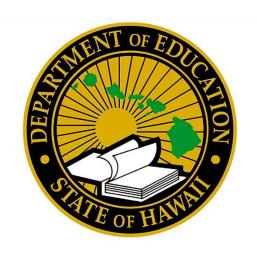
Hawai'i Smarter Balanced Assessments 2023–2024 Technical Report



Submitted to Hawai'i Department of Education by Cambium Assessment, Inc.

TABLE OF CONTENTS

1. (OVERVIEW	1
2. [TEST ADMINISTRATION	3
2.	1 Testing Windows	3
2.	2 Test Options and Administrative Roles	3
	2.2.1 Administrative Roles	4
	2.2.2 Online Administration	6
	2.2.3 Paper-Pencil Test Administration	7
	2.2.4 Braille Test Administration	8
2.	3 Training and Information for Test Coordinators and Administrators	8
	2.3.1 Online Training	9
	2.3.2 Statewide Trainings	12
2.	4 Test Security	13
	2.4.1 Student-Level Testing Confidentiality	13
	2.4.2 System Security	14
	2.4.3 Security of the Testing Environment	14
	2.4.4 Test Security Violations	15
2.	5 Student Participation	16
	2.5.1 Homeschooled Students	16
	2.5.2 Exempt Students	16
2.	6 Online Testing Features and Testing Accommodations	17
	2.6.1 Online Universal Tools for All Students	17
	2.6.2 Designated Supports and Accommodations	19
2.	7 Testing Time	30
2.	8 Data Forensics Program	32
	2.8.1 Changes in Student Performance	32
	2.8.2 Test-Taking Time	
	2.8.3 Inconsistent Item Response Pattern (Person Fit)	
	2.8.4 Item-Response Change	34

	2.9 Prevention and Recovery of Disruptions in the Test Delivery System	34
	2.9.1 High-Level System Architecture	35
	2.9.2 Automated Backup and Recovery	36
	2.9.3 Other Disruption Prevention and Recovery Mechanisms	36
3.	SUMMARY OF 2023–2024 OPERATIONAL TEST ADMINISTRATION	38
	3.1 Student Population	38
	3.2 Summary of Overall Student Performance	40
	3.3 Distribution of Student Ability and Item Difficulty	51
4.	VALIDITY	58
	4.1 Evidence on Test Content	58
	4.2 Evidence on Internal Structure	64
	4.3 Evidence on Relations to Other Variables	66
5.	RELIABILITY	68
	5.1 Marginal Reliability	68
	5.2 Standard Error Curves	69
	5.3 Reliability of Achievement Classification	72
	5.4 Reliability for Subgroups	77
	5.5 Reliability for Claim Scores	80
6.	SCORING	83
	6.1 Estimating Student Ability Using Maximum Likelihood Estimation	83
	6.2 Rules for Transforming Theta to Vertical Scale Scores	84
	6.3 Lowest/Highest Obtainable Scores	85
	6.4 Scoring All Correct and All Incorrect Cases	86
	6.5 Rules for Calculating Strengths and Weaknesses for Claim Scores	86
	6.6 Target Scores	86
	6.6.1 Target Scores Relative to Student's Overall Estimated Ability	87
	6.6.2 Target Scores Relative to Proficiency Standard (Level 3 Cut)	88
	6.7 Handscoring	89
	6.7.1 Rater Selection	89

	6.7.2 Rater Training, Qualification, and Scoring	90
	6.7.3 Rater Monitoring, Feedback, and Evaluation	93
	6.7.4 Rater Agreement	95
	6.8 Automated Scoring	97
	6.8.1 Project Essay Grade	97
	6.8.2 Model Training and Validation	98
	6.8.3 Automated Scoring Processes	103
	6.8.4 Human-Machine Agreement	106
	6.8.5 Recommendations	110
7.	REPORTING AND INTREPRETING SCORES	111
	7.1 Centralized Reporting System	111
	7.1.1 Dashboard	113
	7.1.2 Aggregate Score Reports: Overall Performance	115
	7.1.3 Aggregate Score Reports: Claim and Target Performance	117
	7.1.4 Roster Performance Report	118
	7.1.5 Trend Report	118
	7.1.6 Individual Student Report	119
	7.2 Interpretation of Reported Scores	124
	7.2.1 Scale Score	124
	7.2.2 Standard Error of Measurement	124
	7.2.3 Achievement Level	124
	7.2.4 Performance Category for Claims	125
	7.2.5 Performance Category for Targets	125
	7.2.6 Aggregated Scale Score	125
	7.3 Appropriate Uses of Test Results	126
8.	QUALITY CONTROL PROCEDURES	127
	8.1 Adaptive Test Configuration	127
	8.1.1 Platform Review	127
	8.1.2 User Acceptance Testing and Final Review	
	8.2 Quality Assurance in Document Processing	128

8.3 Quality Assurance in Data Preparation	128
8.4 Quality Assurance in Online Test Delivery System	128
8.4.1 Score Report Quality Check	129
REFERENCES	132

LIST OF TABLES

Table 1. 2023–2024 Testing Windows	3
Table 2. 2023–2024 Testing Options	3
Table 3. Number of Students Who Took Tests Remotely in the 2023–2024 Summative Test Administration	7
Table 4. Number of Students Who Took Paper-Pencil Tests in the 2023–2024 Summative Test Administration	8
Table 5. SY 2023–2024 Universal Tools, Designated Supports, and Accommodations	24
Table 6. ELA/L Total Students with Allowed Embedded and Non-Embedded Accommodations	25
Table 7.ELA/L Total Students with Allowed Embedded Designated Supports	26
Table 8. ELA/L Total Students with Allowed Non-Embedded Designated Supports	26
Table 9. Mathematics Total Students with Allowed Embedded and Non-Embedded Accommodation	
Table 10. Mathematics Total Students with Allowed Embedded Designated Supports	28
Table 11. Mathematics Total Students with Allowed Non-Embedded Designated Supports	29
Table 12. Test-Taking Time: ELA/L	31
Table 13. Test-Taking Time: Mathematics	31
Table 14. Participation Rates by Percentage: ELA/L	38
Table 15. Participation Rates by Percentage: Mathematics	39
Table 16. Number of Students: ELA/L	39
Table 17. Number of Students: Mathematics	40
Table 18. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and Subgroup: ELA/L (Grades 3–5)	•
Table 19. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and Subgroup: ELA/L (Grades 6–8)	
Table 20. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and Subgroup: ELA/L (Grade 11)	•
Table 21. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and Subgroup: Mathematics (Grades 3–5)	•
Table 22. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and Subgroup: Mathematics (Grades 6–8)	•
Table 23. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and Subgroup: Mathematics (Grade 11)	

Table 24. Percentage of Students in Performance Categories by Claim	51
Table 25. Percentage of CAT Delivered Tests Meeting Blueprint Requirements: ELA/L (Grades 3-	
Table 26. Percentage of CAT Delivered Tests Meeting Blueprint Requirements: ELA/L (Grades 6–11)	-8,
Table 27. Percentage of CAT Delivered Tests Meeting Blueprint Requirements for Claims and Targ Mathematics (Grades 3–5)	
Table 28. Percentage of CAT Delivered Tests Meeting Blueprint Requirements for Claims and Targ Mathematics (Grades 6–8)	
Table 29. Percentage of CAT Delivered Tests Meeting Blueprint Requirements for Claims and Targ Mathematics (Grade 11)	-
Table 30. Average and Range of the Number of Unique Targets Assessed Within Each Claim Acros All Delivered CAT Tests	
Table 31. Correlations Among Claims: ELA/L	65
Table 32. Correlations Among Claims: Mathematics	66
Table 33. Relationship Among the Smarter Balanced, Algebra I, and Algebra II Test Scores	67
Table 34. Marginal Reliability: ELA/L and Mathematics	69
Table 35. Average Conditional Standard Error of Measurement by Achievement Level	72
Table 36. Average Conditional Standard Error of Measurement at Each Achievement-Level Cut and Difference of the SEMs Between Two Cuts	
Table 37. Classification Accuracy and Consistency	76
Table 38. Marginal Reliability Coefficients for Overall and by Subgroup: ELA/L (Grades 3–4)	77
Table 39. Marginal Reliability Coefficients for Overall and by Subgroup: ELA/L (Grades 5–6)	77
Table 40. Marginal Reliability Coefficients for Overall and by Subgroup: ELA/L (Grades 7–8)	78
Table 41. Marginal Reliability Coefficients for Overall and by Subgroup: ELA/L (Grade 11)	78
Table 42. Marginal Reliability Coefficients for Overall and by Subgroup: Mathematics (Grades 3–4)	
Table 43. Marginal Reliability Coefficients for Overall and by Subgroup: Mathematics (Grades 5–6	/
Table 44. Marginal Reliability Coefficients for Overall and by Subgroup: Mathematics (Grades 7–8	_
Table 45. Marginal Reliability Coefficients for Overall and by Subgroup: Mathematics (Grade 11).	80
Table 46. Marginal Reliability Coefficients for Claim Scores: ELA/L	81
Table 47. Marginal Reliability Coefficients for Claim Scores: Mathematics	82

Table 48. Vertical Scaling Constants on the Reporting Metric
Table 49. Cut Scores in Scale Scores
Table 50. Extended Lowest and Highest Obtainable Scores
Table 51. Number of Handscored Items in 2023–2024 Smarter Balanced Summative Item Pool, by Grade and Subject
Table 52. Inter-Rater Agreement for ELA/L Short-Answer Items
Table 53. Inter-Rater Agreement for ELA/L Essay Items
Table 54. Inter-Rater Agreement for Mathematics Items
Table 55. Number of Items Eligible for Automated Scoring, by Grade and Subject Area99
Table 56. Initial Model Evaluation Criteria
Table 57. Demographic Variables and Categories
Table 58. Secondary Validation Criteria
Table 59. Summary of Secondary Validation Results, by Grade and Subject Area102
Table 60. Summary of Live Training and Validation Results, by Grade and Subject Area103
Table 61. Flags Currently Established
Table 62. Model Setting
Table 63. Human-Machine Agreement for ELA/L Short-Answer Items on Initial and Secondary Validation Samples, by Grade
Table 64. Human-Machine Agreement for ELA/L Essay Items on Initial and Secondary Validation Samples, by Grade
Table 65. Human-Machine Agreement for Mathematics Items on Initial and Secondary Validation Samples, by Grade
Table 66. Human-Machine Agreement for ELA/L Short-Answer Items on Live Validation Sample, by Grade
Table 67. Human-Machine Agreement for ELA/L Essay Items on Live Validation Sample, by Grade
Table 68. Human-Machine Agreement for Mathematics Items on Live Validation Samples, by Grade
Table 69. Types of Online Score Reports by Level of Aggregation
Table 70. Types of Subgroups
Table 71. Overview of Ouality Assurance Reports

LIST OF FIGURES

Figure 1. Percentage Proficient Across Years: ELA/L	47
Figure 2. Percentage Proficient Across Years: Mathematics	48
Figure 3. Average Scale Score Across Years: ELA/L	49
Figure 4. Average Scale Score Across Years: Mathematics	50
Figure 5. Student Ability—Item Difficulty Distribution: ELA/L	52
Figure 6. Student Ability—Item Difficulty Distribution by Claim: ELA/L (Grades 3–5)	53
Figure 7. Student Ability—Item Difficulty Distribution by Claim: ELA/L (Grades 6–8, and 11)	54
Figure 8. Student Ability—Item Difficulty Distribution: Mathematics	55
Figure 9. Student Ability—Item Difficulty Distribution by Claim: Mathematics (Grades 3–5)	56
Figure 10. Student Ability—Item Difficulty Distribution by Claim: Mathematics (Grades 6–8, 11))57
Figure 11. Conditional Standard Error of Measurement: ELA/L	70
Figure 12. Conditional Standard Error of Measurement: Mathematics	71
Figure 13. PEG Architecture	98
Figure 14. Response Routing Rules	104
LIST OF EXHIBITS	
Exhibit 1. Dashboard: State Level	113
Exhibit 2. Dashboard: Complex Level	114
Exhibit 3. Detailed Dashboard: Complex Level	115
Exhibit 4. Overall Performance Summary Results for Grade 6 Mathematics: Complex Level	116
Exhibit 5. Overall Performance Summary Results for Grade 6 Mathematics by Gender: Complex Level	
Exhibit 6. Claim and Target Level Results for Grade 6 Mathematics: Complex Level	117
Exhibit 7. Roster Performance Report for Grade 6 Mathematics	118
Exhibit 8. Trend Report for ELA/L: Student Level	119
Exhibit 9. Individual Student Report for Grade 5 ELA/L	121

viii

1. OVERVIEW

The Smarter Balanced Assessment Consortium (SBAC) has developed a next-generation assessment system designed to accomplish two goals: first, to measure students' mastery of the *Common Core State Standards* (CCSS) in English language arts/literacy (ELA/L) and mathematics in grades 3–8 and 11, and second, to provide valid, reliable, and fair test scores of students' academic achievement. At the time of development, Hawai'i was one of 18 member states (plus the U.S. Virgin Islands) leading the development of assessments in ELA/L and mathematics. The system includes summative assessments for accountability purposes and optional interim assessments that supply meaningful feedback and actionable data that teachers and educators can use to help students succeed. SBAC, a state-led collaboration, is intended to provide leadership and resources to improve teaching and learning by creating and maintaining a suite of summative and interim assessments and tools aligned to the CCSS in ELA/L and mathematics.

The Hawai'i State Board of Education formally adopted the CCSS in ELA/L and mathematics on June 18, 2010. All students in Hawai'i, including students with significant cognitive disabilities who are eligible to take the Hawai'i State Alternate Assessment (an alternate assessment based on Alternate Academic Achievement Standards), are taught the same academic content standards. The Hawai'i CCSS define the knowledge and skills that students need to succeed in college and careers after graduating from high school. These standards include rigorous content and application of knowledge through higher-order skills and align with college and workforce expectations.

Since the adoption of the CCSS in 2010, the Hawai'i Department of Education (HIDOE) began implementing the CCSS in the 2012–2013 school year with grades K–2 and 11–12. This transition was fully implemented in all grade levels in the 2013–2014 school year. The new Hawai'i statewide assessments in ELA/L and mathematics aligned with the CCSS were administered for the first time in spring 2015 to students in grades 3–8 and 11 in all public elementary and secondary schools.

The Smarter Balanced assessments comprise the end-of-year summative assessment designed for accountability purposes, and the optional interim assessments that support teaching and learning throughout the year. The summative assessments evaluate student achievement based on the CCSS and track student progress toward college and career readiness in ELA/L and mathematics. The summative assessments consist of two parts: a computer-adaptive test (CAT) and a performance task (PT).

- The Computer-Adaptive Test (CAT) provides an individualized assessment for each student.
- The **Performance Task (PT)** challenges students to apply their knowledge and skills to real-world problems. PTs can best be described as collections of items and activities that are coherently connected to a single theme or scenario. They are used to better measure capacities such as depth of understanding, research skills, and complex analysis, which cannot be adequately assessed with selected- or constructed-response items. The computer can score some PT items, but most are handscored.

The optional interim assessments allow teachers to monitor student progress throughout the year and provide information that they can use to improve instruction and learning. These tools are used at the discretion of schools and complex areas, and teachers can employ them to gauge students' progress in mastering specific concepts at strategic points during the school year. There are three types of interim assessments available as fixed-form tests:

1

- The **Interim Comprehensive Assessment (ICA)** tests the same content and reports scores on the same scale as the summative assessments.
- The Interim Assessment Block (IAB) focuses on specific sets of related concepts that measure three to eight assessment targets and provide detailed information about student learning.
- The Focused Interim Assessment Block (FIAB) focuses on specific sets of related concepts that measure no more than three assessment targets and provide more detailed information about student learning than the IAB alone.

In the 2019–2020 school year, the U.S. Department of Education waived testing requirements due to the COVID-19 pandemic (https://www2.ed.gov/policy/gen/guid/secletter/200320.html). For the 2020–2021 school year, the U.S. Department of Education did not grant waivers for standardized testing but did waive certain accountability requirements (e.g., mandatory high participation rates) due to the impacts of the pandemic in many states, resulting in lower participation rates than in previous years. Starting in the 2021–2022 school year, all students were required to take ELA/L and mathematics summative assessments.

Starting with the 2020–2021 Smarter Balanced summative test administration, Hawai'i shortened the full test blueprints for ELA/L and mathematics and allowed schools to administer remote test administrations to individual students.

The American Institutes for Research (AIR) delivered the Hawai'i statewide assessments in ELA/L and mathematics through the 2018–2019 school year. Starting with SY 2020–2021, Cambium Assessment, Inc. (CAI) (formerly a segment of AIR) delivered and scored the Smarter Balanced assessments and produced the score reports. Measurement Incorporated (MI) scored the handscored items.

This report provides a technical summary of Hawai'i's 2023–2024 administration of the Smarter Balanced summative assessments in English language arts/literacy (ELA/L) and mathematics in grades 3–8 and 11. The report is divided into eight chapters: Overview; Test Administration; Summary of the 2023–2024 Operational Test Administration; Validity; Reliability; Scoring; Reporting and Interpreting Scores; and Quality Control Procedures. The data included in this report are based on Hawai'i data for the summative assessment only. For the interim assessments, the number of students who took ICAs and IABs and a summary of their performance are provided in Appendix A.

While this report includes information on all aspects of the technical quality of the Smarter Balanced test administration in Hawai'i, it is an addendum to the 2023–2024 Smarter Balanced technical report. The Smarter Balanced technical report contains information on item and test development, item content review, field-test administration, item-data review, item calibrations, content alignment study, standard setting, and other validity information.

The Smarter Balanced produces a technical report for the Smarter Balanced assessments, including all aspects of the technical qualities for the Smarter Balanced assessments described in the *Standards for Educational and Psychological Testing* (American Educational Research Association [AERA], American Psychological Association [APA], & National Council on Measurement in Education [NCME], 2014) and the requirements of the U.S. Department of Education, *Peer Review of State Assessment Systems: Non-Regulatory Guidance for States* (U.S. Department of Education, 2015). The Smarter Balanced technical report includes information using the data at the consortium level, combining data from the consortium states.

2. TEST ADMINISTRATION

2.1 TESTING WINDOWS

The 2023–2024 Smarter Balanced Assessment (SBA) testing window spanned approximately three months for the summative assessments for most schools and spanned the entire school year for the interim assessments. The paper-pencil fixed forms for the summative assessments were administered concurrently during the three-month online summative window. Table 1 shows the testing windows for both online and paper-pencil assessments.

Tests	Grade	Start Date	End Date	Mode	
		2/20/2024	5/30/2024		
	3–8	3–8 3/11/2024 6		Online Adaptive	
		(Multi-track)	(Multi-track)		
		2/20/2024	5/30/2024		
Summative Assessments	11	11/20/2023	5/30/2024	Online Adaptive	
		(Block Scheduled) (Block Scheduled)			
	3–8, 11	2/20/2023	5/17/2024	Paper Fixed-Form	
	3–8, 11	2/20/2023	6/14/2024	Remote Online Adaptive	
	3–8, 11	2/20/2023	5/17/2024	Braille Paper Fixed-Form	
Interim Comprehensive Assessments	3–8, 11	8/15/2023	7/19/2024	Online Fixed-Form	
Interim Assessment Blocks	3–8, 11	8/15/2023	7/19/2024	Online Fixed-Form	
Focused Interim Assessment Blocks	3–8, 11	8/15/2023	7/19/2024	Online Fixed-Form	

Table 1. 2023–2024 Testing Windows

2.2 TEST OPTIONS AND ADMINISTRATIVE ROLES

The Smarter Balanced Assessment (SBA) is administered primarily online. To ensure that all eligible students in the tested grades were given the opportunity to take the SBA, several assessment options were available to accommodate students' needs. Table 2 lists the testing options offered in 2023–2024. A testing option is selected by content area. Once an option is selected, it is applied to all tests in the content area.

Assessments	Testing Options	Test Mode		
	English	Online		
	Braille	Paper-Pencil/Online		
Summative Assessments	Spanish (mathematics only)	Online		
	Paper-Pencil Fixed-Form	Paper-Pencil		
	Remote	Online		
	English	Online		
Interim Assessments	Braille	Online		
michin Assessments	Spanish (mathematics only)	Online		
	Remote	Online		

Table 2. 2023–2024 Testing Options

To ensure that standardized administration conditions are met, test administrators (TAs) follow procedures outlined in the *Smarter Balanced ELA/L and Mathematics Online, Summative Test Administration Manual* (TAM). TAs must review the TAM before testing to ensure that the testing room is prepared for testing (e.g., removing certain classroom posters, arranging desks). Make-up procedures should be established for students who are absent on the day(s) of testing. TAs follow required administration procedures and directions and read the boxed directions verbatim to students, ensuring standardized administration conditions.

2.2.1 Administrative Roles

The key personnel involved with the test administration are principals (PRs), test coordinators (TCs), and TAs. The main responsibilities of the key personnel are outlined in the following descriptions. More detailed descriptions can be found in the TAM provided online at:

 $\underline{https://smarterbalanced.alohahsap.org/resource-list/en/smarter-balanced-summative-test-administration-manual-2023-2024.}$

Principals

The PR's primary responsibility is to ensure that testing in his or her school is conducted in accordance with the test procedures and security policies established by the Hawai'i State Department of Education (HIDOE).

PRs are responsible for performing the following functions:

- Reviewing all Smarter Balanced policies and test administration documents
- Reviewing scheduling and test requirements with TCs and TAs
- Working with TCs and technology coordinators to ensure that all systems, including the CAI Secure Browser, are properly installed and functioning
- Designating or acting as the TC
- Importing users (TCs) into the Test Information Distribution Engine (TIDE)
- Scheduling and administering training sessions for all TCs, TAs, and technology coordinators (refer to Section 2.3, Training and Information for Test Coordinators and Administrators)
- Ensuring that all personnel understand and are trained on the proper administration of the Smarter Balanced assessments
- Monitoring secure test administration
- Investigating and reporting all testing improprieties, irregularities, and breaches reported by TCs or TAs
- Attending to any secure materials according to state and Smarter Balanced policies

Test Coordinator

The TC's primary responsibility is to coordinate the administration of the Smarter Balanced assessments in the school.

TCs are responsible for performing the following functions:

- Identifying TAs and proctors (if appropriate) and ensuring that TAs complete the TA Certification Course
- Establishing a testing schedule with PRs and TAs based on the testing windows
- Working with technology staff to ensure timely computer setups and installations
- Working with TAs to review student information in TIDE to ensure that student information and test settings for designated supports and accommodations are applied correctly
- Identifying students who may require designated supports and test accommodations and ensuring that procedures for testing these students follow state and Smarter Balanced policies
- Attending all school trainings and reviewing all Smarter Balanced policy and test administration documents
- Ensuring that all TAs attend school trainings and review online training modules posted on the portal
- Establishing secure and separate testing rooms if needed
- Monitoring secure administration of the test
- Monitoring testing progress during the testing window and ensuring that all students participate, as appropriate
- Investigating and reporting all testing improprieties, irregularities, and breaches reported by the TAs in coordination with the PRs
- Attending to any secure materials according to state and Smarter Balanced policies

Test Administrator

The TA's primary responsibility is to administer the Smarter Balanced assessments. The TA's role is designed for test administrators, such as technology staff, who administer tests but should not have access to student results.

TAs are responsible for performing the following functions:

- Completing Smarter Balanced test administration training and reviewing all Smarter Balanced policy and test administration documents before administering any Smarter Balanced assessments
- Reviewing student information for accuracy before testing to ensure that students receive the
 proper test with the appropriate supports and reporting any potential data errors to TCs and PRs,
 as appropriate
- Administering the Smarter Balanced assessments
- Reporting all potential test security incidents to the TCs or PRs in a manner consistent with Smarter Balanced, state, and school policies

2.2.2 Online Administration

Within the state's testing window, schools can set the testing schedule and customize their testing conditions, such as allowing students to test in intervals (i.e., multiple sessions) rather than in one long period and minimizing the interruption of classroom instruction and efficiently using its facility. With online testing, schools do not need to handle test booklets and address the storage and security problems inherent in large shipments of materials to a school site.

Starting with SY 2020–2021, a new feature was developed within the universally used Test Delivery System (TDS) that allowed tests to be administered remotely by a TA to students who remained at home. The decision to allow students to test remotely was made at the school level in cases when a parent or guardian refused to take a student to campus for testing but insisted on the student being tested. This new feature allowed TAs to pre-schedule a testing session, host online video and chat features with a group of students, and video monitor students in a testing session.

To ensure that TAs were able to use these new features, an additional *Remote Testing TA Certification Course* was developed. TAs scheduled to administer remote testing sessions were required to complete this course prior to test administration. In addition, before a student was eligible for remote test administration, a parent or guardian had to provide written consent to the school to administer a remote test that would contain video and audio components allowing the TA to view and monitor the student. The school's TC was responsible for ensuring that these students had positive consent for remote testing within the TIDE system. Additional resources were developed tor TAs to understand the requirements for remote testing and posted to the state portal at https://smarterbalanced.alohahsap.org/resource-list/en/remote-summative-test-administration-2023-2024.

TCs oversee all aspects of testing at their schools and serve as the main point of contact; TAs administer the online assessments only. TAs are trained in the online testing requirements and the mechanics of starting, pausing, and ending a test session. Training materials for the test administration are provided online. All school personnel who serve as TAs must complete an online TA Certification Course. Staff who complete this certification course receive a certificate of completion and are qualified to administer assessments.

To start a test session, the TA must first enter the TA Interface of the online testing system using his or her own computer. A session ID is generated when the test session is created. Students who are taking the assessment with the TA must enter their State Student Identifier (SSID), first name, and session ID into the Student Interface using computers provided by the school. The TA then verifies that the students are taking the appropriate assessments with the appropriate accessibility feature(s) (refer to Section 2.6, Online Testing Features and Testing Accommodations, for a full list of accommodations). Students can begin testing only when the TA confirms the settings. The TA must read the *Directions for Administration* in the *Smarter Balanced Online Summative Test Administration Manual* aloud to the student(s) and walk them through the login process.

Once an assessment is started, the student must answer all the test questions presented on a page before proceeding to the next page. Skipping questions is not permitted. For the CAT, students can review and edit previously answered items as long as these items are in the same test session and this session has not been paused for more than 20 minutes. In addition, students can review and edit only previously answered items before submitting the assessment. During an active CAT session, if a student reviews and changes the response to a previously answered item, all following items to which the student already responded remain the same. No new items are assigned to this student for changing answers. For example, a student

paused for 10 minutes after completing Item 10. After the pause, the student went back to Item 5 and changed the answer. If the updated response to Item 5 changed the item score from wrong to right, the student's overall score would improve; however, there would be no change in Items 6–10. For PTs, there is no pause rule; but the same rules that apply to the CAT for reviews and changes to responses also apply to PTs.

The CAT must be completed within 45 calendar days of the start date, or the assessment opportunity will expire. The ELA/L performance task must be completed within 10 calendar days of the start date.

During a test session, TAs may pause the test for a student or a group of students to take a break. It is up to the TA to determine an appropriate stopping point; however, to ensure the integrity of test scores and testing, the CAT cannot be paused for more than 20 minutes for ELA/L and mathematics. If an assessment is paused for more than 20 minutes, the student must start a new test session and resume the test from the point where he or she paused. Under this circumstance, viewing and editing previous responses is no longer permitted.

The TA must remain in the room when the test is administered in person and be present continuously when using the video feature for remote test administrations to monitor student testing. When the test session ends, the TA must ensure that each student has successfully logged out of the system. The TA must also collect and shred any handouts or scratch paper that students may have used during the CAT session; if handouts or scratch paper were used for the ELA/L PT, the TA must collect and securely store them until the ELA/L PT has been submitted. After the PT's submission, the TA must securely shred all handouts and/or scratch paper.

The number of students who took summative tests remotely in 2023–2024 is presented in Table 3.

Table 3. Number of Students Who Took Tests Remotely in the 2023–2024 Summative Test Administration

Subject	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	Grade 11	Total
ELA/L	8	13	6	13	13	14	4	71
Mathematics	8	13	6	13	13	13	4	70

2.2.3 Paper-Pencil Test Administration

There are two matching versions of the paper-pencil Smarter Balanced ELA/L and mathematics assessments. One version is provided as an accommodation for students who cannot access a computer, and the other is a braille version for students with blindness or visual impairments. Both versions contain the same items and are based on the Smarter Balanced full-length blueprints for ELA/L and mathematics used in SY 2023-24. TCs from schools with any student(s) who require the paper-pencil assessment must submit a request to HIDOE for test materials on behalf of the student(s) before the testing window opens. If the request is approved by HIDOE, the testing contractor will ship the appropriate test booklets and the paper-pencil TAM to the school.

Separate test booklets are used for the ELA/L and mathematics assessments, which are based upon the Smarter Balanced full-length blueprint. The items from the CAT and the PT components are combined into one test booklet, including two sessions for the CAT and one session for the PT in both content areas. Thus, the TA can break up the assessment into separate test sessions. After the student completes the

assessment, the TC will return the test booklets to the testing contractor, and the testing contractor will scan the answer document and score the test, including the handscored items.

The total number of students who took paper-pencil tests is shown in Table 4.

Table 4. Number of Students Who Took Paper-Pencil Tests in the 2023–2024 Summative Test Administration

Subject	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	Grade 11	Total
ELA/L	1		1	1	1	1		5
Mathematics	1		1	1	1	1		5

2.2.4 Braille Test Administration

The adaptive braille test was available with the same test blueprint in both ELA/L and mathematics. In the 2017–2018 test administration, Smarter Balanced added the Braille Hybrid Adaptive Test (Braille HAT) for mathematics. The Braille HAT consists of a fixed-form segment, a computer-adaptive segment, and a fixed-form PT. The fixed-form segment includes items with tactile graphics, which can be embossed at the testing location or received as a package of pre-embossed materials through HIDOE. All items on the Braille HAT can be presented to students using a Refreshable Braille Display (RBD). The blueprints for the Braille HAT follow the Smarter Balanced full-length blueprints for mathematics used in SY 2023-24. This was not an option for administration in Hawai'i in 2023–2024, and no versions of these tests were taken.

The braille interface comprises several formats as follows:

- The braille interface includes a text-to-speech (TTS) component for mathematics consistent with the read-aloud assessment accommodation. The Job Access with Speech (JAWS) screen-reading software provided by Freedom Scientific is an essential component that students use with the braille interface.
- Mathematics items are presented to students in Nemeth Braille Code via a braille embosser through the adaptive online summative test and a fixed-form PT.
- Students taking the summative ELA/L assessment can emboss both reading passages and items as they progress through the assessment. If a student has an RBD, a 40-cell RBD is recommended. The summative ELA/L is presented to the student with items in either contracted or uncontracted literary braille (for items containing only text) and via a braille embosser (for items with tactile or spatial components that cannot be read by an RBD).

Before administering the online summative assessments using the braille interface, TAs must ensure that technical requirements are met. These requirements apply to the student's computer, the TA's computer, and any supporting braille technologies used in conjunction with the braille interface.

2.3 TRAINING AND INFORMATION FOR TEST COORDINATORS AND ADMINISTRATORS

PRs and TCs oversee all aspects of testing at their schools and serve as the main points of contacts; TAs administer the online assessments. The online TA Certification Course, webinars, user guides, manuals,

and training sites are used to train TAs on the online testing requirements and the mechanics of starting, pausing, and ending a test session. Training materials for administration are provided online.

2.3.1 Online Training

Multiple training opportunities are offered to key assessment staff through the state portal.

TA Certification Course

There are three TA Certification Courses that are available for TAs: an Interim Assessment TA Certification Course, a Summative Assessment TA Certification Course, and a Remote Assessment TA Certification Course. TAs must complete an online TA Certification Course every year in order to administer assessments. The Interim Assessment TA Certification Course must be completed to administer Interim Assessments, while the Summative Assessment TA Certification Course must be completed to administer Summative Assessments. For 2023-2024, TAs administering summative tests must complete both the Interim and Summative TA Certification Courses. These web-based courses are each about 30–45 minutes long and cover information on testing policies and the steps for administering Interim and Summative test sessions in the online testing system. The courses are interactive, requiring participants to start test sessions under different scenarios. Participants are required to answer multiple-choice questions about the information provided throughout the training and at the end of the Summative TA course. A third TA Certification Course of about 20 minutes is required for TAs administering tests in a remote format. For 2023–2024, TAs administering remote tests were required to take all courses.

Webinars

The following five webinars were offered to users in the field:

- Accessibility and Accommodations. This webinar provides an overview of the accessibility
 features and supports available to students during testing, including universal tools, designated
 supports, and accommodations.
- Smarter Balanced Test Coordinators Training. This webinar provides information about accessing and using the Interim Assessments, Summative Assessments, Centralized Reporting System, and Digital Library.
- Test Information Distribution Engine. This webinar provides an overview of how to navigate the Test Information Distribution Engine (TIDE), including managing student information and monitoring test progress.
- Centralized Reporting System. This webinar provides information on the Centralized Reporting System (CRS), including an overview of accessing student reports and the distribution of reports to parents and guardians.
- Remote Interim Administration. This webinar provides information about setting up and administering remote interim assessments using the Test Delivery System (TDS) and the CAI Secure Browser.

Each of these webinars is about one hour long. The interactive nature of these training webinars allows the participant to ask questions during and after the presentation. After the live webinar, a streaming video recording of the webinar is made available on the state portal.

Practice and Training Test Site

Starting in August 2022, separate online training sites were opened for TCs, TAs, and students. TAs could practice administering assessments and starting and ending test sessions on the TA Training Site, and students could practice taking an online assessment on the Student Practice and Training Site. The Smarter Balanced assessment practice tests mirror the corresponding summative assessments for ELA/L and mathematics. Each test provides students with a grade-specific testing experience, including a variety of question types and difficulty levels (approximately 30 items each in ELA/L and mathematics) and a performance task in ELA/L.

The training tests are designed to provide students and TAs with opportunities to quickly familiarize themselves with the software and navigational tools that they will use for the Smarter Balanced assessments in ELA/L and mathematics. Training tests are available for both ELA/L and mathematics and are organized by grade bands (grades 3–5, grades 6–8, and grade 11), with each test containing 5–10 questions.

A student can log in to the practice and training test site directly as a "Guest" without a TA-generated test session ID, or the student can log in through a training test session created by the TA in the TA Training Site. Items in the student training test include all item types that are included in the operational item pool, including multiple-choice, grid, and natural language items.

Manuals and User Guides

The following manuals and user guides are available on the Hawai'i Statewide Assessment Program Portal:

The Smarter Balanced Online, Summative, Test Administration Manual provides information for TCs and TAs administering the Smarter Balanced online summative assessments in ELA/L and mathematics. It includes screen captures and step-by-step instructions on how to administer the online tests.

The Smarter Balanced Interim Assessments Test Administration Guide provides an overview of how to prepare for and administer the Smarter Balanced Interim assessments.

The Online Calculators in the Test Delivery System Manual and the Desmos User Guide provide instructions for using the online Desmos Calculators during testing.

The *Braille Requirements and Testing Manual* includes information about the supported operating systems and required hardware and software for braille testing. It also provides information on how to configure JAWS, how to navigate an online test with JAWS, and how to administer a test to a student requiring braille.

The System Requirements for Online Testing document outlines the basic technology requirements for administering an online assessment, including operating system requirements and supported web browsers.

The Secure Browser Installation Manual provides instructions for downloading and installing the CAI Secure Browser on supported operating systems used for online assessments.

The *Technical Specifications Manual for Online Testing* provides technology staff with the technical specifications for online testing, including information on Internet and network requirements, general hardware and software requirements, and the text-to-speech function.

The *Test Information Distribution Engine User Guide* and *Quick Guide to TIDE* are designed to help users navigate TIDE. Users can find information on managing user account information, student account information, student test settings and accommodations, testing incidents, creating and editing rosters, and voice packs.

The Centralized Reporting System User Guide provides information about the CRS, including instructions for viewing score reports, managing test administration, and searching for students. It is also a component of the Smarter Balanced Interim Assessments that allows authorized users to view individual student responses on both the Interim Comprehensive Assessments (ICAs) and the Interim Assessment Blocks (IABs).

The *Guide to Navigating the Online HSAP Administration* is designed to help users navigate the TDS, including the Student Interface and the TA Interface, and to help TAs manage and administer online testing for students.

The Assessment Viewing Application User Guide provides an overview of how to access and use the Assessment Viewing Application (AVA), which allows teachers to view items on the Smarter Balanced interim assessments.

The *Usability, Accessibility, and Accommodations Guidelines* describe the current universal tools, designated supports, and accommodations adopted by the Smarter Balanced states to ensure valid assessment results for all students taking its assessments.

All manuals and user guides pertaining to the 2023–2024 online testing were available on the portal, and PRs and TCs were able to use these manuals and guides when training TAs on test administration policies and procedures.

Training Modules

The following training modules were created to help users in the field understand the overall Smarter Balanced assessments and how each system works. All modules were provided in PowerPoint presentation format; and three modules were also narrated.

The Accessibility and Accommodations Module outlines the designated supports and accommodations available for the online assessments, as described in the Usability, Accessibility, and Accommodations Guidelines available on the Smarter Balanced website.

The Administering a Test Using Speech-to-Text (STT) Software Module provides an overview of key features of the STT accommodation and its functionality during testing.

The *Centralized Reporting Module* provides an overview of the key features of the CRS, which provides teachers with detailed information about their students' performance on the Smarter Balanced Interim Assessments.

The Embedded Universal Tools and Online Features Module acquaints students and teachers with the online universal tools (e.g., types of calculators, expandable text) available in the Smarter Balanced assessments.

The Individual Student Assessment Accessibility Profile (ISAAP) Module offers an overview of the Smarter Balanced Usability, Accessibility, and Accommodations Guidelines, the ISAAP Process, and the ISAAP

Tool. Smarter Balanced suggests a process and tool by which each student's needs can be matched with appropriate universal tools, designated supports, and/or accommodations.

The Performance Task Overview Module provides an introduction to the ELA/L performance task.

The *Read Aloud Module* is designed to help the read-aloud test reader understand the guidelines for the read-aloud designated support and accommodation when administering the Smarter Balanced assessments.

The Scribing Protocol Training Module is designed for test administrators acting as scribes to understand the guidelines for administering this designated support to students with this accommodation for the Smarter Balanced assessments.

The *Student Interface for Online Testing Module* explains how to navigate the Student Interface. The module includes information on how students log in to the testing system, select a test, understand the test layout, and use test tools.

The *Technology Requirements for Online Testing Module* provides current information about technology requirements, site readiness, supported devices, and CAI Secure Browser installation.

The *Test Administrator (TA) Interface for Online Testing Module* presents an overview of how to navigate the TA Interface.

The *Test Information Distribution Engine (TIDE) Module* provides an overview of the TIDE system. It includes information on logging in to TIDE and managing user accounts, student information, rosters, and testing incidents.

The *Testing with Braille Training Module* provides TAs with information on administering online tests to students using braille.

The *What Is a CAT? Module* describes the CAT and how it works when taking ELA/L and mathematics online assessments.

2.3.2 Statewide Trainings

Two series of virtual statewide trainings were held during SY 2023–2024. The first series of virtual statewide trainings was held September 18–19, 2023. The second series of virtual statewide trainings was held November 13–17, 2023. A set of in-person trainings were held January 22–February 1, 2024. These training sessions provided the information necessary for administering the Smarter Balanced assessments in ELA/L and mathematics. New TCs were provided with information on participation guidelines, test security and ethics, accessibility and accommodations, interim assessments, test administration procedures, technology requirements, the CRS, and family reports.

A separate series of trainings was held on August 29, 2023, September 12, 2023, and November 7, 2023. The training sessions held on August 29 and September 12 focused specifically on accessibility and accommodations for all Hawai'i statewide assessments, including the Smarter Balanced summative and interim assessments, while the training held on November 7 focused specifically on the administration of Braille for all Hawai'i statewide assessments.

2.4 TEST SECURITY

The security of assessment instruments and the confidentiality of student information are vital to maintaining the validity, reliability, and fairness of the test results. All test items, test materials, and student-level testing information are classified as secure materials for all assessments. The importance of maintaining test security and the integrity of test items is stressed throughout the webinar trainings and in the user guides, modules, and manuals. Various features of the TDS also protect test security. This section describes student confidentiality, system security, testing environment security, and policies on testing incidents.

2.4.1 Student-Level Testing Confidentiality

All secure websites and software systems enforce role-based security models that protect individual privacy and confidentiality in a manner consistent with the Family Educational Rights and Privacy Act (FERPA) and other federal laws. Secure transmission and password-protected access are basic features of the current system and permit authorized data access only. All aspects of the system, including item development and review, test delivery, and reporting, are secured by password-protected logins. In addition, CAI's systems use role-based security models that ensure that users access only the data to which they are entitled and may edit data according to their user rights only.

Three elements are involved in assuring that students are accessing appropriate test content, including:

- 1. Test eligibility, which refers to the assignment of a test to a particular student
- 2. *Test accommodation*, which refers to the assignment of a test setting to specific students based on student needs
- 3. *Test session*, which refers to the authentication process that TAs must follow when creating a test session, including reviewing and approving a test and its settings for each student, and the student signing on to take the test

FERPA prohibits the public disclosure of student information or test results. The following are examples of prohibited practices:

- Providing login information (usernames and passwords) to other authorized TIDE users or to unauthorized individuals
- Sending a student's name and SSID number together in an email message
- Having a student log in and test under another student's SSID number

Test materials and score reports should not be exposed to reveal student names with test scores except for authorized individuals with an appropriate need to know. If information about a test must be sent via email or fax, only the SSID number should be included, not the student's name.

All students, including homeschooled students, must be enrolled or registered at their testing schools in order to take the online, paper-pencil, or braille assessments. Student enrollment information, including demographic data, is generated using a HIDOE file and uploaded nightly via a secured file transfer site to the online TDS during the testing window.

Students log in to the online assessment using their legal first name, SSID number, and a test session ID. Only students can log in to an online test session. TAs, proctors, or other personnel are not permitted to

log in to the system on behalf of students, although they are permitted to assist students who need help logging in. For the paper-pencil versions of the assessments, TCs and TAs are required to affix the student label to each student's answer document.

After a test session, only staff with the administrative roles of PR, TC, or teacher (TE) can view their students' scores. TAs who are not also teachers do not have access to student scores.

2.4.2 System Security

The objective of system security is to ensure that all data are protected and are accessed only by the appropriate user groups. The end goal of system security entails protecting and maintaining data and system integrity, safeguarding personal information, and ensuring accurate data transfer and appropriate levels of user access.

Hierarchy of Control

As described in Section 2.2.1, Administrative Roles, PRs, TCs, and TAs have well-defined roles and levels of access to the testing system. PRs are responsible for selecting and entering the TC's information into TIDE, and the TC is responsible for entering TAs' and TEs' information into TIDE. Throughout the year, the PR and TC are also expected to delete information in TIDE for any staff members who have transferred to other schools, resigned, or no longer serve as TAs or teachers.

Password Protection

All access points by different roles—at the state, complex area, school principal, and school staff levels—require a password to log in to the system. Newly added TCs, TAs, and TEs receive separate passwords assigned by the school through their personal email addresses.

Secure Browser

A key role of the technology coordinator is to ensure that the CAI Secure Browser is installed correctly on the computers used to administer the online assessments. Developed by the testing contractor, CAI's Secure Browser prevents students from accessing other computers or Internet applications and copying test information. The Secure Browser suppresses access to commonly used browsers such as Internet Explorer and Firefox, and it prevents students from searching for answers on the Internet or communicating with other students. The assessments can be accessed only through the Secure Browser and not by other Internet browsers.

2.4.3 Security of the Testing Environment

The TCs and TAs work together to determine appropriate testing schedules based on the number of computers available, the number of students in each tested grade, and the average amount of time needed to complete each assessment.

Testing personnel are reminded in the online training and user manuals that assessments should be administered in testing rooms that have been set up to prevent students from crowding. Good lighting, ventilation, and protection from noise and other interruptions are also essential factors to consider when selecting testing rooms.

TAs must establish procedures to maintain a quiet environment during each test session, recognizing that some students may finish more quickly than others. If students are allowed to leave the testing room when

they finish their assessments, TAs must explain the procedures for leaving and where students are expected to report once they leave without disrupting others. If students are expected to remain in the testing room until the end of the session, TAs are encouraged to have students read a book after they have completed the assessment.

If a student needs to leave the room for a brief time, the TAs must pause the student's assessment. If a pause lasts longer than 20 minutes during the CAT component, the student can continue the assessment in a new test session. However, the system will not allow the student to return to the items answered before the pause. This measure is implemented to prevent students from using the time spent outside the testing room to look up answers.

Room Preparation

The testing room should be prepared before the start of the test session. Any information displayed on bulletin boards, chalkboards, or charts that students might use to answer test questions should be removed or covered. This rule applies to rubrics, vocabulary charts, student work, posters, graphs, content-area strategy charts, etc. All cell phones belonging to testing personnel and students must be turned off and stored out of sight in the testing room. TAs are encouraged to minimize access to the testing rooms by posting signs in halls and entrances to promote optimal testing conditions; they should also post "TESTING—DO NOT DISTURB" signs on the doors of testing rooms.

Seating Arrangements

TAs should provide adequate spacing between students' seats. Student seating should be arranged to prevent them from looking at other students' answers. Because the online CAT is adaptive, it is unlikely that students will see the same test questions as other students; however, students should be discouraged from communicating through appropriate seating arrangements. For the ELA/L performance task, different forms are distributed throughout the testing room so that students are less likely to receive the same forms as their neighbors.

After the Test

At the end of a test session, TAs must walk through the classroom to pick up any scratch paper that students used and any papers that display students' SSID numbers and names together. These materials should be securely shredded or stored in a locked area immediately. The printed reading passages and questions for any content-area assessment provided for a student allowed to use this accommodation in an individual setting must also be shredded immediately after a test session ends.

For the paper-pencil tests, specific instructions on how to package and secure the test booklets for return to the testing contractor's office are provided in the paper-pencil *Test Administration Manual*.

2.4.4 Test Security Violations

Every individual who administers or proctors the assessments is responsible for understanding the required security procedures associated with administering the assessments. The *Smarter Balanced Online Summative Test Administration Manual* outlines and categorizes prohibited testing practices into three groups, described here.

Impropriety: This is a test security incident that has a low impact on the individual or group of students who are testing and has a low risk of potentially affecting student performance on the test, test security, or test validity (e.g., student[s] leaving the testing room without authorization).

Irregularity: This is a test security incident that affects an individual or group of students who are testing and may potentially affect student performance on the test, test security, or test validity (e.g., a disruption during the test session, such as a fire drill). These circumstances can be contained at the local level.

Breach: This is a test security incident that poses a threat to the validity of the test. Breaches require immediate attention and escalation to the state agency. Examples include exposure of secure materials or a repeatable security/system risk (e.g., administrators modifying student answers, students sharing test items through social media). These circumstances have external implications.

Complex and school personnel are required to document all test security incidents in the test security incident log. This log is the document of record for all test security incidents and should be maintained at the complex level and submitted to HIDOE at the end of testing.

2.5 STUDENT PARTICIPATION

All students enrolled in grades 3–8 and 11 at public or public charter schools in Hawai'i are required to participate in the Smarter Balanced ELA/L and mathematics summative assessments, except the following:

- Students with significant cognitive disabilities who meet the criteria for a state-selected or state-developed ELA/L and mathematics alternate assessment based on the extensions of the Common Core standards (approximately 1% or fewer of the student population)
- Students in the English language learner (ELL) program whose first U.S. school in the past 12 months is a Hawai'i public or public charter school
- Students enrolled in the Hawaiian Language Immersion Program in grades 3–8

Only students in these three categories can be excused from taking the Smarter Balanced ELA/L assessments (all three categories) and/or the Smarter Balanced mathematics assessments (categories one and three). Students must be tested in the enrolled grade assessment; out-of-grade-level testing is not allowed for the administration of Smarter Balanced assessments.

2.5.1 Homeschooled Students

Students who are homeschooled may participate in the Smarter Balanced assessments at the request of their parent or guardian. If requested, schools must provide these students with one testing opportunity for each relevant content area.

2.5.2 Exempt Students

The following categories of students are exempt from participating in the Smarter Balanced assessments based on required documentation:

- A student who has a significant medical emergency
- A student who is receiving services at an out-of-state residential program
- An ELL who has moved to the country within the year (ELA/L exemption only)

 A student who meets the requirements of Regulation 4140, Exceptions to Compulsory School Attendance

2.6 ONLINE TESTING FEATURES AND TESTING ACCOMMODATIONS

The Smarter Balanced Assessment Consortium's *Usability, Accessibility, and Accommodations Guidelines (Guidelines)* are intended for school-level personnel and decision-making teams, including Individualized Education Program (IEP) and Section 504 Plan teams, as they prepare for and implement the Smarter Balanced assessments. The *Guidelines* provide information for classroom teachers, English language development educators, special education teachers, and instructional assistants to select and administer universal tools, designated supports, and accommodations for students who need them. The *Guidelines* are also intended for assessment staff and administrators who oversee the decisions made in instruction and assessment.

The *Guidelines* apply to all students. They emphasize an individualized approach to the implementation of assessment practices for students who have diverse needs and participate in large-scale content assessments. The *Guidelines* focus on universal tools, designated supports, and accommodations for the Smarter Balanced assessments of ELA/L and mathematics. At the same time, the *Guidelines* support important instructional decisions about accessibility and accommodations for students who participate in the Smarter Balanced assessments.

The summative assessments contain universal tools, designated supports, and accommodations in both embedded and non-embedded formats. Embedded resources are part of the computer administration system, whereas non-embedded resources are provided outside of that system.

State-level users, TCs, and teachers can set embedded and non-embedded designated supports and accommodations based on their user role in TIDE. Designated supports and accommodations must be set in TIDE prior to starting a test session.

All the embedded and non-embedded universal tools will be activated for use by all students during a test session. Before students begin testing, one or more of the preselected universal tools can be deactivated by a TC in TIDE or a TA in the TA Interface of the testing system for a student who may be distracted by the ability to access a specific tool during a test session.

For additional information about the availability of designated supports and accommodations, refer to the Smarter Balanced *Usability, Accessibility, and Accommodations Guidelines* at: https://smarterbalanced.alohahsap.org/resource-item/en/usability-accessibility-and-accommodations-guidelines-2023-2024.

2.6.1 Online Universal Tools for All Students

Universal tools are access features of an assessment or exam that are embedded or non-embedded components of the test administration system. Universal tools are available to all students based on their preference and selection and have been preset in TIDE. In the 2023–2024 test administration, the following universal tools were available for *all* students to access. For specific information on how to access and use these features, refer to the *Smarter Balanced Online, Summative, Test Administration Manual* at: https://smarterbalanced.alohahsap.org/resource-list/en/smarter-balanced-summative-test-administration-manual-2023-2024.

Embedded Universal Tools

Breaks (Pause). A student can pause the assessment and return to the test question that he or she was working on. However, if an assessment is paused for more than 20 minutes, students will not be allowed to return to previously attempted test questions.

Calculator (for calculator-allowed mathematics items only in grades 6–8, 11). This is an embedded on-screen digital calculator for calculator-allowed items that students can access by clicking the calculator button. This tool is available only with specific items that the Smarter Balanced item specifications have indicated as appropriate.

Digital Notepad. This tool is used for making notes about an item. The digital notepad is item-specific and is available through the end of the test segment. Notes are not saved when the student moves on to the next segment or after a break of more than 20 minutes.

English Dictionary. An English dictionary is available for the full-write portion of an ELA/L performance task. A full-write is the second component of a performance task.

English Glossary. This feature displays grade- and context-appropriate definitions of specific construct-irrelevant terms in English on the screen via a pop-up. The student can access the embedded glossary by clicking any of the pre-selected terms.

Expandable Passages and/or Stimuli. Each passage or stimulus can be expanded to take up a larger portion of the screen.

Global Notes. Global notes is a notepad that is available for the ELA/L performance task in which students complete a full-write. Students click the notepad icon for the notepad to appear. During the ELA/L performance task, the notes are retained from segment to segment and allow a student to return to the notes even though he or she cannot go back to specific items in the previous segment.

Highlighter. This tool is used to mark desired text, test questions, item answers, or parts of these with color. An enhanced highlighting feature allows multiple color options. Highlighted text remains available throughout each test segment. This tool is not available while the Line Reader tool is in use.

Keyboard Navigation. This tool allows students to navigate text using a keyboard.

Line Reader. Students use an onscreen universal tool to assist in reading by raising and lowering the tool for each line of text on the screen. If the enhanced line reader mode is enabled, all content except for the line in focus is grayed out for greater emphasis. This tool is not available while the Highlighter tool is in use.

Mark for Review. Students can mark a question for review in order to return to it later. However, for the CAT, if the assessment is paused for more than 20 minutes, students are not allowed to return to marked test questions.

Math Tools. These digital tools (e.g., embedded ruler, embedded protractor) are used for measurements related to mathematics items. They are available only with the specific items that the Smarter Balanced item specifications have indicated that one or more of these tools are appropriate.

Spellcheck. This is a writing tool for checking the spelling of words in student-generated responses. Spellcheck indicates only that a word is misspelled; it does not provide the correct spelling. This tool is

available only with the specific items that the Smarter Balanced item specifications have indicated as appropriate. Spellcheck is bundled with other embedded writing tools for all performance task full-write items: planning, drafting, revising, and editing.

Strikethrough. This feature allows the student to cross out answer options. If an answer option is an image, a strikethrough line will not appear, but the image will be grayed out.

Thesaurus. A thesaurus is available for the full-write portion of an ELA/L performance task. A full-write is the second part of a performance task.

Writing Tools. Selected writing tools (e.g., bold, italic, bullets, undo, redo) are available for all student-generated responses. (Also, refer to spellcheck.)

Zoom. Students can zoom in on test questions, text, or graphics. This tool makes these features appear larger on the screen.

Non-Embedded Universal Tools

Breaks. Breaks may be given at predetermined intervals or after completion of sections of the assessment for students taking a paper-pencil test. Sometimes students can take breaks when individually needed to reduce cognitive fatigue when they experience heavy assessment demands. The use of this universal tool may result in the student needing additional overall time to complete the assessment.

English Dictionary. An English dictionary can be provided for the full-write portion of an ELA/L performance task. A full-write is the second part of a performance task. The use of this universal tool may result in the student needing additional time to complete the assessment.

Scratch Paper. Scratch paper to make notes, write computations, or record responses may be made available. Only plain paper or lined paper is appropriate for ELA/L. Graph paper is required beginning in grade 6 and can be used on all mathematics assessments. A student may use an assistive technology device for scratch paper as long as the device is consistent with the child's IEP and acceptable to the State.

Thesaurus. A thesaurus provides synonyms of terms while a student interacts with text included in the assessment. This tool is available for the full-write portion of an ELA/L performance task. A full-write is the second part of a performance task. The use of this universal tool may result in the student needing additional time to complete the assessment.

2.6.2 Designated Supports and Accommodations

Designated supports for the Smarter Balanced assessments are features available for use by any student for whom the need has been indicated by an educator (or team of educators with the parent or guardian and student). Scores achieved by students using designated supports will be included for federal accountability purposes. It is recommended that a consistent process be used to determine which supports should be designated for individual students. All educators making these decisions should be trained to use this process and should be made aware of the range of available designated supports. Smarter Balanced members have identified digitally embedded and non-embedded designated supports for students for whom an adult or team has indicated a need for the support.

Accommodations are modifications in procedures or materials that increase equitable access during the Smarter Balanced assessments. Assessment accommodations generate valid assessment results for students who need them; they allow these students to show what they know and can do. Accommodations

are available only for students with documented IEPs or Section 504 Plans. Consortium-approved accommodations do not compromise the learning expectations, construct, grade-level standard, or intended outcome of the assessments.

Embedded Designated Supports

Color Contrast. Students can adjust the screen background or font color based on their needs or preferences. This may include reversing the colors for the entire interface or choosing the color of the font and background. Black on white, reverse contrast, black on rose, medium gray on light gray, and yellow on blue were offered for the online assessments.

Illustration Glossaries (for mathematics items). Illustration glossaries are provided for selected construct-irrelevant terms for mathematics. Illustrations for these terms appear on the computer screen when students select them. Students can also adjust the size of the illustration and move it around the screen. Only students with the illustration glossary setting enabled can use this accommodation.

Masking. Masking involves blocking off content that is not of immediate need or that may be distracting to the student. This tool allows students to focus their attention on a specific part of a test item.

Mouse Pointer. This support allows the mouse pointer to be set to a larger size and for the color to be changed. A TA sets the size and color of the mouse pointer prior to testing.

Streamline. This accommodation provides a streamlined interface of the test in an alternative, simplified format in which the items are displayed below the stimuli.

Text-to-Speech (for mathematics stimuli and items and ELA/L items; not for ELA/L reading passages). Text is read aloud to the student via embedded text-to-speech technology. The student can control the speed and raise or lower the volume of the voice via a volume control. This support is also available in Spanish for mathematics tests when students have a Spanish language support selected.

Text-to-Speech in Spanish (for mathematics stimuli and items). Text is read aloud to the student via embedded text-to-speech technology in Spanish. The student can control the speed and raise or lower the volume of the voice via a volume control.

Translated Student Interface Messages (for mathematics tests in Spanish). Translation of the student interface messages is a language support available prior to beginning the actual test items. Students can see test directions in Spanish. As an embedded designated support, translated test directions are automatically a part of the Spanish language translations designated support.

Translations (Glossaries) (for mathematics items). Translated glossaries are a language support. The translated glossaries are provided for selected construct-irrelevant terms in mathematics. Translations for these terms appear on the computer screen when students click them. The following language glossaries were offered: Arabic, Burmese, Cantonese, Filipino, Hmong, Korean, Mandarin, Punjabi, Russian, Somali, Spanish, Ukrainian, and Vietnamese.

Translations (Spanish) (for mathematics items). Dual language translations are a linguistic support available for some students; dual language translations provide the full translation of each test item above the original English language version of the item.

Turn Off Any Universal Tools. A TA may disable any universal tools that might be distracting, that students do not need to use, or that students are unable to use.

Non-Embedded Designated Supports

Amplification. Students may adjust the volume control beyond the computer's built-in settings using headphones or other non-embedded devices.

Bilingual Dictionary. The bilingual/dual-language word-to-word dictionary is a language support that can be provided for the full-write portion of an ELA/L performance task.

Color Contrast. Test content of online items may be printed with different colors.

Color Overlays. Color transparencies may be placed over a paper-pencil assessment.

Illustration Glossaries (for mathematics paper-pencil tests). The illustration glossaries are a language support provided for selected construct-irrelevant terms for mathematics. Illustrations for these terms appear in a supplement to the paper-pencil test and are identified by item number.

Magnification. The size of specific areas of the screen (e.g., text, formulas, tables, graphics, navigation buttons) may be adjusted by the student with an assistive technology device. Magnification allows students to increase the size of images and text on the screen to a level not allowed by the universal Zoom tool.

Medical Supports. Students may have access to an electronic device for medical purposes (e.g., glucose monitor). The device may include a cell phone and should support the student for medical reasons only during testing.

Noise Buffers. Ear mufflers, white noise, and/or other equipment that reduces environmental noises may be used.

Printed Test Directions in English. Available as a supplement to the TAM, a printed copy of oral test directions in English may be provided to the student. The use of this support may result in the student needing additional overall time to complete the assessment.

Read-Aloud (for mathematics stimuli and items and ELA/L items; not for ELA/L reading passages). The text is read aloud to the student by a trained and qualified human reader who follows the administration guidelines provided in the Smarter Balanced Online Summative Test Administration Manual and the Guidelines for Read Aloud, Test Reader. All or portions of the content may be read aloud.

Read-Aloud in Spanish (for mathematics, all grades). Spanish text is read aloud to the student by a trained and qualified human reader who follows the administration guidelines provided in the Smarter Balanced Online Summative Test Administration Manual and the Guidelines for Read-Aloud, Test Reader. All or portions of the content may be read aloud.

Scribe (for all items except ELA/L PT full-writes). Students dictate their responses to a human who records verbatim what they dictate. The scribe must be trained and qualified and must follow the administration guidelines provided in the Smarter Balanced Online Summative Test Administration Manual.

Separate Setting. The test location is altered so that the student is tested in a setting different from that made available to most students.

Simplified Test Directions. The TA simplifies or paraphrases the test directions found in the test administration manual according to the Simplified Test Directions guidelines.

Translated Student Interface Messages. A bilingual adult may read aloud a PDF file of directions translated in each of the languages currently supported.

Translated Test Directions in American Sign Language (ASL). Test directions that include test administration scripts are translated into ASL video. The ASL human signer and the signed test content are viewed at the same time. Students may view portions of the ASL video as often as needed.

Translations (Glossaries) (for mathematics paper-pencil tests). Translated glossaries are a language support provided for selected construct-irrelevant terms for mathematics. Glossary terms are listed by item and include the English term and its translated equivalent.

Embedded Accommodations

American Sign Language (ASL) (for ELA/L listening items and mathematics items). This accommodation allows test content to be translated into an ASL video. An ASL human signer and the signed test content are viewed on the same screen. Students may view portions of the ASL video as often as needed.

Braille. This is a raised-dot code that individuals read with their fingertips. Graphic material (e.g., maps, charts, graphs, diagrams, illustrations) is presented in a raised format (paper or thermoform). Contracted and non-contracted braille is available; Nemeth Braille Code is available for mathematics.

Braille Transcript (for ELA/L listening passages). This is a braille transcript of the closed captioning created for the listening passages. The braille transcripts are available in uncontracted and contracted English Braille American Edition (EBAE).

Closed Captioning (for ELA/L listening items). Printed text may appear on the computer screen as audio materials are presented.

Text-to-Speech (for ELA/L reading passages). Text is read aloud to the student via embedded text-to-speech technology. The student can control the speed and raise or lower the volume of the voice via a volume control.

Word Prediction. This allows students to begin writing a word and choose from a list of words that have been predicted from word frequency and syntax rules. Word prediction is delivered via an embedded software program. The program must use only single-word prediction. Functionality such as phrase prediction, predict ahead, or next word must be deactivated. The program must have settings that allow only a basic dictionary. Expanded dictionaries, such as topic dictionaries and word banks, must be deactivated. Phonetic spelling functionality and programs with built-in speech output that reads back the information the student has written may also be used. Students who use word prediction in conjunction with speech output will need headphones unless tested individually in a separate setting.

Non-Embedded Accommodations

100s Number Table. A paper-based table listing numbers 1–100 is available for reference.

Abacus. This tool may be used in place of scratch paper for students who typically use an abacus.

Alternate Response Options. Alternate response options include but are not limited to adapted keyboards, large keyboards, Sticky Keys, MouseKeys, FilterKeys, adapted mouse, touch screen, head wand, and switches.

Braille (paper-pencil assessment). This is a raised-dot code that individuals read with their fingertips. Graphic material (e.g., maps, charts, graphs, diagrams, illustrations) is presented in a raised format (paper or thermoform). The following codes are available for the ELA/L paper-pencil assessment: EBAE uncontracted, EBAE contracted, Unified English Braille (UEB) uncontracted, and UEB contracted. The following codes are available for the mathematics paper-pencil assessment: EBAE uncontracted with Nemeth Braille Code, EBAE contracted with Nemeth, UEB uncontracted with Nemeth, UEB contracted with Nemeth, UEB uncontracted with UEB mathematics.

Calculator (for calculator-allowed items mathematics items only in grades 6–8, 11). This is a non-embedded calculator for students needing a special calculator, such as a braille calculator or a talking calculator, currently unavailable in the assessment platform.

Math Manipulatives. This accommodation allows eligible students with IEPs and Section 504 Plans to represent their understanding of mathematical concepts using visual and tactile concrete materials. This list of approved mathematics manipulatives that may be provided on-site includes Algebra Tiles (recommended for grade 6 and above), Base Ten Blocks, Colored Tiles, Geoblocks Set, Geoboards and Geobands, Multi-Link Cubes, Pop Cubes, or Similar Cubes, Multi-Sensory Learning (MSL) Kit, One-Inch Blocks, Pattern Blocks, Transparent Sheets, and Two-Color Counters. Up to four manipulatives may be selected for a student; other accommodations not listed can be requested for verification.

Multiplication Table. A paper-based single digit (1–9) multiplication table is available for reference.

Print-on-Demand. This accommodation allows TAs to print paper copies of either passages/stimuli and/or items for students. For students needing a paper copy of a passage or stimulus, permission for the students to request printing must first be set in TIDE. The TC must fill out a Verification of Student Need Form and contact HIDOE to have the accommodation set for the student.

Read-Aloud (for ELA/L reading passages). Text is read aloud to the student via an external screen reader or by a trained and qualified human reader who follows the administration guidelines provided in the Smarter Balanced Online Summative Test Administration Manual and Read-Aloud Guidelines. All or portions of the content may be read aloud. Refer to the Guidelines for Choosing the Read-Aloud Accommodation when deciding if this accommodation is appropriate for a student.

Scribe (for ELA/L PT full-write items). Students dictate their responses to a human who records verbatim what they dictate. The scribe must be trained and qualified and must follow the administration guidelines provided in the Smarter Balanced Online Summative Test Administration Manual.

Speech-to-Text. Voice recognition allows students to use their voices as input devices to the computer in order to dictate responses or give commands (e.g., opening application programs, pulling down menus, saving work). Voice recognition software generally can recognize speech up to 160 words per minute. Students may use their own assistive technology devices.

Word Prediction. This allows students to begin writing a word and choose from a list of words that have been predicted from word frequency and syntax rules. Word prediction is delivered via a non-embedded software program. The program must use only single-word prediction. Functionality such as phrase prediction, predict ahead, or next word must be deactivated. The program must have settings that allow only a basic dictionary. Expanded dictionaries, such as topic dictionaries and word banks, must be deactivated. Phonetic spelling functionality and programs with built-in speech output that reads back the information the student has written may also be used. Students who use word prediction in conjunction

with speech output will need headphones unless tested individually in a separate setting. Students may use their own assistive technology devices.

Table 5 presents a list of universal tools, designated supports, and accommodations that were offered in the 2023–2024 administration. Tables 6–11 provide the numbers of students who utilized any of the offered accommodations and designated supports. Note that the overall count in the designated support tables may not match the sum of students in ELL and students with disabilities because some students are counted in both categories or because these features were approved for some students other than ELL and students with disabilities.

Table 5. SY 2023–2024 Universal Tools, Designated Supports, and Accommodations

Universal Tools	Designated Supports	Accommodations
	Embedded	
Breaks (Pause) Calculator ¹ Digital Notepad English Dictionary ² English Glossary Expandable Passages and/or Stimuli Global Notes ³ Highlighter Keyboard Navigation Line Reader Mark for Review Math Tools ⁴	Color Contrast Illustration Glossaries ⁶ Masking Mouse Pointer Streamline Text-to-Speech ⁷ Translated Student Interface Messages ⁶ Translated Test Directions ⁶ Translations (Glossaries) ⁶ Translations (Spanish) ⁶ Turn Off Any Universal Tools	American Sign Language ⁸ Braille Braille Transcript ⁹ Closed Captioning ⁹ Text-to-Speech ¹⁰ Word Prediction
Spellcheck Strikethrough Thesaurus ² Writing Tools ⁵ Zoom	Non-Embedded	
Breaks	Amplification	100s Number Table
English Dictionary ²	Bilingual Dictionary ²	Abacus
Scratch Paper	Color Contrast	Alternate Response Options ¹⁴
Thesaurus ²	Color Overlay	Braille ¹⁵
	Illustration Glossaries ¹¹	Calculator ¹
	Magnification	Math Manipulatives ¹⁶
	Medical Supports	Multiplication Table
	Noise Buffers	Print-on-Demand
	Printed Test Directions in English	Read-Aloud ¹⁷
	Read-Aloud ¹²	Scribe ²
	Read-Aloud in Spanish ⁶	Speech-to-Text
	Scribe ¹³	Word Prediction
	Separate Setting	
	Simplified Test Directions	
	Translated Student Interface	
	Messages	
	Translated Test Directions in ASL	
	Translations (Glossaries) ¹¹	

^{*} Items shown are available for ELA/L and mathematics unless otherwise noted.

¹ For calculator-allowed mathematics items only in grades 6–8 and 11

² For ELA/L performance task full-write items

Table 6. ELA/L Total Students with Allowed Embedded and Non-Embedded Accommodations

Accommodations				Grade			
Accommodations	3	4	5	6	7	8	11
	Emb	edded Acco	mmodation	ıs			
American Sign Language	8	4	2	5	2	6	6
Braille			1	1			
Braille Transcript	5		1	1			5
Closed Captioning	9	10	14	12	10	13	17
Text-to-Speech: Reading Passages and Items	6	4	7	5	5	7	4
Word Prediction		1	1				
	Non-E	mbedded A	ccommodat	ions			
Alternate Response Options		1					
Print-on-Demand: Stimuli & Items		2		1	1		
Read-Aloud Passages	6	3	7	2	1		4
Scribe (Full-Write)	7	7	2	3	4		1
Speech-to-Text	4	12	9	7	7	7	

³ For ELA/L performance tasks

⁴ Includes embedded ruler, embedded protractor

⁵ Includes bold, italic, underline, indent, cut, paste, spellcheck, bullets, undo, redo

⁶ For mathematics items

⁷ For mathematics stimuli and items and ELA/L items (not for ELA/L reading passages): must be set in TIDE before test begins. Available in both English and Spanish for the mathematics tests.

⁸ For ELA/L listening items and mathematics items

⁹ For ELA/L listening items

¹⁰ For ELA/L reading passages. Must be set in TIDE by state-level user. TCs must submit a student's Verification of Need form to the Assessment Section for review and approval or disapproval.

¹¹ For mathematics paper-pencil tests

¹² For mathematics stimuli and items and ELA/L items (not for ELA/L reading passages)

¹³ For all items except for ELA/L performance task full-writes

¹⁴ Includes adapted keyboards, large keyboard, Sticky Keys, MouseKeys, FilterKeys, adapted mouse, touch screen, head wand, and switches

¹⁵ For paper-pencil assessments

¹⁶ Includes Algebra Tiles (recommended for grade 6 and above), Base Ten Blocks, Colored Tiles, Geoblocks Set, Geoboards and Geobands, Multi-Link Cubes, Pop Cubes, or Similar Cubes, Multi-Sensory Learning (MSL) Kit, One-Inch Blocks, Pattern Blocks, Transparent Sheets, and Two-Color Counters

¹⁷ For ELA/L reading passages, all grades

Table 7.ELA/L Total Students with Allowed Embedded Designated Supports

Designated Commonts	Ch	Grade						
Designated Supports	Subgroup	3	4	5	6	7	8	11
	Overall	1	1	21	8	8	1	1
Color Contrast	ELL			1		2		
	Disability	3 4 5 6 7 1 1 21 8 8 1 8 1 8 100 26 185 44 20 6 1 5 4 6 24 16 38 34 15 3 53 2 6 3 1 2 3 9 2 6 101 29 40 134 27 5 3 5 28 7 25 16 19 82 23 3,003 3,022 3,161 2,209 1,385 694 764 764 579 401 768 929 949 681 487 27 22 31 18 5 1 2 3 5 3 3 1 2 6 5	1	1				
	Overall	100	26	185	44	20	19	_
Masking	ELL	6	1	5	4	6	2	
	Disability	24	16	38	34	15	17	
	Overall		3	53	2	6		
Mouse Pointer	ELL			3	1	2		
	Disability		3	9	2	6	1 1 1 1 1 1 1 2 5 1 7 2 4 4 3 1 8 5 1,444 1 484 7 467 5 5 3 1 1 90 1,450 5 486	
	Overall	101	29	40	134	27	24	3
Streamline	ELL	5	3	5	28	7	4	
	Disability	25	16	19	82	23	18	3
	Overall	3,003	3,022	3,161	2,209	1,385	1,444	47
Text-to-Speech: CAT Items	ELL	694	764	764	579	401	484	13
	Disability	768	929	949	681	487	467	37
	Overall	27	22	31	18	5	5	8
Text-to-Speech: PT Items	ELL	1	2	3		2		4
	Disability	6	5	9	8	4	5	5
	Overall	8	5	2	3	5	3	
Text-to-Speech: PT Stimuli	ELL	3	3		1		1	
	Disability	2			2	4	1	
The state of the part of the state of the st	Overall	3,169	3,057	3,178	2,198	1,390	1,450	45
Text-to-Speech: PT Stimuli and	ELL	726	773	769	581	405	486	1 1 3 3 47 13 37 8 4 5
Items	Disability	794	947	953	682	492	473	37

Table 8. ELA/L Total Students with Allowed Non-Embedded Designated Supports

Designated Supports	Ch				Grade	rade			
	Subgroup	3	4	5	6	7	8	11	
	Overall	1	2		2			4	
Amplification	ELL		1		1			1	
-	Disability	1	1					4	
	Overall	24	10	17	15	11	4	21	
Bilingual Dictionary	ELL	22	10	14	12	11	4	17	
	Disability	1		1		4		6	
	Overall			1	1				
Color Contrast	ELL								
	Disability				1				
	Overall	1	2	2	3	3	3	2	
Magnification	ELL				1		1		
	Disability		2	2	3	2	1	1	
	Overall	2	1	2	2	1	4	1	
Medical Supports	ELL			1					
	Disability	1		2		1			

Designated Supports	C 1	-			Grade			
	Subgroup	3	4	5	6	7	8	11
	Overall	3		1	19	3	3	
Noise Buffers	ELL				3			
	Disability	3		1	5	2	2	
D: (1T (D: ())	Overall		1	1	15	9	6	
Printed Test Directions in English	ELL			1	12	9	6	
Engrish	Disability			1		5		
	Overall	130	105	107	19	6	9	18
Read-Aloud Items	ELL	34	19	17	4	2	1	7
	Disability	57	63	55	17	6	9	16
	Overall	124	98	101	13	3	9	13
Read-Aloud Stimuli	ELL	32	17	17	3	1	1	6
	Disability	52	58	50	12	3	9	12
	Overall	8	13	9	4	4		1
Scribe (Not Full-Write)	ELL	1	2	1	1	1		
	Disability	6	11	6	3	4		1
	Overall	385	415	382	231	143	139	52
Separate Setting	ELL	65	67	60	31	26	18	5
	Disability	252	288	294	182	2 1 6 9 3 9 1 1 3 9 4 1 4 143 139 26 18 121 119 24 32 10 7 19 25	119	43
	Overall	223	214	145	70	24	32	24
Simplified Test Directions	ELL	46	25	25	9	10	7	2
_	Disability	86	90	64	47	19	25	22
Translated Student Interface	Overall	8	4	3	3	1		5
	ELL	7	3	3	2	1		2
Messages	Disability		1		1			4
m 1 · 1m · D: · · ·	Overall	5	2	1				5
Translated Test Directions in	ELL	2						1
ASL	Disability	5	2	1				5

Table 9. Mathematics Total Students with Allowed Embedded and Non-Embedded Accommodations

Assemmedations	Grade								
Accommodations	3	4	5	6	7	8	11		
	Emb	edded Acco	mmodation	18					
American Sign Language	3	5	2	5	2	5	5		
Braille			1	1					
	Non-Ei	nbedded A	ccommodat	ions					
100s Number Table	29	27	29	23	12	2			
Abacus		1		1	1				
Alternate Response Options		1							
Calculator			1	2	5	2			
Math Manipulatives	27	34	26	17	7	2			
Multiplication Table		2	12	11	11	3			
Print-on-Demand: Stimuli & Items		2		1	1				
Speech-to-Text	4	14	9	7	7	7	1		

Table 10. Mathematics Total Students with Allowed Embedded Designated Supports

D : 10	G 1				Grade			
Designated Supports	Subgroup	3	4	5	6	7	8	11
	Overall	1	1	14	1	6		1
Color Contrast	ELL					1		
	Disability		1	6	1	6		1
	Overall	48	75	125	238	206	219	
Illustration Glossaries	ELL	31	47	46	143	150	163	
	Disability	5	11	19	60	52	57	
	Overall	96	30	186	42	18	15	
Masking	ELL	6	1	6	5	5	2	
-	Disability	23	20	39	34	14	12	
	Overall		3	53	2	6		
Mouse Pointer	ELL			3	1	2		
	Disability		3	9	2	6		
	Overall	99	30	42	134	22	23	3
Streamline	ELL	5	3	6	28	2	3	
	Disability	28	18	20	81	22	17	3
	Overall	10	4	8	5	4	5	1
Text-to-Speech: Items	ELL	2		1		1		
-	Disability	5	2	6	5	4	5	1
	Overall	1	7	2	3	1	1	
Text-to-Speech: Stimuli	ELL				1	1		
-	Disability		2		1	1	1	
	Overall	3,283	3,205	3,283	2,283	1,405	1,433	51
Text-to-Speech: Stimuli and	ELL	766	788	787	614	417	481	15
Items	Disability	806	972	975	695	487	458	41
	Overall	7	4	3	9	17	18	1
Translations (Glossaries):	ELL	6	4	3	9	15	18	1
Spanish	Disability	2			1		1	
T. 1.4. (Cl. 1.) 2.3	Overall	9	20	31	52	59	62	
Translations (Glossaries): Other	ELL	7	19	24	51	59	60	
Languages	Disability			1	3	5	3	
	Overall	10	5	4	7	9	11	2
Translations (Spanish)	ELL	10	5	4	7	9	11	2
` -	Disability							

Table 11. Mathematics Total Students with Allowed Non-Embedded Designated Supports

Cub				Grade			
Subgroup	3	4	5	6	7	8	11
Overall	1	2		1			2
ELL		1					1
Disability	1	1					2
Overall				1			
ELL							
Disability				1			
Overall	13	19	13	22	10	2	
ELL	11	15	10	6	2	2	
Disability	4	7	9	15	9	1	
Overall	1	2	1	3	3	2	1
ELL				1			
Disability		2	1	3	2		
Overall	2	1	2	2	1	4	1
ELL			1				
Disability	1		2		1		
Overall	3		1	19	3	3	
ELL				3			
Disability	3		1	5	2	2	
Overall		1	1	13	11	6	
ELL			1	10	11	6	
Disability			1		5		
Overall	123	98	117	19	7	7	12
ELL	30	20	21	4	1		5
Disability	53	58	62	17	7	7	12
Overall	1		2	1	1		
ELL	1			1	1		
Disability	1				1		
Overall	120	96	104	13	4	7	10
ELL	29	19	19	3	1		5
Disability				12	4	7	10
Overall	2		2	1	2		
	1			1			
Disability	2		1		1		
		11	7	4	4		1
			1	1	1		
	6		6	3	4		1
•						135	51
							5
							42
							22
ELL	39	26	25	10	10	6	3
						24	22
Disability	84	90	62	40	20	Z 4	2.7.
Disability Overall	84	90	62	3	20		
Disability Overall ELL	4 1	6 4	3 3	3 2	1 1	24	4 2
	ELL Disability Overall	Overall ELL Disability Overall Sell Disability Overall Overall Sell Disability Overall Sell Disability Overall Overall Sell Disability Overall Overall Overall Sell Disability Overall Overall Overall Sell Disability Overall Overall Overall Overall Overall Overall	Overall 1 2 ELL 1 1 Disability 1 1 Overall 13 19 ELL 11 15 Disability 4 7 Overall 1 2 ELL 1 2 Disability 1 2 Overall 3 3 ELL 1 3 Disability 3 3 Overall 1 1 ELL 30 20 Disability 53 58 Overall 1 1 ELL 1 1 Disability 53 58 Overall 1 2 ELL 1 1 Disability 50 57 Overall 2 2 Overall 8 11 ELL 1 2 Disability 6 11 <td>Overall ELL Disability 1 2 Doverall ELL Disability 1 1 Overall ELL Disability 13 19 13 ELL Disability 4 7 9 Overall 1 1 2 1 ELL Disability 2 1 2 ELL Disability 1 2 1 Overall 2 1 2 1 ELL Disability 3 1 1 EVAR Disability 3 1 1 Overall 1 1 1 1 1 ELL 30 20 21 21 2 2 1 Disability 53 58 62 62 2 2 1 Overall 1 2 96 104 1 2 2 1 Disability 1 00 1 0 <td< td=""><td>Subgroup 3 4 5 6 Overall ELL Disability 1 1 1 Overall ELL Disability 1 1 1 Overall ELL Disability 13 19 13 22 ELL Disability 4 7 9 15 Overall Disability 1 2 1 3 22 1 3 1 1 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 1 <t< td=""><td> Subgroup 3</td><td>Subgroup 3 4 5 6 7 8 Overall 1 2 1 2 1 3 2 2 2 1 3 3 2 2 1 3 3 2 1 2 2 2 2 2 2 1 3 2 2 1 4 4 4 7 9 1 5 9 1 1 1 1 1 1 1 1 1 1 1 1</td></t<></td></td<></td>	Overall ELL Disability 1 2 Doverall ELL Disability 1 1 Overall ELL Disability 13 19 13 ELL Disability 4 7 9 Overall 1 1 2 1 ELL Disability 2 1 2 ELL Disability 1 2 1 Overall 2 1 2 1 ELL Disability 3 1 1 EVAR Disability 3 1 1 Overall 1 1 1 1 1 ELL 30 20 21 21 2 2 1 Disability 53 58 62 62 2 2 1 Overall 1 2 96 104 1 2 2 1 Disability 1 00 1 0 <td< td=""><td>Subgroup 3 4 5 6 Overall ELL Disability 1 1 1 Overall ELL Disability 1 1 1 Overall ELL Disability 13 19 13 22 ELL Disability 4 7 9 15 Overall Disability 1 2 1 3 22 1 3 1 1 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 1 <t< td=""><td> Subgroup 3</td><td>Subgroup 3 4 5 6 7 8 Overall 1 2 1 2 1 3 2 2 2 1 3 3 2 2 1 3 3 2 1 2 2 2 2 2 2 1 3 2 2 1 4 4 4 7 9 1 5 9 1 1 1 1 1 1 1 1 1 1 1 1</td></t<></td></td<>	Subgroup 3 4 5 6 Overall ELL Disability 1 1 1 Overall ELL Disability 1 1 1 Overall ELL Disability 13 19 13 22 ELL Disability 4 7 9 15 Overall Disability 1 2 1 3 22 1 3 1 1 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 2 2 1 3 1 <t< td=""><td> Subgroup 3</td><td>Subgroup 3 4 5 6 7 8 Overall 1 2 1 2 1 3 2 2 2 1 3 3 2 2 1 3 3 2 1 2 2 2 2 2 2 1 3 2 2 1 4 4 4 7 9 1 5 9 1 1 1 1 1 1 1 1 1 1 1 1</td></t<>	Subgroup 3	Subgroup 3 4 5 6 7 8 Overall 1 2 1 2 1 3 2 2 2 1 3 3 2 2 1 3 3 2 1 2 2 2 2 2 2 1 3 2 2 1 4 4 4 7 9 1 5 9 1 1 1 1 1 1 1 1 1 1 1 1

Designated Supports	Sub avaus	_			Grade			
Designated Supports	Subgroup	3	4	5	6	7	8	11
T 1. 1T (D) (; ;	Overall	4	1	1				4
Translated Test Directions in ASL	ELL	1						1
ASL	Disability	4	1	1				4
T. 1. (Cl)	Overall	6	2	5	1			1
Translations (Glossaries): Spanish	ELL	5	2	3	1			1
Spanish	Disability	1						
T. 1 (Cl) O	Overall	1	1	4	1	1	1	1
Translations (Glossaries): Other	ELL	1	1	3	1	1	1	1
Languages	Disability							1

2.7 TESTING TIME

The online environment allows item response time to be captured as the item page time (i.e., the time each item page is presented on the screen) in milliseconds. For discrete items, each item appears on the screen one item at a time, whereas stimulus-based items appear on the screen together. For discrete items, the page time is the time spent on one item; and, for stimulus-based items, it is the time spent on all items associated with a stimulus. For each student, the total time taken to complete the test is computed by adding up the page time for all items and item groups (stimulus-based items).

The Smarter Balanced summative assessments are not timed, and an individual student may need more or less time than average overall. The length of a test session is determined by PRs or TCs who are knowledgeable about the class periods in the school's instructional schedule and the timing needs associated with the assessments. Students should be allowed extra time if they need it, but TAs must use their best professional judgment when allowing students extra time.

Tables 12 and 13 present the average testing time and the testing time at percentiles for the overall test, the computer-adaptive test (CAT) component, and the performance task (PT) component.

Table 12. Test-Taking Time: ELA/L

	Average	SD of	Median		Testing Time	e in Percenti	iles (hh:mm)	
Grade	Testing Time (hh:mm)	Testing Time (hh:mm)	Testing Time (hh:mm)	75th	80th	85th	90th	95th
				Overall Te	st			
3	2:35	1:42	2:12	3:17	3:36	4:01	4:41	5:42
4	3:07	2:07	2:37	3:53	4:20	4:52	5:40	7:19
5	3:18	2:06	2:48	4:08	4:33	5:05	5:55	7:24
6	3:00	1:47	2:35	3:44	4:05	4:33	5:13	6:20
7	2:40	1:28	2:22	3:20	3:37	4:00	4:31	5:24
8	2:41	1:27	2:23	3:24	3:44	4:07	4:38	5:27
11	1:57	0:59	1:50	2:26	2:36	2:51	3:08	3:40
			(CAT Compo	nent			
3	0:52	0:32	0:45	1:03	1:09	1:16	1:27	1:49
4	0:59	0:40	0:49	1:10	1:17	1:26	1:40	2:08
5	1:03	0:39	0:53	1:15	1:22	1:33	1:48	2:14
6	1:04	0:36	0:58	1:17	1:23	1:31	1:44	2:07
7	0:57	0:29	0:52	1:09	1:14	1:20	1:30	1:46
8	0:56	0:28	0:52	1:09	1:14	1:20	1:29	1:45
11	0:46	0:21	0:44	0:56	1:00	1:04	1:11	1:21
				PT Compon	ent			
3	1:43	1:21	1:24	2:15	2:31	2:53	3:23	4:15
4	2:08	1:40	1:44	2:45	3:06	3:35	4:12	5:27
5	2:15	1:39	1:52	2:53	3:15	3:40	4:16	5:27
6	1:55	1:22	1:34	2:29	2:45	3:09	3:39	4:33
7	1:43	1:11	1:27	2:13	2:27	2:48	3:12	3:58
8	1:45	1:09	1:29	2:15	2:32	2:52	3:17	4:01
11	1:12	0:45	1:04	1:33	1:41	1:51	2:06	2:32

Table 13. Test-Taking Time: Mathematics

-	Average	SD of	Median		Testing Time	e in Percenti	iles (hh:mm)	
Grade	Testing Time (hh:mm)	Testing Time (hh:mm)	Testing Time (hh:mm)	75th	80th	85th	90th	95th
			Overall T	Test (CAT C	Component)			
3	0:49	0:30	0:42	1:00	1:06	1:14	1:26	1:48
4	0:58	0:38	0:48	1:11	1:20	1:29	1:44	2:14
5	1:05	0:42	0:55	1:21	1:28	1:39	1:54	2:23
6	1:02	0:35	0:54	1:15	1:22	1:31	1:44	2:08
7	1:00	0:33	0:54	1:14	1:20	1:27	1:39	1:59
8	1:04	0:34	0:59	1:21	1:28	1:36	1:47	2:06
11	0:44	0:22	0:41	0:55	0:58	1:03	1:11	1:24

2.8 DATA FORENSICS PROGRAM

The validity of test scores depends on the integrity of the test administration. Any irregularities in test administration could cast doubt on the validity of the inferences based on those test scores. Multiple facets ensure that tests are administered properly, including clear test administration policies, effective TA training, and tools to identify possible irregularities in test administrations.

For online administrations, a set of quality assurance (QA) reports is generated during and after the testing window. One of the QA reports focuses on flagging possible testing anomalies. Testing anomalies are analyzed by examining changes in student performance from year to year, test-taking time, item response patterns using a person-fit index, and item response change analyses.

Analyses are performed at the student level and summarized for each aggregate unit, including the testing session, TA, and school. The flagging criteria used for these analyses are described in the following section and are configurable by an authorized user. When the aggregate unit size is small, the aggregate unit is flagged if the percentage of flagged students is greater than 50% in the analysis. The default small aggregate unit size is five or fewer students, but this value is configurable. For each aggregate unit, small groups are identified based on the number of tests included in the aggregate unit from that analysis. Thus, a small unit identified in one analysis may not be a small unit in another analysis. The QA reports are provided to state clients to monitor testing anomalies after the testing window closes.

2.8.1 Changes in Student Performance

Changes in student scores between administration years are examined using a regression model to check for outliers. For these between-year comparisons, students' current-year scores are regressed on their test scores from the previous year and on the number of days between the two years' test-end dates (to control for the instruction time between the two test scores).

A large score gain or loss in student scores between administration years is detected by examining the residuals for outliers. The residuals are computed as the observed value minus the regression model's predicted value. The studentized residuals are computed to detect unusual residuals. An unusual increase or decrease in student scores between administration years is flagged when the absolute value of the studentized residual is greater than 3.

The residuals of students are also aggregated for a testing session, TA, and school. The system flags any unusual changes in an aggregate performance between administrations and/or years based on the average of the residuals in the aggregate unit (e.g., testing session, TA, school). For each aggregate unit, a t value is computed and flagged when |t| is greater than 3,

$$t = \frac{\sum_{i=1}^{n} \hat{e}_i / n}{\sqrt{\frac{s^2}{n} + \frac{\sum_{i=1}^{n} \sigma^2 (1 - h_{ii})}{n^2}},$$

where s is the standard deviation of residuals in an aggregate unit; n is the number of students in an aggregate unit (e.g., testing session, TA, school), σ^2 is the MSE from the regression, h_{ii} is the leverage from the regression for the ith student, and \hat{e}_i is the residual for the ith student.

The variance of average residuals in the denominator is estimated in two components, conditioning on the true residual e_i , $var(E(\hat{e}_i|e_i)) = s^2$ and $E(var(\hat{e}_i|e_i)) = \sigma^2(1 - h_{ii})$. Following the law of total variance (Billingsley, 1995, p. 456),

$$var(\hat{e}_i) = var(E(\hat{e}_i|e_i)) + E(var(\hat{e}_i|e_i)) = s^2 + \sigma^2(1 - h_{ii}), \text{ hence,}$$
$$var(\frac{\sum_{i=1}^n \hat{e}_i}{n}) = \frac{\sum_{i=1}^n (s^2 + \sigma^2(1 - h_{ii}))}{n^2} = \frac{s^2}{n} + \frac{\sum_{i=1}^n (\sigma^2(1 - h_{ii}))}{n^2}.$$

2.8.2 Test-Taking Time

The summative assessments are not timed, and thus, individual test-taking times may vary across students. However, unusual test-taking times such as excessively shorter or longer test-taking times may indicate irregularities in test administration. An example of an unusual test-taking time is a test record for an individual who scores very well on the test even though the average time spent is far less than that required of students statewide. If students already know the answers to the questions, the test-taking time may be much shorter than the test-taking time for those who have no prior knowledge of the item content. Conversely, if a TA helps students by coaching them to change their responses during the test, the testing time could be longer than expected.

The state average testing time and standard deviation are computed based on all students available when the analysis was performed. Students and aggregate units are flagged if the test-taking time is different from the state average by three standard deviations or more, although the flagging criteria can be adjusted by an authorized user.

2.8.3 Inconsistent Item Response Pattern (Person Fit)

In item response theory (IRT) models, person-fit measurement is used to identify test takers whose response patterns are improbable given an IRT model. If a test has psychometric integrity, little irregularity will be seen in the item responses of the individual who responds to the items fairly and honestly.

If a test taker has prior knowledge of some test items (or is provided answers during the exam), he or she will respond correctly to those items at a higher probability than indicated by his or her ability as estimated across all items. In this case, the person-fit index will be large for the student. However, if a student has prior knowledge of the entire test content, this will not be detected based on the person-fit index, although the item response time index might flag such a student.

The person-fit index is based on all item responses in a test. An unlikely response to a single test question may not result in a flagged person-fit index. Of course, not all unlikely patterns indicate cheating, as in the case of a student who is able to guess a significant number of correct answers. Therefore, the evidence of person-fit index should be evaluated along with other testing irregularities to determine possible testing irregularities. The number of flagged students is summarized for every testing session, TA, and school.

The person-fit index is computed using a standardized log-likelihood statistic. Following Drasgow, Levine, and Williams (1985) and Sotaridona, Pornell, and Vallejo (2003), an aberrant response pattern is defined as a deviation from the expected item score model. Snijders (2001) showed that the distribution of l_z is asymptotically normal (i.e., with an increasing number of administered items). Even at shorter test lengths of 8 or 15 items, the "asymptotic error probabilities are quite reasonable for nominal Type I error probabilities of 0.10 and 0.05" (Snijders, 2001).

Sotaridona et al. (2003) report promising results of using l_z for systematic flagging of aberrant response patterns. Students with l_z values less than -3 are flagged. Aggregate units are flagged with t less than -3,

$$t = \frac{\text{Average } l_z \text{ values}}{\sqrt{s^2/n}},$$

where s = standard deviation of l_z values in an aggregate unit and n = number of students in an aggregate unit.

2.8.4 Item-Response Change

Students are allowed to revisit items as many times as they wish within a session and may also mark items to be revisited prior to completing the session. However, excessively high rates of response change, especially high rates of item score increases (i.e., response changes from wrong to right), may indicate irregularities in test administration. For example, TAs could review students' responses and either coach them to modify their responses or keep the session active and change responses themselves.

To identify irregular patterns of response change, the item score for the final response to each item and the penultimate response if one exists are examined, and the number of instances in which the item score increases are counted.

The average and standard deviation of positive item score changes are computed based on all students available when the analysis was performed. Students and aggregate units are flagged if the number of positive item score changes is larger than the state average by three standard deviations or more, although the flagging criteria can be adjusted by an authorized user.

2.9 Prevention and Recovery of Disruptions in the Test Delivery System

CAI is continuously improving its ability to protect testing systems from interruptions. CAI's TDS is designed to ensure that student responses are captured accurately and stored on more than one server in case of a failure. The CAI architecture, described in the following section, is designed to recover from a failure of any component with little interruption. Each system is redundant, and crucial student response data are transferred to a different data center each night.

CAI has developed a unique monitoring system that is extremely sensitive to changes in server performance. Most monitoring systems provide warnings when something is going wrong. The CAI system does, too, but it also provides warnings when any given server performs differently from its performance over the few hours prior or differently than the other servers performing the same jobs. Subtle changes in performance often precede actual failure by hours or days, allowing CAI to detect potential problems, investigate them, and mitigate them. This system has enabled CAI to make adjustments and replace equipment on multiple occasions before any problems occurred.

CAI has also implemented an escalation procedure to alert clients within minutes of any disruption. The emergency alert system notifies CAI's executive and technical staff by text message, who then immediately join a call to identify and address the problem.

The following section describes CAI's system architecture and how it recovers from device failures, Internet interruptions, and other problems.

2.9.1 High-Level System Architecture

Our architecture provides the redundancy, robustness, and reliability required by a large-scale, high-stakes testing program. The general approach, which Smarter Balanced has adopted as standard policy, is pragmatic and well supported by the system architecture.

CAI posits that any system built around an expectation of the flawless performance of computers or networks within schools and complex areas is bound to fail. Therefore, the system is designed to ensure that the testing results and experience respond robustly to such inevitable failures. CAI's TDS is designed to protect data integrity and prevent student data loss at every point throughout the test administration process. Fault tolerance and automated recovery are built into every component of the system.

The key elements of the testing system, including the data integrity processes, are described in the following paragraphs.

Student Machine

Student responses are conveyed to CAI's servers in real time as students respond. Long responses, such as essays, are saved automatically at configurable intervals (usually set to one minute) so that student work is not at risk of being unrecorded during testing.

Responses are saved asynchronously, with a background process on the student machine waiting to confirm that the data has been successfully stored on the server. If confirmation is not received within the designated time (usually 30–90 seconds), the system will prevent the student from completing more work until connectivity is restored. The student is offered the choice of asking the system to try again or pausing the test and completing it at another time. For example:

- If connectivity is lost and restored within the designated time, the student may be unaware of the momentary interruption.
- If connectivity cannot be silently restored, the student is prevented from testing and given the option of logging out or retrying the save.
- If the system fails completely, upon logging back into the system, the student returns to the item at which the failure occurred.

In short, data integrity is preserved by confirmed saves to CAI servers and the prevention of further testing if confirmation is not received.

Test Delivery Satellites

The test delivery satellites communicate with the student machines to deliver items and receive responses. Each satellite is a collection of web and database servers. Each satellite is equipped with a redundant array of independent disks (RAID) systems to mitigate the risk of disk failure. Each response is stored on multiple independent disks.

One server operates as a backup hub for every four satellites. This server continually monitors and stores all changed student response data from the satellites, creating an additional copy of the real-time data. In the unlikely event of failure, data are completely protected. Satellites are automatically monitored, and they are removed from service upon failure. Real-time student data are immediately recoverable from the satellite, backup hub, or hub (as described in the following paragraphs), with backup copies remaining on the drive arrays of the disabled satellite.

If a satellite fails, students will exit the system. The automatic recovery system enables students to log in again within seconds or minutes of the failure without data loss. The hub manages this process. Data will remain on the satellites until the satellite receives notice from the demographic and history servers that the data are safely stored on those disks.

Hub

Hub servers are redundant clusters of database servers with RAID drive systems. Hub servers continuously gather data from the test delivery satellites and their mini-hubs and store that data as described earlier. This real-time backup copy remains on the hub until the hub receives a notification from the demographic and history servers that the data have reached the designated storage location.

Demographic and History Servers

The demographic and history servers store student data for the duration of the testing window. They are clustered database servers, also equipped with RAID subsystems, providing the redundant capability to prevent data loss in the event of server or disk failure. At the normal conclusion of a test, these servers receive completed tests from the test delivery satellites. Once the data are successfully stored, these servers notify the hub and satellites that it is safe to delete student data.

Quality Assurance System

The QA system gathers data that detect cheating, monitor real-time item function, and evaluate test integrity. Every completed test runs through the QA system, and any anomalies (such as unscored or missing items, unexpected test lengths, or other unlikely issues) are flagged. A notification then goes out to CAI's psychometricians and project team immediately.

Database of Record

The Database of Record (DOR) is the final storage location for the student data. These clustered database servers equipped with RAID systems hold the completed student data.

2.9.2 Automated Backup and Recovery

Industry-standard backup and recovery procedures are in place to ensure the safety, security, and integrity of all data, and every system is backed up nightly. This set of systems and processes is designed to provide complete data integrity and prevent the loss of student data. Redundant systems at every point, real-time data integrity protection and checks, and well-considered real-time backup processes prevent the loss of student data, even in the unlikely event of system failure.

2.9.3 Other Disruption Prevention and Recovery Mechanisms

These testing systems are designed to be extremely fault-tolerant. The systems can withstand the failure of any component with little or no service interruption. This robustness is archived through redundancy. Key redundant systems are as follows:

- The system's hosting provider has redundant power generators that operate for up to 60 hours without refueling. In addition, with multiple refueling contracts in place, these generators can operate indefinitely.
- The hosting provider has multiple redundancies in the flow of information to and from the system's

data centers through their partnership with nine different network providers. Each fiber carrier must enter the data center at separate physical points, protecting the data center from a complete service failure caused by an unlikely network cable cut.

- At the network level, there are redundant firewalls and load balancers throughout the environment.
- The system uses redundant power and switching in all server cabinets.
- Data are protected by nightly backups. A full weekly backup and incremental nightly backups protect data. Should a catastrophic event occur, CAI can reconstruct real-time data using the data retained on the TDS satellites and hubs.
- The server backup agents send alerts to notify system administration staff in the event of a backup error, at which time they will inspect the error to determine whether the backup was successful or if they need to rerun the backup.

To summarize, the system's TDS is hosted in an industry-leading facility with redundant power, cooling systems, state-of-the-art security, and other features that protect the system from failure. The system is redundant at every component, and in the event of failure, the unique design ensures that data are always stored in at least two locations. The engineering that led to this system protects student responses from loss.

3. SUMMARY OF 2023–2024 OPERATIONAL TEST ADMINISTRATION

3.1 STUDENT POPULATION

All students enrolled in grades 3–8 and 11 in all public elementary and secondary schools must participate in the Smarter Balanced English language arts/literacy (ELA/L) and mathematics assessments. Before the testing window opened for the 2023–2024 test administration, the state or complex area sends CAI a student enrollment file to load to the Test Information Distribution Engine (TIDE). Using this enrollment file, the participation rates were calculated as the percentage of students who attempted the test. Tables 14 and 15 present the participation rates and the percentage of students who attempted the test by subgroups. Tables 16 and 17 present the number of Hawai'i students who met attemptedness requirements for scoring and reporting the results of the Smarter Balanced summative assessments.

Table 14. Participation Rates by Percentage: ELA/L

Group	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	Grade 11
All Students	94.7	95.5	95.3	94.6	93.8	94.4	91.2
Female	95.2	95.8	95.4	94.6	93.9	94.6	91.1
Male	94.3	95.2	95.3	94.7	93.6	94.2	91.3
African American	97.7	96.8	97.8	95.4	98.1	99.4	93.9
AmerIndian/Alaskan	100.0	72.7	90.9	90.0	96.0	88.9	84.2
Asian/Pacific Islander	97.0	97.8	97.5	96.6	97.2	97.3	95.0
Hispanic	95.4	95.3	95.3	94.4	93.1	93.4	91.0
Hawai'i Pacific Islander	89.1	91.4	90.8	90.6	88.7	90.6	85.0
White	96.8	97.9	97.6	96.5	96.5	96.1	91.0
Multi-Racial	96.9	96.7	97.0	96.7	94.8	95.2	93.0
ELL	94.1	93.6	93.2	91.5	92.0	92.3	84.0
Disadvantaged	94.6	94.8	94.8	93.9	92.1	92.8	88.6
Migrant	98.7	93.6	97.5	97.0	97.7	95.7	90.1
Disability	85.7	87.7	89.3	86.9	83.9	85.4	76.1

Note. AmerIndian/Alaskan = American Indian/Alaskan Native; ELL = English Language Learner; Disadvantaged = Economic Disadvantage Status

Table 15. Participation Rates by Percentage: Mathematics

Group	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	Grade 11
All Students	95.2	95.8	95.7	95.3	94.5	94.9	91.3
Female	95.7	96.0	95.8	95.3	94.6	94.9	91.2
Male	94.7	95.6	95.7	95.3	94.3	94.9	91.4
African American	97.7	97.3	97.8	95.4	98.7	99.4	94.5
AmerIndian/Alaskan	100.0	72.7	90.9	90.0	100.0	94.4	73.7
Asian/Pacific Islander	98.3	98.6	98.6	98.0	98.2	98.1	95.6
Hispanic	95.6	95.5	95.4	94.9	93.6	93.8	90.7
Hawai'i Pacific Islander	89.5	91.7	91.3	91.0	89.4	91.0	85.2
White	97.1	98.1	97.5	96.8	96.9	96.6	90.6
Multi-Racial	96.9	96.8	97.2	97.0	95.4	95.8	92.8
ELL	97.9	96.4	96.7	96.2	95.0	95.1	86.3
Disadvantaged	95.1	95.1	95.3	94.6	93.1	93.3	88.3
Migrant	98.7	93.6	97.5	96.4	98.4	96.8	88.7
Disability	86.0	88.0	89.4	87.4	84.6	86.2	76.1

Table 16. Number of Students: ELA/L

Group	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	Grade 11
All Students	12,256	12,785	13,141	12,400	12,167	12,202	10,884
Female	5,983	6,165	6,320	5,992	5,859	5,925	5,240
Male	6,273	6,620	6,821	6,408	6,308	6,277	5,644
African American	172	186	139	151	158	167	154
AmerIndian/Alaskan	16	8	21	18	24	16	17
Asian/Pacific Islander	2,695	2,895	3,050	2,928	3,104	3,235	3,497
Hispanic	2,363	2,431	2,626	2,403	2,347	2,293	1,868
Hawai'i Pacific Islander	2,782	3,008	3,037	2,971	2,878	2,943	2,384
White	1,479	1,515	1,427	1,315	1,322	1,284	1,061
Multi-Racial	2,749	2,742	2,841	2,614	2,334	2,263	1,899
ELL	1,549	1,569	1,451	1,248	1,314	1,344	717
Disadvantaged	5,634	5,877	5,838	5,477	5,438	5,168	3,940
Migrant	153	133	154	160	126	177	135
Disability	1,195	1,356	1,443	1,323	1,316	1,274	910

Table 17. Number of Students: Mathematics

Group	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	Grade 11
All Students	12,317	12,831	13,189	12,479	12,257	12,270	10,893
Female	6,016	6,182	6,344	6,033	5,907	5,946	5,248
Male	6,301	6,649	6,845	6,446	6,350	6,324	5,645
African American	172	186	139	150	159	167	156
AmerIndian/Alaskan	16	8	21	18	25	17	14
Asian/Pacific Islander	2,730	2,919	3,084	2,971	3,135	3,260	3,521
Hispanic	2,369	2,436	2,627	2,417	2,359	2,303	1,860
Hawai'i Pacific Islander	2,798	3,018	3,050	2,984	2,902	2,957	2,390
White	1,483	1,517	1,423	1,317	1,327	1,292	1,056
Multi-Racial	2,749	2,747	2,845	2,622	2,350	2,273	1,892
ELL	1,594	1,563	1,473	1,311	1,347	1,358	732
Disadvantaged	5,670	5,899	5,860	5,519	5,497	5,204	3,930
Migrant	154	134	154	158	127	180	133
Disability	1,202	1,365	1,445	1,333	1,326	1,289	909

3.2 SUMMARY OF OVERALL STUDENT PERFORMANCE

Tables 18–23 present a summary of the 2023–2024 summative test results for all students and by subgroup, including the average and the standard deviation of scale scores, the percentage of students in each achievement level, and the percentage of proficient students.

Figures 1 and 2 present the percentage of proficient students over the past five test administrations for all students (cohort comparisons). Figures 3 and 4 present the average scale scores in five test administrations for all students. In Figures 1–4, the 2019–2020 performance is not included because the testing was canceled due to the COVID-19 pandemic.

Appendix B, Student Performance Across Four Years for All Students and by Subgroup, provides the average and standard deviations of scale scores and the percentage of proficient students by subgroup for each test administration across four years.

Table 18. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and by Subgroup: ELA/L (Grades 3–5)

	Number	Scale Score	Scale	%	%	%	%	%
Group	Tested	Mean	Score SD	Level 1	Level 2	Level 3	Level 4	Proficient
			Grade 3					
All Students	12,256	2425.29	102.63	29	22	22	27	49
Female	5,983	2434.85	100.69	25	22	23	30	53
Male	6,273	2416.17	103.64	32	22	21	25	46
African American	172	2424.33	84.78	22	33	26	20	45
AmerIndian/Alaskan	16	2420.70	82.71	19	38	25	19	44
Asian/Pacific Islander	2,695	2456.05	99.77	19	20	23	38	61
Hispanic	2,363	2411.51	96.72	33	23	22	21	44
Hawai'i Pacific Islander	2,782	2375.20	95.72	47	24	17	11	29
White	1,479	2455.17	95.50	18	19	25	38	63
Multi-Racial	2,749	2441.68	100.64	23	21	24	32	56
ELL	1,549	2372.04	94.47	50	23	17	11	28
Disadvantaged	5,634	2394.73	98.52	40	24	19	17	36
Migrant	153	2375.08	91.00	47	27	14	12	26
Disability	1,195	2325.57	85.55 Grade 4	71	19	6	4	11
All Students	12,785	2466.17	106.19	31	19	21	28	49
Female	6,165	2476.94	102.51	28	19	23	30	53
Male	6,620	2476.34	102.51	35	19	20	26	46
African American	186	2464.50	91.60	32	21	23	24	47
AmerIndian/Alaskan	8*	2404.30	91.00	32	21	23	24	47
Asian/Pacific Islander	2,895	2495.89	103.40	21	18	22	40	62
Hispanic	2,431	2448.27	100.60	37	22	19	22	41
Hawai'i Pacific Islander	3,008	2418.83	101.70	49	20	17	13	31
White	1,515	2501.59	99.57	19	17	27	37	64
Multi-Racial	2,742	2483.27	101.32	25	18	24	32	56
ELL	1,569	2403.47	96.92	54	21	16	10	26
Disadvantaged	5,877	2431.68	101.30	44	21	19	17	35
Migrant	133	2415.69	101.61	50	19	21	10	31
Disability	1,356	2353.52	87.44	77	15	6	3	9
			Grade 5					
All Students	13,141	2510.89	109.69	26	19	28	27	56
Female	6,320	2524.85	105.32	21	19	29	31	60
Male	6,821	2497.95	112.06	30	19	28	24	51
African American	139	2511.78	103.89	22	21	31	26	57
AmerIndian/Alaskan	21	2519.89	83.18	19	10	52	19	71
Asian/Pacific Islander	3,050	2544.12	105.67	16	16	31	38	69
Hispanic	2,626	2495.44	105.02	30	21	29	21	50
Hawai'i Pacific Islander	3,037	2458.03	103.92	44	20	23	12	35
White	1,427	2541.22	100.90	15	19	31	36	67
Multi-Racial	2,841	2530.64	104.92	19	18	30	33	63
ELL	1,451	2426.45	95.04	56	21	18	5	23
Disadvantaged	5,838	2475.26	106.68	37	21	25	16	42
Migrant	154	2442.84	103.26	47	23	23	7	31
Disability	1,443	2386.92	92.35	73	16	8	3	11

^{*} Suppressed the data due to the small sample size, n < 10.

Table 19. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and by Subgroup: ELA/L (Grades 6–8)

	Number	Scale Score	Scale	%	%	%	%	%
Group	Tested	Mean	Score SD	Level 1	Level 2	Level 3	Level 4	Proficient
			Grade 6					
All Students	12,400	2530.55	104.77	25	24	30	21	52
Female	5,992	2545.92	101.89	20	23	32	26	58
Male	6,408	2516.18	105.39	29	25	29	18	46
African American	151	2519.26	97.67	24	30	29	17	46
AmerIndian/Alaskan	18	2493.33	102.67	39	17	33	11	44
Asian/Pacific Islander	2,928	2560.63	98.30	15	21	34	30	64
Hispanic	2,403	2516.78	102.40	29	25	28	17	46
Hawai'i Pacific Islander	2,971	2480.42	98.74	42	26	23	9	32
White	1,315	2568.72	98.07	13	18	37	32	69
Multi-Racial	2,614	2548.21	99.91	18	24	34	25	58
ELL	1,248	2441.21	85.45	58	26	14	2	16
Disadvantaged	5,477	2499.88	101.23	34	27	26	13	39
Migrant	160	2482.20	99.20	43	22	24	11	35
Disability	1,323	2413.54	83.58	73	19	7	2	9
4.11 G . 1	10.167	0545.50	Grade 7	26	22	22	10	50
All Students	12,167	2547.52	111.67	26	22	33	19	52
Female	5,859	2564.85	104.85	20	22	36	22	58
Male	6,308	2531.42	115.34	32	22	31	16	46
African American	158	2555.94	105.96	23	20	37 50	19	56
AmerIndian/Alaskan	24	2571.05	89.64	13	21	50	17	67
Asian/Pacific Islander	3,104	2584.08	105.50	15	18	38	28	67
Hispanic	2,347	2528.98	106.09	30	25	32	12	45
Hawai'i Pacific Islander	2,878	2491.71	104.92	44	27	23	7	29
White	1,322	2589.33	100.07	14	18 21	39	29 22	68
Multi-Racial	2,334	2561.88	107.94	22		35		57
ELL Di la 1	1,314	2457.29	97.65	58	25 25	15	2	17
Disadvantaged	5,438 126	2511.91	108.06	37 44	25 25	28	10	38
Migrant Disability	1,316	2490.38 2422.65	94.06 92.23	72	23 18	27 9	3 1	30 10
Disability	1,310	2422.03	92.23 Grade 8	12	10	9	1	10
All Students	12,202	2558.84	110.32	26	24	32	17	49
Female	5,925	2575.61	103.73	20	24	35	20	55
Male	6,277	2543.01	113.96	32	25	30	14	44
African American	167	2568.21	100.14	20	28	38	15	53
AmerIndian/Alaskan	16	2532.70	98.67	38	19	38	6	44
Asian/Pacific Islander	3,235	2593.57	102.53	16	21	38	25	63
Hispanic	2,293	2537.35	107.39	32	27	29	11	41
Hawai'i Pacific Islander	2,943	2505.58	102.11	43	29	22	6	28
White	1,284	2600.16	102.36	14	20	40	26	66
Multi-Racial	2,263	2576.24	107.01	20	23	36	21	56
ELL	1,344	2478.62	90.54	52	30	16	1	17
Disadvantaged	5,168	2523.88	106.65	37	27	27	9	36
Migrant	177	2498.86	98.84	47	27	23	3	27
Disability	1,274	2435.73	91.95	73	19	7	1	8

Table 20. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and by Subgroup: ELA/L (Grade 11)

Group	Number Tested	Scale Score Mean	Scale Score SD	% Level 1	% Level 2	% Level 3	% Level 4	% Proficient
	Testeu	Ivican	Grade 11	LCVCII	LCVCI 2	<u> Level 5</u>	LCVCI 4	Troncicit
All Students	10,884	2596.91	116.21	19	23	34	25	58
Female	5,240	2616.49	106.51	13	22	37	28	65
Male	5,644	2578.74	121.75	25	23	31	21	52
African American	154	2591.83	104.89	14	29	39	18	57
AmerIndian/Alaskan	17	2598.88	116.38	12	29	35	24	59
Asian/Pacific Islander	3,497	2623.77	109.31	13	19	36	32	68
Hispanic	1,868	2580.23	111.15	22	25	35	19	53
Hawai'i Pacific Islander	2,384	2543.64	108.70	33	29	29	9	38
White	1,061	2633.81	118.78	12	18	31	39	70
Multi-Racial	1,899	2610.28	115.81	16	21	34	29	63
ELL	717	2489.49	91.08	51	33	15	1	16
Disadvantaged	3,940	2561.19	114.48	28	27	31	15	45
Migrant	135	2537.09	106.36	34	27	30	8	39
Disability	910	2467.06	99.26	62	26	11	1	12

Table 21. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and by Subgroup: Mathematics (Grades 3–5)

	Number	Scale Score	Scale	%	%	%	%	%
Group	Tested	Mean	Score SD	Level 1	Level 2	Level 3	Level 4	Proficient
			Grade 3					
All Students	12,317	2439.48	94.71	26	21	27	26	53
Female	6,016	2436.23	90.44	26	22	28	24	51
Male	6,301	2442.58	98.51	25	20	26	28	54
African American	172	2423.73	83.76	28	30	25	17	42
AmerIndian/Alaskan	16	2426.62	85.90	38	13	31	19	50
Asian/Pacific Islander	2,730	2472.77	91.90	15	18	28	39	67
Hispanic	2,369	2423.46	90.15	30	24	26	20	46
Hawai'i Pacific Islander	2,798	2393.06	88.25	45	24	20	11	31
White	1,483	2465.96	85.55	15	19	31	35	66
Multi-Racial	2,749	2454.24	91.25	19	20	30	30	60
ELL	1,594	2395.94	93.90	44	22	21	13	34
Disadvantaged	5,670	2410.58	92.40	37	24	24	16	40
Migrant	154	2394.85	95.35	44	25	19	12	31
Disability	1,202	2347.27	95.95 Grade 4	65	19	10	6	16
All Students	12,831	2478.81	94.09	23	29	25	23	48
Female	6,182	2474.76	88.09	23	31	26	20	46
Male	6,649	2474.70	99.20	23	27	24	26	50
African American	186	2472.58	81.72	21	34	30	15	45
AmerIndian/Alaskan	8*	2472.36	01.72	21	37	30	13	43
Asian/Pacific Islander	2,919	2512.74	92.34	12	24	28	35	63
Hispanic	2,436	2461.81	88.25	27	34	22	17	39
Hawai'i Pacific Islander	3,018	2432.94	88.21	39	33	18	10	28
White	1,517	2507.02	87.09	13	25	29	33	62
Multi-Racial	2,747	2493.14	88.07	18	28	28	27	55
ELL	1,563	2426.08	92.22	42	33	15	9	25
Disadvantaged	5,899	2448.33	88.70	33	33	21	13	34
Migrant	134	2446.58	86.52	34	36	16	14	30
Disability	1,365	2384.50	85.49	65	24	8	3	11
			Grade 5					
All Students	13,189	2507.27	103.11	31	25	18	26	44
Female	6,344	2505.12	98.20	31	28	18	23	41
Male	6,845	2509.25	107.42	31	23	19	28	46
African American	139	2500.10	94.14	27	33	19	21	40
AmerIndian/Alaskan	21	2506.26	88.99	29	29	29	14	43
Asian/Pacific Islander	3,084	2548.24	99.42	17	22	21	39	61
Hispanic	2,627	2489.51	98.19	36	28	17	18	35
Hawai'i Pacific Islander	3,050	2455.28	93.95	50	26	14	10	23
White	1,423	2529.55	93.87	21	25	20	33	53
Multi-Racial	2,845	2524.19	98.88	24	24	21	31	51
ELL	1,473	2436.60	92.05	59	26	9	7	16
Disadvantaged	5,860	2475.09	98.86	42	27	15	15	30
Migrant	154	2439.82	92.83	56	25	10	8	19
Disability	1,445	2399.83	90.42	76	15	6	3	9

^{*} Suppressed the data due to the small sample size, n < 10.

Table 22. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and by Subgroup: Mathematics (Grades 6–8)

Group	Number	Scale Score	Scale	% L 1.1	% L 1.2	%	% L 1.4	% Day 6 a day
•	Tested	Mean	Score SD Grade 6	Level 1	Level 2	Level 3	Level 4	Proficient
All Students	12,479	2516.17	115.28	33	28	18	21	39
Female	6,033	2517.26	111.55	33	28	18	21	39
Male	6,446	2517.20	111.55	34	27	18	21	39
African American	150	2517.86	100.54	35	31	18	17	35
AmerIndian/Alaskan	18	2449.04	122.24	61	17	17	6	22
Asian/Pacific Islander	2,971	2554.69	112.27	21	26	21	32	53
Hispanic	2,417	2496.10	112.97	41	27	16	16	32
Hawai'i Pacific Islander	2,984	2459.30	109.02	53	28	12	7	19
White	1,317	2558.28	103.10	18	27	23	32	55
Multi-Racial	2,622	2534.96	104.92	26	29	21	24	45
ELL	1,311	2425.00	105.09	68	21	7	4	11
Disadvantaged	5,519	2482.11	111.47	45	28	15	12	27
Migrant	158	2473.66	104.87	46	30	12	11	23
Disability	1,333	2390.75	101.15	80	14	4	2	6
	-,		Grade 7				_	
All Students	12,257	2519.63	121.11	38	26	19	17	36
Female	5,907	2518.06	116.87	38	27	18	16	35
Male	6,350	2521.10	124.91	38	24	19	18	37
African American	159	2519.32	111.98	31	35	16	17	33
AmerIndian/Alaskan	25	2541.13	66.75	20	44	24	12	36
Asian/Pacific Islander	3,135	2567.40	120.21	24	24	22	30	53
Hispanic	2,359	2496.63	112.22	44	28	17	11	28
Hawai'i Pacific Islander	2,902	2456.50	106.37	60	25	10	5	15
White	1,327	2560.00	109.00	24	27	26	24	50
Multi-Racial	2,350	2533.97	115.76	33	27	21	20	40
ELL	1,347	2428.47	107.65	71	19	7	3	10
Disadvantaged	5,497	2480.15	113.05	51	27	14	8	23
Migrant	127	2442.71	116.28	65	24	8	3	11
Disability	1,326	2395.98	95.85	82	14	3	1	4
			Grade 8					
All Students	12,270	2527.55	125.97	44	24	16	16	32
Female	5,946	2528.97	121.25	43	26	16	16	32
Male	6,324	2526.21	130.25	45	23	15	17	32
African American	167	2527.95	109.00	40	33	12	15	27
AmerIndian/Alaskan	17	2442.45	120.36	65	24	12	0	12
Asian/Pacific Islander	3,260	2574.02	126.28	29	25	20	26	46
Hispanic	2,303	2500.39	116.52	52	26	13	10	22
Hawai'i Pacific Islander	2,957	2464.64	106.96	65	21	9	5	13
White	1,292	2572.39	115.18	28	25 25	23	24	47
Multi-Racial	2,273	2545.29	123.08	38	25	17	19	37
ELL	1,358	2442.04	105.07	75	17	6	3	9
Disadvantaged	5,204	2488.03	116.49	57	23	11	8	20
Migrant	180	2472.71	99.29	61	27	8	4	12
Disability Note The record of soft	1,289	2401.14	100.33	86	11	2	1	3

Table 23. Descriptive Statistics and Percentage of Students in Achievement Levels for Overall and by Subgroup: Mathematics (Grade 11)

Group	Number Tested	Scale Score Mean	Scale Score SD	% Level 1	% Level 2	% Level 3	% Level 4	% Proficient
	Testeu	1,10411	Grade 11	Beveri	Ecverz	Ecvere	Ec ver i	Troncicii
All Students	10,893	2544.39	123.26	51	24	16	8	25
Female	5,248	2548.87	115.92	49	26	17	8	25
Male	5,645	2540.22	129.59	52	23	16	9	25
African American	156	2520.01	101.89	61	26	9	4	13
AmerIndian/Alaskan	14	2569.61	130.99	36	29	21	14	36
Asian/Pacific Islander	3,521	2580.80	123.97	38	27	22	13	35
Hispanic	1,860	2522.76	110.83	59	24	12	4	17
Hawai'i Pacific Islander	2,390	2486.07	103.62	72	20	7	2	8
White	1,056	2578.06	130.88	40	24	21	15	36
Multi-Racial	1,892	2554.33	121.24	46	26	20	8	28
ELL	732	2456.93	96.93	84	13	3	1	4
Disadvantaged	3,930	2507.60	113.94	64	21	11	4	14
Migrant	133	2495.31	96.47	72	19	7	2	9
Disability	909	2426.14	90.13	90	8	1	0	2

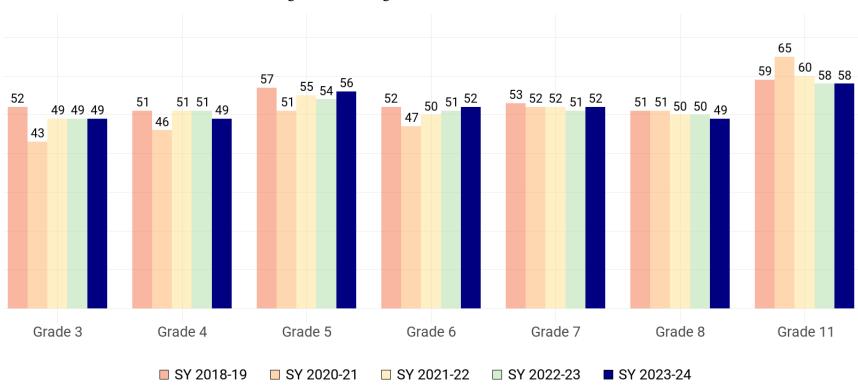
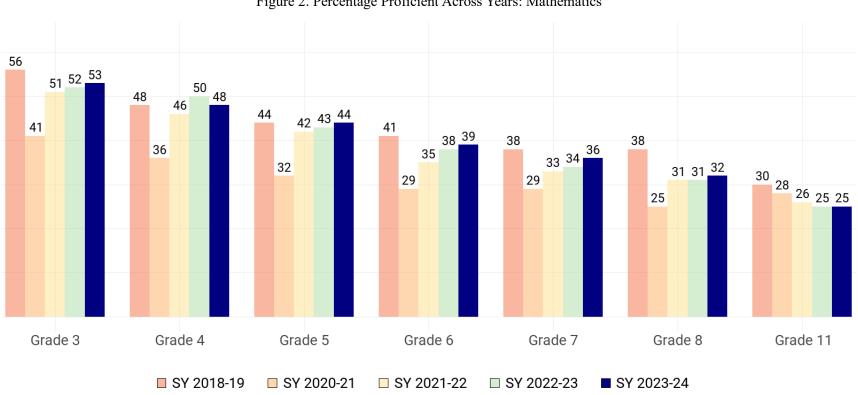


Figure 1. Percentage Proficient Across Years: ELA/L



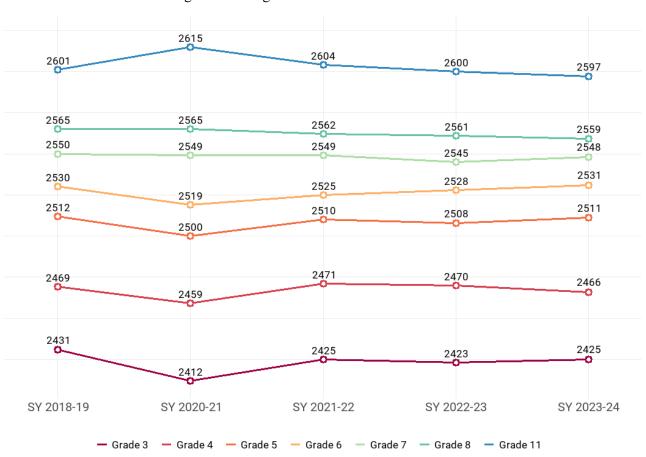


Figure 3. Average Scale Score Across Years: ELA/L

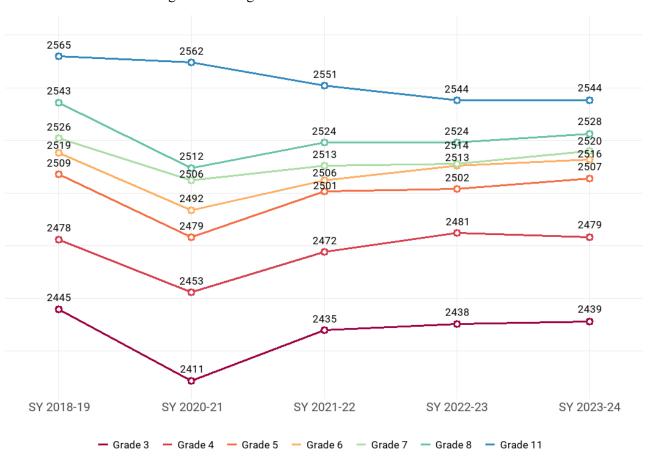


Figure 4. Average Scale Score Across Years: Mathematics

Because the precision of scores in each claim is not sufficient to report scores, given a small number of items, the scores on each claim are reported using one of the three performance categories, taking into account the standard error of measurement (SEM) of the claim score: (1) Below Standard, (2) At/Near Standard, or (3) Above Standard (see Section 7.5, Rules for Calculating Strengths and Weaknesses for Claim Scores, for the rules). Given the reduction in the number of items in Hawai'i's shortened blueprints, the reliabilities for claim scores are low, especially for Claim 3 and Claim 4 in ELA/L and Claims 2 and 4 combined and Claim 3 in mathematics. Therefore, starting with 2021–2022, the performance category for claim scores were reported only for Claims 1 and 2 in ELA/L and Claim 1 in mathematics at individual student level. Table 24 presents the distribution of performance categories for the reported claims.

Table 24. Percentage of Students in Performance Categories by Claim

	Performance	EI	A/L	Mathematics
Grade	Category	Claim 1 Reading	Claim 2 Writing	Claim 1 Concepts and Procedures
	Below	22	28	27
3	At/Near	60	51	40
	Above	18	21	33
	Below	19	28	29
4	At/Near	61	54	40
	Above	20	18	30
	Below	18	24	32
5	At/Near	61	51	40
	Above	21	25	27
	Below	27	26	39
6	At/Near	54	54	38
	Above	20	20	23
	Below	24	25	42
7	At/Near	58	51	38
	Above	18	24	21
	Below	28	27	42
8	At/Near	54	55	40
	Above	18	19	17
	Below	19	19	55
11	At/Near	58	53	33
	Above	23	29	12

3.3 DISTRIBUTION OF STUDENT ABILITY AND ITEM DIFFICULTY

Figures 5–10 display the empirical distribution of the Hawai'i student scale scores in the 2023–2024 test administration and the distribution of the administered summative item-difficulty parameters for each grade for overall and by claim. For overall, the student ability distribution shifted to the left in all grades and subjects, a pattern more pronounced in the mathematics upper grades, indicating that the pool includes more difficult items than the ability of students in the tested population. The pool includes difficult items to accurately measure high-performing students but needs additional easy items to better measure low-performing students.

At the claim level, the student ability distribution shifted to the left for all claims except for Claim2 grades 4-7 in ELA/L. In mathematics, the student ability distribution shifted to the left for all claims except for Claim 1 in grades 3–5. The Smarter Balanced Assessment Consortium plans to add additional easy items to the pool and to augment the pool in proportion to the test blueprint constraints (e.g., content, Depth of Knowledge [DOK], item type, item difficulties) to better measure low-performing students.

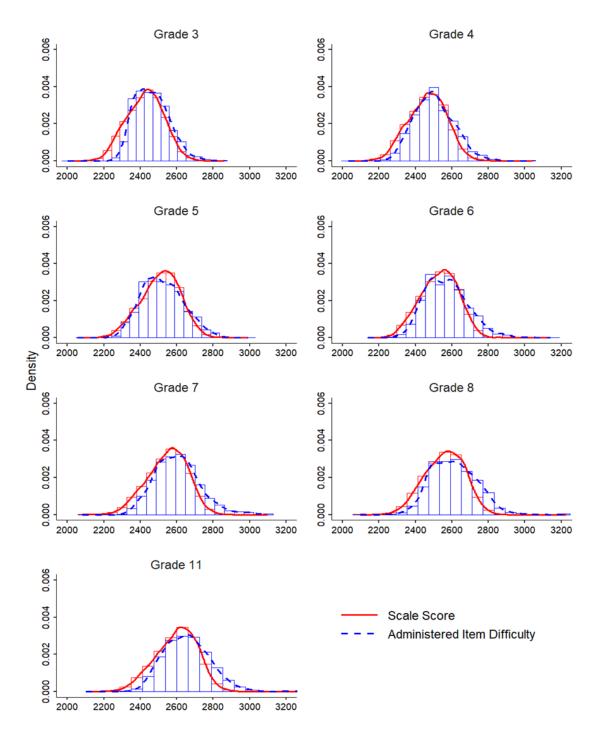


Figure 5. Student Ability—Item Difficulty Distribution: ELA/L

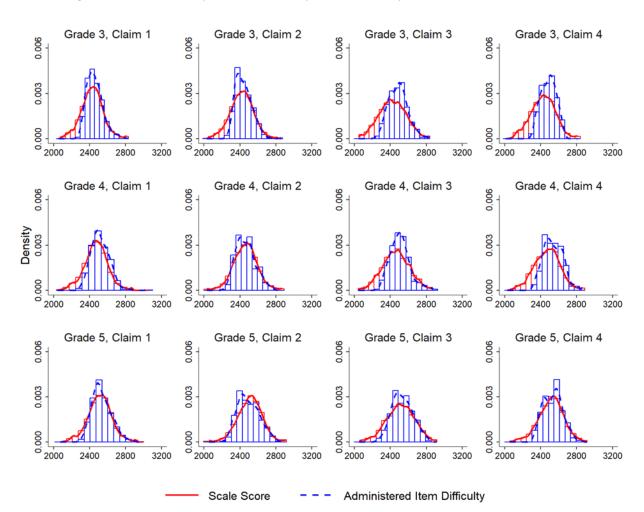


Figure 6. Student Ability—Item Difficulty Distribution by Claim: ELA/L (Grades 3–5)

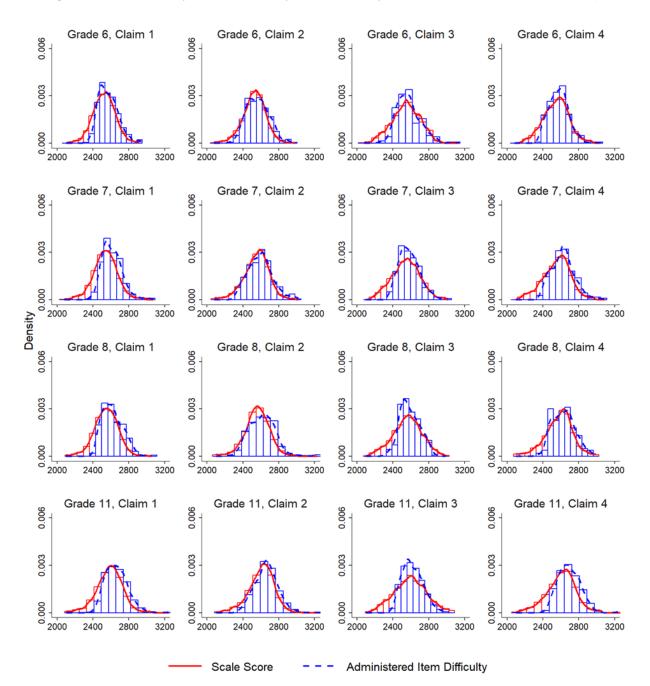


Figure 7. Student Ability—Item Difficulty Distribution by Claim: ELA/L (Grades 6–8, and 11)

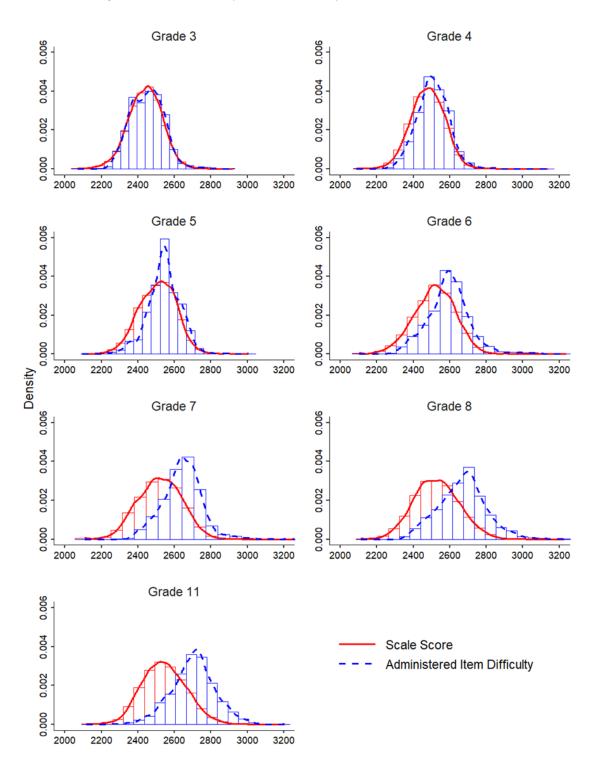


Figure 8. Student Ability—Item Difficulty Distribution: Mathematics

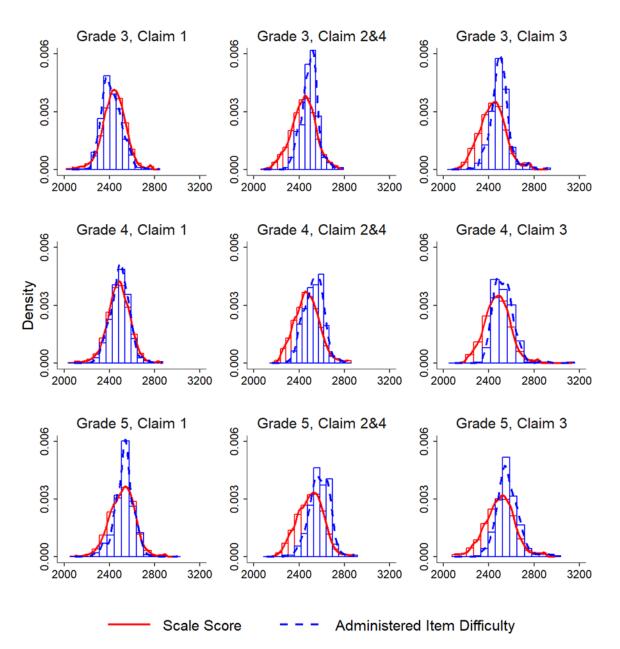
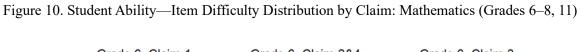
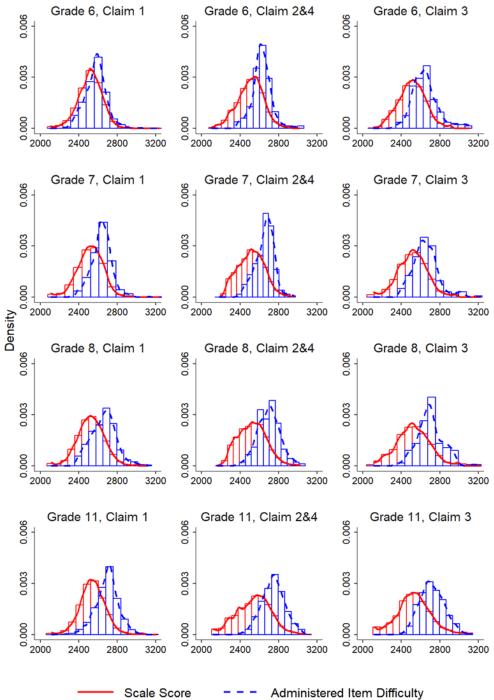


Figure 9. Student Ability—Item Difficulty Distribution by Claim: Mathematics (Grades 3–5)





4. VALIDITY

According to the *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014), validity refers to the degree to which evidence and theory support the interpretations of test scores as described by the intended uses of assessments. The validity of an intended interpretation of test scores relies on all the evidence accrued about the technical quality of a testing system, including test development and construction procedures, test score reliability, accurate scaling and equating, procedures for setting meaningful achievement standards, standardized test administration and scoring procedures, and attention to fairness for all test takers. The appropriateness and usefulness of the Smarter Balanced summative assessments depends on the assessments meeting the relevant standards of validity.

Validity evidence provided in this chapter is as follows:

- Test Content
- Internal Structure
- Relations to Other Variables (External Structure)

Evidence on test content validity is provided with the blueprint match rates for the delivered tests. Evidence on internal structure is examined in the results of intercorrelations among claim scores.

Some of the evidence on standardized test administration, scoring procedures, and attention to fairness for all test takers is provided in other chapters.

4.1 EVIDENCE ON TEST CONTENT

The Smarter Balanced summative assessment includes two components: the computer-adaptive test (CAT) and the performance task (PT). For the CAT, each student receives a different set of items adapted to his or her ability. For the PT, each student is administered a fixed-form test. The content coverage in all PT forms is the same. The test blueprint constraints for CAT and PT can be found at:

https://smarterbalanced.alohahsap.org/resource-list/en/hawaii-shortened-summative-assessment-final-blueprints.

In the adaptive item-selection algorithm, item selection takes place in two discrete stages: blueprint satisfaction and match-to-ability. The blueprints specify a range of items to be administered in each claim, content domain/standard, and target. Moreover, blueprints constrain the Depth of Knowledge (DOK) and item and passage types. For DOK constraints, the Smarter Balanced blueprint specifies either the minimum or maximum number of items, not *both* the minimum and maximum. In blueprints, all content blueprint elements are configured to obtain a strictly enforced range of items administered. The algorithm also seeks to satisfy target-level constraints, but these ranges are not strictly enforced. In English language arts/literacy (ELA/L), the blueprints also specify the number of passages in reading (Claim 1) and listening (Claim 3) claims.

For the Smarter Balanced item pool, all items are developed in English. A portion of the English item pool was transcribed in braille or translated into Spanish to accommodate students who use braille and students who require tests administered in Spanish. The ELA/L pool is available in English and braille. The mathematics pool is available in English, braille, and Spanish. For each of these pools, a portion of items in each pool was further divided to accommodate American sign language (ASL), translations glossaries, and illustration glossaries. The translations glossaries and illustration glossaries were for mathematics

items while the ASL was for mathematics items and listening items in ELA/L. Since the accommodated pools are small, few tests have violations in a few blueprint constraints.

Tables 25 and 26 present the percentages of tests aligned with the ELA/L CAT test blueprint constraints for items in claims, targets, DOK, and number of passage requirement. All tests met the blueprint requirements except for a few targets and DOK in grade 6 due to the application of pool filters limiting item pool.

Tables 27–29 provide the percentages of tests aligned with the test blueprint constraints for the mathematics CAT, the blueprint match rates for claims, DOK, and target constraints. All tests met all blueprint constraints, except for a few tests in grades 4, 5, 7, and 8. The violations appeared on tests due to the application of pool filters limiting the item pool. Pool filters, such as using only items with illustration or language glossaries, can result in an accommodated CAT item pool that is too limited to meet all test blueprint requirements, especially if multiple pool filters are employed on the same test.

Table 25. Percentage of CAT Delivered Tests Meeting Blueprint Requirements: ELA/L (Grades 3–5)

- CI - I		Required	% BP Match			
Claim	Content Category/Target	Items/Passages	Grade 3	Grade 4	Grade 5	
1	Literary Text	4	100.00	100.00	100.00	
	Target 2: Central Ideas Target 4: Reasoning and Evaluation	1–3	100.00	100.00	100.00	
	Targets 1, 3, 5, 6, and 7	1–3	100.00	100.00	100.00	
	Long Literary Text Passage Short Literary Text Passage	1	100.00	100.00	100.00	
	Informational Text	4	100.00	100.00	100.00	
	Target 9: Central Ideas Target 11: Reasoning and Evaluation	1–3	100.00	100.00	100.00	
	Targets 8, 10, 12, 13, and 14	1–3	100.00	100.00	100.00	
	Long Informational Text Passage Short Informational Text Passage	1	100.00	100.00	100.00	
	DOK 2	≥ 4	100.00	100.00	100.00	
	DOK 3 or 4	≥ 1	100.00	100.00	100.00	
2	Writing	5	100.00	100.00	100.00	
	Target 1, 3, or 6: Organization/Purpose	1	100.00	100.00	100.00	
	Target 1, 3, or 6: Evidence/Elaboration	1	100.00	100.00	100.00	
	Target 8: Language and Vocabulary Use	1	100.00	100.00	100.00	
	Target 9: Edit/Clarify	2	100.00	100.00	100.00	
	DOK 2	≥ 2	100.00	100.00	100.00	
3	Listening	4	100.00	100.00	100.00	
	Target 4: Listen/Interpret	4	100.00	100.00	100.00	
	DOK 2 or Higher	≥ 2	100.00	100.00	100.00	
	Listening Passage	2	100.00	100.00	100.00	
4	Research	5	100.00	100.00	100.00	
	Target 2: Interpret and Integrate Information	1–2	100.00	100.00	100.00	
	Target 3: Analyze Information/Sources	1–2	100.00	100.00	100.00	
	Target 4: Use Evidence	1–2	100.00	100.00	100.00	

Table 26. Percentage of CAT Delivered Tests Meeting Blueprint Requirements: ELA/L (Grades 6-8, 11)

		Required	Required	% BP Match			
Claim	Content Category/Target	Items/ Passages in Grades 6–8	Items/ Passages in Grade 11	Grade 6	Grade 7	Grade 8	Grade 11
1	Literary Text	4	4	100.00	100.00	100.00	100.00
	Target 2: Central Ideas Target 4: Reasoning and Evaluation	1–3	1–3	100.00	100.00	100.00	100.00
	Targets 1, 3, 5, 6, and 7	1–3	1–3	100.00	100.00	100.00	100.00
	Long Literary Text Passage	1	1	100.00	100.00	100.00	100.00
	Informational Text	6	6	100.00	100.00	100.00	100.00
	Target 9: Central Ideas Target 11: Reasoning and Evaluation	2–4	2–4	100.00	100.00	100.00	100.00
	Targets 8, 10, 12, 13, and 14	2–4	2–4	100.00	100.00	100.00	100.00
	Long Informational Text Passage	1	1	100.00	100.00	100.00	100.00
	Short Informational Text Passage	1	1	100.00	100.00	100.00	100.00
	DOK 1	≤ 3	≤ 2	100.00	100.00	100.00	100.00
	DOK 3 or Higher	≥ 1	≥ 2	100.00	100.00	100.00	100.00
2	Writing	5	5	100.00	100.00	100.00	100.00
	Target 1, 3, or 6: Organization/Purpose	1	1	100.00	100.00	100.00	100.00
	Target 1, 3, or 6: Evidence/Elaboration	1	1	99.99	100.00	100.00	100.00
	Target 8: Language and Vocabulary Use	1	1	100.00	100.00	100.00	100.00
	Target 9: Edit/Clarify	2	2	99.99	100.00	100.00	100.00
	DOK 2	≥ 2	≥ 2	99.99	100.00	100.00	100.00
3	Listening	4	4	100.00	100.00	100.00	100.00
	Target 4: Listen/Interpret	4	4	100.00	100.00	100.00	100.00
	DOK 2 or Higher	≥ 2	≥ 2	100.00	100.00	100.00	100.00
	Listening Passage	2	2	100.00	100.00	100.00	100.00
4	Research	5	5	100.00	100.00	100.00	100.00
	Target 2: Analyze/Integrate Information	1–2	1–2	100.00	100.00	100.00	100.00
	Target 3: Evaluate Information/Sources	1–2	1–2	100.00	100.00	100.00	100.00
	Target 4: Use Evidence	1–2	1–2	100.00	100.00	100.00	100.00

Table 27. Percentage of CAT Delivered Tests Meeting Blueprint Requirements for Claims and Targets: Mathematics (Grades 3–5)

		Grad	de 3	Grae	de 4	Grade 5		
Claim	Content Domain	Required Items	% BP Match	Required Items	% BP Match	Required Items	% BP Match	
1	Overall	12	100.00	12	100.00	12	100.00	
	DOK 2 or Higher	≥ 4	100.00	≥ 4	100.00	≥ 4	100.00	
	Priority Cluster	9	100.00					
	Targets B, C, G, I	4	100.00					
	Targets D, F	4	100.00					
	Target A	1	100.00					
	Supporting Cluster	3	100.00					
	Targets E, J, K	2	100.00					
	Target H	1	100.00					
	Priority Cluster			9	100.00			
	Targets A, E, F			5	100.00			
	Target G			2	100.00			
	Target D			1	100.00			
	Target H			1	100.00			
	Supporting Cluster			3	100.00			
	Targets I, K			1	100.00			
	Targets B, C, J			1	100.00			
	Target L			1	100.00			
	Priority Cluster					9	100.00	
	Targets E, I					4	100.00	
	Target F					3	100.00	
	Targets C, D					2	100.00	
	Supporting Cluster					3	100.00	
	Targets J, K					2	100.00	
	Targets A, B, G, H					1	100.00	
2 and 4	Overall	5	100.00	5	100.00	5	100.00	
	DOK 3 or Higher	≥ 2	100.00	≥ 2	99.96	≥ 2	99.88	
	2. Target A	1	100.00	1	99.98	1	100.00	
	2. Targets B, C, D	1	100.00	1	99.98	1	100.00	
	4. Targets A, D	1	100.00	1	100.00	1	100.00	
	4. Targets B, E	1	100.00	1	100.00	1	100.00	
	4. Targets C, F	1	100.00	1	100.00	1	100.00	
3	Overall	5	100.00	5	100.00	5	100.00	
	DOK 3 or Higher	≥ 2	100.00	≥ 2	100.00	≥ 2	100.00	
	Targets A, D	2	100.00	2	100.00	2	100.00	
	Targets B, E	2	100.00	2	100.00	2	100.00	
	Targets C, F	1	100.00	1	100.00	1	100.00	

Table 28. Percentage of CAT Delivered Tests Meeting Blueprint Requirements for Claims and Targets: Mathematics (Grades 6–8)

		Grad	de 6	Grad	de 7	Grade 8		
Claim	Content Domain	Required	% BP	Required	% BP	Required	% BP	
		Items	Match	Items	Match	Items	Match	
1	Overall	12	100.00	12	100.00	12	100.00	
	DOK 2 or Higher	≥ 4	100.00	≥ 4	100.00	≥ 4	100.00	
	Priority Cluster	9	100.00					
	Targets E, F	4	100.00					
	Target A	2	100.00					
	Targets G, B	2	100.00					
	Target D	1	100.00					
	Supporting Cluster	3	100.00					
	Targets C, H, I, J	3	100.00					
	Priority Cluster			9	99.54			
	Targets A, D			5	100.00			
	Targets B, C			4	99.54			
	Supporting Cluster			3	99.54			
	Targets E, F			2	99.54			
	Targets G, H, I			1	100.00			
	Priority Cluster					9	100.00	
	Targets C, D					3	99.97	
	Targets B, E, G					3	99.97	
	Targets F, H					3	100.00	
	Supporting Cluster					3	100.00	
	Targets A, I, J					3	100.00	
2 and 4	Overall	5	100.00	5	100.00	5	100.00	
	DOK 3 or Higher	≥ 2	100.00	≥ 2	99.98	≥ 2	100.00	
	2. Target A	1	100.00	1	100.00	1	100.00	
	2. Targets B, C, D	1	100.00	1	100.00	1	100.00	
	4. Targets A, D	1	100.00	1	100.00	1	100.00	
	4. Targets B, E	1	100.00	1	100.00	1	100.00	
	4. Targets C, F	1	100.00	1	100.00	1	100.00	
3	Overall	5	100.00	5	100.00	5	100.00	
="	DOK 3 or Higher	≥ 2	100.00	≥ 2	99.99	≥ 2	100.00	
	Targets A, D	2	100.00	2	100.00	2	100.00	
	Targets B, E	2	100.00	2	100.00	2	100.00	
	Targets C, F, G	1	100.00	1	100.00	1	100.00	

Table 29. Percentage of CAT Delivered Tests Meeting Blueprint Requirements for Claims and Targets: Mathematics (Grade 11)

Claire	Contant Danis	Grade 11			
Claim	Content Domain	Required Items	% BP Match		
1	Overall	14	100.00		
	DOK 2 or Higher	≥ 4	100.00		
	Priority Cluster	10	100.00		
	Targets D, E	1–2	100.00		
	Target F	1	100.00		
	Targets G, H, I	3	100.00		
	Target J	1–2	100.00		
	Target K	1–2	100.00		
	Targets L, M, N	2	100.00		
	Supporting Cluster	4	100.00		
	Target O	0–2	100.00		
	Target P	0–2	100.00		
	Targets A, B	0–1	100.00		
	Target C	0–1	100.00		
2 and 4	Overall	5	100.00		
	DOK 3 or Higher	≥ 2	100.00		
	2. Target A	1	100.00		
	2. Targets B, C, D	1	100.00		
	4. Targets A, D	1	100.00		
	4. Targets B, E	1	100.00		
	4. Targets C, F	1	100.00		
3	Overall	5	100.00		
	DOK 3 or Higher	≥ 2	100.00		
	Targets A, D	2	100.00		
	Targets B, E	2	100.00		
	Targets C, F, G	1	100.00		

Table 30 summarizes target coverage by claim and includes the average and range of the number of unique targets administered in each delivered CAT component. The Smarter Balanced blueprints for ELA/L did not require every target to be covered in a claim; therefore, all targets listed in the blueprint are not expected to be covered in every test. Although the target coverage varies somewhat across individual tests, all targets are covered at an aggregate level across all tests combined.

Table 30. Average and Range of the Number of Unique Targets Assessed Within Each Claim Across All Delivered CAT Tests

Crada	T	otal Tar	gets in B	P		Ave	rage		Range	(Minim	um–Maxi	mum)
Grade	C1	C2	C3	C4	C1	C2	C3	C4	C1	C2	C3	C4
ELA/L												
3	14	5	1	3	7.5	4.0	1.0	3.0	4–8	4–4	1-1	3-3
4	14	5	1	3	7.8	4.0	1.0	3.0	6–8	4–4	1-1	3-3
5	14	5	1	3	7.5	4.0	1.0	3.0	5–8	4–4	1-1	3-3
6	14	5	1	3	9.1	4.0	1.0	3.0	6-10	3-4	1-1	3-3
7	14	5	1	3	9.3	4.0	1.0	3.0	7–10	4–4	1-1	3-3
8	14	5	1	3	9.1	4.0	1.0	3.0	7–10	4–4	1-1	3-3
11	14	5	1	3	8.4	4.0	1.0	3.0	6–10	4–4	1-1	3-3
					N	Aathem	atics					
3	11	4	6	6	10.0	2.0	4.1	3.0	9–10	2-2	3–5	3–3
4	12	4	6	6	9.0	2.0	4.0	3.0	9_9	1-2	3-5	3-3
5	11	4	6	6	8.0	2.0	3.9	3.0	7–8	2-2	3-5	3-3
6	10	4	7	6	9.0	2.0	3.8	3.0	8–9	2-2	3-5	3-3
7	9	4	7	6	6.9	2.0	3.9	3.0	6–7	2-2	3-5	3-3
8	10	4	7	6	10.0	2.0	4.3	3.0	8-10	2-2	3-5	3-3
11	16	4	7	6	12.7	2.0	3.9	3.0	10–14	2-2	3-5	3-3

An adaptive-testing algorithm constructs a test form unique to each student, targeting the student's level of ability and meeting the test blueprints. Consequently, the test forms will not be statistically parallel (e.g., equal test difficulty) across individual students, but test scores from the individual tests are comparable since all test forms measure the same content, albeit with a different set of test items. Although each form is unique with respect to its items, all forms align with the same curricular expectations outlined in the test blueprints.

4.2 EVIDENCE ON INTERNAL STRUCTURE

The measurement model used in the Smarter Balanced assessments assumes a single underlying latent trait in student ability estimates, which supports the reporting of a single total ability score. During the test construction phase, the test blueprint was designed to cover multiple distinct claims under each subject. The item selection algorithm prioritizes blueprint matching to ensure each test contains an appropriate combination of items from each claim. Assessing the relationship between these different claim scores is a measure of internal validity according to the *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014). The presence of high correlations among claim scores is evidence that the Smarter Balanced assessments measure a single underlying ability, and that the claim scores are related to each other.

The correlations among claim scores, both observed (below diagonal) and corrected for attenuation (above diagonal), are presented in Tables 31 and 32. The correction for attenuation indicates what the correlation would be if claim scores could be measured with perfect reliability and corrected (adjusted) for measurement error estimates.

The observed correlation between two claim scores with measurement errors can be corrected for attenuation $r_{x'y'} = \frac{r_{xy}}{\sqrt{r_{xx} \times r_{yy}}}$, where $r_{x'y'}$ is the correlation between x and y corrected for attenuation, r_{xy} is

the observed correlation between x and y, r_{xx} is the reliability coefficient for x, and r_{yy} is the reliability coefficient for y.

When corrected for attenuation (above diagonal), the correlations among claim scores are higher than observed correlations. The disattenuated correlations are quite high in both subjects, showing evidence of unidimensional tests. The correction for attenuation is large in both ELA/L and mathematics because the marginal reliabilities of claim scores are low due to the reduction in the test length.

Table 31. Correlations Among Claims: ELA/L

Condo	Claire	Obse	erved and Disatt	tenuated Correl	lation
Grade	Claim	Claim 1	Claim 2	Claim 3	Claim 4
	Claim 1: Reading		0.93	1	0.95
2	Claim 2: Writing	0.62		1	0.93
3	Claim 3: Listening	0.51	0.50		1
	Claim 4: Research	0.58	0.61	0.49	
	Claim 1: Reading		0.92	1	0.94
4	Claim 2: Writing	0.60		1	0.91
4	Claim 3: Listening	0.51	0.51		1
	Claim 4: Research	0.56	0.59	0.48	
	Claim 1: Reading		0.91	1	0.97
5	Claim 2: Writing	0.61		1	0.92
5	Claim 3: Listening	0.53	0.51		1
	Claim 4: Research	0.60	0.63	0.52	
	Claim 1: Reading		0.88	1	0.93
6	Claim 2: Writing	0.62		1	0.91
0	Claim 3: Listening	0.54	0.50		1
	Claim 4: Research	0.60	0.59	0.48	
	Claim 1: Reading		0.89	1	0.95
7	Claim 2: Writing	0.61		1	0.92
/	Claim 3: Listening	0.52	0.50		1
	Claim 4: Research	0.60	0.62	0.49	
	Claim 1: Reading		0.89	1	0.92
8	Claim 2: Writing	0.61		1	0.91
8	Claim 3: Listening	0.53	0.50		1
	Claim 4: Research	0.58	0.59	0.48	
	Claim 1: Reading		0.87	1	0.92
11	Claim 2: Writing	0.59		0.96	0.91
11	Claim 3: Listening	0.49	0.46		1
	Claim 4: Research	0.57	0.59	0.45	

Table 32. Correlations Among Claims: Mathematics

C	Clair	Observed	d and Disattenuated Co	rrelation
Grade	Claim	Claim 1	Claims 2 & 4	Claim 3
	Claim 1		1	1
3	Claims 2 & 4	0.74		1
	Claim 3	0.71	0.64	
	Claim 1		1	1
4	Claims 2 & 4	0.71		1
	Claim 3	0.73	0.64	
	Claim 1		1	1
5	Claims 2 & 4	0.71		1
	Claim 3	0.67	0.6	
	Claim 1		1	1
6	Claims 2 & 4	0.70		1
	Claim 3	0.68	0.59	
	Claim 1		1	1
7	Claims 2 & 4	0.68		1
	Claim 3	0.67	0.56	
	Claim 1		1	0.96
8	Claims 2 & 4	0.68		1
	Claim 3	0.57	0.51	
	Claim 1		1	0.93
11	Claims 2 & 4	0.63		0.97
	Claim 3	0.59	0.48	

Legend. Claim 1: Concepts and Procedures; Claims 2 & 4: Problem Solving / Modeling and Data Analysis; Claim 3: Communicating Reasoning

4.3 EVIDENCE ON RELATIONS TO OTHER VARIABLES

Validity evidence based on relations to other variables can address a variety of questions. At its core, this type of validity addresses the relationship between test scores and variables of interest that are derived outside the testing system. One type of validity evidence based on relations to other variables is evidence for convergent and discriminant validity. Evidence for convergent validity is based on the degree to which test scores correlate with other measures of the same attribute—scores from two tests measuring the same attribute should be correlated. Conversely, evidence for discriminant validity is obtained when test scores are not correlated with measures of construct-irrelevant attributes.

Evidence for convergent and discriminant validity is determined by examining the patterns of correlations between Smarter Balanced assessments and performance on other tests. Observed correlations should be limited only by the unreliability of the measures.

When both assessments measure student achievement in common subject areas, as with, for example, test scores based on mathematics in the Smarter Balanced summative test and the Algebra I and Algebra II End-of-Course (EOC) tests, we expect test scores between the common subject-area assessments to be substantially correlated. In addition, we expect that the magnitude of observed correlations between test scores in different subject areas will be lower than correlations between test scores in a common subject area.

The relationship between the Smarter Balanced scores and the Algebra I and II scores was examined to evaluate the convergent and discriminant aspects of validity using grade 8 and grade 11 assessment data—Smarter Balanced mathematics and Hawai'i Algebra I and II EOC test scores for two different traits (contents) and the Smarter Balanced ELA/L. In examining the convergent and discriminant aspects of validity, Algebra I (grade 8) and II (grade 11) EOC test scores were considered.

It was expected that the correlation between the Smarter Balanced mathematics scores and the Algebra I and II scores for the same subject (convergent validity) would be moderate and higher than the correlation between Smarter Balanced ELA/L and Smarter Balanced mathematics (discriminant validity). That is, the correlation between two tests measuring the same content would be higher than the correlation between tests measuring different contents. For Algebra I and II EOC test, the scores would show a higher correlation with the Smarter Balanced mathematics scores than with the Smarter Balanced ELA/L scores (discriminant validity).

The results are provided in Table 33. In most scenarios, the results are as would be expected given the criteria set forth by Campbell and Fiske (1959), providing the validity evidence.

First, the reliability coefficients (numbers in boldface) were higher than the convergent and discriminant coefficients for all tests.

Second, the scores between similar traits measured by the different methods correlated more highly with each other than they did with different traits measured by the same method. This is the evidence needed for convergent validity (numbers underlined). For example, the correlation between the Smarter Balanced mathematics and Algebra I in grade 8 scores is 0.81. This is higher than the correlation between the Smarter Balanced ELA/L and Smarter Balanced mathematics scores (r = 0.58) and between the Smarter Balanced ELA/L and Hawai'i Algebra I EOC test scores (r = 0.59). The same pattern is shown in grade 11 Algebra II EOC scores. The correlation between the Smarter Balanced mathematics and Algebra II score is 0.68 which is higher than the correlation between the Smarter Balanced ELA/L and Smarter Balanced mathematics scores (r = 0.57) and between the Smarter Balanced ELA/L and Hawai'i Algebra II EOC test scores (r = 0.44).

Last, the correlations of scores between different traits are lower than the correlations between similar traits. This is the evidence needed for discriminant validity (numbers in a rectangle). The correlations between the Smarter Balanced ELA/L scores and the Smarter Balanced mathematics and Algebra I and II EOC test scores in a rectangle are lower than the underlined correlations.

Overall, the observed pattern of correlations in each multitrait-multimethod matrix conforms to the criteria expected for convergent and discriminant validity.

Table 33. Relationship Among the Smarter Balanced, Algebra I, and Algebra II Test Scores

Test/Subject	SB ELA/L	SB Mathematics	EOC Algebra									
Grade 8 (N = 1,640)												
SB ELA/L	0.77											
SB Mathematics	0.58	0.86										
Algebra I	0.59	<u>0.81</u>	0.91									
	Grade 11	(N = 1,194)										
SB ELA/L	0.82											
SB Mathematics	0.57	0.81										
Algebra II	0.44	0.68	0.83									

5. RELIABILITY

According to the *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014), reliability refers to the consistency of test scores across replications of a testing procedure. Reliability is related to the precision of measurement for a test and is evaluated, in part, in terms of the scores' standard error of measurement (SEM). In classical test theory, reliability is defined as the ratio of the true score variance to the observed score variance, assuming the error variance is the same for all scores, and reliability coefficients are the correlation between scores on two equivalent forms of the test.

Within the item response theory (IRT) framework, measurement error is conditional on ability and varies across the ability scale. The amount of precision in estimating achievement can be determined by the test information function, which describes the amount of information provided by the test at each score point along the ability continuum. Test information is the inverse of measurement error; the larger the measurement error, the less test information is being provided. In computer-adaptive testing, items administered vary among students, so the amount of measurement error differs from one test to another, which yields conditional standard errors of measurement (CSEM).

The reliability evidence of the Smarter Balanced summative tests is provided with marginal reliability, CSEM, and classification accuracy and consistency in each achievement level.

5.1 MARGINAL RELIABILITY

For reliability, the marginal reliability was computed for the scale scores, taking into account the varying measurement errors across the ability range. Marginal reliability is a measure of the overall reliability of an assessment based on the average CSEM, estimated at different points on the ability scale, for all students.

The marginal reliability $(\bar{\rho})$ is defined as

$$\bar{\rho} = [\sigma^2 - \left(\frac{\sum_{i=1}^{N} CSEM_i^2}{N}\right)]/\sigma^2,$$

where N is the number of students, $CSEM_i$ is the CSEM of the scale score for student i, and σ^2 is the variance of the scale score. The higher the reliability coefficient, the greater the precision of the test.

Another way to examine test reliability is with the SEM. In the IRT, SEM is estimated as a function of test information provided by a given set of items that make up the test. In computer-adaptive testing (CAT), items administered vary among all students, so the SEM also can vary among students, which yields CSEM. The average CSEM can be computed as

Average
$$CSEM = \sigma \sqrt{1 - \bar{\rho}} = \sqrt{\sum_{i=1}^{N} CSEM_i^2 / N}$$
.

The smaller the value of average CSEM, the greater the accuracy of test scores.

Table 34 presents the marginal reliability coefficients and the average CSEM for the total scale scores.

Table 34. Marginal Reliability: ELA/L and Mathematics

Grade	N	Number of Items Specified in Test Blueprint	Marginal Reliability	Scale Score Mean	Scale Score SD	Average CSEM
			ELA/L			
3	12,256	24	0.89	2425.29	102.63	34.16
4	12,785	24	0.88	2466.17	106.19	36.49
5	13,141	24	0.89	2510.89	109.69	36.25
6	12,400	26	0.89	2530.55	104.77	35.36
7	12,167	26	0.89	2547.52	111.67	37.28
8	12,202	26	0.88	2558.84	110.32	37.48
11	10,884	26	0.87	2596.91	116.21	41.23
			Mathematics			
3	12,317	22	0.92	2439.48	94.71	27.36
4	12,831	22	0.92	2478.81	94.09	27.20
5	13,189	22	0.90	2507.27	103.11	31.82
6	12,479	22	0.91	2516.17	115.28	35.25
7	12,257	22	0.89	2519.63	121.11	40.30
8	12,270	22	0.88	2527.55	125.97	44.08
11	10,893	24	0.87	2544.39	123.26	45.11

5.2 STANDARD ERROR CURVES

Figures 11 and 12 present plots of the CSEM of scale scores across the range of ability. The vertical lines indicate the three cut scores for the four achievement levels. For most of the ability range, the selection algorithm matched items to each student's ability and to the test blueprints with similar precision. Because the item pool is finite and has fewer items located at the extremes of the ability scale, the selection algorithm had to prioritize meeting blueprint requirements over matching items to ability level for those students with very high or very low abilities. This results in higher standard errors for students with very high or very low abilities compared to students with abilities around and between the three cut scores.

Given that classifying students into achievement levels, especially into proficient or not proficient levels based on the Level 3 cut score, is a high-stakes decision for schools, it is important that ability levels near and between the cut scores are measured with as much precision as possible. This increased precision near and between the cut scores is achieved by having more items in the item pool for abilities across the middle of the scale, where the cut scores are located.

A consequence of the selection algorithm's prioritization of meeting blueprint requirements is that student ability near the low and high extremes of the scale is measured with relatively less precision. This produces the expected *u-curve* shape for the CSEM plots shown in Figures 11 and 12. An adaptive test with an infinitely large item pool and a selection algorithm that focused on maximizing information over blueprint requirements would produce CSEM curves that are flatter. The Smarter Balanced assessments focus on increasing precision where it is most needed, i.e., the ability scores near and in between the cut scores. It is worth noting that larger standard errors are observed at the lower ends of the score distribution, relative to the higher ends. This occurs because the item pools currently have a shortage of easy items that are better targeted toward these lower-achieving students. Content experts use this information to consider how to further target and populate item pools.

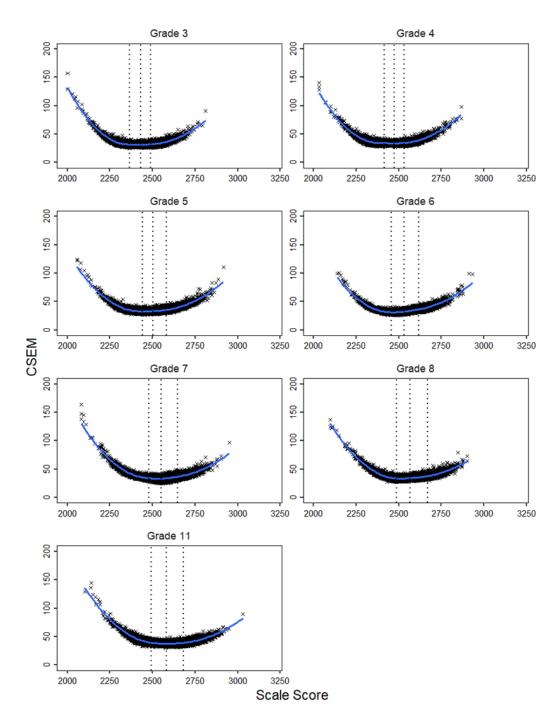


Figure 11. Conditional Standard Error of Measurement: ELA/L

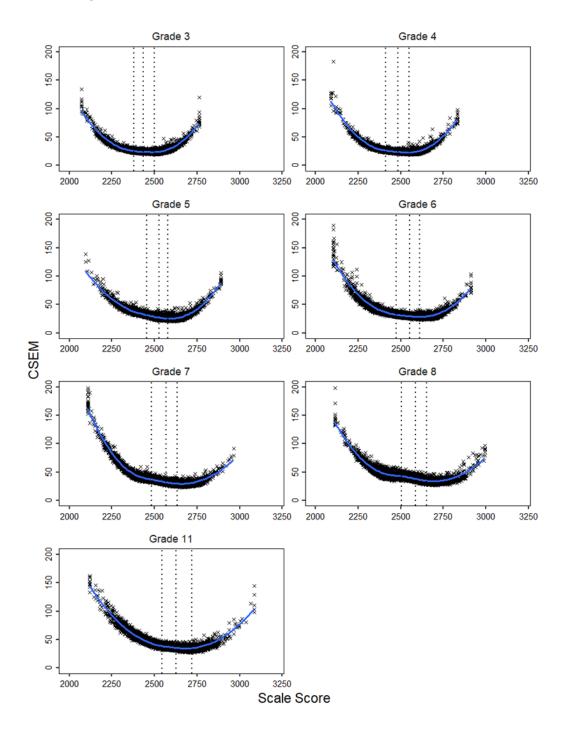


Figure 12. Conditional Standard Error of Measurement: Mathematics

The CSEMs presented in Figures 11 and 12 are summarized in Tables 35 and 36. Table 35 provides the average CSEM for all scale scores and by achievement level. Table 36 presents the average CSEMs at each cut score and the difference in average CSEMs between two cut scores. As shown in Figures 11 and 12, the greatest average CSEM is in Level 1 in both ELA/L and mathematics. Average CSEMs at all cut scores are similar in ELA/L, but larger in Level 2 cut scores in mathematics.

Table 35. Average Conditional Standard Error of Measurement by Achievement Level

Grade	Level 1	Level 2	Level 3	Level 4	Average CSEM
			ELA/L		
3	37.03	31.42	31.63	35.04	34.16
4	37.95	33.97	33.97	38.31	36.49
5	36.80	32.99	34.17	39.78	36.25
6	35.72	31.48	34.31	40.11	35.36
7	41.31	33.33	34.64	40.09	37.28
8	40.85	33.59	35.54	40.85	37.48
11	48.21	38.08	38.11	42.18	41.23
		M	athematics		
3	32.27	24.43	23.70	27.82	27.36
4	33.38	24.86	23.47	26.94	27.20
5	38.72	29.11	26.32	28.65	31.82
6	43.23	30.62	29.07	31.40	35.25
7	50.38	34.89	30.68	31.00	40.30
8	50.67	40.85	36.24	35.62	44.08
11	51.85	37.27	34.66	39.67	45.11

Table 36. Average Conditional Standard Error of Measurement at Each Achievement-Level Cut and Difference of the SEMs Between Two Cuts

Grade	L2 Cut	L3 Cut	L4 Cut	L2-L3	L3-L4	L2-L4
			ELA/L			
3	32.00	31.92	32.35	0.08	0.43	0.35
4	33.72	33.94	34.66	0.21	0.72	0.93
5	33.36	33.35	35.15	0.01	1.80	1.79
6	30.25	31.81	36.83	1.56	5.02	6.58
7	34.22	32.99	36.01	1.24	3.02	1.79
8	32.64	33.57	37.28	0.92	3.72	4.64
11	40.20	37.57	39.02	2.63	1.45	1.18
		I	Mathematics			
3	25.51	23.76	23.68	1.75	0.07	1.83
4	25.98	23.52	23.28	2.46	0.24	2.70
5	31.83	27.39	25.88	4.44	1.51	5.95
6	31.87	29.93	28.96	1.94	0.97	2.91
7	37.12	31.61	29.79	5.52	1.82	7.33
8	42.46	38.94	34.53	3.52	4.41	7.93
11	38.58	35.15	32.86	3.43	2.29	5.71

5.3 RELIABILITY OF ACHIEVEMENT CLASSIFICATION

When student performance is reported in terms of achievement levels, the reliability of achievement classification is computed in terms of the probabilities of accurate and consistent classification of students as specified in Standard 2.16 in *The Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014). The indexes consider the accuracy and consistency of classifications.

For a fixed-form test, the accuracy and consistency of classifications are estimated on a single form's test scores from a single test administration based on the true-score distribution estimated by fitting a bivariate beta-binomial model or a four-parameter beta model (Huynh, 1976; Livingston & Wingersky, 1979; Subkoviak, 1976; Livingston & Lewis, 1995). For the CAT, because the adaptive testing algorithm constructs a test form unique to each student, the classification indexes are computed based on all sets of items administered across students using an IRT-based method (Guo, 2006).

The classification index can be examined in terms of the classification accuracy and the classification consistency. The term *classification accuracy* refers to the agreement between classifications that were made based on the form actually taken and classifications that would be made based on the test takers' true scores if their true scores could somehow be known. Classification consistency refers to the agreement between the classifications based on the form (adaptively administered items) actually taken and the classifications that would be made based on an alternative form (another set of adaptively administered items given the same ability), that is, the percentages of students who are consistently classified in the same achievement levels on two equivalent test forms.

In reality, the true ability is unknown, and students do not take an alternate, equivalent form; therefore, the classification accuracy and the classification consistency are estimated based on students' item scores, item parameters, and assumed underlying latent ability distribution as described in this section. The true score is an expected value of the test score with a measurement error.

For the *i*th student, the student's estimated ability is $\hat{\theta}_i$ with SEM of $se(\hat{\theta}_i)$, and the estimated ability is distributed as $\hat{\theta}_i \sim N\left(\theta_i, se^2(\hat{\theta}_i)\right)$, assuming a normal distribution, where θ_i is the unknown true ability of the *i*th student. The probability of the true score at achievement level *l* based on the cut scores c_{l-1} and c_l is estimated as

$$\begin{split} p_{il} &= p(c_{l-1} \leq \theta_i < c_l) = p\left(\frac{c_{l-1} - \hat{\theta}_i}{se(\hat{\theta}_i)} \leq \frac{\theta_i - \hat{\theta}_i}{se(\hat{\theta}_i)} < \frac{c_l - \hat{\theta}_i}{se(\hat{\theta}_i)}\right) = p\left(\frac{\hat{\theta}_i - c_l}{se(\hat{\theta}_i)} < \frac{\hat{\theta}_i - \theta_i}{se(\hat{\theta}_i)} \leq \frac{\hat{\theta}_i - c_{l-1}}{se(\hat{\theta}_i)}\right) \\ &= \Phi\left(\frac{\hat{\theta}_i - c_{l-1}}{se(\hat{\theta}_i)}\right) - \Phi\left(\frac{\hat{\theta}_i - c_l}{se(\hat{\theta}_i)}\right). \end{split}$$

Instead of assuming a normal distribution of $\hat{\theta}_i \sim N\left(\theta_i, se^2(\hat{\theta}_i)\right)$, the above probabilities can be estimated directly using the likelihood function.

The likelihood function of theta given a student's item scores represents the likelihood of the student's ability at that theta value. Integrating the likelihood values over the range of theta at and above the cut point (with proper normalization) represents the probability of the student's latent ability or the true score being at or above that cut point. If a student with estimated theta is below the cut point, a probability of being at or above the cut point is an estimate of the chance that this student is misclassified as below the cut, and that probability subtracted from 1 is the estimate of the chance that the student is correctly classified as being below the cut score. Using this logic, the various classification probabilities can be defined.

The probability of the *i*th student being classified at achievement level l ($l=1,2,\cdots,L$) based on the cut scores cut_{l-1} and cut_l , given the student's item scores $\mathbf{z}_i = (z_{i1},\cdots,z_{iJ})$ and item parameters $\mathbf{b} = (\mathbf{b}_1,\cdots,\mathbf{b}_J)$ and using the J administered items, can be estimated as

$$\begin{aligned} p_{il} &= P(cut_{l-1} \leq \theta_i < cut_l | \mathbf{z}, \mathbf{b}) = \frac{\int_{cut_{l-1}}^{cut_l} L(\theta | \mathbf{z}, \mathbf{b}) d\theta}{\int_{-\infty}^{+\infty} L(\theta | \mathbf{z}, \mathbf{b}) d\theta} \text{ for } l = 2, \cdots, L-1, \\ p_{i1} &= P(-\infty < \theta_i < cut_1 | \mathbf{z}, \mathbf{b}) = \frac{\int_{-\infty}^{cut_1} L(\theta | \mathbf{z}, \mathbf{b}) d\theta}{\int_{-\infty}^{+\infty} L(\theta | \mathbf{z}, \mathbf{b}) d\theta}, \\ p_{iL} &= P(cut_{L-1} \leq \theta_i < \infty | \mathbf{z}, \mathbf{b}) = \frac{\int_{cut_{L-1}}^{\infty} L(\theta | \mathbf{z}, \mathbf{b}) d\theta}{\int_{-\infty}^{+\infty} L(\theta | \mathbf{z}, \mathbf{b}) d\theta}, \end{aligned}$$

where the likelihood function, based on general IRT models, is

$$L(\theta|\mathbf{z}_i,\mathbf{b}) = \prod_{j \in \mathbf{d}} \left(z_{ij} c_j + \frac{(1-c_j) exp\left(z_{ij} Da_j(\theta-b_j)\right)}{1+exp\left(Da_j(\theta-b_j)\right)} \right) \prod_{j \in \mathbf{p}} \left(\frac{exp\left(Da_j\left(z_{ij}\theta - \sum_{k=1}^{z_{ij}} b_{ik}\right)\right)}{1+\sum_{m=1}^{K_j} exp\left(Da_j\left(\sum_{k=1}^{m} (\theta-b_{jk})\right)\right)} \right),$$

where d stands for dichotomous and p stands for polytomous items; $\mathbf{b}_j = (a_j, b_j, c_j)$ if the jth item is a dichotomous item, and $\mathbf{b}_j = (a_j, b_{j1}, ..., b_{jK_i})$ if the jth item is a polytomous item; a_j is the item's discrimination parameter (for Rasch model, $a_j = 1$), c_j is the guessing parameter (for Rasch and 2PL models, $c_j = 0$), and D is 1.7 for non-Rasch models and 1 for Rasch model.

Classification Accuracy

Using p_{il} , a $L \times L$ table can be constructed as

$$\begin{pmatrix} n_{a11} & \cdots & n_{a1L} \\ \vdots & \vdots & \vdots \\ n_{aL1} & \cdots & n_{aLL} \end{pmatrix},$$

where $n_{alm} = \sum_{pl_i=l} p_{im}$. n_{alm} is the expected number of students at achievement level lm, pl_i is the ith student's achievement level, and p_{im} is the probability of the ith student being classified at achievement level m. In the above table, the row represents the observed level, and the column represents the expected level.

The classification accuracy (CA) at level $l(l = 1, \dots, L)$ is estimated by

$$CA_l = \frac{n_{all}}{\sum_{m=1}^{L} n_{alm}},$$

and the overall classification accuracy is estimated by

$$CA = \frac{\sum_{l=1}^{L} n_{all}}{N},$$

where *N* is the total number of students. Because classifying students as proficient or not proficient is such a high-stakes decision, classification accuracy is also considered at the proficiency level by repeating the process for overall classification accuracy of achievement levels but with the four achievement levels collapsed into two proficiency categories: proficient (achievement levels 3 and 4) and not proficient (achievement levels 1 and 2).

Classification Consistency

Using p_{il} , which is similar to accuracy, another $L \times L$ table can be constructed by assuming the test is administered twice independently to the same student group

$$\begin{pmatrix} n_{c11} & \cdots & n_{c1L} \\ \vdots & \vdots & \vdots \\ n_{cL1} & \cdots & n_{cLL} \end{pmatrix},$$

where $n_{clm} = \sum_{i=1}^{N} p_{il} p_{im}$. p_{il} and p_{im} are the probabilities of the *i*th student being classified at achievement level l and m, respectively, based on observed scores and hypothetical scores from an equivalent test form.

The classification consistency (CC) at level l ($l = 1, \dots, L$) is estimated by

$$CC_l = \frac{n_{cll}}{\sum_{m=1}^{L} n_{clm}},$$

and the overall classification consistency is

$$CC = \frac{\sum_{l=1}^{L} n_{cll}}{N}.$$

As with classification accuracy, classification consistency is also considered at the proficiency level by repeating the process for overall classification consistency of achievement levels but with the four achievement levels collapsed into two proficiency categories: proficient (achievement levels 3 and 4) and not proficient (achievement levels 1 and 2).

The analysis of the classification index is performed based on the overall scale scores. Table 37 provides the percentages of classification accuracy and consistency for overall, by achievement level, and at proficiency cut score.

The overall classification index ranged from 74% to 80% for accuracy and from 66% to 73% for the consistency across all grades and subjects. For achievement levels, the classification index is higher in L1 and L4 than in L2 and L3. The higher accuracy at L1 and L4 is due to the fact that the intervals used to compute the classification probabilities for students in L1 and L4 [$-\infty$, L2 cut; L4 cut, ∞] are wider than the intervals used to compute the classification probabilities for students in L2 and L3 [L2 cut, L3 cut; L3 cut, L4 cut]. The misclassification probability tends to be higher for narrower intervals. Classification accuracy and classification consistency at the proficiency cut scores were high, ranging from 90% to 93% for accuracy and from 87% to 90% for consistency.

The accuracy of classifications is higher than the consistency of classifications in all achievement levels. The accuracy is higher than the consistency because the accuracy is based on one test with a measurement error and the true score while the consistency is based on two tests with measurement errors. The classification indexes by subgroup are provided in Appendix C, Classification Accuracy and Consistency Index by Subgroup.

Table 37. Classification Accuracy and Consistency

Cuada	Achievement	EI	LA/L	Math	ematics
Grade	Level	% Accuracy	% Consistency	% Accuracy	% Consistency
	Overall	75	67	78	69
	L1	89	82	85	79
2	L2	62	50	64	52
3	L3	58	47	71	61
	L4	86	79	88	82
	Proficiency Cut	91	87	92	89
	Overall	74	66	79	71
	L1	89	83	87	79
4	L2	55	43	73	63
4	L3	56	45	71	60
	L4	85	78	88	82
	Proficiency Cut	90	87	92	89
	Overall	75	67	77	69
	L1	89	82	88	82
_	L2	57	46	68	57
5	L3	66	56	59	48
	L4	85	78	87	81
	Proficiency Cut	91	87	92	89
	Overall	76	67	78	70
	L1	88	81	90	84
_	L2	66	55	69	59
6	L3	69	59	60	49
	L4	83	74	87	80
	Proficiency Cut	91	87	92	88
	Overall	76	68	78	70
	L1	89	82	89	84
_	L2	64	52	66	56
7	L3	71	63	64	52
	L4	82	72	87	79
	Proficiency Cut	91	87	91	88
	Overall	76	67	76	68
	L1	88	82	88	83
	L2	66	55	61	49
8	L3	72	63	58	45
	L4	81	70	87	78
	Proficiency Cut	91	87	92	88
	Overall	75	66	80	73
	L1	86	78	90	86
	L2	66	54	65	54
11	L3	69	60	69	57
	L4	83	75	84	73
	Proficiency Cut	91	87	93	90

5.4 RELIABILITY FOR SUBGROUPS

The reliability of test scores is also computed by subgroup. Tables 38–45 present the marginal reliability coefficients by the subgroup: gender, ethnicity groups, ELLs, disadvantaged (free or reduced lunch), migrant, and students with disabilities. The reliability coefficients are similar across subgroups but somewhat lower for the ELL and students with disabilities subgroups. A large percentage of students in these subgroups received Level 1 with large CSEMs.

Table 38. Marginal Reliability Coefficients for Overall and by Subgroup: ELA/L (Grades 3–4)

C-harren			Grade 3					Grade 4		
Subgroup	N	MR	SS	SD	CSEM	N	MR	SS	SD	CSEM
All Students	12,256	0.89	2425.29	102.63	34.16	12,785	0.88	2466.17	106.19	36.49
Female	5,983	0.89	2434.85	100.69	34.02	6,165	0.88	2476.94	102.51	36.24
Male	6,273	0.89	2416.17	103.64	34.28	6,620	0.89	2456.13	108.56	36.73
African American	172	0.85	2424.33	84.78	33.06	186	0.85	2464.50	91.60	35.54
AmerIndian/Alaskan	16	0.84	2420.70	82.71	32.62	8*				
Asian/Pacific Islander	2,695	0.88	2456.05	99.77	33.96	2,895	0.87	2495.89	103.40	36.69
Hispanic	2,363	0.88	2411.51	96.72	33.74	2,431	0.87	2448.27	100.60	36.12
Hawai'i Pacific Islander	2,782	0.87	2375.20	95.72	35.12	3,008	0.87	2418.83	101.70	36.87
White	1,479	0.88	2455.17	95.50	33.68	1,515	0.87	2501.59	99.57	36.47
Multi-Racial	2,749	0.89	2441.68	100.64	34.04	2,742	0.87	2483.27	101.32	36.26
ELL	1,549	0.86	2372.04	94.47	35.22	1,569	0.85	2403.47	96.92	37.01
Disadvantaged	5,634	0.88	2394.73	98.52	34.51	5,877	0.87	2431.68	101.30	36.50
Migrant	153	0.86	2375.08	91.00	34.47	133	0.87	2415.69	101.61	36.76
Disability	1,195	0.81	2325.57	85.55	37.17	1,356	0.80	2353.52	87.44	39.35

Note. * Suppressed the data due to the small sample size, n < 10.

Table 39. Marginal Reliability Coefficients for Overall and by Subgroup: ELA/L (Grades 5–6)

Cubanaun			Grade 5					Grade 6		
Subgroup	N	MR	SS	SD	CSEM	N	MR	SS	SD	CSEM
All Students	13,141	0.89	2510.89	109.69	36.25	12,400	0.89	2530.55	104.77	35.36
Female	6,320	0.88	2524.85	105.32	36.21	5,992	0.88	2545.92	101.89	35.59
Male	6,821	0.90	2497.95	112.06	36.29	6,408	0.89	2516.18	105.39	35.14
African American	139	0.88	2511.78	103.89	35.63	151	0.87	2519.26	97.67	34.69
AmerIndian/Alaskan	21	0.82	2519.89	83.18	35.08	18	0.88	2493.33	102.67	35.10
Asian/Pacific Islander	3,050	0.88	2544.12	105.67	36.91	2,928	0.87	2560.63	98.30	35.77
Hispanic	2,626	0.88	2495.44	105.02	35.65	2,403	0.88	2516.78	102.40	34.85
Hawai'i Pacific Islander	3,037	0.88	2458.03	103.92	35.88	2,971	0.88	2480.42	98.74	34.72
White	1,427	0.87	2541.22	100.90	36.56	1,315	0.86	2568.72	98.07	36.24
Multi-Racial	2,841	0.88	2530.64	104.92	36.36	2,614	0.87	2548.21	99.91	35.67
ELL	1,451	0.86	2426.45	95.04	36.17	1,248	0.83	2441.21	85.45	34.83
Disadvantaged	5,838	0.89	2475.26	106.68	35.82	5,477	0.88	2499.88	101.23	34.92
Migrant	154	0.88	2442.84	103.26	36.32	160	0.88	2482.20	99.20	34.63
Disability	1,443	0.83	2386.92	92.35	38.04	1,323	0.81	2413.54	83.58	36.34

Table 40. Marginal Reliability Coefficients for Overall and by Subgroup: ELA/L (Grades 7–8)

Subgroup			Grade 7					Grade 8		
Subgroup	N	MR	SS	SD	CSEM	N	MR	SS	SD	CSEM
All Students	12,167	0.89	2547.52	111.67	37.28	12,202	0.88	2558.84	110.32	37.48
Female	5,859	0.88	2564.85	104.85	36.86	5,925	0.87	2575.61	103.73	37.09
Male	6,308	0.89	2531.42	115.34	37.66	6,277	0.89	2543.01	113.96	37.85
African American	158	0.88	2555.94	105.96	36.58	167	0.87	2568.21	100.14	36.71
AmerIndian/Alaskan	24	0.84	2571.05	89.64	35.66	16	0.87	2532.70	98.67	35.78
Asian/Pacific Islander	3,104	0.87	2584.08	105.50	37.35	3,235	0.87	2593.57	102.53	37.25
Hispanic	2,347	0.88	2528.98	106.09	37.22	2,293	0.88	2537.35	107.39	37.61
Hawai'i Pacific Islander	2,878	0.87	2491.71	104.92	37.91	2,943	0.86	2505.58	102.11	37.86
White	1,322	0.87	2589.33	100.07	36.75	1,284	0.87	2600.16	102.36	37.35
Multi-Racial	2,334	0.88	2561.88	107.94	36.81	2,263	0.88	2576.24	107.01	37.34
ELL	1,314	0.84	2457.29	97.65	39.20	1,344	0.82	2478.62	90.54	38.20
Disadvantaged	5,438	0.88	2511.91	108.06	37.50	5,168	0.87	2523.88	106.65	37.76
Migrant	126	0.85	2490.38	94.06	36.91	177	0.86	2498.86	98.84	37.14
Disability	1,316	0.80	2422.65	92.23	41.33	1,274	0.79	2435.73	91.95	42.06

Table 41. Marginal Reliability Coefficients for Overall and by Subgroup: ELA/L (Grade 11)

Cubaraun			Grade 11		
Subgroup	N	MR	SS	SD	CSEM
All Students	10,884	0.87	2596.91	116.21	41.23
Female	5,240	0.85	2616.49	106.51	40.58
Male	5,644	0.88	2578.74	121.75	41.83
African American	154	0.86	2591.83	104.89	39.88
AmerIndian/Alaskan	17	0.87	2598.88	116.38	41.69
Asian/Pacific Islander	3,497	0.86	2623.77	109.31	40.71
Hispanic	1,868	0.86	2580.23	111.15	41.13
Hawai'i Pacific Islander	2,384	0.85	2543.64	108.70	42.13
White	1,061	0.88	2633.81	118.78	41.86
Multi-Racial	1,899	0.88	2610.28	115.81	40.89
ELL	717	0.77	2489.49	91.08	43.65
Disadvantaged	3,940	0.87	2561.19	114.48	42.00
Migrant	135	0.85	2537.09	106.36	41.34
Disability	910	0.77	2467.06	99.26	47.78

Table 42. Marginal Reliability Coefficients for Overall and by Subgroup: Mathematics (Grades 3-4)

Cuhanoun			Grade 3					Grade 4		
Subgroup	N	MR	SS	SD	CSEM	N	MR	SS	SD	CSEM
All Students	12,317	0.92	2439.48	94.71	27.36	12,831	0.92	2478.81	94.09	27.20
Female	6,016	0.91	2436.23	90.44	26.97	6,182	0.91	2474.76	88.09	26.61
Male	6,301	0.92	2442.58	98.51	27.72	6,649	0.92	2482.57	99.20	27.75
African American	172	0.90	2423.73	83.76	26.66	186	0.89	2472.58	81.72	26.71
AmerIndian/Alaskan	16	0.91	2426.62	85.90	25.82	8*				
Asian/Pacific Islander	2,730	0.91	2472.77	91.90	27.05	2,919	0.92	2512.74	92.34	26.76
Hispanic	2,369	0.91	2423.46	90.15	27.26	2,436	0.91	2461.81	88.25	27.12
Hawai'i Pacific Islander	2,798	0.90	2393.06	88.25	28.44	3,018	0.89	2432.94	88.21	28.88
White	1,483	0.91	2465.96	85.55	26.35	1,517	0.91	2507.02	87.09	26.01
Multi-Racial	2,749	0.91	2454.24	91.25	27.21	2,747	0.91	2493.14	88.07	26.52
ELL	1,594	0.91	2395.94	93.90	28.64	1,563	0.89	2426.08	92.22	30.02
Disadvantaged	5,670	0.91	2410.58	92.40	28.18	5,899	0.90	2448.33	88.70	28.17
Migrant	154	0.91	2394.85	95.35	27.95	134	0.90	2446.58	86.52	27.27
Disability	1,202	0.88	2347.27	95.95	33.05	1,365	0.85	2384.50	85.49	32.82

Note. * Suppressed the data due to the small sample size, n < 10.

Table 43. Marginal Reliability Coefficients for Overall and by Subgroup: Mathematics (Grades 5–6)

Ch. c			Grade 5					Grade 6		
Subgroup	N	MR	SS	SD	CSEM	N	MR	SS	SD	CSEM
All Students	13,189	0.90	2507.27	103.11	31.82	12,479	0.91	2516.17	115.28	35.25
Female	6,344	0.90	2505.12	98.20	31.55	6,033	0.90	2517.26	111.55	34.65
Male	6,845	0.91	2509.25	107.42	32.06	6,446	0.91	2515.15	118.65	35.80
African American	139	0.89	2500.10	94.14	31.78	150	0.89	2517.86	100.54	33.36
AmerIndian/Alaskan	21	0.88	2506.26	88.99	30.92	18	0.87	2449.04	122.24	43.93
Asian/Pacific Islander	3,084	0.91	2548.24	99.42	30.54	2,971	0.91	2554.69	112.97	33.74
Hispanic	2,627	0.89	2489.51	98.19	32.39	2,417	0.90	2496.10	110.46	35.48
Hawai'i Pacific Islander	3,050	0.87	2455.28	93.95	34.28	2,984	0.87	2459.30	109.02	39.23
White	1,423	0.90	2529.55	93.87	30.11	1,317	0.90	2558.28	103.10	32.38
Multi-Racial	2,845	0.90	2524.19	98.88	30.72	2,622	0.90	2534.96	104.92	33.29
ELL	1,473	0.85	2436.60	92.05	35.83	1,311	0.83	2425.00	105.09	42.98
Disadvantaged	5,860	0.89	2475.09	98.86	33.21	5,519	0.89	2482.11	111.47	37.05
Migrant	154	0.85	2439.82	92.83	35.39	158	0.88	2473.66	104.87	36.98
Disability	1,445	0.81	2399.83	90.42	39.36	1,333	0.80	2390.75	101.15	45.77

Table 44. Marginal Reliability Coefficients for Overall and by Subgroup: Mathematics (Grades 7–8)

Subgroup			Grade 7					Grade 8		
Subgroup	N	MR	SS	SD	CSEM	N	MR	SS	SD	CSEM
All Students	12,257	0.89	2519.63	121.11	40.30	12,270	0.88	2527.55	125.97	44.08
Female	5,907	0.88	2518.06	116.87	39.81	5,946	0.87	2528.97	121.25	43.54
Male	6,350	0.89	2521.10	124.91	40.75	6,324	0.88	2526.21	130.25	44.58
African American	159	0.86	2519.32	111.98	41.65	167	0.85	2527.95	109.00	42.90
AmerIndian/Alaskan	25	0.73	2541.13	66.75	34.42	17	0.84	2442.45	120.36	48.33
Asian/Pacific Islander	3,135	0.91	2567.40	120.21	36.57	3,260	0.89	2574.02	126.28	41.29
Hispanic	2,359	0.86	2496.63	112.22	41.52	2,303	0.84	2500.39	116.52	45.89
Hawai'i Pacific Islander	2,902	0.82	2456.50	106.37	45.57	2,957	0.80	2464.64	106.96	47.72
White	1,327	0.89	2560.00	109.00	36.73	1,292	0.87	2572.39	115.18	41.02
Multi-Racial	2,350	0.89	2533.97	115.76	38.72	2,273	0.88	2545.29	123.08	42.87
ELL	1,347	0.79	2428.47	107.65	49.69	1,358	0.77	2442.04	105.07	50.62
Disadvantaged	5,497	0.85	2480.15	113.05	43.64	5,204	0.84	2488.03	116.49	46.62
Migrant	127	0.81	2442.71	116.28	50.33	180	0.77	2472.71	99.29	47.27
Disability	1,326	0.68	2395.98	95.85	54.20	1,289	0.70	2401.14	100.33	54.54

Table 45. Marginal Reliability Coefficients for Overall and by Subgroup: Mathematics (Grade 11)

Cubarana			Grade 11		
Subgroup	N	MR	SS	SD	CSEM
All Students	10,893	0.87	2544.39	123.26	45.11
Female	5,248	0.86	2548.87	115.92	43.87
Male	5,645	0.87	2540.22	129.59	46.23
African American	156	0.81	2520.01	101.89	44.86
AmerIndian/Alaskan	14	0.89	2569.61	130.99	43.30
Asian/Pacific Islander	3,521	0.88	2580.80	123.97	42.46
Hispanic	1,860	0.83	2522.76	110.83	45.58
Hawai'i Pacific Islander	2,390	0.77	2486.07	103.62	49.83
White	1,056	0.89	2578.06	130.88	43.53
Multi-Racial	1,892	0.87	2554.33	121.24	44.04
ELL	732	0.70	2456.93	96.93	52.98
Disadvantaged	3,930	0.82	2507.60	113.94	47.97
Migrant	133	0.77	2495.31	96.47	46.25
Disability	909	0.58	2426.14	90.13	58.14

5.5 RELIABILITY FOR CLAIM SCORES

The marginal reliability, average and standard deviation of scale scores, and average of CSEM are also computed for claim scores by test and grade. In mathematics, Claims 2 and 4 are combined to have enough items to generate a score. Given the reduction in the small number of items in the Hawai'i shortened blueprint, the reliabilities for claim scores are low, especially for Claim 3 and Claim 4 in ELA/L and Claims 2 and 4 combined and Claim 3 in mathematics. In 2023–2024, the performance category for claim scores was reported at the individual student level for only Claims 1 and 2 in ELA/L and Claim 1 in mathematics.

Tables 46 and 47 present the marginal reliability coefficients and descriptive statistics by claim in ELA/L and mathematics, respectively.

Table 46. Marginal Reliability Coefficients for Claim Scores: ELA/L

Grade	Claim	Number of Items Specified in Test Blueprint	Marginal Reliability	Scale Score Mean	Scale Score SD	Average CSEM
	Claim 1: Reading	8	0.62	2428.04	122.94	76.26
3	Claim 2: Writing	6	0.73	2419.15	128.45	67.00
3	Claim 3: Listening	4	0.30	2423.50	148.12	123.69
	Claim 4: Research	6	0.61	2427.61	137.01	85.92
	Claim 1: Reading	8	0.60	2471.90	132.11	83.09
4	Claim 2: Writing	6	0.72	2459.01	134.10	71.49
4	Claim 3: Listening	4	0.33	2462.12	150.62	123.59
	Claim 4: Research	6	0.59	2470.22	144.34	92.77
	Claim 1: Reading	8	0.61	2512.79	134.23	84.05
5	Claim 2: Writing	6	0.74	2509.11	138.31	70.88
3	Claim 3: Listening	4	0.34	2511.19	157.19	128.07
	Claim 4: Research	6	0.63	2514.99	140.72	85.90
	Claim 1: Reading	10	0.69	2525.75	127.55	70.51
6	Claim 2: Writing	6	0.72	2526.07	129.79	69.09
O	Claim 3: Listening	4	0.27	2543.67	157.48	134.74
	Claim 4: Research	6	0.59	2545.54	142.65	91.90
	Claim 1: Reading	10	0.65	2539.64	135.85	80.16
7	Claim 2: Writing	6	0.73	2549.28	141.44	73.02
7	Claim 3: Listening	4	0.30	2542.81	155.29	129.58
	Claim 4: Research	6	0.61	2554.21	154.50	96.62
	Claim 1: Reading	10	0.67	2547.89	133.39	76.09
8	Claim 2: Writing	6	0.71	2557.70	137.61	74.53
8	Claim 3: Listening	4	0.31	2561.85	160.92	133.23
	Claim 4: Research	6	0.60	2577.04	155.38	98.56
	Claim 1: Reading	10	0.65	2584.04	145.28	86.23
	Claim 2: Writing	6	0.71	2604.01	144.09	77.95
11	Claim 3: Listening	4	0.33	2590.62	179.45	147.29
	Claim 4: Research	6	0.59	2602.57	161.91	104.20

Table 47. Marginal Reliability Coefficients for Claim Scores: Mathematics

Grade	Claim	Number of Items Specified in Test Blueprint	Marginal Reliability	Scale Score Mean	Scale Score SD	Average CSEM
	Claim 1	12	0.85	2442.31	103.82	40.50
3	Claims 2 & 4	5	0.61	2438.16	109.88	68.59
	Claim 3	5	0.59	2432.86	115.96	74.53
	Claim 1	12	0.86	2482.31	102.71	38.74
4	Claims 2 & 4	5	0.55	2469.46	110.61	74.52
	Claim 3	5	0.63	2475.57	114.47	69.90
	Claim 1	12	0.84	2513.01	113.93	46.05
5	Claims 2 & 4	5	0.51	2499.13	118.96	83.53
	Claim 3	5	0.54	2497.92	134.91	91.69
	Claim 1	12	0.85	2518.87	127.31	49.24
6	Claims 2 & 4	5	0.54	2508.11	137.01	93.42
	Claim 3	5	0.50	2510.53	142.37	100.75
	Claim 1	12	0.82	2517.62	135.39	57.31
7	Claims 2 & 4	5	0.39	2516.08	136.82	106.53
	Claim 3	5	0.53	2516.39	153.73	105.62
	Claim 1	12	0.81	2526.18	139.35	60.93
8	Claims 2 & 4	5	0.44	2525.76	145.45	108.96
	Claim 3	5	0.43	2521.78	169.41	128.17
	Claim 1	14	0.80	2543.07	131.99	59.24
11	Claims 2 & 4	5	0.48	2541.42	176.77	128.04
	Claim 3	5	0.50	2529.72	178.33	125.73

Legend. Claim 1: Concepts and Procedures; Claims 2 & 4: Problem Solving / Modeling and Data Analysis; Claim 3: Communicating Reasoning

6. SCORING

The Smarter Balanced Assessment Consortium (SBAC) provided the vertically scaled item parameters by linking across all grades using common items in adjacent grades. All scores are estimated based on these item parameters. Each student received an overall scale score, an overall achievement level, and a performance category for Claims 1 and 2 in English language arts/literacy (ELA/L) and Claim 1 in mathematics. This section describes the rules used to generate the scores and the handscoring procedure. The rules and procedures for generating scores are the same in all operational administration years.

6.1 ESTIMATING STUDENT ABILITY USING MAXIMUM LIKELIHOOD ESTIMATION

The Smarter Balanced tests are scored using maximum likelihood estimation (MLE). The likelihood function for generating the MLEs is based on a mixture of item types.

Indexing items by i, the likelihood function based on the jth person's score pattern for I items is

$$L_i(\theta_i|\mathbf{z}_i, \mathbf{a}_i b_1, \dots b_k) = \prod_{i=1}^{I} p_{ij}(z_{ij}|\theta_i, a_{i,} b_{i,1}, \dots b_{i,m_i}),$$

where $b'_i = (b_{i,1}, ..., b_{i,m_i})$ for the *i*th item's step parameters, m_i is the maximum possible score of this item, a_i is the discrimination parameter for item i, z_{ij} is the observed item score for person j, and k indexes the step of item i.

Depending on the item score points, the probability $p_{ij}(z_{ij}|\theta_j, a_i, b_{i,1}, ..., b_{i,m_i})$ takes either the form of a two-parameter logistic (2PL) model for items with one point or the form based on the generalized partial-credit model (GPCM) for items with two or more points.

In the case of items with one score point, $m_i = 1$,

$$p_{ij}(z_{ij}|\theta_{j},a_{i,}b_{i,1},\dots b_{i,m_{i}}) = \begin{cases} \frac{exp\left(Da_{i}(\theta_{j}-b_{i,1})\right)}{1+exp\left(Da_{i}(\theta_{j}-b_{i,1})\right)} = p_{ij}, if \ z_{ij} = 1\\ \frac{1}{1+exp\left(Da_{i}(\theta_{j}-b_{i,1})\right)} = 1-p_{ij}, if \ z_{ij} = 0 \end{cases};$$

in the case of items with two or more points,

$$p_{ij}(z_{ij}|\theta_{j},a_{i,}b_{i,1},...b_{i,m_{i}}) = \begin{cases} \frac{exp(\sum_{k=1}^{z_{ij}} Da_{i}(\theta_{j} - b_{i,k}))}{s_{ij}(\theta_{j},a_{i,}b_{i,1,...}b_{i,m_{i}})}, if \ z_{ij} > 0 \\ \frac{1}{s_{ij}(\theta_{j},a_{i,}b_{i,1,...}b_{i,m_{i}})}, if \ z_{ij} = 0 \end{cases},$$

where $s_{ij}(\theta_j, a_{i,b_{i,1,...}}b_{i,m_i}) = 1 + \sum_{l=1}^{m_i} exp(\sum_{k=1}^l Da_i(\theta_j - b_{i,k}))$, and D = 1.7.

Standard Error of Measurement

With MLE, the standard error (SE) for student *j* is

$$SE(\theta_j) = \frac{1}{\sqrt{I(\theta_j)}},$$

where $I(\theta_i)$ is the test information for student j, calculated as

$$I(\theta_{j}) = \sum_{i=1}^{l} D^{2} a_{i}^{2} \left(\frac{\sum_{l=1}^{m_{i}} l^{2} exp(\sum_{k=1}^{l} Da_{i}(\theta_{j} - b_{ik}))}{1 + \sum_{l=1}^{m_{i}} exp(\sum_{k=1}^{l} Da_{i}(\theta_{j} - b_{ik}))} - \left(\frac{\sum_{l=1}^{m_{i}} lexp(\sum_{k=1}^{l} Da_{i}(\theta_{j} - b_{ik}))}{1 + \sum_{l=1}^{m_{j}} exp(\sum_{k=1}^{l} Da_{i}(\theta_{j} - b_{ik}))} \right)^{2} \right),$$

where m_i is the maximum possible score point (starting from 0) for the *i*th item, and *D* is the scale factor, 1.7. The SE is calculated based on the answered item(s) only for both complete and incomplete tests. The upper bound of the SE is set to 2.5 on the θ metric. Any value larger than 2.5 is truncated at 2.5 on the θ metric.

The algorithm allows previously answered items to be changed; however, it does not allow items to be skipped. Item selection requires iteratively updating the estimate of the overall ability estimates after each item is answered. When a previously answered item is changed, the proficiency estimate is adjusted to account for the changed responses when the next new item is selected. Although the update of the ability estimates is performed at each iteration, the overall scores are recalculated using all data at the end of the assessment for the final score.

6.2 RULES FOR TRANSFORMING THETA TO VERTICAL SCALE SCORES

The student's performance in each subject is summarized in an overall test score referred to as a *scale score*. The scale scores represent a linear transformation of the ability estimates (theta scores) using the formula $SS = a * \theta + b$. The scaling constants a and b are provided by SBAC. Table 48 presents the scaling constants for each subject for the theta-to-scale score linear transformation. Scale scores are rounded to an integer.

Table 48. Vertical Scaling Constants on the Reporting Metric

Subject	Grade	Slope (a)	Intercept (b)
ELA/L	3–8, 11	85.8	2508.2
Mathematics	3–8, 11	79.3	2514.9

Standard errors of the MLEs are transformed to be placed onto the reporting scale. This transformation is

$$SE_{ss} = a * SE_{\theta}$$

where SE_{SS} is the standard error of the ability estimate on the reporting scale, SS_{θ} is the standard error of the ability estimate on the θ scale, and a is the slope of the scaling constant that transforms θ into the reporting scale.

The scale scores are mapped into four achievement levels using three achievement standards (i.e., cut scores). Table 49 provides three achievement standards for each grade and content area.

Table 49. Cut Scores in Scale Scores

Grade		ELA/L		Mathematics				
Graue	Level 2	Level 3	Level 4	Level 2	Level 3	Level 4		
3	2367	2432	2490	2381	2436	2501		
4	2416	2473	2533	2411	2485	2549		
5	2442	2502	2582	2455	2528	2579		
6	2457	2531	2618	2473	2552	2610		
7	2479	2552	2649	2484	2567	2635		
8	2487	2567	2668	2504	2586	2653		
11	2493	2583	2682	2543	2628	2718		

6.3 LOWEST/HIGHEST OBTAINABLE SCORES

Although the observed score is measured more precisely in an adaptive test than in a fixed-form test, especially for high- and low-performing students, if the item pool does not include enough easy or difficult items to measure low- and high-performing students, the standard error could be large in the low and high ends of the ability range. SBAC decided to truncate extreme, unreliable student ability estimates. Table 50 presents the lowest obtainable theta (LOT) and scale score (LOSS) and the highest obtainable theta (HOT) and scale score (HOSS) in both theta and scale score metrics. Estimated thetas lower than LOT or higher than HOT are truncated to the LOT and HOT values and are assigned LOSS and HOSS associated with the LOT and HOT. LOT and HOT were applied to all tests and total scores. The standard error for the LOT and HOT is computed using the LOT and HOT ability estimates given the administered items.

Table 50. Extended Lowest and Highest Obtainable Scores

Condo	Theta N	Metric	Scale Sco	re Metric
Grade	LOT	НОТ	LOSS	HOSS
		ELA/L		
3	-5.9110	3.5332	2001	2811
4	-5.5500	4.1826	2032	2867
5	-5.2670	4.7546	2056	2916
6	-5.0000	5.0000	2079	2937
7	-4.9660	5.3119	2082	2964
8	-4.7925	5.6063	2097	2989
11	-4.7305	6.1096	2102	3032
		Mathematics		
3	-5.6030	3.1219	2071	2762
4	-5.3601	4.0264	2090	2834
5	-5.3012	4.7426	2095	2891
6	-5.1942	5.0000	2103	2911
7	-5.1311	5.6630	2108	2964
8	-5.0681	6.0272	2113	2993
11	-5.0000	7.1896	2118	3085

6.4 SCORING ALL CORRECT AND ALL INCORRECT CASES

In the item response theory (IRT) maximum likelihood ability estimation methods, zero and perfect scores are assigned the ability of minus and plus infinity. For all correct and all incorrect cases, the highest obtainable scores (HOT and HOSS) and the lowest obtainable scores (LOT and LOSS) were assigned in the 2014–2015 administration. Since the 2015–2016 administration, all incorrect and correct cases were scored by either adding 0.5 to or subtracting 0.5 from an item score with the smallest item discrimination parameter among the administered operational items (computer-adaptive testing [CAT] and performance tasks [PTs]) for a student.

6.5 RULES FOR CALCULATING STRENGTHS AND WEAKNESSES FOR CLAIM SCORES

In ELA/L, claim scores are computed and reported for Claims 1 and 2 at the individual student level; in mathematics, claim scores are computed and reported for Claim 1 only. For the claim, three performance categories, indicating relative strength and weakness, are produced.

The difference between the proficiency cut score and the claim score plus or minus 1.5 times standard error of the claim is used to determine the relative strengths and weaknesses. For summative tests, the specific rules are as follows:

- Below Standard (Code = 1): if $round(SS_{rc} + 1.5 * SE(SS_{rc}), 0) < SS_{v}$
- At/Near Standard (Code = 2): if $round(SS_{rc} + 1.5 * SE(SS_{rc}), 0) \ge SS_p$ and $round(SS_{rc} 1.5 * SE(SS), 0) < SS_p$, a strength or weakness is indeterminable
- Above Standard (Code = 3): if $round(SS_{rc} 1.5 * SE(SS_{rc}), 0) \ge SS_p$

where SS_{rc} is the student's scale score on a claim, SS_p is the proficiency scale score cut (Level 3 cut), and $SE(SS_{rc})$ is the standard error of the student's scale score on the claim.

6.6 TARGET SCORES

The target-level reports are impossible to produce for a fixed-form test because the number of items included per target (i.e., benchmark) is too small to produce a reliable score at the target level. A typical fixed-form test includes only one or two items per target. Even when aggregated, these data narrowly reflect the benchmark because they reflect only one or two ways of measuring the target. An adaptive test, however, offers a tremendous opportunity for target-level data at the class, school, and complex-area level. With an adequate item pool, a class of 20 students might respond to 10 or 15 different items measuring any given target. Target scores are computed for attempted tests based on the responded items. Target scores are computed in each claim (four claims) for ELA/L and in Claim 1 only for mathematics. Target scores can be computed for any aggregate group of students, and Chapter 7: Reporting and Interpreting Scores provides details on which aggregate groups of students have target scores computed and who has access to the reports.

Target scores are computed in two ways: (1) target scores relative to a student's overall estimated ability (θ) , and (2) target scores relative to the proficiency standard (Level 3 cut).

6.6.1 Target Scores Relative to Student's Overall Estimated Ability

By defining $p_{ij} = p(z_{ij} = 1)$, indicating the probability that student j responds correctly to item i, z_{ij} represents the jth student's score on the ith item. For items with one score point, the 2PL IRT model is used to calculate the expected score on item i for student j with estimated ability $\hat{\theta}_j$ as:

$$E(z_{ij}) = \frac{\exp(Da_i(\widehat{\theta}_j - b_i))}{1 + \exp(Da_i(\widehat{\theta}_j - b_i))}.$$

For items with two or more score points, using the generalized partial credit model (GPCM), the expected score for student j with estimated ability $\hat{\theta}_j$ on an item i with a maximum possible score of m_i is calculated as

$$E(z_{ij}) = \sum_{l=1}^{m_i} \frac{lexp(\sum_{k=1}^{l} Da_i(\widehat{\theta}_j - b_{i,k}))}{1 + \sum_{l=1}^{m_i} exp(\sum_{k=1}^{l} Da_i(\widehat{\theta}_j - b_{i,k}))}.$$

For each item i, the residual between observed and expected score for each student is defined as

$$\delta_{ij} = z_{ij} - E(z_{ij}).$$

Residuals are summed for items within a target. The sum of residuals is divided by the total number of points possible for items within the target, *T*:

$$\delta_{jT} = \frac{\sum_{i \in T} \delta_{ji}}{\sum_{i \in T} m_i}.$$

For an aggregate unit, a target score is computed by averaging the individual student target scores for the target across all students in the aggregate unit.

$$\bar{\delta}_{Tg} = \frac{1}{n_g} \sum_{j \in g} \delta_{jT}$$
, and $se(\bar{\delta}_{Tg}) = \sqrt{\frac{1}{n_g(n_g-1)} \sum_{j \in g} (\delta_{jT} - \bar{\delta}_{Tg})^2}$,

where n_g is the number of students who responded to any of the items that belong to the target T for an aggregate unit g. If a student did not happen to see any items on a particular target, the student is *not* included in the n_g count for the aggregate.

A difference from zero in these aggregates may indicate that a roster, teacher, school, complex, or complex area is more effective (if $\bar{\delta}_{Tg}$ is positive) or less effective (negative $\bar{\delta}_{Tg}$) in teaching a given target.

Direct reporting of the statistic $\bar{\delta}_{Tg}$ is not suggested. Instead, reporting whether, in the aggregate, a group of students performs better, worse, or as expected on this target is recommended. In some cases, insufficient information will be available, and that will be indicated, as well. For a target within an aggregate group, a minimum amount of precision is required to report target performance for the group. There are no requirements for a minimum number of items or students.

For target-level strengths/weaknesses, the following are reported:

- If $\bar{\delta}_{Tg} \ge +1 * se(\bar{\delta}_{Tg})$, then performance is *better* than on the overall test.
- If $\bar{\delta}_{Tg} \leq -1 * se(\bar{\delta}_{Tg})$, then performance is *worse* than on the overall test.

- Otherwise, performance is *similar to* performance on the test as a whole.
- If $se(\bar{\delta}_{Tg}) > 0.2$, data are insufficient.

6.6.2 Target Scores Relative to Proficiency Standard (Level 3 Cut)

By defining $p_{ij} = p(z_{ij} = 1)$, indicating the probability that student j responds correctly to item i, z_{ij} represents the jth student's score on the ith item. For items with one score point, the 2PL IRT model is used to calculate the expected score on item i for student j with $\theta_{Level\ 3\ cut}$ as:

$$E(z_{ij}) = \frac{exp(Da_i(\theta_{Level\ 3\ cut} - b_i))}{1 + exp(Da_i(\theta_{Level\ 3\ cut} - b_i))}.$$

For items with two or more score points, using the GPCM, the expected score for student j with a Level 3 cut on an item i with a maximum possible score of m_i is calculated as

$$E(z_{ij}) = \sum_{l=1}^{m_i} \frac{lexp(\sum_{k=1}^{l} Da_i(\theta_{Level\ 3\ cut} - b_{i,k}))}{1 + \sum_{l=1}^{m_i} exp(\sum_{k=1}^{l} Da_i(\theta_{Level\ 3\ cut} - b_{i,k}))}.$$

For each item i, the residual between observed and expected score for each student is defined as

$$\delta_{ij} = z_{ij} - E(z_{ij}).$$

Residuals are summed for items within a target. The sum of residuals is divided by the total number of points possible for items within the target, T:

$$\delta_{jT} = \frac{\sum_{i \in T} \delta_{ji}}{\sum_{i \in T} m_i}.$$

For an aggregate unit, a target score is computed by averaging the individual student target scores for the target across all students in the aggregate unit.

$$\bar{\delta}_{Tg} = \frac{1}{n_g} \sum_{j \in g} \delta_{jT}$$
, and $se(\bar{\delta}_{Tg}) = \sqrt{\frac{1}{n_g(n_g-1)} \sum_{j \in g} (\delta_{jT} - \bar{\delta}_{Tg})^2}$,

where n_g is the number of students who responded to any of the items that belong to the target T for an aggregate unit g. If a student did not happen to see any items on a particular target, the student is NOT included in the n_g count for the aggregate.

A difference from zero in these aggregates may indicate that a class, teacher, school, complex, or complex area is more effective (if $\bar{\delta}_{Tg}$ is positive) or less effective (negative $\bar{\delta}_{Tg}$) in teaching a given target.

Direct reporting of the statistic $\bar{\delta}_{Tg}$ is not suggested. Instead, reporting whether, in the aggregate, a group of students performs better, worse, or as expected on this target is recommended. In some cases, insufficient information will be available, and that will be indicated, as well.

For target-level strengths/weaknesses, the following are reported:

• If $\bar{\delta}_{Tg} \ge +1 * se(\bar{\delta}_{Tg})$, then performance is *above* the Proficiency Standard.

- If $\bar{\delta}_{Tg} \leq -1 * se(\bar{\delta}_{Tg})$, then performance is *below* the Proficiency Standard.
- Otherwise, performance is *near* the Proficiency Standard.
- If $se(\bar{\delta}_{Tg}) > 0.2$, data are insufficient.

6.7 HANDSCORING

Constructed-response short-answer (SA) items and essay (i.e., full write) items in English language arts/literacy (ELA/L) and SA items in mathematics for the summative assessments administered by Cambium Assessment Inc. (CAI) are routed to Measurement Incorporated (MI) for scoring. MI provides handscoring using human raters and automated scoring using the Project Essay Grade (PEG) engine. Some Smarter Balanced member states have elected to use handscoring exclusively, while others have elected to use a hybrid automated scoring/handscoring approach. The methods and results for hand scoring and hybrid automated scoring are described in the following sections.

For handscoring items, the total number of items and the summary of rater agreements were calculated based on across all states and territories that participated in the 2023–2024 summative assessments in grades 3–8 and 10–11. Grade 11 data are based on the students in grades 10 and 11.

For the 2023–2024 summative operational item pool, there were a total of 669 ELA/L SA items, 186 ELA/L essay items, and 334 mathematics items. Table 51 shows the number of items by grade and subject.

Table 51. Number of Handscored Items in 2023–2024 Smarter Balanced Summative Item Pool, by Grade and Subject

Cuada	ELA	/L	Mathamatian
Grade	Short Answer	Essay	Mathematics
3	67	25	54
4	77	27	49
5	83	27	86
6	85	20	51
7	86	29	22
8	78	29	30
11	193	29	42
Total	669	186	334

All guidelines for handscoring responses were specified by Smarter Balanced. Outlined below is the handscoring process MI followed in spring 2024 in accordance with the Smarter Balanced guidelines. This process applied to the scoring of all students constructed responses for ELA/L SA and essay items and mathematics items.

6.7.1 Rater Selection

MI has developed a pool of approximately five thousand raters experienced in scoring the Smarter Balanced assessments. MI first recruited qualified raters who had experience scoring these assessments. Rater accuracy data, collected during prior administration scoring, was used to prioritize recruitment of the most accurate, experienced raters. Once recruited, experienced raters were assigned to the content area and grade band(s) with which they were most experienced.

To supplement this pool, MI also recruited raters with experience successfully scoring other large-scale assessments. MI assigned those raters to the grade level, subject area, and item type for which they were most qualified based on their performance on similar projects. Returning raters were selected based on experience and performance, as well as attendance, and cooperation with work procedures and MI policies. MI maintains evaluations and performance data for all staff who work on each scoring project in order to determine employment eligibility for future projects. Finally, MI targeted recruitment of new raters as needed, in an effort to continue to identify talent across the country that will best fulfill the handscoring requirements.

All raters possessed, at a minimum, a four-year college degree. MI collected proof of degree for all raters as a condition of employment. All raters resided in the United States, and properly completed Form I-9 to verify their identity and employment authorization. Raters' I-9 forms are retained on file as required by law and made available for inspection by authorized government officers as needed. MI is an equal-opportunity employer, and believes that a diverse work force is of the utmost importance. When hiring, MI strives to ensure the work force is diverse across age, ethnicity, gender, and other demographic groups.

In selecting team leaders to monitor the raters, MI scoring leadership reviewed records of all returning staff. They looked for people who were experienced team leaders with a record of good performance on previous projects, and they also considered raters who had been recommended for promotion to the team leader position or otherwise displayed exemplary performance.

MI requires all handscoring project staff (scoring directors, team leaders, raters, and clerical staff) to sign a confidentiality/nondisclosure agreement before receiving any training or viewing any secure project materials. The employment agreement indicates that no participant in training and/or scoring may reveal information about the test, the scoring criteria, or the scoring methods to any person.

6.7.2 Rater Training, Qualification, and Scoring

All raters hired to score the Smarter Balanced assessments were trained using the rubric(s), anchor sets, and training/qualifying sets provided by Smarter Balanced. Many of these sets were created during the original field-test scoring in 2014 and approved by Smarter Balanced. Additional sets were created as new items were field-tested. The same anchor sets are used each year. Additionally, MI conducts an annual review of the rater agreement and scoring materials to inform the development of item-specific, supplemental training materials. Supplemental materials are developed each summer and implemented in the subsequent operational administration. These additional materials are developed with a focus on challenging areas identified during the previous operational administration, as indicated by suboptimal rater accuracy (based on validity responses) and/or rater agreement. Supplemental materials may address item- or response-specific concerns. Supplemental materials are also created for newly operational items for which MI identifies a need for additional examples. For instance, MI may find an approach to a mathematics item that was not encountered during field testing but appears frequently during operational scoring, or an uncommon but valid way to address a Research prompt that is not reflected in the existing rubric. In these cases, MI provides examples of these specific approaches along with guidance on how to score them correctly. MI also supplement materials to provide raters with additional guidance for contentwide challenging spots—such as full write conventions—or to help them more accurately identify responses that should be flagged as non-scorable.

Once hired, raters were assigned to a scoring group corresponding to the subject/grade that they were deemed best suited to score. Raters were trained to score a specific item group of either SA (research, brief write, reading, and mathematics) or essay (i.e., full-write) items. Within each item group, raters were

divided into teams supervised by team leaders and a scoring director. Each scoring director, team leader, and rater was assigned a unique ID used to track their scoring work throughout the scoring effort. The number of items an individual rater scored was minimized to allow the rater to more quickly develop experience scoring responses to a small number of items.

All raters, regardless of experience, were required to train on all anchor and training sets. Following training and practice, all raters were required to pass a qualification to prove that they understood and could apply the criteria accurately. The scoring director and team leaders had access to all practice and qualification results, which were reviewed to identify frequently mis-scored responses and inform initial monitoring and feedback needs.

Until a rater had trained and qualified successfully, the rater was not permitted to score operational student responses. Training was structured so that raters understood that all scoring decisions must be grounded in the training materials. In addition, raters learned how to navigate the anchor set, developed the knowledge and flexibility needed to evaluate or escalate a variety of responses, and retained the necessary consistency to score all responses accurately.

When beginning working, all scoring personnel logged in to MI's secure Scoring Resource Center (SRC). SRC includes all online training modules, serves as the portal to MI's Virtual Scoring Center (VSC) interface, and host scoring reports used for rater monitoring. MI's training system (VSC Train) provides a remote, secure application for training both team leaders and raters. VSC Train provided each trainee with a training lesson for each item that allowed the trainee to complete the following steps:

- 1) Review the anchor set(s)
- 2) Score the practice set(s)
- 3) Review an annotated version of the practice set(s) after submitting scores
- 4) Score the qualification sets

Training and qualification design varied slightly depending on Smarter Balanced item type:

- ELA/L full write: Raters trained and qualified on a baseline training lesson for a grade and writing purpose (e.g., grade 3 narrative, grade 6 argumentative, etc.). After qualifying on the baseline, raters then completed qualifying sets for each item associated with that grade and purpose. Raters could only score those items for which they have passed the qualifying set.
- ELA/L brief write, reading, and research SA: Raters trained and qualified on a baseline lesson within a specific grade band and target. Qualification on the baseline lesson permitted the rater to score all items in that grade band and target.
- Mathematics SA: Raters trained and qualified on baseline lessons within a specific grade band. Qualification on a baseline lesson permitted the rater to score that item and all items associated with it; for items with no associated items, training was for the specific item.

An additional validation stage supplemented full write, brief write, reading, and research rater qualification. Following the training and qualification steps described above, all prospective full write, brief write, reading, and research raters were required to score, for most items, a 20-response set of prescored student responses sourced from the prior test administration. Like the qualification step, raters were required to meet accuracy standards during this validation to score operational responses for a given item. Any raters who failed to meet validation accuracy standards were automatically disqualified from scoring

the item despite having passed qualification. This additional validation matches the full write qualification methods that have been in place since the start of Smarter Balanced scoring in 2015 and adds an additional level of quality assurance.

Rater training time varied by grade and content area. Training for SA brief write, reading, research, and mathematics items could typically be accomplished in one day, while training for essay items took up to five days to complete. Raters generally worked 3-7 hours per day. The hours worked per day were flexible, based on the raters' shift preference and item(s) being scored. At a minimum, most raters scored 15 hours per week (day shift) or 10 hours per week (evening shift), with many scoring over 30 hours per week (day shift) or 20 hours per week (evening shift).

In addition to item-specific scoring expectations, a variety of substantive procedural and policy information was provided to each trainee during training. These included instructions for how to identify and flag particular types of responses as well as how to communicate with leadership during handscoring.

Raters were trained to recognize non-scorable responses, and these responses were systematically routed to scoring supervisors for final condition-code assignment per Smarter Balanced requirements. For some item types, such as essays, condition-code responses were scored by scoring leaders trained to specialize in the scoring of these types of responses.

An "alerts" procedure was explained to raters during training sessions, where raters are trained to recognize "alerts" in their various forms, including those for suicide, criminal activity, alcohol or drug use, extreme depression, violence, rape, sexual or physical abuse, self-harm, intent to harm others, and neglect.

The training process, including this additional information, ensured that raters were fully prepared to hand score responses and understood all responsibilities and scoring requirements before they began operational scoring.

Following training, all training materials remained available to raters throughout scoring via the VSC Score Resource Library. This library included the item and rubric, the annotated anchor and practice sets, and any associated supplemental materials.

When scoring, raters had access only to those items for which they had successfully trained and qualified. The handscoring system sorts individual student responses into small sets of 5-10, grouped by item. When a rater is qualified to score multiple items, this approach eases cognitive load by presenting the rater with a scoring set in which all responses relate to the same item.

Multiple strategies were employed to minimize rater bias during scoring. First, raters did not have access to any student identifiers. Unless the students signed their names, wrote about their hometowns, or in some way provided other identifying information as part of their response, the raters had no knowledge of student characteristics. Second, all raters were trained using Smarter Balanced–provided materials, which were approved as unbiased examples of responses at the various score points. Training involved constant comparisons with the rubric and anchor papers so that raters' judgments were based solely on the scoring criteria. Finally, following training, a cycle of diagnosis and feedback was maintained to identify any issues. Specifically, raters were closely monitored during scoring, and any instances of raters making scoring decisions based on anything except the criteria were discussed with the raters. After this feedback had been provided, raters were further monitored, and if any continue to exhibit bias after receiving a reasonable amount of feedback, they were dismissed.

A series of automated score verifications were implemented to further ensure the accuracy of scores. For example, a blank check was conducted, which reset scores when a condition code of "blank" was assigned to a response that had one or more characters in the response string (e.g., a response comprised of spaces or tabs). In this case, only after three independent raters had assigned a condition code of "blank" to a response that appeared blank, but which included characters in the response string, was the score recorded. A similar check was run when a score or condition code other than "blank" was assigned to a response that included no characters in the response string. Automatic resetting of double-scored responses when two raters assign non-adjacent scores, mismatched condition codes, or a combination of a condition code and a numeric score provided an additional score verification. In addition to automatically resetting and rescoring these responses, the raters' information was captured in a report and reviewed by scoring directors, one of many tools used to determine retraining needs.

6.7.3 Rater Monitoring, Feedback, and Evaluation

During operational scoring, five percent of the responses scored comprised pre-approved validity responses. Validity responses serve as benchmark responses as the most appropriate score for each validity response is predetermined by key stakeholders. A small set of validity responses is provided by Smarter Balanced for all vendors to use, and these are supplemented with responses selected and approved by MI scoring management. The validity pool includes anchor validity responses originating from the field test administration. The pool of validity responses is selected to be generally representative of operational responses, while ensuring sufficient examples of each score point. Validity results compare the score assigned by a rater to a validity response with the benchmark score of the same response. Validity responses provide a more direct measurement of rating quality than measures of inter-rater reliability (Raczynski et al., 2015).

MI calibrates validity responses to fit a unidimensional Item Response Theory (IRT) model for each content area/item type. This approach involves transforming raters' validity response scores into accuracy scores. Specifically, if the rater's score matches the "true" score of the validity response, an accuracy score of 2 is assigned. If the rater's score is adjacent to the score of the validity response, an accuracy score of 1 is assigned. Otherwise, for scores that are non-adjacent, an accuracy score of 0 is assigned. All accuracy score data for validity responses and raters are then fitted to a Generalized Partial Credit Model (GPCM) IRT model. Utilizing the resulting IRT parameters, MI calculates accuracy values for each rater based on a given set of validity responses. This calculation is conducted several times each day during scoring, providing real-time measures of rater accuracy.

In addition to validity responses, 15% of handscored responses received blind second reads, the results of which were used to calculate inter-rater reliability. To support interpretability, second reads were conducted exclusively by expert (i.e., highly-accurate) raters, described further below.

The VSC system automatically and randomly routed the requisite number of responses to raters for second reads and validity in an inconspicuous manner. In this way raters had no means of discerning whether they

-

¹ Responses and results of the 2014-15 Smarter Balanced field test administration were used to derive the base scale to which subsequent item parameters are aligned.

were scoring a first read, a second read, or a validity response. This system also prohibited raters from being eligible to score second reads for responses they had already scored.

Scoring accuracy during handscoring was maintained by continuously assessing rater performance using validity responses. MI specifically evaluated how closely raters' scores aligned with the benchmark scores of these validity responses. Key performance measures included the agreement between rater and benchmark scores, quantified using Quadratic Weighted Kappa (QWK)², and the comparison of mean score differences between the distributions of benchmark and rater-assigned scores.

The system automatically generated performance metrics several times a day based on the most recent data, providing raters and scoring managers with daily, automated summaries of rater performance. This ensured that all handscoring staff were kept informed of their current performance and any issues that needed attention. In addition to these daily summaries, detailed manager-level reports were produced to identify raters who required retraining or, if necessary, removal due to accuracy or productivity concerns. These reports enabled scoring management to direct scoring leaders to specific VSC reports, allowing them to pinpoint the areas where individual raters needed improvement.

The monitoring system afforded the objective, dynamic identification of the most accurate raters, referred to as "expert raters." Specifically, expert raters are those who demonstrate highly accurate and consistent scoring of validity responses. Rater status changed daily based on current rater performance to ensure that any rater drift did not negatively impact scoring accuracy. Expert rater status was a precondition for conducting second readings.

During scoring, raters received automated feedback system based on recent performance. The automated feedback system identifies raters who require additional feedback—based on accuracy metrics—and automatically generates a custom set of responses for the rater to review. The system functions at the item level, thus providing feedback even to those raters with relatively high accuracy when the data identifies there are one or more items on which they can improve.

VSC provided real-time reports throughout the scoring effort. These reports were available for access by handscoring management and clients. Inter-rater reliability reports provide the percentage of exact, adjacent, and non-adjacent agreement for scorable responses. Score point frequency distribution reports provide the percentage per score point and include the mean and standard deviation for each item. Validity performance reports provide the percentage of exact, adjacent, and non-adjacent agreement for validity responses and were used to monitor drift. Validity performance reports are typically used to monitor and correct drift at the group level. If the data indicate that raters as a group are scoring validity responses either consistently high or consistently low, leadership will recalibrate the group by having raters review key training responses that reflect the types of responses being missed in validity. Leadership may also provide raters with a supplemental set of responses that help reinforce the lines for the various score-points and re-anchor the raters to the proper position, arresting groupwide drift.

Reports using item-level accuracy expectations identified any items not meeting the expected levels of agreement. Specifically, these reports indicated the difference between expected accuracy and current accuracy for each item. Expected accuracy was defined based on historical data; in some cases (e.g., most

² QWK is a measure used to assess the agreement between two raters, accounting for the possibility of agreement occurring by chance and giving more weight to larger discrepancies between ratings.

Mathematics items) expected accuracy exceeded Smarter Balanced's minimum accuracy thresholds. In this way, reports informed improvements to the scoring accuracy of all items.

Automated removal of raters and score resets were performed when item and rater performance failed to meet accuracy expectations. In these cases, all responses scored by a rater during a period of poor performance were reset and redistributed to other qualified raters for rescoring By limiting raters to scoring relatively fewer items, this approach also maximized accuracy across items.

In addition to the automated feedback, scoring leadership provided individualized feedback to raters based on their performance. Specifically, leadership reviewed the rater's mis-scored validity responses and associated data and looked for a trend that suggests the rater has drifted from the anchored responses. If such a trend is present, leadership can tailor feedback specific to that rater, typically by presenting them with live responses they have mis-scored in a way that is reflective of their overall drift from the anchor set criteria and providing targeted, thoughtful rationales for the "correct" scores.

Finally, as a supplement to automated assessments, team leaders spot-checked (i.e., read behind) raters' scoring to ensure that the raters were on target, and conducted one-on-one retraining sessions to address any problems found. At the beginning of the project, team leaders read behind every rater every day; they became more selective about the frequency and number of read-behinds as raters became more proficient at scoring.

6.7.4 Rater Agreement

Rater inter-rater reliability (IRR) was computed based only on scorable responses (numeric scores) scored by two independent raters. Non-scorable responses (e.g., off-topic, off-purpose, or foreign-language responses) were scored by scoring leadership per the handscoring rules—and not by one expert and one random rater—and were thus excluded from IRR computations. For the handscored items, the human-human agreement was computed based on the combined data across all states and territories that participated in the 2023–2024 summative assessment.

In ELA/L essay (i.e., full writes) item responses were scored in three dimensions: conventions (0–2 rubric), evidence/elaboration (1–4 rubric), and organization/purpose (1–4 rubric). All ELA/L SA items were scored using a 0–2 rubric. Mathematics SA items were scored using 0–1, 0–2, or 0–3 rubrics.

Tables 52 through 54 provide a summary of the human-human IRR based on items with a sample size greater than or equal to 50. For Mathematics and ELA/L essay items, the tables show the majority of the items administered. For ELA/L SA items, relatively fewer items reached a sample size greater than or equal to 50, and thus a subset of the items administered are represented in the tables. The IRR is presented with mean of percent exact agreement, minimum and maximum percent exact agreements, combined percent exact and adjacent agreement, and the mean, minimum and maximum QWK. The average number of responses, as well as minimum and maximum number of responses to a given item are presented as well.

Table 52. Inter-Rater Agreement for ELA/L Short-Answer Items

Grade	Number	Number of Responses			%Exact			%(Exact+	QWK		
	of Items	Mean	Min	Max	Mean	Min	Max	Adjacent)	Mean	Min	Max
3	30	483.0	61	885	77.1	61.2	89.0	100.0	0.69	0.48	0.87
4	41	442.3	67	1103	74.9	56.1	86.5	100.0	0.70	0.33	0.90
5	38	454.7	76	1095	71.7	57.0	84.6	100.0	0.69	0.43	0.87
6	54	687.0	50	3178	73.2	57.1	86.7	100.0	0.65	0.29	0.82
7	56	656.9	54	2339	74.2	61.5	85.2	100.0	0.70	0.46	0.82
8	60	655.6	56	2898	73.0	64.1	89.8	100.0	0.70	0.43	0.88
11	84	435.0	53	1092	73.3	60.1	92.6	100.0	0.71	0.46	0.94

Table 53. Inter-Rater Agreement for ELA/L Essay Items

Grade	Dimension	Number		mber o			%Exact		%(Exact+		QWK	
		of Items	Mean	Min	Max	Mean	Min	Max	Adjacent)	Mean	Min	Max
	Conventions	25	693.6	330	924	70.6	64.2	77.7	100.0	0.67	0.63	0.74
3	Evid/Elab	25	693.6	330	924	69.7	55.8	79.2	100.0	0.72	0.64	0.80
	Org/Purp	25	693.6	330	924	69.7	56.9	79.4	100.0	0.72	0.64	0.80
	Conventions	27	750.4	443	1019	67.6	62.7	73.8	100.0	0.72	0.62	0.81
4	Evid/Elab	27	750.4	443	1019	69.3	60.9	79.7	100.0	0.74	0.66	0.83
	Org/Purp	27	750.4	443	1019	69.2	60.9	79.4	100.0	0.74	0.64	0.83
	Conventions	27	818.9	554	1028	69.1	58.6	77.9	100.0	0.66	0.59	0.74
5	Evid/Elab	27	818.9	554	1028	68.0	59.8	73.4	100.0	0.75	0.71	0.80
	Org/Purp	27	818.9	554	1028	68.1	60.6	73.8	100.0	0.76	0.71	0.80
	Conventions	20	999.1	708	1217	70.9	67.0	74.2	100.0	0.68	0.63	0.72
6	Evid/Elab	20	999.1	708	1217	69.3	63.6	74.5	100.0	0.73	0.67	0.81
	Org/Purp	20	999.1	708	1217	68.9	62.8	74.2	100.0	0.73	0.67	0.81
	Conventions	29	700.1	444	884	70.5	66.6	74.3	100.0	0.68	0.62	0.74
7	Evid/Elab	29	700.1	444	884	71.5	65.1	77.7	100.0	0.77	0.72	0.81
	Org/Purp	29	700.1	444	884	71.5	64.8	77.3	100.0	0.77	0.71	0.81
	Conventions	29	742.3	524	910	74.6	70.2	80.6	100.0	0.68	0.60	0.75
8	Evid/Elab	29	742.3	524	910	72.3	64.9	81.5	100.0	0.77	0.72	0.84
	Org/Purp	29	742.3	524	910	72.0	63.7	81.4	100.0	0.77	0.72	0.83
	Conventions	29	711.2	549	843	72.5	68.3	76.1	100.0	0.69	0.63	0.75
11	Evid/Elab	29	711.2	549	843	73.3	68.2	78.0	100.0	0.79	0.73	0.82
	Org/Purp	29	711.2	549	843	73.4	68.4	78.0	100.0	0.79	0.73	0.82

Note. Evid/Elab: Evidence/Elaboration, Org/Purp: Organization/Purpose

Number of Score OWK^a Number %Exact %(Exact+ Responses Grade **Point** Adjacent) of Items Mean Mean Range Min Max Mean Min Max Min Max 3 0 - 113 1086.0 763 1312 93.1 85.6 97.7 100.0 NA NA NA 4 0 - 110 1324.3 1113 1564 89.2 85.2 95.5 100.0 NA NA NA 5 0 - 1951.9 818 91.5 82.0 98.1 12 1101 100.0 NA NA NA 0 - 197.1 99.8 6 10 1315.7 693 2047 89.4 100.0 NA NA NA 7 0 - 18 1989.2 1641 2575 95.9 87.3 99.0 100.0 NA NA NA 8 0 - 19 2542.2 2023 2736 89.2 85.7 96.9 100.0 NA NA NA 0 - 11300.7 142 1901 94.2 89.0 100 11 16 100.0 NA NA NA 3 0-2 35 1121.7 337 1611 91.0 78.0 100 0.92 1.00 100.0 0.58 4 0-21270.0 439 91.0 99.2 0.90 0.99 35 1819 81.2 100.0 0.72 5 0-265 1002.7 646 1300 88.3 74.2 96.8 100.0 0.88 0.67 0.97 6 0-241 1691.4 1291 1975 88.5 78.2 98.9 100.0 0.85 0.75 0.99 2214.1 95.4 0.96 7 0-213 1634 2509 90.7 82.7 100.0 0.86 0.69 8 0-219 2287.6 1848 2821 89.4 78.0 97.7 100.0 0.87 0.75 0.98 11 0-220 1492.5 614 1937 92.0 80.6 99.2 100.0 0.90 0.74 0.99 3 0-36 934.0 611 1629 91.1 86.9 94.8 100.0 0.96 0.94 0.98 4 0 - 34 1163.5 1095 1353 88.8 88.3 90.0 100.0 0.95 0.95 0.96 5 0-39 993.7 605 1218 87.5 83.0 97.2 100.0 0.90 0.83 0.95 7 0-31 2405.0 2405 2405 92.6 92.6 92.6 100.0 0.93 0.93 0.93 8 0-32 2572.5 2529 2616 84.6 84.5 84.7 100.0 0.95 0.94 0.95 0 - 31668 1847 88.8 82.2 93.8 100.0 0.91 0.89 0.94 11 6 1765.3

Table 54. Inter-Rater Agreement for Mathematics Items

Note. ^a QWK is not presented for 0-1 items due to the binary score scale.

6.8 AUTOMATED SCORING

MI's PEG automated scoring technology was used to score eligible SA and essay items in ELA/L and SA items in mathematics. This section describes PEG, the training and validation sample and process, and the automated scoring process, concluding with the human-machine (HM) agreement statistics.

6.8.1 Project Essay Grade

Figure 13 presents the architecture of MI's PEG engine. During engine training, this architecture allows PEG to generate hundreds of custom linguistic (rule-based) features, which are determined by codified English linguistic rules such as syntax and semantics and extracted from representative student responses. In addition to rule-based features, PEG also includes features extracted by Latent Semantic Analysis (LSA) and Latent Dirichlet Allocation (LDA) procedures.

PEG's item and trait specific scoring models use computed features from the training responses along with the scores assigned to them by expert human raters. Using hundreds of parameterizations across several machine-learning algorithms, via cross-validation and optimization, PEG determines which algorithms best predict the expert-assigned scores. These algorithms draw on many of the latest advances in the field of machine learning to generate linear and non-linear classification and regression models. These approaches typically result in 100 candidate models for a single item or trait. PEG then uses an ensembling procedure to combine the best models into a robust final model. The ensembling procedure utilizes a linear

regression, where the objective is to maximize a continuous relaxation of the quadratic-weighted-kappa (QWK) metric, thus maximizing PEG's agreement with the expert human raters.

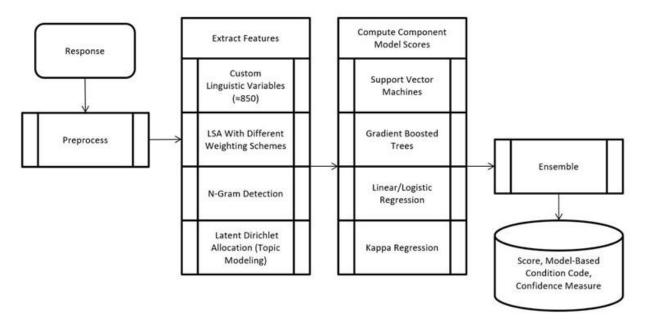


Figure 13. PEG Architecture

The sections that follow describe the process used to train and validate the engine, followed by a description and results of the hybrid human-automated scoring process.

6.8.2 Model Training and Validation

Sample

Automated scoring models were not created for items that had an insufficient quantity of training responses. This was this case for items with low exposure to students, as dictated by the adaptive testing algorithm. Additionally, mathematics performance task items that had multiple parts with scoring dependencies were not considered for automated scoring. Table 55 shows that pretrained models existed for 595 items, thus, no additional training was conducted in preparation for the spring 2024 administration. The remainder of this section describes the process used to train and validate the 595 existing models.

Table 55. Number of Items Eligible for Automated Scoring, by Grade and Subject Area

Grade	Items With Existing Models			Items Without Models		
	ELA/L		Mathamatics	ELA/L		Mathamatics
	Short-Answer	Essay	Mathematics	Short-Answer	Essay	Mathematics
3	12	13	44	0	0	0
4	13	16	42	0	0	0
5	13	10	50	0	0	0
6	32	10	41	0	0	0
7	45	17	15	0	0	0
8	49	14	24	0	0	0
11	80	17	38	0	0	0
Total	244	97	254	0	0	0

Training Data

Student responses used for training and validation were sourced from the 2018–2019, 2020–2021, 2021–2022, and 2022–2023 Smarter Balanced operational test administrations. Responses were randomly sampled from available on-grade responses in the operational population. For all items, the sample included 1,500–2,000 responses, stratified by score point. The score of record used to train the engine was the score assigned to each response by an expert rater.

For each item, the sample was divided as follows:

- Approximately 85% of the responses were assigned to a training set used to build the model.
- Approximately 15% of the responses were assigned to a validation set used to evaluate the accuracy of the model.

Model Training

Component model training requires inputs of response "features." For items that assess writing quality (e.g., essays), PEG processes the responses and calculates approximately 850 linguistic variables that describe the responses in mathematical terms. These variables range in complexity from simple to highly complex. Examples of simple variables are measures such as word count or sentence length, word choice and spelling errors, and the number and severity of grammatical errors. The most complex variables measure patterns that represent style, fluidity, smoothness of transitions, clarity of communication, and other sophisticated concepts.

For content-based items (e.g., SA mathematics items), the number of variables is unknown until the models are built. Because the content varies significantly from item to item, and therefore from model to model, PEG examines training responses and identifies the variables that most accurately capture the content in question. To do this, MI uses techniques like LSA, N-Gram Detection, and LDA. To further refine the variable generation process, MI built a computer language to perform a simultaneous search over semantic, lexicographic and syntactic features of responses.

To build an essay scoring model, PEG examines the variables and text features of responses, correlates them with the human scores previously assigned, and identifies those variables that have high predictive value.

To build a content scoring model, PEG analyzes training responses and calculates features that pertain to the content in question. PEG then sends the features to hundreds of different algorithms that compete to see which algorithms best associate the features with the human-assigned scores. These algorithms draw on many of the latest advances in the field of machine learning to generate both linear and non-linear models. Examples of approaches used include Support Vector Machines, Gradient Boosted Trees, and various regression approaches.

Note that building component models for each item—and for multi-dimensional items, each trait or dimension—prevents variables from being generalized across items or traits, allowing PEG to faithfully reproduce humans' application of the scoring rubrics. This means that the resultant models are reasonably robust to gaming attempts, as each represents a unique valuation of the item- (or trait-) specific text features similarly valued by expert professional raters.

The approaches just described typically result in 100 models for a single item or essay trait. Ensembling is the process of selecting the "best of the best" models, to result in a small set of strong, yet dissimilar component models. A linear-kappa regression is used to determine the model ensembling weights. The more accurate a given model is, the more weight it carries in the final score decision.

Scoring a response involves first preprocessing the response. The purpose of preprocessing is twofold: (1) create raw and canonical representations of the response from which features can be extracted, and (2) filter out responses for which the scoring model does not apply (e.g., blank or insufficient responses). The response is then scored with the associated component models. A final score is produced performing a weighted sum using the ensembling weights.

Model Validation

Model validation involved a two-phase approach: an initial validation using held-out training data and a secondary validation using operational data from the current administration.

Initial Validation

Initial validation was conducted by applying each model to score a respective validation set of responses. The validation set is independent of the training set, in that none of the responses it contains have been used to build the model. Two or more professional raters will not always agree on what score to give a student's response; therefore, modeling is considered successful when the engine produces scores that agree with professional raters to the same or greater extent than the raters agree with each other. The initial evaluation was made using the criteria shown in Table 56, based on criteria proposed by Williamson, Xi, and Breyer (2012). While Williamson et al. (2012) recommend an agreement between human and machine scores of 0.70 quadratic weighted kappa (QWK) for normally distributed data, a QWK threshold of 0.65 was adopted due to the prevalence of skewed distributions in response data. The degradation (QWK) criterion of .07 is slightly more stringent than proposed by Williamson et al. (2012). The evaluation process was used for both the item-specific scoring models and the condition code models.

Table 56. Initial Model Evaluation Criteria

Criterion	Threshold
Agreement of automated scores with human scores	$QWK_{H:M} \ge 0.65$
Degradation from the human-human score agreement	$QWK_{H:H} - QWK_{H:M} < 0.07$
Standardized mean score difference between human and automated scores	$ SMD_{H:M} < 0.15$

Note. QWK = Quadratic weighted kappa. SMD = Standardized mean difference. H:H = human:human. H:M = human:machine.

Bias Considerations. Subgroup differences in responses to constructed response items can introduce construct-irrelevant variance in scores, in turn threatening valid score interpretations. MI investigated potential sources of bias annually, for newly modeled items, as part of the initial validation process using available data from previous summative administration. Table 57 shows the demographic variables and categories considered. MI received separate datafiles containing (1) hand score data and (2) student demographic data associated with responses.

Demographic Variable	Categories
Gender	Male
Gender	Female
	American Indian or Alaska Native
	Asian
	Native Hawaiian or Pacific Islander
D /Ed : '4-	Filipino
Race/Ethnicity	Hispanic or Latino
	Black or African American
	White
	Two or More Races
LEDG	LEP
LEP Status	Non LEP

Table 57. Demographic Variables and Categories

For each new item being modeled, analysis was performed on a subgroup if the number of observations (i.e., human-machine scores) was at least 10. A subgroup was flagged for bias if $|SMD| \ge 0.125$ and if the SMD was significant at an overall significance level of 95%. A Bonferroni correction was used to adjust the significance level for each subgroup comparison. An item was flagged for bias, excluded from automated scoring, and handscored if any subgroup comparison associated with the item was flagged.

Secondary Validation

All models associated with items that passed initial validation were subject to a secondary validation at the start of the spring 2024 administration using an early sample of operational responses from that administration. This sample was comprised of the first available 500 responses/item across states, at a minimum. Responses from this sample were scored by both the automated scoring engine and an expert rater. During this interval the human score was reported as the score of record. If the PEG scores were found to be consistent with the scores assigned by the expert raters, subsequent student responses for a given item were scored by PEG using a hybrid human-automated scoring approach. If not, the item was handscored. Table 58 presents the secondary validation criteria. Note that since expert raters are the only humans that score the secondary validation sample, a second human score is not collected and thus QWK degradation is not part of the criteria.

Table 58. Secondary Validation Criteria

Criterion	Threshold
Agreement of automated scores with human scores	$QWK_{H:M} \ge 0.65$
Standardized mean score difference between human and automated scores	$ SMD_{H:M} \leq 0.15$

Note. QWK = Quadratic weighted kappa. SMD = Standardized mean difference. H:M = human:machine.

Table 59 presents the secondary validation results. Of the 595 items with models subject to secondary validation, models associated with 454 of the items (76.3%) passed all secondary evaluation criteria.

Table 59. Summary of Secondary Validation Results, by Grade and Subject Area

G 1		All Models Passi lidation Criteria	0	g Secondary a			
Grade	ELA	/L	Mathamatics	ELA/L		Mathamatics	
	Short-Answer	Essay	Mathematics	Short-Answer	Essay	Mathematics	
3	12	13	44	12	3	44	
4	13	16	42	13	6	40	
5	13	10	50	13	5	47	
6	32	10	41	19	5	40	
7	45	17	15	27	9	15	
8	49	14	24	31	9	22	
11	80	17	38	46	10	38	
Total	244	97	254	161	47	246	

Live Training and Validation

Additionally, in April-May 2024 when operational scoring was underway, a live training and validation effort was undertaken for those handscored items lacking validated models from prior efforts but having sufficient 2024 operational responses to train and validate new models. In general, these items were associated with models that had previously failed an initial and/or secondary validation. In such cases, training with 2024 operational responses offered potential to improve model performance. All models associated with these items were thus trained using either exclusively 2024 responses (when a minimum of 1,400 2024 responses/item existed) or 2024 responses supplemented with 2023 responses. In either case, the validation sets consisted exclusively of 2024 responses. Because live validation involved operational data, it was unnecessary to conduct a secondary validation.

Table 60 summarizes the results of the live training and validation. Of the 356 items associated with models that underwent live training and validation, models associated with 211 of the items (59.3%) passed all evaluation criteria. While this pass rate is considerably lower than the pass rates during secondary (76.3%) validation efforts, it is most likely explained by the nature of the items modeled. Specifically, since all item models in this sample had failed a prior validation, by design the sample consisted of difficult-to-model items.

12

89

Grade	I	tems Trained		Items with All Valid	l Models Pas lation Criter	
Graue	ELA/L		3.7.11	ELA/L	1	3.4 d
	Short-Answer	Essay	Mathematics	Short-Answer	Essay	Mathematics
3	1	25	9	1	16	4
4	3	24	9	3	19	1
5	1	25	33	1	14	19
6	24	16	10	15	10	4
7	28	20	7	18	12	4
8	26	25	9	17	6	7

Table 60. Summary of Live Training and Validation Results, by Grade and Subject Area

Following initial validation, secondary validation, and live training and validation, a total of 665 items, comprised of 240 ELA/L SA, 136 essay, and 289 mathematics SA, were scored using a hybrid process, described next.

4

81

24

79

21

156

6.8.3 Automated Scoring Processes

36

119

Hybrid Scoring Process

11

Total

As all models associated with a given item passed secondary validation (or live validation), subsequent student responses were scored using a hybrid human-automated scoring approach. If all models associated with a given item did not pass secondary validation, responses associated with the item continued to be handscored by the larger pool of raters. These raters were monitored and evaluated as described in the handscoring section above.

Figure 14 shows the response routing rules under the hybrid scoring process. In the hybrid model, responses with associated scoring models were first pre-processed for automated scoring and to filter alert responses and certain non-scorable cases (e.g., insufficient text to score or high proportion of copied prompt text). Flags were used to indicate condition codes as defined in the handscoring criteria (see Table 61 and Table 62). For example, PEG flags responses that lack proper development, lack enough content to be scored, are written in an unsupported language, or contain vulgar language or other alert words or phrases that indicate that the response should be reviewed by the client. Responses were then sent to the automated scoring engine, where text features were extracted, the scoring model(s) applied, and responses assigned a score and measure of score confidence. Low-confidence responses straddle the lines between score point values on a rubric and are difficult to score accurately because they exhibit characteristics of multiple score points Higher-confidence responses received the engine score as the score of record, while lower-confidence responses were routed directly to expert raters, who assigned the score of record. Note that the expert rater pool was dynamic, and raters were added or removed several times each day based on their current performance. Overall, approximately 15% of responses to engine-scored items were flagged as low confidence and scored by expert raters.

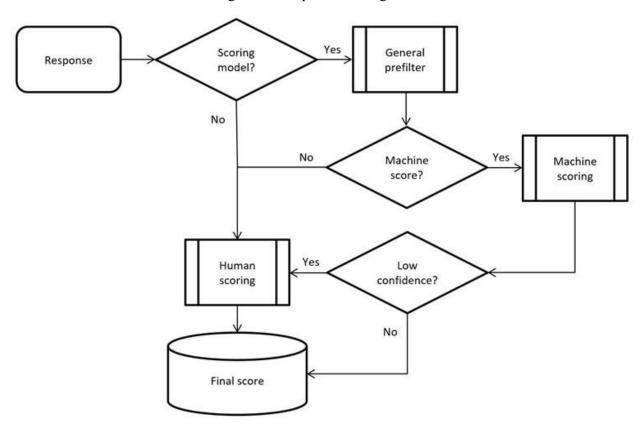


Figure 14. Response Routing Rules

Upon receipt and validation of each response, MI routed responses for those items eligible for automated scoring to PEG and the remainder of the responses to the VSC handscoring system.

Table 61. Flags Currently Established

FLAG	USAGE DESCRIPTION	*SCORABLE
0	Standard scoring	YES
200	Too few words (i.e., blank, or extremely short response)	NO
240	Too long (i.e., too many characters submitted; 30,000 characters is the current limit)	NO
250	Expected essay fields are null or empty; set when nulls are discovered within the processing pipeline. Not client configurable.	NO
400	Unexpected item_id (i.e., the item_id is not one of the items PEG AI has modeled)	NO
500	Scorable alert (i.e., an essay which seems perfectly scorable, but happens to contain alert language); client may configure alert scanning to "on" or "off", but other changes are not recommended.	YES
501-599	Non-scorable alert (i.e., alert language was detected, and the essay could not be scored). If alert scanning is "on", then any code in the 500-599 range is possible. Not client configurable.	NO
620	Applies when the ratio of copied characters exceeds specified threshold (e.g.; 0.5 means 50%). Can be used for all Smarter items for which prompt content was provided.	YES
650	Insufficient Condition Code (I): Response holds strong general resemblance to those marked 'Insufficient' by human readers, but is nonetheless PEG scorable (and, so scores are provided). <u>PEG Configuration</u> : Item agnostic; but for 2021 onwards, applicable to ELA/L items only.	YES

FLAG	USAGE DESCRIPTION	*SCORABLE
660	Language Non-English Condition Code (L): Response holds strong general resemblance to those marked 'Non-English' by human readers, but is nonetheless PEG scorable (and, so scores are provided). *PEG Configuration*: Item agnostic; but for 2021 onwards, applicable to ELA/L items only.	YES
670	Off-Topic: Applicable to ELA/L essays only and is item specific in the PEG environment.	YES
680	Off-Mode: Applicable to ELA/L essays only and is item specific in the PEG environment.	YES
900	Timeout (i.e., unable to complete essay score prediction within time limits). Not client configurable.	NO
950	System error processing essay (i.e., internal PEG error). Not client configurable.	NO

Note. Scorable flags indicate instances where PEG will return both the applicable flag and a score.

Table 62. Model Setting

ТҮРЕ	ASSOCIATED FLAG(S)	DESCRIPTION	VALUES
Minimum Words	200	Triggers if there are fewer than the associated value of word-tokens in a response. The flag may also appear regardless of setting if the response is blank.	0-15
Alert	500	Current setting (PREDC1) is for the	Standard settings in
Aleit	501-599	standard alert scan.	place
Plagiarism	620	Prompt and source material text is included in model configuration.	50% of prompt and source material characters triggers flag

Scoring Infrastructure

During the automated scoring process, response data are transferred from CAI to MI's IT project team. Data are then passed to PEG from the IT project team via an internal server, at which point they are processed through the PEG Streaming Scoring Service—a cloud-deployed, horizontally scalable, distributed parallel computing application. Scored batches were typically completed within one day. All data are then transferred from PEG to the IT project team, who ultimately sends the data/scores back to CAI.

Quality Assurance

MI's hybrid scoring approach included numerous quality assurance steps. First, models were trained using exclusively scores assigned by expert raters and the associated responses. Second, each automated scoring model was subjected to an evaluation process, as described in the model validation section. This involved evaluating the quality of the human-scored training data, as well as comparing the performance of the engine to the performance of expert raters. Third, for models trained using responses from prior administrations, the generalizability of each model to the 2023-24 operational responses was confirmed via a secondary validation. Finally, quality was further assured during scoring by routing a minimum of 15% of the responses that were most different from the training responses to expert raters and assigning the human score.

"Alert" Procedures

MI implemented a formal process for informing clients when student responses reflect a possibly dangerous situation for the test-taker. Specifically, MI employed a set of alert procedures to notify the client of responses indicating endangerment, abuse, or psychological and/or emotional difficulties. PEG employed a rule-based detection system to flag responses that are indicative of potentially dangerous situations. Responses flagged by PEG as possible alerts were reviewed by scoring leadership, who decided whether each response should be forwarded to the client. Once vetted, all alerts were provided to CAI, who associated the pertinent student information with the response(s) and contacts the state. In addition, CAI separately evaluates all responses and student-generated text for possible alerts.

Score Delivery

As scores were assigned by PEG, MI verified and delivered them to CAI. MI received confirmation from CAI that each response had been received and had passed data validation.

6.8.4 Human-Machine Agreement

This section summarizes the human-machine agreement for all items scored using a hybrid process in spring 2024, including (1) items passing initial model validation, (2) items passing secondary validation, and (3) items passing live validation.

Tables 63 through 65 present the human-machine agreement on the initial and secondary validation samples for ELA/L SA items, ELA/L essay items, and mathematics SA items, respectively. For the PEG-scored items, the human-machine agreement was computed based on the combined data across all states with hybrid scoring in the 2023–2024 summative assessment.

Table 63. Human-Machine Agreement for ELA/L Short-Answer Items on Initial and Secondary Validation Samples, by Grade

		Initial V	alidation		Secondary Validation			
Grade	Number of Items	% Exact	% (Exact + Adjacent)	QWK	Number of Items	% Exact	% (Exact + Adjacent)	QWK
3	12	79.6	99.6	0.81	12	82.3	99.5	0.77
4	13	80.1	99.8	0.84	13	80.9	99.8	0.80
5	13	75.4	99.6	0.81	13	77.4	99.8	0.78
6	19	78.7	99.5	0.81	19	79.1	99.6	0.77
7	27	76.3	99.4	0.79	27	76.4	99.4	0.75
8	31	76.2	99.5	0.78	31	75.8	99.4	0.75
11	46	77.2	99.5	0.79	46	76.1	99.5	0.77

Table 64. Human-Machine Agreement for ELA/L Essay Items on Initial and Secondary Validation Samples, by Grade

			Initial V	alidation		\$	Secondar	y Validation	
Grade	Trait	Number of Items	% Exact	% (Exact+ Adjacent)	QWK	Number of Items	% Exact	% (Exact+ Adjacent)	QWK
3	Conventions	3	71.6	99.7	0.72	3	72.5	99.5	0.70
3	Evid/Elab	3	77.9	99.2	0.82	3	78.2	99.7	0.77
3	Org/Purp	3	75.0	99.7	0.8	3	79.1	99.6	0.78
4	Conventions	6	69.2	99.0	0.74	6	69.7	99.3	0.74
4	Evid/Elab	6	73.6	99.5	0.84	6	73.5	99.1	0.79
4	Org/Purp	6	72.2	99.2	0.82	6	74.2	99.2	0.79
5	Conventions	5	72.5	99.6	0.71	5	73.0	99.6	0.72
5	Evid/Elab	5	73.0	99.0	0.82	5	72.6	99.6	0.80
5	Org/Purp	5	72.2	99.6	0.83	5	72.7	99.6	0.80
6	Conventions	5	75.5	99.0	0.72	5	73.5	99.5	0.74
6	Evid/Elab	5	71.4	98.7	0.78	5	76.2	99.6	0.78
6	Org/Purp	5	69.8	98.9	0.78	5	76.2	99.6	0.78
7	Conventions	9	76.1	99.7	0.70	9	75.5	99.8	0.74
7	Evid/Elab	9	75.6	99.7	0.83	9	81.7	99.8	0.84
7	Org/Purp	9	75.6	99.6	0.84	9	81.6	99.9	0.84
8	Conventions	9	77.0	99.1	0.71	9	76.1	99.7	0.74
8	Evid/Elab	9	73.7	99.1	0.82	9	76.9	99.6	0.80
8	Org/Purp	9	75.1	99.7	0.84	9	77.2	99.6	0.80
11	Conventions	10	79.1	99.7	0.75	10	77.1	99.6	0.73
11	Evid/Elab	10	76.5	99.7	0.86	10	75.6	99.9	0.84
11	Org/Purp	10	76.4	99.7	0.86	10	75.8	99.9	0.83

Table 65. Human-Machine Agreement for Mathematics Items on Initial and Secondary Validation Samples, by Grade

	Score		Initial Va	alidation		Secondary Validation			
Grade	Point Range	Number of Items	% Exact	%(Exact+ Adjacent)	QWK	Number of Items	% Exact	%(Exact+ Adjacent)	QWK ^a
3	0-1	10	94.2	100	0.86	10	94.1	100.0	NA
4	0-1	7	91.0	100	0.79	7	92.3	100.0	NA
5	0-1	7	92.6	100	0.81	7	93.5	100.0	NA
6	0-1	8	96.6	100	0.81	8	95.8	100.0	NA
7	0-1	7	96.9	100	0.85	7	96.8	100.0	NA
8	0-1	5	90.2	100	0.75	5	90.5	100.0	NA
11	0-1	16	95.6	100	0.87	16	94.2	100.0	NA
3	0-2	28	90.8	99.3	0.91	28	90.6	99.4	0.89
4	0-2	29	91.0	99.7	0.91	29	91.6	99.7	0.89
5	0-2	38	88.3	99.6	0.88	38	87.9	99.5	0.84
6	0-2	32	88.9	99.6	0.86	32	89.1	99.5	0.84
7	0-2	8	87.0	99.4	0.80	8	88.9	99.9	0.8
8	0-2	16	89.1	99.8	0.89	16	90.3	99.7	0.86
11	0-2	17	89.1	99.4	0.88	17	88.1	99.4	0.87
3	0-3	6	91.1	99.8	0.96	6	92.5	99.9	0.96
4	0-3	4	87.9	99.8	0.94	4	86.8	99.6	0.93
5	0-3	2	90.8	98.4	0.94	2	89.4	98.3	0.90
8	0-3	1	78.2	98.0	0.88	1	86.1	98.4	0.92
11	0-3	5	85.5	99.0	0.89	5	83.7	99.0	0.88

Note. ^a QWK is not presented for 0-1 items due to the binary score scale.

Tables 66 through 68 present the human-machine agreement on the live validation samples for ELA/L SA items, ELA/L essay items, and mathematics SA items, respectively. Recall live training did not involve a secondary validation since 2023-24 operational data were used to build the models.

Table 66. Human-Machine Agreement for ELA/L Short-Answer Items on Live Validation Sample, by Grade

	Live Validation							
Grade	Number of Items	% Exact	%(Exact+ Adjacent)	QWK				
3	1	73.8	99.3	0.66				
4	3	79.7	99.7	0.81				
5	1	70.4	98.0	0.73				
6	15	77.6	99.5	0.73				
7	18	78.5	99.7	0.74				
8	17	76.1	99.6	0.74				
11	24	76.5	99.6	0.77				

Table 67. Human-Machine Agreement for ELA/L Essay Items on Live Validation Sample, by Grade

		Live Validation			
Grade	Trait	Number	%	%(Exact+	QWK
		of Items	Exact	Adjacent)	QWK
3	Conventions	16	70.5	99.6	0.71
3	Evid/Elab	16	73.4	98.8	0.77
3	Org/Purp	16	72.8	99.0	0.77
4	Conventions	19	69.4	99.2	0.73
4	Evid/Elab	19	72.2	98.9	0.78
4	Org/Purp	19	73.0	99.2	0.79
5	Conventions	14	70.8	99.5	0.70
5	Evid/Elab	14	70.1	99.0	0.78
5	Org/Purp	14	70.2	99.1	0.79
6	Conventions	10	73.2	99.4	0.72
6	Evid/Elab	10	73.6	99.3	0.79
6	Org/Purp	10	74.0	99.4	0.79
7	Conventions	12	71.5	99.6	0.72
7	Evid/Elab	12	74.6	99.4	0.80
7	Org/Purp	12	74.8	99.4	0.81
8	Conventions	6	76.7	99.6	0.72
8	Evid/Elab	6	76.9	99.8	0.84
8	Org/Purp	6	74.8	99.8	0.83
11	Conventions	12	75.8	99.5	0.73
11	Evid/Elab	12	76.0	99.7	0.84
11	Org/Purp	12	76.2	99.8	0.84

Table 68. Human-Machine Agreement for Mathematics Items on Live Validation Samples, by Grade

	Score	Live Validation			
Grade	Point Range	Number of Items	% Exact	%(Exact+ Adjacent)	QWK ^a
3	0-1	3	94.4	100.0	NA
4	0-1	1	88.7	100.0	NA
5	0-1	4	95.4	100.0	NA
6	0-1	1	91.4	100.0	NA
7	0-1	1	100	100.0	NA
8	0-1	3	87.8	100.0	NA
3	0-2	1	100	100.0	1.00
5	0-2	14	84.1	99.4	0.82
6	0-2	3	87.3	99.2	0.81
7	0-2	3	90.1	99.1	0.88
8	0-2	3	92.3	100.0	0.92
11	0-2	3	97.6	100.0	0.98
5	0-3	1	88.3	98.7	0.91
8	0-3	1	72.2	97.0	0.89
11	0-3	1	90.2	98.8	0.89

Note. QWK^a is not presented for 0-1 items due to the binary score scale.

6.8.5 Recommendations

The 2023 administrations highlighted the importance of expanding automated monitoring and implementing further interventions to maximize score quality. Building on this, the 2024 administration successfully broadened the additional rater validation stage—originally introduced in 2023 for brief write and research rater qualification—to encompass all ELA/L item types. Furthermore, validity-based measures of scoring accuracy were refined in 2024 to include a comparison of mean score differences between the distributions of benchmark and rater-assigned scores in addition to the previously utilized agreement (QWK). This enhancement provided a more nuanced and sensitive measure of rater quality, ensuring that scoring accuracy is maintained at a high standard.

Despite these improvements, the primary challenge faced during the spring 2024 administration was related to rater productivity, with raters not meeting the expected number of working hours projected from 2023. This issue became particularly evident in April and May, leading to bottlenecks, especially in the scoring of full write and brief write responses, which are time-consuming to train for and score accurately. In response, additional raters were recruited, and pay incentives were offered in key production bottleneck areas. However, some responses still experienced delays in scoring. To address these challenges for the 2025 administration, it is recommended to develop a core pool of full-time raters, establish a minimum work commitment for part-time raters, and collect a measure of rater quality earlier, ideally during qualification. Additionally, surveying raters on their availability and work preferences, as well as enhancing the rater management system, will be crucial steps in improving rater productivity and maintaining the quality and timeliness of scoring.

Furthermore, a review of the scoring outcomes revealed that while the mean QWK values for inter-rater agreement generally met expectations, there were concerns regarding the relatively low minimum QWKs observed for some ELA/L short-answer items, as indicated by the minimum QWK values in Table 52. These low QWK values suggest variability in rater agreement for certain items, which could undermine the overall reliability of the scoring process. To address this issue, it is recommended that additional targeted training and calibration sessions be conducted for raters assigned to items with historically low QWK values. This could include additional focused trainings on interpreting and applying scoring rubrics for those items, the development of supplemental materials, as well as implementing more frequent monitoring and feedback loops during the scoring process.

7. REPORTING AND INTREPRETING SCORES

The Centralized Reporting System (CRS) generates a set of online score reports that includes the information describing student performance for students, parents, educators, and other stakeholders. The online score reports are produced immediately after students complete tests and handscored items are scored. Because the score reports on students' performance are updated every time students complete tests and handscored items are scored, authorized users (e.g., school principals, teachers) can readily access information on students' test performance and use it to improve student learning. In addition to individual student's score reports, the CRS also produces aggregate score reports by class, school, complex, complex area, and state. The timely accessibility of aggregate score reports helps users monitor students' performance in each subject by grade area, evaluate the effectiveness of instructional strategies, and inform the adoption of strategies to improve student learning and teaching during the school year.

This section contains a detailed description of the types of scores reported in the CRS and how to interpret and use these scores.

7.1 CENTRALIZED REPORTING SYSTEM

The CRS is designed to help educators and students answer questions about how well students have performed on the English language arts/literacy (ELA/L) and mathematics assessments. The CRS is an online tool that provides all stakeholders with timely, relevant score reports. The CRS for the Smarter Balanced assessments was designed such that score reports are easy to read and understand for all stakeholders. This is achieved by using plain, non-technical language to facilitate review by parents and the general public. The CRS is also designed to present student performance in a uniform format. For example, similar colors are used for groups of similar elements, such as achievement levels, throughout the design. This design strategy allows readers to compare similar elements and avoid comparing dissimilar elements.

Generally, the CRS provides two categories of online score reports: (1) aggregate score reports, and (2) student score reports. Table 69 summarizes the types of online score reports available at the aggregate level and the individual student level. Detailed information about the online score reports and instructions on how to navigate the online score reporting system can be found in the *Centralized Reporting System User Guide*, located via a help button in the CRS.

Table 69. Types of Online Score Reports by Level of Aggregation

Level of Aggregation	Types of Online Score Reports		
State	 Number of students tested and percentage of proficient students (for overall students and by subgroup) 		
Complex Area Complex	 Average scale score and standard error of average scale score on the overall test and claim (for overall students and by subgroup) 		
School Teacher	 Percentage of students at each achievement level on the overall test (for overall students and by subgroup) 		
Roster	Performance category in each target (for overall students)On-demand student roster report		
	Total scale score and standard error of measurement		
C4-14	 Achievement level for the overall score and claim scores with achievement-level descriptors 		
Student	 Average scale scores and standard errors of average scale scores for individual complex, complex areas, and states 		
	 Writing performance descriptors and scores by dimensions 		

Aggregate score reports at a selected aggregate level are provided for overall students and by subgroup. Users can see student assessment results by any of the subgroups. Table 70 presents the types of subgroups and subgroup categories provided in the CRS.

Table 70. Types of Subgroups

Subgroup	Subgroup Category	
Gender	Male	
	Female	
ELL	Yes	
	No	
Disability	01 - Autism	
	02 - Deaf-Blindness	
	03 - Deafness	
	04 - Developmental Delay (Age 3-5)	
	05 - Developmental Delay (Age 6-8)	
	06 - Emotional Disturbance	
	07 - Hearing Impaired	
	08 - Mental Retardation	
	09 - Multiple Disability	
	10 - Orthopedic Impairment	
	11 - Other Health Impairment	
	12 - Specific Learning Disability	
	13 - Speech/Language Impairment	
	14 - Traumatic Brain Injury	
	15 - Visual Impairment including Blindness	
	16 - Autism Spectrum Disorder	
	17 - Other Health Disability	

Subgroup	Subgroup Category
	18 - Speech or Language Disability
	19 - Intellectual Disability
	20 - Visual Disability Including Blindness
	21 - Hard of Hearing
	22 - Orthopedic Disability
Migrant Status	Yes
	No
Disadvantaged	C, D, E, F, R, 1, 2, 3
Ethnicity	American Indian/Alaskan Native
	Asian/Pacific Islander
	African American
	Hispanic
	Hawai'i Pacific Islander
	White
	Multi-Racial

7.1.1 Dashboard

The CRS provides a state dashboard for authorized state-level users to track student performance for a test across the entire state. The dashboard summarizes students' performance for both ELA/L and mathematics in each grade, including (1) student count, (2) average score and standard error of the average score, (3) percentage and counts of students at each achievement level, and (4) test date last taken.

Exhibit 1 presents a sample state dashboard page.

Q Dashboard Selector > Dashboard Generator > State Dashboard 다. Fines Average Score and Performance Distribution, by Assessment: Hawaii Department of Education, 2023-2024 Filtered By School: All Schools | Test Reasons: All Test Rea Test Groups Student Count Average Score Performance Distribution Date Last Taken Spring 2024 (Smarter Summative) Grade 6 Math 2516 ± 1 🕕 06/12/2024 Spring 2024 (Smarter Grade 6 ELA Spring 2024 (Smarter 2479 ± 1 0 06/03/2024 Grade 4 Math Spring 2024 (Smarter 05/31/2024 Spring 2024 (Smarter Grade 11 ELA Spring 2024 (Smarter Grade 7 Math 2520 ± 1 📵 Spring 2024 (Smarter 05/30/2024 Spring 2024 (Smarter Grade 11 Math 05/30/2024 Spring 2024 (Smarter Grade 3 ELA 2425±1 0 05/30/2024 Spring 2024 (Smarter Grade 8 Math 05/29/2024 Click here to view more tests in this test group (10 of 14 Total Tests)

Exhibit 1. Dashboard: State Level

When authorized users at the complex area, complex, school, and teacher level log in to the CRS, the dashboard page shows the overall test results for all tests that the students have taken grouped by test family (i.e., Smarter Balanced Summative ELA/L). The dashboard summarizes students' performance by test family for both ELA/L and mathematics across all grades, including (1) the grades of the students who have tested, (2) the number of tests taken, (3) the test date last taken, and (4) the percentage and counts of students at each achievement level. State personnel and complex area personnel would select a specific complex to view the aggregate results.

Exhibit 2 presents a sample dashboard page at the complex level.

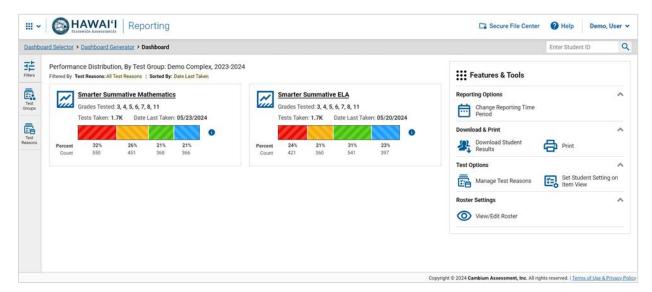


Exhibit 2. Dashboard: Complex Level

When a user clicks on a test family for further exploration, he or she will be taken to a detailed dashboard, where the results will be displayed by test (e.g., grade 3 ELA/L). The detailed dashboard page will appear by test in each grade. The detailed dashboard summarizes students' performance by test in each grade, including (1) the number of students tested, (2) average score and standard error of the means, and (3) percentage and counts of students at each performance level.

Exhibit 3 presents a sample detailed dashboard page for Smarter Balanced summative mathematics at the complex level.

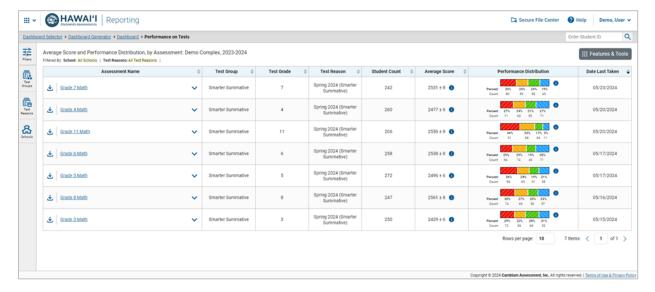


Exhibit 3. Detailed Dashboard: Complex Level

7.1.2 Aggregate Score Reports: Overall Performance

Student performance for each grade in a subject area for a selected aggregate level is presented when users select a specific assessment name. On each aggregate report, the summary report presents the summary results for the selected aggregate unit and the summary results for the state and the aggregate unit both above and below the selected aggregate. For example, if a complex is selected, the summary results of the state and individual schools within the complex are provided as well as the complex summary results so that complex performance can be compared with the other aggregate levels.

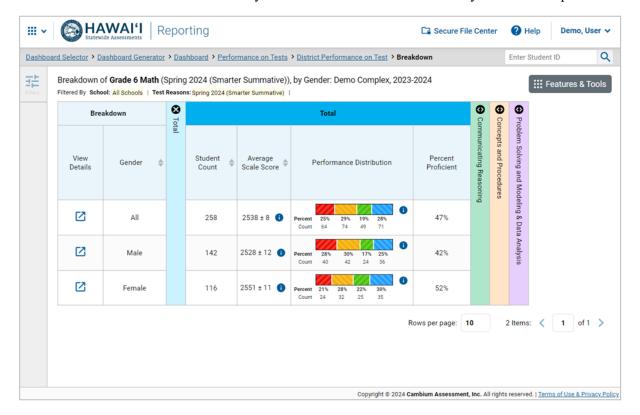
The aggregated summary report provides the summaries on a specific grade in a subject, including (1) the student count, (2) the average scale score and standard error of the average scale score, (3) the percentage and counts of students in each achievement level, and (4) the percentage of proficient students. The summaries are also presented for students overall and by subgroup.

Exhibit 4 presents a sample overall performance summary results page for grade 6 mathematics at the complex level, and Exhibit 5 presents an example summary for grade 6 mathematics by gender.

HAWAI'I | Reporting Secure File Center Help Q <u>Dashboard Selector</u> > <u>Dashboard Generator</u> > <u>Dashboard</u> > <u>Performance on Tests</u> > <u>District Performance on Tests</u> Enter Student ID Average Score and Performance Distribution for Grade 6 Math (Spring 2024 (Smarter Summative)), by School and Reporting ::: Features & Tools Category: Demo Complex, 2023-2024 Filtered By School: All Schools | Test Reasons: Spring 2024 (Smarter Summative) | 8 0 0 School 0 Total Problem Solving and Modeling & Data Analysis Student Average Performance Distribution Scale Score Count Proficient 12393 State 2516 ± 1 1 39% Complex Area 1001 56% Complex 258 2538 ± 8 1 47% 40 Demo School 1 Rows per page: 1 5 Items: < 1 of 5 > Copyright © 2024 Cambium Assessment, Inc. All rights reserved. | Terms of Use & Privacy Policy

Exhibit 4. Overall Performance Summary Results for Grade 6 Mathematics: Complex Level

Exhibit 5. Overall Performance Summary Results for Grade 6 Mathematics by Gender: Complex Level



7.1.3 Aggregate Score Reports: Claim and Target Performance

Detailed summaries on aggregated claim and target results are also available on the same report page when a claim on the right side of the page is selected. For the claim result, both the average scale score and standard error of the average scale score are presented. For the target result, the strength or weakness indicators on each target within a claim are presented. These strength or weakness indicators are presented in two ways. The "Proficient?" measure indicates whether the group's performance on each target is better than (checkmark), less than (x mark), or not different from (half-filled circle) the proficiency standard for the selected test. The "Weak or Strong?" measure presents whether the group's performance on each target is lower than (minus sign), higher than (plus sign), or not different from (equal sign) the group's overall performance. If there is insufficient information in the "Proficient?" measure or "Weak or Strong?" measure, this is indicated with a star sign (*).

Like the overall performance summary results, the summary report presents results for the selected aggregate unit, for the state, and for the aggregate unit both above and below the selected aggregate unit. Also, the summaries on claim and target-level performance can be presented for overall students and by subgroup.

Exhibit 6 presents a sample claim and target-level results page for grade 6 mathematics at the complex level.

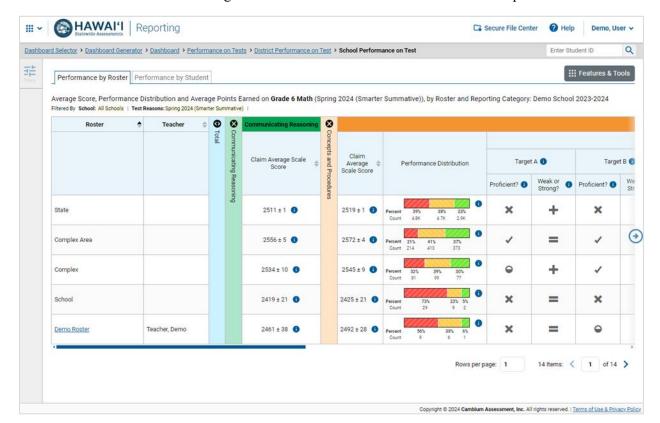


Exhibit 6. Claim and Target Level Results for Grade 6 Mathematics: Complex Level

7.1.4 Roster Performance Report

Class, teacher, and school performance rosters provide users with performance data for a group of students belonging to a system-defined or user-defined class. The report includes (1) the student's overall subject scale scores with standard error of measurement, and (2) the performance level.

Exhibit 7 shows a sample roster performance report page for the grade 6 mathematics summative assessment.

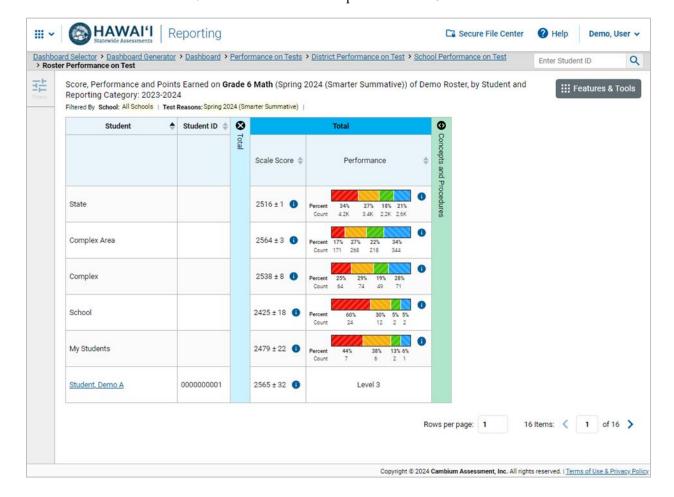


Exhibit 7. Roster Performance Report for Grade 6 Mathematics

7.1.5 Trend Report

The trend (i.e., longitudinal) page provides the trend of student performance for individual level and aggregate level over time. The trend report can be set to plot either average scale scores or percentage of students in each achievement level on the graph for the selected aggregate unit. The trend report is also available at the individual student level. Exhibit 8 presents an example trend report page for ELA/L at the individual student level.

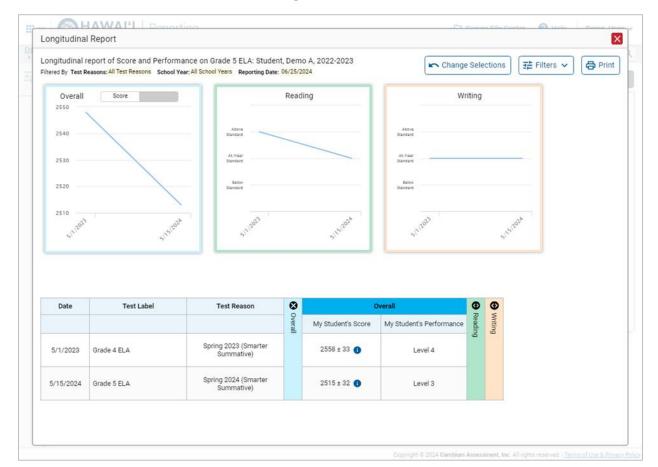


Exhibit 8. Trend Report for ELA/L: Student Level

7.1.6 Individual Student Report

An individual student report (ISR) can be generated and exported as a PDF. The ISR shows the student's overall performance on the test with detailed information on multiple pages. In each subject area, the ISR provides (1) the scale score and SEM; (2) achievement level for the overall test; (3) average scale scores for student's state, complex area, complex, and school; and (4) writing performance descriptors in each dimension (ELA/L only).

On the first page of the ISR, the student's name, scale score with the SEM, and achievement level for ELA are shown at the top of the page. In the middle section, the student's performance is described in detail using a barrel chart. In the barrel chart, the student's scale score is presented with the SEM using a "±" sign. The SEM represents the precision of the scale score, or the range in which the student would likely score if a similar test were administered multiple times. Furthermore, in the barrel chart, achievement-level descriptors with cut scores at each achievement level are provided. These define the content-area knowledge, skills, and processes that test takers at the achievement level are expected to possess.

Average scale scores and standard errors of the average scale scores for the student's state, complex area, complex, and school are displayed at the bottom of the page so the student's achievement can be compared with the above-aggregate levels. It should be noted that the "±" next to the student's scale score is the

standard error of measurement of the scale score, whereas the "±" next to the average scale scores for aggregate levels represents the standard error of the average scale scores.

The second page shows the student's performance on claims (i.e., Claims 1 and 2 for ELA and Claim 1 only for mathematics) which is displayed alongside a description of his or her performance on the claim. At the bottom of the page, the student's performance on the different writing dimensions is displayed alongside a detailed description. The last page provides the trend of the student's performance over time. Student scale scores and achievement levels over time are graphed, showing how the student's scale scores changed over time and whether the student met the standards each year.

Exhibit 9 presents a sample ISR for grade 5 ELA/L.

Exhibit 9. Individual Student Report for Grade 5 ELA/L

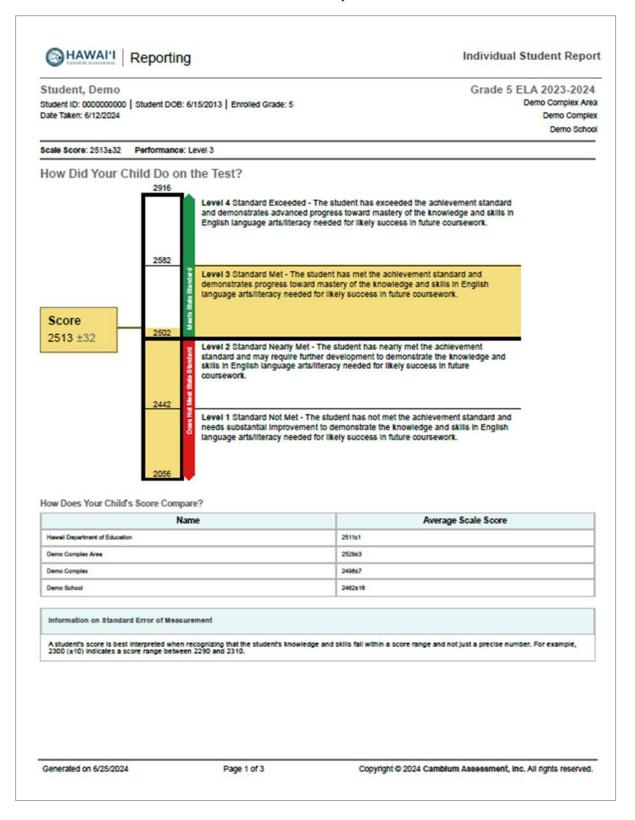


Exhibit 9. Individual Student Report for Grade 5 ELA/L (Continued)

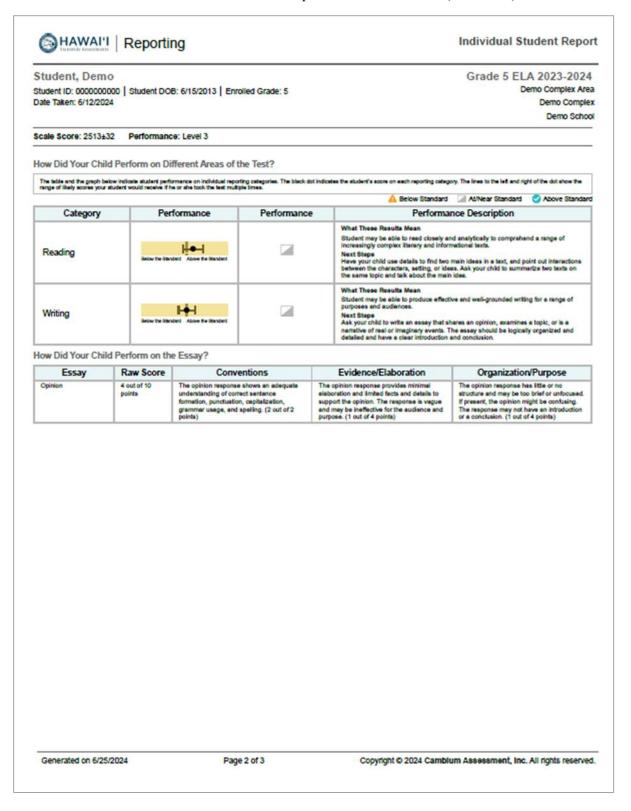
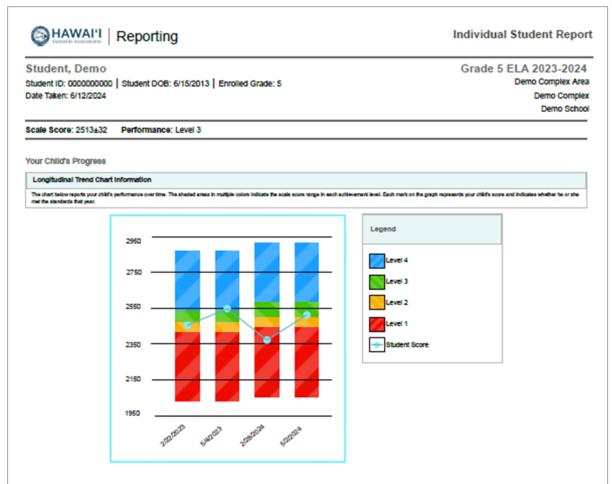


Exhibit 9. Individual Student Report for Grade 5 ELA/L (Continued)



Your Child's Progress

Date	Test Reason	Test Label	Scale Score	Performance
2/22/2023	Opportunity 1	Grade 4 ELA - Interim (ICA)	2455 ± 24	Level 2
5/4/2023	Spring 2023 (Smarter Summative)	Grade 4 ELA	2548 ± 33	Level 4
2/28/2024	Opportunity 1	Grade 5 ELA - Interim (ICA)	2373 ± 31	Level 1
5/2/2024	Spring 2024 (Smarter Summative)	Grade 5 ELA	2513 ± 32	Level 3

Generated on 6/25/2024

Page 3 of 3

Copyright © 2024 Cambium Assessment, Inc. All rights reserved.

7.2 Interpretation of Reported Scores

A student's performance on a test is reported as a scale score and an achievement level for the overall test. Students' scores and achievement levels are also summarized at the aggregate levels. The next section provides a description of how to interpret these scores.

7.2.1 Scale Score

A scale score is used to describe how well a student performed on a test and can be interpreted as an estimate of the student's knowledge and skills measured. The scale score is the transformed score from a theta score, which is estimated based on mathematical models. Low scale scores can be interpreted to mean that the student does not possess sufficient knowledge and skills measured by the test. Conversely, high scale scores can be interpreted to mean that the student has proficient knowledge and skills measured by the test. Scale scores can be used to measure student growth across school years. The interpretation of scale scores is more meaningful when the scale scores are used along with achievement levels and achievement-level descriptors.

7.2.2 Standard Error of Measurement

A scale score (observed score on any test) is an estimate of the true score. If a student takes a similar test multiple times, the resulting scale score will vary across administrations, sometimes being a little higher, a little lower, or the same. The standard error of measurement (SEM) represents the precision of the scale score, or the range in which the student would likely score if a similar test was administered multiple times. When interpreting scale scores, it is recommended to consider the range of scale scores incorporating the SEM of the scale score.

The " \pm " next to the student's scale score provides information about the certainty, or confidence, of the score's interpretation. The boundaries of the score band are one SEM above and below the student's observed scale score, representing a range of score values that is likely to contain the true score. For example, 2680 ± 10 indicates that if a student was tested again, it is likely that the student would receive a score between 2670 and 2690. The SEM can be different for the same scale score, depending on how closely the administered items match the student's ability.

7.2.3 Achievement Level

Achievement levels are proficiency categories on a test that students fall into based on their scale scores. For the Smarter Balanced assessments, scale scores are mapped into four achievement levels (i.e., Level 1, Level 2, Level 3, and Level 4) using three achievement standards (i.e., cut scores). Achievement-level descriptors (ALDs) are a description of content-area knowledge and skills that test takers at each achievement level are expected to possess. Thus, achievement levels can be interpreted based on ALDs. For the achievement level in ELA/L, for instance, ALDs are described for grade 6 Level 3 as: "The student has met the achievement standard and demonstrates progress toward mastery of the knowledge and skills in English language arts/literacy needed for likely success in entry-level credit-bearing college coursework after high school." Generally, students performing at Levels 3 and 4 on Smarter Balanced tests are on track to demonstrate progress toward mastery of the knowledge and skills necessary for college and career readiness.

7.2.4 Performance Category for Claims

Students' performance on each claim is reported in three categories: (1) Below Standard, (2) At/Near Standard, and (3) Above Standard. Unlike the achievement level for the overall test, student performance on each claim is evaluated with respect to the "Meets Standard" achievement standard. For students performing at "Below Standard" or "Above Standard," this can be interpreted to mean that their performance is clearly below or above the "Meets Standard" cut score for a specific claim. For students performing at "At/Near Standard," this can be interpreted to mean that their performance does not provide enough information to tell whether they reached the "Meets Standard" mark for the specific claim.

7.2.5 Performance Category for Targets

Teachers and educators sometimes need more detailed reports on student performance for instructional purposes. The target report provides information on student performance about relative strength and weakness scores for each target within a claim. The strengths and weaknesses reports are generated for aggregate units of classroom, school, and complex and provide information about how a group of students in a class, school, or complex performed on each target, either relative to the proficiency standard (i.e., "Proficient?" target measure) or relative to their overall performance on the test (i.e., "Weak or Strong?" target measure). Target-level reports are produced for the aggregate units only, not for individual students, because each student is administered too few items in a target to produce a reliable score for each target.

For the "Proficient?" target measure, students' observed performance on items within the reporting element is compared to the expected performance on those items of someone who has an ability equal to the proficiency cut score (i.e., the Achievement Level 3 cut). At the aggregate level, when the observed performance within a target is greater than the proficiency cut, the reporting unit shows relative strength in that target compared to the proficiency standard. Conversely, when observed performance within a target is below the proficiency cut, the reporting unit shows relative weakness in that target.

For the "Weak or Strong?" target measure, students' observed performance on items within the reporting element is compared with the expected performance based on the overall ability estimate. At the aggregate level, when the observed performance within a target is greater than the expected performance, the reporting unit (e.g., roster, teacher, school, complex) shows relative strength in that target. Conversely, when the observed performance within a target is below the level expected based on overall achievement, the reporting unit shows relative weakness in that target.

Although performance categories for targets provide some evidence to help address students' strengths and weaknesses, they should not be over-interpreted because student performance on some targets may be based on relatively few items, especially for a small group.

7.2.6 Aggregated Scale Score

Students' scale scores are aggregated at roster, teacher, school, complex, complex area, and state levels to represent how a group of students performs on a test. When students' scale scores are aggregated, the average scale scores can be interpreted as an estimate of the knowledge and skills that a group of students possesses. Given that student scale scores are estimates, the average scale scores are also estimates and are subject to measures of uncertainty. In addition to the average scale scores, the percentage of students in each achievement level for overall are reported at the aggregate level to represent how well a group of students performs.

7.3 APPROPRIATE USES OF TEST RESULTS

Assessment results can provide information about individual students' achievements on the test. Overall, assessment results show what students know and are able to do in certain subject areas and provide further information on whether students are on track to demonstrate the knowledge and skills necessary for college and career readiness. Additionally, assessment results can be used to identify students' relative strengths and weaknesses in certain content areas. For example, performance categories for targets can be used to identify a group's relative strengths and weaknesses among targets within a claim.

Assessment results on student achievement on the test can be used to help teachers or schools make decisions on how best to support students' learning. Aggregate score reports at the teacher and school level provide information regarding the strengths and weaknesses of their students and can be used to improve teaching and student learning. For example, a group of students may perform very well overall on the test but potentially not perform as well in several targets compared to their overall performance. In this case, teachers and schools would be able to identify the strengths and weaknesses of their students through the group performance by claim and target. They could then promote instruction in the specific claim or target areas in which their students perform relatively lower. Further, by narrowing the student performance results by subgroup, teachers and schools can determine which strategies may be best suited to improving student learning, particularly for students from disadvantaged subgroups. For example, teachers can examine student assessment results by limited English proficiency (LEP) status and may observe that LEP students need help particularly in a certain specific area, such as reading literary responses and analysis. Teachers can then provide additional focused instruction for these students to enhance their achievement in any specific target or claim in which they are struggling.

In addition, assessment results can be used to compare performance among different students and among different groups. Teachers can evaluate how their students perform compared with other students in their school, complex, and complex area for overall scores and by claim. Although all students are administered different sets of items in each computer-adaptive test, scale scores are comparable across students. Furthermore, scale scores can be used to measure the growth of individual students over time when data are available. In the Smarter Balanced assessments, the scale scores across grades are on the same scale because the scores are vertically linked across grades. Therefore, scale scores from one grade can be compared with the next grade, i.e., measuring the growth.

While assessment results provide valuable information to understand students' performance, these scores and reports should be used with caution. It is important to note that scale scores reported are estimates of true scores and hence do not represent the precise measure for student performance. A student's scale score is associated with measurement error and thus users need to consider measurement error when using student scores to make decisions about student achievement. Moreover, although student scores may be used to help make important decisions about students' placement and retention, or teachers' instructional planning and implementation, the assessment results should not be used as the only source of information. Given that assessment results measured by a test provide limited information, other sources on student achievement such as classroom assessment and teacher evaluation should be considered when making decisions on student learning. Finally, when student performance is compared across groups, users need to consider the group size. The smaller the group size, the larger the measurement error related to these aggregate data, thus requiring interpretation with more caution.

8. QUALITY CONTROL PROCEDURES

Quality assurance (QA) procedures are enforced throughout all stages of the Smarter Balanced assessment development, administration, scoring, and reporting of results. CAI uses a series of quality control (QC) steps to ensure the error-free production of score reports in both online and paper-pencil formats. The quality of the information produced in the Test Delivery System (TDS) is tested thoroughly before, during, and after the testing window opens.

8.1 ADAPTIVE TEST CONFIGURATION

For the computer-adaptive testing (CAT) component, a test configuration file is the key resource that contains all specifications for the item-selection algorithm and the scoring algorithm, such as the test blueprint, cut scores, item information (i.e., answer keys, item attributes, item parameters, and passage information), and slopes and intercepts for theta-to-scale score transformation. The accuracy of the information in the configuration file is independently checked and confirmed before the testing window opens.

CAI uses simulated test administrations along with the test configuration file to configure the adaptive algorithm in order to optimize item selection to meet blueprint specifications while targeting test information to student ability. First, the simulator generates a sample of students with an ability distribution that matches that of the population in the previous year's data. The ability of each simulated student is used to generate a sequence of item-response scores while matching the blueprint and minimizing measurement error. These simulations provide a rigorous test of the adaptive algorithm. The results of these simulations are used to configure and evaluate the adequacy of the item-selection algorithm used to administer the Smarter Balanced summative assessments.

After the adaptive testing simulations, another set of simulations for the combined tests (CAT and performance task [PT] components) are performed for scoring engine verification. The simulated data are generated such that verification of the scoring engine is based on a wide range of student response patterns. CAI rigorously checks whether the scoring rules specified in scoring specifications were applied accurately. The scores in the simulated data file are checked independently.

8.1.1 Platform Review

CAI's TDS supports a variety of item layouts. Each item goes through an extensive platform review on different operating systems such as Windows, Linux, and iOS to ensure that the item looks consistent in all of them. Some of the layouts have the stimulus and item response options/response area displayed side by side. In each of these layouts, both stimulus and response options have independent scroll bars.

Platform review is a process during which each item is checked to ensure that it is displayed appropriately on each tested platform. A platform is a combination of a hardware device and an operating system. In recent years, the number of platforms has proliferated, and platform review now takes place on various platforms that are significantly different from one another.

Platform review is conducted by a team. The team leader projects the item as it was web approved in the Item Tracking System (ITS), and team members, each using a different platform, view the same item to ensure that it renders as expected.

8.1.2 User Acceptance Testing and Final Review

Before deployment, the testing system and content are deployed to a staging server, where they are subject to user acceptance testing (UAT). UAT of the TDS serves as both a software evaluation and a content approval role. The UAT period provides HIDOE with an opportunity to interact with the exact test that the students will use.

8.2 QUALITY ASSURANCE IN DOCUMENT PROCESSING

The Smarter Balanced assessments are administered primarily online; however, a few students take paper-pencil assessments. When test documents are scanned, a QC sample of documents consisting of 10 test cases per document type (normally between 500 and 600 documents) is created so that all possible responses and all demographic grids are verified, including various typical errors that required editing via Measurement Incorporated's (MI) Data Inspection, Correction, and Entry (DICE) application. This structured testing method provides exact test parameters and a methodical way of determining that the output received from the scanner(s) is correct. MI staff carefully compare the documents and the data file created from them to further ensure that the results from the scanner, the editing process (validation and data correction), and the transfer to the CAI database are correct.

8.3 QUALITY ASSURANCE IN DATA PREPARATION

CAI's TDS has a real-time quality-monitoring component built in. After a test is administered to a student, the TDS passes the resulting data to CAI's QA system. The QA system conducts a series of data integrity checks, ensuring, for example, that the record for each test contains information for each item, keys for multiple-choice items, score points for each item, and the total number of field-test items and operational items. It also ensures that the test record contains no data from items that have been invalidated.

Data pass directly from the Quality Monitor System (QM) to the Database of Record (DOR), which serves as the repository for all test information from which all test information for reporting is pulled. The Data Extract Generator is the tool that is used to pull data from the DOR for delivery to HIDOE. CAI staff ensure that data in the extract files match the DOR before it is delivered.

8.4 QUALITY ASSURANCE IN ONLINE TEST DELIVERY SYSTEM

To monitor the performance of the TDS during the test administration window, CAI statisticians examine the delivery demands, including the number of tests to be delivered, the length of the testing window, and the historic, state-specific behaviors, to model the likely peak loads. Using data from the load tests, these calculations indicate the number of each type of server necessary to provide continuous, responsive service, and CAI contracts for service in excess of this amount. Once deployed, the servers are monitored at the hardware, operating system, and software platform levels with monitoring software that alerts CAI's engineers at the first signs that trouble may arise. The applications log not only errors and exceptions, but also latency (timing) information for crucial database calls. This information enables CAI to know instantly whether the system is performing as designed or if it is starting to slow down or experience a problem. In addition, latency data, such as data about how long it takes to load, view, or respond to an item, are captured for each assessed student. All this information is logged, enabling CAI to automatically identify schools or complex areas experiencing unusual slowdowns, often before they even notice.

A series of quality assurance reports, such as blueprint match rate, item exposure rate, and item statistics, can also be generated at any time during the online assessment window for the early detection of any

unexpected issues. Any deviations from the expected outcome are flagged, investigated, and resolved. In addition to these statistics, a cheating analysis report is produced to flag any unlikely patterns of behavior in a testing session, as discussed in Section 2.8, Data Forensics Program.

For example, an item statistics analysis report allows psychometricians to ensure that items are performing as intended and serves as an empirical key check throughout the operational testing window. The item statistics analysis report is used to monitor the performance of test items throughout the testing window and serves as a key check for the early detection of potential problems with item scoring, including the incorrect designation of a keyed response or other scoring errors and potential breaches of test security that may be indicated by changes in the difficulty of test items. This report generates classical item analysis indicators including item *p*-value and item discrimination index and item response theory item-fit statistics. The report is configurable and can be produced so that only items with statistics falling outside of a specified range are flagged for reporting or to generate reports based on all items in the pool.

For the CAT component, other reports, such as blueprint match and item exposure reports, allow psychometricians to verify that test administrations conform to the simulation results. The QA reports can be generated on any desired schedule. Item analysis and blueprint match reports are evaluated frequently at the opening of the testing window to ensure that test administrations conform to the blueprint and that items are performing as anticipated.

Table 71 presents an overview of the QA reports.

QA Reports	Purpose	Rationale	
Item Statistics	To confirm whether items work as expected	Early detection of errors (key errors for selected-response items and scoring errors for constructed-response, performance, or technology-enhanced items)	
Blueprint Match Rates	To monitor unexpectedly low blueprint match rates	Early detection of unexpected blueprint match issue	
Item Exposure Rates	To monitor unlikely high exposure rates of items or passages or unusually low item pool usage (highly unused items/passages)	Early detection of any oversight in the blueprint specification	
Cheating Analysis	To monitor testing irregularities	Early detection of testing irregularities	

Table 71. Overview of Quality Assurance Reports

8.4.1 Score Report Quality Check

Two types of score reports were produced in the Smarter Balanced summative assessments: (1) online reports, and (2) printed reports (family reports only). In Hawaii, printed reports are not generated.

8.4.1.1 Online Report Quality Assurance

The system automatically assigns scores for the online assessments in real time. Every test undergoes a series of validation checks. Once the QA system signs off, data are passed to the DOR, which serves as the central location for all student scores and responses, ensuring that there is only one place where the official record is stored. Only after scores have passed the QA checks and are uploaded to the DOR are they passed to the Centralized Reporting System (CRS), which is responsible for presenting individual-level results and calculating and presenting aggregate results. Absolutely no score is reported in the CRS until it passes all the QA system's validation checks. All of these processes take milliseconds to complete,

with CAI receiving handscores and passing them through QA validation checks in less than one second and making the composite score available in the CRS immediately.

8.4.1.2 Paper Report Quality Assurance

Statistical Programming

The family reports contain custom programming and require rigorous QA processes to ensure their accuracy. All custom programming is guided by the detailed and precise specifications outlined in CAI's reporting specifications document. Analytic rules are programmed upon approval of the specifications, and each program is extensively tested on test decks and real data from other programs. The final programs are reviewed by two senior statisticians and one senior programmer to ensure that they implemented agreed-on procedures. Custom programming is implemented independently by two statistical programming teams working from the specifications. The scripts are released for production only when the output from both teams matches precisely.

Much of the statistical processing is repeated, and CAI has implemented a structured software development process to ensure that the repeated tasks are implemented correctly and identically each time. Small programs (called *macros*) are written to take specified data as input and produce data sets containing derived variables as output. Approximately 30 such macros reside in CAI's library for score reports. Each macro is extensively tested and stored in a central development server. Once a macro is tested and stored, changes to the macro must be approved by the director of score reporting, the director of psychometrics, and the project directors for affected projects.

Each change is followed by a complete retesting with the entire collection of scenarios on which the macro was originally tested. The main statistical program is mostly made up of calls to various macros, including macros that read in and verify the data and conversion tables and the macros that do the many complex calculations. This program is developed and tested using artificial data generated to test both typical and extreme cases. Additionally, the program goes through a rigorous code review by a senior statistician.

Display Programming

The paper report development process uses graphical programming, which takes place in a Xerox-developed programming language called Variable Data Intelligent PostScript Printware (VIPP) and allows virtually infinite control of the visual appearance of the reports. After our designers create backgrounds, CAI's VIPP programmers write code that indicates where to place all variable information (data, graphics, and text) on the reports. The VIPP code is tested using both artificial and real data. CAI's data generation utilities can read the output layout specifications and generate artificial data for direct input into the VIPP programs. This allows testing of these programs to begin before the statistical programming is complete. In later stages, artificial data are generated according to the input layout and are run through the psychometric process and the score reporting statistical programs, and the output is formatted as VIPP input. This process enables CAI to test the entire system.

Programmed output goes through multiple stages of review and revision by graphics editors and the CAI score reporting team to ensure that design elements are accurately reproduced and data are correctly displayed. Once CAI receives the final data and VIPP programs, the CAI score reporting team reviews proofs that contain actual data based on CAI's standard quality assurance documentation. Several CAI staff members review a large sample of the reports to ensure that all data are correctly placed on reports. This rigorous review is conducted over several days and takes place in a secure location in a CAI building. All reports containing actual data are stored in a locked storage area. Before the reports are printed, CAI

provides a live data file and individual student reports with sample complex areas for HIDOE staff review. CAI will work closely with the Hawai'i to resolve questions and correct any problems. The reports will not be delivered unless the Department approves the sample reports and data file.

REFERENCES

- American Educational Research Association (AERA), American Psychological Association (APA), & National Council on Measurement in Education (NCME). (2014). *Standards for educational and psychological testing*. Washington, DC.
- Billingsley, P. (1995). Probability and Measure (3rd ed.). New York, NY: John Wiley & Sons.
- Drasgow, F., Levine, M. V., & Williams, E. A. (1985). Appropriateness measurement with polychotomous item response models and standardized indices. *British Journal of Mathematical and Statistical Psychology*, 38(1), 67–86.
- Guo, F. (2006). Expected classification accuracy using the latent distribution. *Practical Assessment, Research & Evaluation, 11*(6).
- Huynh, H. (1976). On the reliability of decisions in domain-referenced testing. *Journal of Educational Measurement*, 13(4), 253–264.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32(2), 179–197.
- Livingston, S. A., & Wingersky, M. S. (1979). Assessing the reliability of tests used to make pass/fail decisions. *Journal of Educational Measurement*, 16(4), 247–260.
- Raczynski, K. R., Cohen, A. S., Engelhard, G., Jr., & Lu, Z. (2015). Comparing the effectiveness of self-paced and collaborative frame-of-reference training on rater accuracy in a large-scale writing assessment. *Journal of Educational Measurement*, 52(3), 301–318.
- Snijders, T. A. B. (2001). Asymptotic null distribution of person fit statistics with estimated person parameter. *Psychometrika*, 66(3), 331–342.
- Sotaridona, L. S., Pornel, J. B., & Vallejo, A. (2003). Some applications of item response theory to testing. *The Philippine Statistician*, *52*(1–4), 81–92.
- Subkoviak, M. J. (1976). Estimating reliability from a single administration of a criterion-referenced test. *Journal of Educational Measurement*, 13(4), 265–276.
- U.S. Department of Education. (2015). Peer Review of State Assessment Systems: Non-Regulatory Guidance for States. Washington, D.C. Retrieved from https://www2.ed.gov/policy/elsec/guid/assessguid15.pdf
- Williamson, D. M., Xi, X., & Breyer, F. J. (2012). A framework for evaluation and use of automated scoring. *Educational Measurement: Issues and Practice*, 31(1), 2–13.