensionality study, Smarter Balanced found that the use of the unidimensional item response theory (IRT) model and test design was appropriate. A detailed discussion and the results of the dimensionality study can be found in the online [Smarter Balanced 2013–2014 Technical Report \(2016\)](https://portal.smarterbalanced.org/library/en/2013-14-technical-report.pdf).¹

12.3.2 Social Studies

For M-STEP social studies, MDE conducted two analyses to evaluate the unidimensionality assumption with OP items only. The first set was an exploratory factor analysis (EFA) using the Mplus software with the WLSMV² estimator. Barendse, Oort, and Timmerman (2015) found that WLSMV is the preferred estimation method and is recommended to rely on the Root Mean Squared Error of Approximation (RMSEA) index (in which values less than 0.05 are desired) if the primary interest is in major factors. The second set of analyses is a principle component analysis (PCA) using *MATLAB* (2018). For PCA results, the magnitude of the first and second eigenvalues are examined. Both the eigenvalues-greater-than-one rule and the scree plot approach are considered. The RMSEA values for one-factor EFA models and the first two eigenvalues from

¹ <https://portal.smarterbalanced.org/library/en/2013-14-technical-report.pdf>

² “WLSMV-weighted least square parameter estimates using a diagonal weight matrix with standard errors and mean- and variance-adjusted chi-square test statistic that use a full weight matrix” (Muthén and Muthén, 2012, p. 603)

each PCA model are reported in Table 12-17.

As shown in Table 12-17, the dimensionality assessment for social studies is evaluated by administration mode at each grade level.³ Both the EFA and PCA results failed to reject the unidimensionality assumption, which is a supporting piece of evidence for the use of unidimensional IRT models at each grade for social studies.

Table 12-17. RMSEA from 1-Factor EFA and the First Two Eigenvalues from PCA

Content Area	Grade	Form	RMSEA (1-Factor EFA)	PCA First Eigenvalue	PCA Second Eigenvalue
Social Studies	5	1–3	0.013	1.4607	0.2675
Social Studies	5	4	0.011	1.2058	0.3183
Social Studies	8	1–3	0.016	1.7482	0.3006
Social Studies	8	4	0.014	1.3544	0.3461
Social Studies	11	1–3	0.016	1.7725	0.2870
Social Studies	11	4	0.018	1.7578	0.3458

12.4 Validity Evidence

The *Standards for Educational and Psychological Testing* defines validity as “the degree to which evidence and theory support the interpretations of test scores for proposed uses of tests. Validity is, therefore, the most fundamental consideration in developing tests and evaluating tests” (AERA, APA, & NCME, 2014, p. 11). The purpose of test score validation is not to validate the test itself but to validate interpretations of the test scores for particular purposes or uses. Test score validation is not a quantifiable property but an ongoing process, beginning at initial conceptualization and continuing throughout the entire assessment process. Every aspect of an assessment provides evidence that either supports or challenges its validity, including design, content specifications, item development, psychometric quality, and inferences made from the results.

The validity of score interpretations for M-STEP is supported by multiple sources of evidence. Chapter 1 of the *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014) specifies the following sources of validity evidence that are important to gather and document in order to support validity claims for an assessment:

- Test content
- Response processes
- Internal test structure
- Relation to other variables
- Consequences of test use

³ Note that for each grade, forms 1–3 are online forms and form 4 is a paper/pencil form. All OP items are the same across forms 1–3 for social studies at each grade. Form 4, however, has somewhat different OP items from the online forms because technology-enhanced items cannot be put on a paper/pencil form.

It is important to note that these categories are not mutually exclusive. One source of validity evidence often falls into more than one category, as discussed in more detail in this section. The process of gathering evidence of the validity of score interpretations is best characterized as ongoing throughout test development, administration, scoring, reporting, and beyond. As the technical report has progressed, it has covered the different phases of the testing cycle. Each part of the technical report has detailed the procedures and processes applied in Michigan and the corresponding results. Each part has also highlighted the meaning and significance of the procedures, processes, and results in terms of validity and their relationship to specific sections of the *Standards*. The current section now addresses these final issues in validity: test content, response processes, internal test structure, relation to other variables, and consequences of test use.

12.4.1 Minimization of Construct-Irrelevant Variance and Construct Underrepresentation

Minimization of construct-irrelevant variance and construct underrepresentation is addressed in the following steps of the test development process:

1) specification, 2) item writing, 3) review, 4) field-testing, 5) test construction, and 6) item calibration (see Chapter 3 for more information on steps 1 through 5 and Chapter 8 for more information on calibration).

Construct-irrelevant variance refers to error variance that is caused by factors unrelated to the constructs measured by the test. For example, when tests are not administered under standardized conditions (e.g., one administration may be timed, but another administration may not be timed), differences in student performance may be partially associated with the different administration conditions. Careful specification of content and review of the items representing that content are the first steps in minimizing construct-irrelevant variance. Then, empirical evidence, especially item-level data, is used to infer construct irrelevance. For additional details with respect to ELA and mathematics, please see the *Smarter Balanced 2017–2018 Technical Report* (2018).

Construct underrepresentation occurs when the content of the assessment does not reflect the full range of content that the assessment is expected to cover. Specification and review, in which test blueprints are developed and reviewed, are primary steps in the development process and are designed to ensure that content is appropriately represented.

12.4.2 Evidence Based on Test Content

According to the *Standards*, evidence based on test content “can include logical or empirical analyses of the adequacy with which the test content represents the content domain and of the relevance of the content domain to the proposed interpretation of test scores” (AERA, APA, & NCME, 2014, p. 14). Documentation of the content domains, how the content is sampled and represented, and alignment of items to the content were discussed in Chapter 3. The documentation showed how test specification documents derived from earlier developmental activities guided the final phases of test development and ultimately yielded the test forms that were administered to students.

Chapter 3 also showed that the participation of Michigan educators in that process provided a solid rationale for having confidence in the content and design of Michigan M-STEP as a tool from which to derive valid inferences about Michigan student performance. Particularly for social studies, use of classroom teachers brought into the process the enacted curriculum perspective and the written curriculum perspective. The test development process and the involvement of Michigan educators in that process formed an important part of the validity of the entire Michigan M-STEP assessment.

12.4.3 Evidence Based on Response Processes

According to the *Standards*, evidence based on response processes “generally comes from analyses of individual responses” (AERA, APA, & NCME, 2014, p. 15). Hence, the best opportunity for detecting and eliminating potential sources of invalidity occurs during the test development process (U.S. Department of Education, 2015). As indicated in Chapter 3, all items for M-STEP were carefully reviewed through multiple cycles of the item development process for issues related to ambiguity, bias, sensitivity, irrelevance, and inaccuracy to ensure a fit between the construct and the nature of the actual performance.

12.4.4 Evidence Based on Internal Test Structure

According to the *Standards*, evidence based on internal structure reflects “the degree to which the relationships among test items and test components conform to the construct on which the proposed test score interpretations are based” (AERA, APA, & NCME, 2014, p. 16). Three important sources of internal structure evidence have been addressed within this technical document: measurement invariance, dimensionality, and reliability. Evidence of measurement invariance is provided in Chapter 11 by using DIF. Moreover, Appendix F provides support for measurement invariance when discussing the mode comparability for social studies. Additional support for measurement invariance can be found in Section 12.2.5, which reports the subgroup reliability estimates. The dimensionality investigation mentioned in Section 12.3 also provides supporting evidence of the internal test structure.

12.4.5 Evidence Based on Relations to Other Variables

Convergent validity is a subtype of construct validity that can be estimated by the extent to which measures of constructs that theoretically should be related to each other are, in fact, observed as related to each other. Analyses of the internal structure of a test can indicate the extent to which the relationships among test items conform to the construct the test purports to measure. For example, M-STEP mathematics test is designed to measure a single overall construct—mathematics achievement. Therefore, the items comprising the M-STEP mathematics test should measure only mathematics.

For M-STEP assessments,⁴ this technical report summarizes additional statistics that contribute to construct validity, reliability—as reported previously in this chapter and Chapter 8—and item fit. The internal consistency coefficient reported above is a measure of item homogeneity. For a group of items to be homogeneous, they must measure the same construct (construct validity)

⁴ For ELA and mathematics, not all psychometric characteristics are provided in this report. Additional details can be found in the Smarter Balanced Technical Reports (2016 & 2017).

or represent the same content domain (content validity). Because IRT models were used to calibrate test items and to report student scores, item fit is also relevant to construct validity. The extent to which test items function as the IRT model prescribes is relevant to the validation of test scores. Additional evidence to support construct validity is examined by the correlations between the claim scores for ELA and mathematics in the next section.

12.4.6 Correlations among Claims as Evidence of Convergent Validity

In this section, the strength of the interrelationships among the claims are reported by computing the correlations between them. Two types of correlations are reported here: the uncorrected Pearson product-moment (PPM) correlation coefficients and the PPM corrected for attenuation (CAPPM).

AERA, APA, & NCME (2014) Standard 1.21 states the following:

When statistical adjustments, such as those for restriction of range or attenuation, are made, both adjusted and unadjusted coefficients, as well as the specific procedure used, and all statistics used in the adjustment, should be reported. Estimates of the construct-criterion relationship that remove the effects of measurement error on the test should be clearly reported as adjusted estimates. (p. 29)

We can correct for the attenuation of the PPM statistically using Spearman's formula:

$$CAPPM = \frac{r_{xy}}{\sqrt{r_{xx}r_{yy}}} \quad (12.9)$$

where r_{xy} is the PPM between two claims, r_{xx} is the reliability of one of those claims, and r_{yy} is the reliability for the other claim.

Tables 12-18 and 12-19 report the PPM and CAPPM described above. The PPM among the claim scores is presented below the diagonal portion of the matrix, and the CAPPM is presented above the diagonal portion of the matrix in each table.

The uncorrected PPM in Tables 12-18 and 12-19 should be interpreted in the context of the reliability coefficient. In general, it is expected to see lower PPM coefficients between variables that are less reliable. In most cases, the PPM coefficients show that performance on one claim is moderately related to performance on another claim within the same grade and content area. In cases where there is a limited number of items per claim, caution should be used when comparing the PPM coefficients measuring the relationships between claims to those measuring the relationships between content areas.). We expect to see a more modest relationship (smaller correlation coefficients) reported between the claims as a consequence of the lower number of items measuring each of the reporting categories. The PPM between two claim scores may be artificially low because of measurement error.

Across all tables, the CAPPM indicates strong relationships between the claims. In some cases, the CAPPM is greater than 1.00. "Disattenuated values greater than 1.00 indicate that measurement error is not randomly distributed" (Schumacker, 1996). The strong relationships suggested by the CAPPM in Tables 12-18 and 12-19 are further evidence of the validity of the

test construct. Since the overall content area is composed of the claim scores, and the content area is expected to measure a single dimension, it is expected that these claim scores are also highly related.

Table 12-18. Uncorrected Correlation Coefficient (below Diagonal) and Corrected Correlation Coefficient (above Diagonal) among Claims: English Language Arts

Grade	Claim No.	Claim	Number of Items	1	2	3	4
3	1	Reading	15–16		0.94	0.97	0.91
3	2	Writing	13	0.78		0.93	0.91
3	3	Listening	8–9	0.69	0.66		0.91
3	4	Research	8	0.72	0.71	0.62	
4	1	Reading	15–16		0.92	0.94	0.91
4	2	Writing	13	0.77		0.90	0.89
4	3	Listening	8–9	0.67	0.65		0.90
4	4	Research	8	0.71	0.70	0.61	
5	1	Reading	15–16		0.93	0.94	0.92
5	2	Writing	13	0.77		0.91	0.91
5	3	Listening	8–9	0.68	0.66		0.91
5	4	Research	8	0.74	0.73	0.64	
6	1	Reading	15–16		0.93	0.96	0.94
6	2	Writing	13	0.75		0.93	0.93
6	3	Listening	8–9	0.68	0.66		0.94
6	4	Research	8	0.71	0.70	0.62	
7	1	Reading	15–16		0.94	0.97	0.95
7	2	Writing	13	0.75		0.92	0.93
7	3	Listening	8–9	0.69	0.66		0.95
7	4	Research	8	0.72	0.70	0.63	
8	1	Reading	15–16		0.94	0.97	0.96
8	2	Writing	13	0.74		0.94	0.95
8	3	Listening	8–9	0.67	0.67		0.96
8	4	Research	8	0.74	0.75	0.67	

Table 12-19. Uncorrected Correlation Coefficient (below Diagonal) and Corrected Correlation Coefficient (above Diagonal) among Claims: Mathematics

Grade	Claim No.	Claim	Number of Items	1	3	2 & 4
3	1	Concepts and Procedures	20		0.90	0.95
3	3	Communicating Reasoning	8	0.73		0.95
3	2 & 4	Problem Solving and Modeling and Data Analysis	8	0.73	0.64	
4	1	Concepts and Procedures	20		0.91	0.91
4	3	Communicating Reasoning	8	0.73		0.94
4	2 & 4	Problem Solving and Modeling and Data Analysis	8	0.73	0.65	
5	1	Concepts and Procedures	20		0.94	1.01
5	3	Communicating Reasoning	8	0.75		1.03
5	2 & 4	Problem Solving and Modeling and Data Analysis	8	0.70	0.63	
6	1	Concepts and Procedures	20		0.94	1.10
6	3	Communicating Reasoning	8	0.70		1.09
6	2 & 4	Problem Solving and Modeling and Data Analysis	8	0.73	0.58	
7	1	Concepts and Procedures	20		1.00	1.07
7	3	Communicating Reasoning	8	0.68		1.13
7	2 & 4	Problem Solving and Modeling and Data Analysis	8	0.69	0.54	
8	1	Concepts and Procedures	20		1.04	1.01
8	3	Communicating Reasoning	8	0.71		1.09
8	2 & 4	Problem Solving and Modeling and Data Analysis	8	0.69	0.56	

12.4.7 Divergent (Discriminant) Validity

Measures of different constructs should not be highly correlated with each other. Divergent validity is a subtype of construct validity that can be assessed by the extent to which measures of constructs that theoretically should not be related to each other are, in fact, observed as not related to each other. Typically, correlation coefficients among measures of unrelated or distantly related constructs are examined in support of divergent validity.

To assess the divergent validity of M-STEP, correlations were computed between the ELA and mathematics scale scores for students who took both assessments. These correlation results are shown in Table 12-20. The correlation coefficients ranged from 0.77 (between ELA and mathematics in grade 5) to 0.79 (between ELA and mathematics in grades 6, 7, and 8). The correlation coefficients suggest that individual student scores for ELA and mathematics are highly related. Despite high correlations, the tests are not perfectly related to each other, suggesting that different constructs are being tapped; however, the test scores do appear as highly related to one another, suggesting they may be tapping into a similar knowledge base or general underlying ability.

Table 12-20. Inter-correlation of ELA and Mathematics Scale Scores

Grade	Inter-Correlation
3	0.78
4	0.78
5	0.77
6	0.79
7	0.79
8	0.79

12.4.8 Evaluation of Item Exposure for CAT ELA and Mathematics

Controlling item exposure is of concern with CAT administrations, which impacts the validity of the interpretation of the test scores. Overexposed items could be a threat to validity because students may become familiar with the items over time and, thus, decrease the difficulty of the item, which would impact the ability estimate (Georgiadou, Triantafillou, & Economides, 2007). Item exposure rates were obtained using all completed, online, adaptive tests for which item data were available. The exposure rate for a given item is the proportion of tests (in the grade and content area) on which the item appeared.

Table 12-21 presents a summary of the item exposure results for ELA and mathematics. Within each grade, the table presents the number of items in the OP pool (N) and various descriptive statistics, including the mean, standard deviation (SD), range (Min, Max), and median of the observed exposure rates. Table 12-21 shows that, on average, the same item appeared in 5% of the grade 3 tests; in other words, 5% of grade 3 examinees saw the same item. As a rule of thumb, Smarter Balanced attempts to maintain a maximum exposure rate of 25% (meaning that 25% of examinees will see the same item). Table 12-21 shows that the mean and median exposure rates for ELA and mathematics CAT items are well below 25%.

Table 12-21. Summary of ELA Item Exposure Rates by Grade and Component

Content Area	Grade	<i>N</i>	Mean	SD	Min	Max	Median
ELA	3	865	0.05	0.09	0.00	0.47	0.01
ELA	4	827	0.05	0.09	0.00	0.52	0.01
ELA	5	788	0.06	0.09	0.00	0.49	0.02
ELA	6	750	0.06	0.10	0.00	0.68	0.01
ELA	7	655	0.07	0.11	0.00	0.61	0.02
ELA	8	735	0.06	0.10	0.00	0.48	0.01
Mathematics	3	1242	0.03	0.04	0.00	0.21	0.01
Mathematics	4	1275	0.03	0.04	0.00	0.21	0.01
Mathematics	5	1205	0.03	0.04	0.00	0.20	0.01
Mathematics	6	1127	0.03	0.05	0.00	0.23	0.01
Mathematics	7	1010	0.04	0.06	0.00	0.27	0.00
Mathematics	8	895	0.04	0.06	0.00	0.26	0.01

Table 12-22 provides further information about the exposure rates by showing the number of items in the OP pool (*N*) and proportion of items with exposure rates falling into certain ranges (bins with a width of 20%), including those that were completely unexposed (Unused). The majority of CAT items, for both ELA and mathematics, had item exposure rates between 0% and 20%.

There were a handful of items in ELA with higher-than-desirable exposure rates. This occurred when there were few items measuring elements in the blueprint. There were also items in both content areas that were unused. There is a trade-off between blueprint fidelity and exposure, with the adaptive CAT engine weighting blueprint fidelity more heavily. In addition, for ELA, it was requested to use all or almost all items with a passage so students were not given numerous passages to read to meet the blueprint.

Table 12-22. Percentage of CAT Items by Exposure Rate

Content Area	Grade	Total Number of Items	Unused*	0%–20%	21%–40%	41%–60%	61%–80%	81%–100%
ELA	3	865	10.06	92.37	6.47	1.16	0.00	0.00
ELA	4	827	6.17	91.05	8.34	0.60	0.00	0.00
ELA	5	789	8.37	91.89	7.22	0.89	0.00	0.00
ELA	6	750	8.40	89.47	8.93	1.33	0.27	0.00
ELA	7	655	10.08	87.79	10.08	1.98	0.15	0.00
ELA	8	735	9.25	88.84	9.93	1.22	0.00	0.00
Mathematics	3	1242	11.84	99.92	0.08	0.00	0.00	0.00
Mathematics	4	1275	7.06	99.69	0.31	0.00	0.00	0.00
Mathematics	5	1205	7.22	100.00	0.00	0.00	0.00	0.00
Mathematics	6	1127	8.25	99.65	0.35	0.00	0.00	0.00
Mathematics	7	1011	4.06	98.12	1.88	0.00	0.00	0.00
Mathematics	8	895	10.39	96.76	3.24	0.00	0.00	0.00

*Note: "Unused" is also included in the 0% to 20% range.

12.4.9 Evidence Based on Consequences of Test Use

The *Standards* incorporate the intended and unintended consequences of test use into the concept of validity. It indicates that information about the consequences of testing does not in and of itself detract from the validity of intended test interpretations (AERA, APA, & NCME, 2014, p. 19). Rather, according to the *Standards*, a more searching inquiry into the sources of those consequences given the intended purposes of an assessment is a basis for evaluating the quality of the validity evidence. The test data alone do not provide sufficient verification of this type of evidence. For this reason, it is not straightforward to measure and collect evidence on the consequential aspects of validity.

To address the intended consequences of M-STEP, the purposes of M-STEP must be specified. MDE has carefully articulated the intended purposes of M STEP as driving features of the selection of Smarter Balanced items, the development of social studies tests, and the implementation of the testing program. The specific purposes associated with M-STEP include the following:

- M-STEP accurately describes both student achievement (i.e., how much students know at the end of the year) and student growth (i.e., how much students have improved since the previous year) to inform program evaluation and school-, district-, and state-accountability systems and to provide valid, reliable, and fair measures of students' progress toward, and attainment of, the knowledge and skills required to be college- and career-ready.
- M-STEP informs state and federal accountability.
- M-STEP assessments are fair for all students, including those with disabilities or limited English proficiency, at all levels of achievement.

12.5 Summary

In summary, Chapter 12 of this report demonstrates M-STEP's adherence to the AERA, APA, & NCME (2014) *Standards* regarding reliability and construct-related validity. The analyses described above address multiple best practices of the testing industry, particularly the following standards:

- Standard 2.0—Appropriate evidence of reliability/precision should be provided for the interpretation for each intended score use.
- Standard 2.1—The range of replications over which reliability/precision is being evaluated should be clearly stated, along with a rationale for the choice of this definition, given the testing situation.
- Standard 2.3—For each total score, subscore, or combination of scores that is to be interpreted, estimates of relevant indices of reliability/precision should be reported.
- Standard 2.13—The standard error of measurement, both overall and conditional (if reported), should be provided in units of each reported score.
- Standard 2.14—When possible and appropriate, conditional standard errors of measurement should be reported at several score levels unless there is evidence that the standard error is constant across score levels. Where cut scores are specified for selection or classification, the standard errors of measurement should be reported in the vicinity of each cut score.
- Standard 2.16—When a test or combination of measures is used to make classification decisions, estimates should be provided of the percentage of test takers who would be classified in the same way on two replications of the procedure.
- Standard 2.19—Each method of quantifying the reliability/precision of scores should be described clearly and expressed in terms of statistics appropriate to the method. The sampling procedures used to select test takers for reliability/precision analyses and the descriptive statistics on these samples, subject to privacy obligations where applicable, should be reported.
- Standard 4.3—Test developers should document the rationale and supporting evidence for the administration, scoring, and reporting rules used in computer-adaptive, multistage-adaptive, or other tests delivered using computer algorithms to select items. This documentation should include procedures used in selecting items or sets of items for administration, in determining the starting point and termination conditions for the test, in scoring the test, and in controlling item exposure.

References

- American Institutes for Research (2016). Smarter Balanced scoring specification: Summative and Interim Assessments: ELA/Literacy Grades 3-8;11 and Mathematics Grade 3-8;11. Los Angeles, CA: Smarter Balanced Assessment Consortium.
- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for educational and psychological testing*. Washington, DC: American Educational Research Association.
- Barendse, M. T., Oort, F. J., & Timmerman, M. E. (2015). Using exploratory factor analysis to determine the dimensionality of discrete responses. *Structural Equation Modeling*, 22(1), 87–101.
- Birnbaum, A. (1968). *Some latent trait models and their use in inferring an examinee's ability*. In F. M. Lord & M. R. Novick (Eds.), *Statistical Theories of Mental Test Scores* (pp. 374–472). Reading, MA: Addison-Wesley Publishing.
- Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: An application of an EM algorithm. *Psychometrika*, 46, 443–459.
- Brennan, R. L. (2004). *Manual for BB-CLASS: A computer program that uses the beta-binomial model for classification consistency and accuracy (Version 1.1)*. Iowa City, IA: Center for Advanced Studies in Measurement and Assessment.
- Cai, L. (2017). flexMIRT (Version 3.51) [Computer software]. Chapel Hill, NC: Vector Psychometric Group.
- Cai, L., & Hansen, M. (2013). Limited-information goodness-of-fit testing of hierarchical item factor models. *British Journal of Mathematical and Statistical Psychology*, 66, 245–276.
- Camilli, G. (2006). Test fairness. In R. L. Brennan (Ed.), *Educational Measurement* (4th ed., pp. 221–256). Westport, CT: American Council on Education and Praeger.
- Camilli, G., & Shepard, L. A. (1994). *Methods for identifying biased test items*. Thousand Oaks, CA: SAGE Publications.
- Candell, G. L. & Drasgow, F. (1988). An iterative procedure for linking metrics bias in item response theory. *Applied Psychological Measurement*, 12(3), 253–260.
- Cattell, R. B. (1966). *The scree test for the number of factors*. *Multivariate Behavioral Research*, 1, 245–276.
- Cizek, G. J. & Bunch, M. B. (2007). *Standard setting: A guide to establishing and evaluating performance standards on tests*. Thousand Oaks, CA: Sage.
- Crocker, L., & Algina, J. (1986). *Introduction to classical and modern test theory*. Belmont, CA: Wadsworth.

References

- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16, 297–334.
- Cronbach, L. J. (1971). Test validation. In R. L. Thorndike (Ed.), *Educational measurement* (2nd ed.). Washington, DC: American Council on Education.
- Cudeck, R., & Henly, S. J. (1991). Model selection in covariance structures analysis and the “problem” of sample size: A clarification. *Psychological Bulletin*, 109(3), 512–519.
- Darling-Hammond, L., & Pecheone, R. (2010). Developing an Internationally Comparable Balanced Assessment System that Supports High-Quality Learning. Retrieved from <http://www.k12center.org/publications.html>.
- Data Recognition Corporation (2017). *Technology User Guide*. Maple Grove, MN: Author.
- Dorans, N. J., & Schmitt, M. P. (1991). *Constructed response and differential item functioning: A pragmatic approach*. Princeton: Educational Testing Service.
- Embretson, S., & Reise, S. (2000). *Item response theory for psychologists*. Mahwah, NJ: Lawrence Erlbaum Associates.
- ETS. (2012). Smarter Balanced Assessment Consortium: Bias and sensitivity guidelines. Princeton, NJ: ETS.
- Georgiadou, E., Triantafillou, E., Economides, A. (2007). A review of item exposure control strategies for computerized adaptive testing developed from 1983 to 2005. *Journal of Technology, Learning, and Assessment*, 5(8).
- Green, D. R. (1975). *Procedures for assessing bias in achievement tests*. Paper presented at the National Institute of Education Conference on Test Bias, Annapolis, MD.
- Haertel, E. H. (2006). Reliability. In R. L. Brennan (Ed.), *Educational Measurement* (4th ed., pp. 65–110). Westport, CT: American Council on Education and Praeger.
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Hingham, MA: Kluwer-Nijhoff Publishing.
- Hansen, E.G. & Mislevy, R.J. (2008). *Design Patterns for Improving Accessibility for Test Takers With Disabilities*. Princeton, NJ, ETS Research Report No. RR-08-49.
- Hanson, B. A., & Brennan, R. L. (1990). An investigation of classification consistency indexes estimated under alternative strong true score models. *Journal of Educational Measurement*, 27(4), 345–359.
- Hays, W. L. (1994). *Statistics* (5th ed.). Belmont, CA: Wadsworth.
- Kolen, M. J., & Brennan, R. L. (2004). *Test equating, scaling, and linking: Methods and practices* (2nd ed.). New York, NY: Springer-Verlag.

References

- Leung, C.-K., Chang, H.-H., & Hau, K.-T. (2003). Computerized adaptive testing: A comparison of three content balancing methods. *Journal of Technology, Learning, and Assessment*, 2(5). Available from <http://www.jtla.org>.
- Lewis D. M., Mitzel, H. C., & Green, D. R. (1996). Standard setting: A bookmark approach. In D. R. Green (Chair), IRT-based standard-setting procedures utilizing behavioral anchoring. Symposium conducted at the 1996 Council of Chief State School Officers 1996 National Conference on Large Scale Assessment, Phoenix, AZ.
- Lewis, D. M., Mitzel, H. C., Mercado, R. L., & Schultz, E. M. (2012). *The bookmark standard setting procedure*. In G. J. Cizek (Ed.), *Setting performance standards: Foundations, methods, and innovations*. New York, NY: Routledge.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32, 179–197.
- Lord, F. M. (1965). A strong true score theory, with applications. *Psychometrika*, 30, 239–270.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theories of mental test scores*. Reading, MA: Addison-Wesley.
- Mantel, N. (1963) Chi-Square Tests with One Degree of Freedom: Extensions of the Mantel-Haenszel Procedure. *Journal of the American Statistical Association*, 58, 690-700.
- Mantel, N., & Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute*, 22, 719–748.
- MATLAB and Statistics Toolbox Release 2018b, The MathWorks, Inc., Natick, Massachusetts, United States.
- McDonald, R. P. (1999). *Test theory: A unified treatment*. Hillsdale, NJ: Erlbaum.
- Michigan Department of Education (2016). *2017 Scribing Protocol for the M-STEP, MI-Access, SAT, ACT, and WIDA Assessments*. Retrieved from https://www.michigan.gov/documents/mde/M-STEP_Scribing_Protocol_477116_7.pdf
- Michigan Department of Education (2016). *2016–2017 Guide to State Assessments*. Retrieved from https://www.michigan.gov/documents/mde/Guide_to_State_Assessments_622260_7.pdf
- Michigan Department of Education (2016). *Arabic Read-Aloud Guidelines: M-STEP Mathematics, Spring 2017*. Retrieved from https://www.michigan.gov/documents/mde/Arabic_Read-Aloud_Guidelines_M-STEP_Mathematics_536798_7.pdf
- Michigan Department of Education (2016). *Assessment Integrity Guide*. Retrieved from https://www.michigan.gov/documents/mde/Assessment_Integrity_Guide_291950_7.pdf
- Michigan Department of Education (2017). *M-STEP Guide to Reports*. Retrieved from https://www.michigan.gov/documents/mde/2017_M-STEP_GTR_598970_7.pdf

References

- Michigan Department of Education (2016). *M-STEP, MI-Access, SAT, ACT WorkKeys, and WIDA Student Supports and Accommodations Tables*. Retrieved from https://www.michigan.gov/documents/mde/M-STEP_Supports_and_Accommodations_Table_477120_7.pdf
- Michigan Department of Education (2016). *Read-Aloud Guidelines: M-STEP Mathematics and English Language Arts, Spring 2017*. Retrieved from https://www.michigan.gov/documents/mde/Math_and_ELA_Read-Aloud_Guidelines_512366_7.pdf
- Michigan Department of Education (2016). *Spanish Read-Aloud Guidelines: M-STEP Mathematics, Spring 2017*. Retrieved from https://www.michigan.gov/documents/mde/Math_Read-Aloud_Spanish_Guidelines_512131_7.pdf
- Michigan Department of Education and Michigan Virtual University (2016). *MDE Assessment Security*. Retrieved from <http://bit.ly/MDEAssessmentSecurity>
- Mislevy, R. J., & Haertel, G. D. (2006). Implications of evidence-centered design for educational testing. *Educational Measurement: Issues and Practice*, 25(4), 6–20.
- Mislevy, R. J., Steinberg, L. S., & Almond, R. G. (2003). On the structure of educational assessments. *Measurement: Interdisciplinary Research and Perspectives*, 1(1), 3–62.
- Muraki, E. (1992). A generalized partial credit model: Application of an EM algorithm. *Applied Psychological Measurement*, 16, 159–176.
- Muraki, E., & Bock, R. D. (2003). PARSCALE 4 for Windows: IRT based test scoring and item analysis for graded items and rating scales [Computer software].
- Muthén, B. O., & Muthén, L. K. (2012). *Mplus user's guide: Statistical analysis with latent variables*. Los Angeles, CA: Muthén & Muthén.
- Schumacker, R. E. (1996). *Disattenuating correlation coefficients*. *Rasch Measurement Transactions*, 10, 479.
- Penfield, R. D. (2005). DIFAS: Differential item functioning analysis system. *Applied Psychological Measurement*, 29, 150–151.
- Penfield, R. D. (2007). An approach for categorizing DIF in polytomous items. *Applied Measurement in Education*, 20, 335–355.
- Smarter Balanced Assessment Consortium. (2014a). *Accessibility and accommodations framework*. Retrieved from *Smarter Balanced Accessibility and Accommodations Framework*. Los Angeles, CA: Author.
- Smarter Balanced Assessment Consortium (2014b). *Usability, accessibility, and accommodations guidelines*. Los Angeles, CA: Author.
- Smarter Balanced Assessment Consortium (2014f). *Interpretation and use of scores and achievement levels*. Los Angeles, CA: Author.

References

- Smarter Balanced Assessment Consortium. (2014g). *Reporting system user guide*. Los Angeles, CA: Author.
- Smarter Balanced Assessment Consortium. (2015a). *Content specifications for the summative assessment of the common core state standards for English language arts and literacy in history/social studies, science, and technical subjects*. Los Angeles, CA: Author.
- Smarter Balanced Assessment Consortium. (2015b). *Content specifications for the summative assessment of the common core state standards for mathematics*. Los Angeles, CA: Author.
- Smarter Balanced Assessment Consortium. (2015c). *Item and task specifications*. Los Angeles, CA: Author.
- Smarter Balanced Assessment Consortium. (2015d). *The Smarter Balanced Assessment Consortium: Achievement level setting final report*. Los Angeles, CA: Author.
- Smarter Balanced Assessment Consortium. (2016). *2014–2015 technical report*. Los Angeles, CA: Author.
- Smarter Balanced Assessment Consortium. (2017). *2016–2017 technical report*. Los Angeles, CA: Author.
- U.S. Department of Education. (2007). Standards and assessments peer review guidance: Information and examples for meeting requirements of the No Child Left Behind act of 2001. Retrieved from US Department of Education Policy and Guidance
- Yen, W. M. (1981). Using simulation results to choose a latent trait model. *Applied Psychological Measurement*, 5, 245–262.
- Yen, W. M., & Fitzpatrick, A. R. (2006). *Item response theory*. In R. L. Brennan (Ed.), *Educational measurement* (4th ed., pp. 111–153). Westport, CT: Praeger
- Young, M. J., & Yoon, B. (1998, April). Estimating the consistency and accuracy of classifications in a standards-referenced assessment. (CSE Technical Report 475). Center for the Study of Evaluation, Standards, and Student Testing. Los Angeles, CA: University of California, Los Angeles.
- Zhang, T., Haertel, G., Javitz, H., Mislevy, R., Murray, E., & Wasson, J. (2009). *A design pattern for a spelling bee assessment for students with disabilities*. A paper presented at the annual conference of the American Psychological Association, Montreal, Canada.
- Zwick, R., Donoghue, J. R., & Grima, A. (1993). *Assessment of differential item functioning for performance tasks*. *Journal of Educational Measurement*, 30, 233–251.