

South Dakota Science Assessment

2022–2023

Volume 1: Annual Technical Report



south dakota
DEPARTMENT OF EDUCATION
Learning. Leadership. Service.

TABLE OF CONTENTS

1. INTRODUCTION 1

1.1 Background and Historical Context of Tests 1

1.2 Purpose and Intended Uses of South Dakota Science Assessment 2

1.3 Participants in the Development and Analysis of the South Dakota Science Assessment 3

 1.3.1 South Dakota Department of Education 3

 1.3.2 South Dakota Educators 3

 1.3.3 Technical Advisory Committee 4

 1.3.4 Cambium Assessment, Inc. 4

 1.3.5 Caveon Test Security 4

1.4 Available Test Formats and Special Versions 4

1.5 Student Participation 4

2. OPERATIONAL PRACTICES AND PROCEDURES 5

2.1 Test Administration 5

2.2 Test Administrators 6

2.3 Testing Environment 6

2.4 Simulations 6

2.5 Universal Features, Designated Supports, and Accommodations 7

3. ITEM BANK AND TEST DESIGN 9

3.1 Shared Science Assessment Item Bank 9

3.2 Field Testing 11

 3.2.1 2023 Field Tests 11

3.3 Test Design 21

4. FIELD-TEST CLASSICAL ANALYSIS 22

4.1 Item Discrimination 23

4.2 Item Difficulty 23

4.3 Response Time 23

4.4 Differential Item Functioning 24

4.5 Classical Analysis Results 26

5. ITEM CALIBRATION 30

5.1 Model Description 30

 5.1.1 Latent Structure 30

 5.1.2 Item Response Function 32

 5.1.3 Multigroup Model 32

5.2 Estimation 33

 5.2.1 Estimation 33

 5.2.2 Overview of the Operational Item Bank 34

6.	SCORING	36
6.1	Marginal Maximum Likelihood Function	36
6.2	Derivative	37
6.3	Extreme Case Handling	39
6.4	Standard Error of Measurement	39
6.5	Scoring Incomplete Tests	39
6.6	Student-Level Scale Score	40
6.7	Rules for Calculating Achievement Levels	41
	6.7.1 <i>Strengths and Weaknesses for Disciplines Relative to Proficiency Cut Score</i> ..	42
6.8	Disciplinary Core Idea-Level Reporting	42
	6.8.1 <i>Relative to Overall Achievement</i>	42
	6.8.2 <i>Relative to Proficiency Cut Score</i>	43
7.	QUALITY CONTROL PROCEDURES	44
7.1	Quality Assurance Reports	44
	7.1.1 <i>Item Analysis</i>	44
	7.1.2 <i>Blueprint Match</i>	44
	7.1.3 <i>Item Exposure Rates</i>	45
	7.1.4 <i>Cheating Detection Analysis</i>	45
7.2	Scoring Quality Check	45
8.	REFERENCES	46

LIST OF TABLES

Table 1. Required Uses and Citations for the SDSA	3
Table 2. Number of Students Participating in the SDSA in Spring 2023.....	5
Table 3. Distribution of Demographic Characteristics of Student Population	5
Table 4. SDSA Testing Windows.....	6
Table 5. Number of Testing Sessions with Allowed Designated Supports	8
Table 6. Number of Testing Sessions with Allowed Accommodations	9
Table 7. Number of Field-Test Items Administered, Spring 2023	12
Table 8. Number of Common Elementary School Field-Test Items Administered and Calibrated, Spring 2023.....	14
Table 9. Number of Common Middle School Field-Test Items Administered and Calibrated, Spring 2023.....	16
Table 10. Number of Common High School Field-Test Items Administered and Calibrated, Spring 2023.....	18
Table 11. Number of Field-Test Item Administration, Rubric Validation, and Item Data Review, Spring 2023.....	20
Table 12. Summary of Shared Science Assessment Item Bank, Spring 2023.....	21
Table 13. Thresholds for Flagging in Classical Item Analysis.....	22
Table 14. DIF Classification Rules.....	26
Table 15. Distribution of p-Values for Field-Test Items in Spring 2023.....	27
Table 16. Distribution of Item Biserial Correlations for Field-Test Items in Spring 2023	27
Table 17. Summary of Response Times for Field-Test Items Administered in Spring 2023.....	28
Table 18. Differential Item Functioning Classifications for Field-Test Items Administered in Spring 2023.....	29
Table 19. Groups Per Grade Band for the Spring 2023 Calibration of Field-Test Items	34
Table 20. Science Reporting Scale Linear Transformation Constants, Theta, and Corresponding Scaled-Score Limits for Extreme Ability Estimates (for 2021 θ Scale).....	41
Table 21. Achievement-Level Cut Scores	42

LIST OF FIGURES

Figure 1. Directed Graph of the Science IRT Model.....	32
Figure 2. SDSA Assertion Difficulty and Student Proficiency Distributions, Grade 5.....	35
Figure 3. SDSA Assertion Difficulty and Student Proficiency Distributions, Grade 8.....	35
Figure 4. SDSA Assertion Difficulty and Student Proficiency Distributions, Grade 11.....	36

LIST OF APPENDICES

Appendix 1-A. Caveon Test Security Overview	
Appendix 1-B. Shared Science Assessment Item Bank: Field Testing	
Appendix 1-C. Calibration of the Shared Science Assessment Item Bank	
Appendix 1-D. Distribution of Scale Scores and Achievement Levels	
Appendix 1-E. Distribution of Scale Scores by Science Discipline	
Appendix 1-F. Distribution of Scale Scores and Achievement Levels by Subgroup	

1. INTRODUCTION

The South Dakota Science Assessment (SDSA) is an assessment for grades 5, 8, and 11. The *South Dakota Science Assessment 2022–2023 Technical Report* is provided to document and make transparent all methods used in item development, test construction, psychometrics, standard setting, test administration, and score reporting, including summaries of student results and evidence and support for the intended uses and interpretations of the test scores. The technical reports are reported as six separate, self-contained volumes as described in the following list:

- 1) **Annual Technical Report.** This volume is updated each year and provides a global overview of the tests administered to students annually.
- 2) **Test Development.** This volume summarizes the procedures used to construct test forms and provides summaries of the item bank and development process.
- 3) **Confirming Achievement Standards.** This addendum is to document the background, processes, methods, and results of the SDSA standards confirmation workshop.
- 4) **Evidence of Reliability and Validity.** This volume provides technical summaries of the test quality and special studies conducted to support the intended uses and interpretations of the test scores.
- 5) **Test Administration.** This volume describes the security protocols, accessibility features (including accommodations), methods used, and system characteristics developed to administer tests.
- 6) **Score Interpretation Guide.** This volume describes the score types reported and details the appropriate inferences that can be drawn from each score reported.

The South Dakota Department of Education (SDDOE) communicates the quality of the SDSA by making these technical reports accessible to the public on the state’s website.

1.1 BACKGROUND AND HISTORICAL CONTEXT OF TESTS

South Dakota adopted three-dimensional science standards as the new South Dakota Science Standards in May of 2015, based on *A Framework for K–12 Science Education* (National Research Council, 2012). The SDDOE and its assessment vendor, CAI, developed and administered a new online assessment to measure these new standards. The SDSA was administered operationally for the first time in 2020–2021 and measures the science knowledge and skills of South Dakota students in grades 5, 8, and 11. The SDSA continued its administration to these grades in 2022–2023.

The SDDOE provides an overview of the SDSA at: <https://doe.sd.gov/Assessment/science.aspx>. Information about the South Dakota Science Standards is available at: <https://doe.sd.gov/contentstandards/>.

1.2 PURPOSE AND INTENDED USES OF SOUTH DAKOTA SCIENCE ASSESSMENT

The SDSA is a criterion-referenced test established using principles of evidence-centered design to yield overall and discipline-level test scores at the student level and other levels of aggregation that reflect student achievement of the SDSA. The three-dimensional science standards (i.e., the South Dakota Science Standards) establish a set of knowledge and skills that all students need to be prepared for a wide range of high-quality post-secondary opportunities, including higher education and entering the workplace. The three-dimensional South Dakota Science Standards reflect the latest research and advances in modern science and differ from previous science standards in multiple ways.

First, rather than describing general knowledge and skills that students should know and be able to do, they describe specific performances that demonstrate what students know and can do. The South Dakota Science Standards refers to such performed knowledge and skills as *performance expectations (PEs)*. Second, while unidimensionality is a typical goal of standards (and the items that measure them), the South Dakota Science Standards are intentionally multi-dimensional. Each PE incorporates all three dimensions from *A Framework for K–12 Science Education* (National Research Council, 2012)—a science or engineering practice, a disciplinary core idea, and a crosscutting concept. Another unique feature of the South Dakota Science Standards is the assumption that students should learn all science disciplines rather than a select few, as is traditionally done in many high schools, where students may elect, for example, to take biology and chemistry but not physics or astronomy.

The SDSA supports instruction and student learning by providing valuable feedback to educators and parents, which can be used to form instructional strategies to remediate or enrich instruction. An array of reporting metrics is provided to evaluate achievement at the student and aggregate levels and to monitor improvement at the student and group levels over time.

The SDSA tests draw all items from the Independent College and Career Readiness (ICCR) item pool, which is part of the Shared Science Assessment Item Bank that consists of items owned by several other states and one U.S. territory. Each of these states and the territory has signed a Memorandum of Understanding (MOU) to share content, leadership, and new ideas and methods (see Volume 2, Test Development, for more information). In 2023, the full members of the MOU were Connecticut, Hawaii, Idaho, Montana, Oregon, Rhode Island, Utah, West Virginia, and Wyoming. New Hampshire, North Dakota, South Dakota, and U.S. Virgin Islands observed and participated in some activities. CAI had a supporting and coordinating role for the MOU collaborative efforts. CAI also worked with the SDDOE to ensure that the items in the test forms which were constructed for all grades within the state uniquely measured the three-dimensional South Dakota Science Standards.

Table 1 outlines the required uses and citations for the SDSA based on the South Dakota Codified Law (SDCL), the Administrative Rule of South Dakota (Article 24:55), and the federal *Every Student Succeeds Act* (ESSA) plan. The SDSA fulfills all the requirements described in Table 1.

Table 1. Required Uses and Citations for the SDSA

Required Use	Required Use Citation
Required Use	Required Use Citation
Indicator of academic achievement and progress	ESSA section 1111(b)(2)(B)(ii) SDCL §13-3-55 SD Administrative Rule 24:55:01:04
Test administration frequency and grade levels	ESSA section 1111(b)(2)(B)(v)(II) SDCL §13-3-55
Disaggregation of test scores	ESSA section 1111(b)(2)(B)(v) ESSA section 1111(b)(3)(C)(xiii)
Publication of test scores	ESSA section 1111(b)(3)(C)(xii) SD Administrative Rule 24:55
Requirement of the alignment of test to academic content standards	SDCL §13-3-55

1.3 PARTICIPANTS IN THE DEVELOPMENT AND ANALYSIS OF THE SOUTH DAKOTA SCIENCE ASSESSMENT

The SDDOE manages the SDSA with the assistance of several participants, including South Dakota educators, a Technical Advisory Committee (TAC), and vendors. The SDDOE fulfills the diverse requirements of implementing South Dakota’s statewide assessments while meeting or exceeding the guidelines established 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). To comply with the *Standards*, scale development, scoring, linking, and evaluation of differential item functioning are addressed in the current volume; item development, test design, and test blueprints are documented in Volume 2, Test Development; development of cut scores is summarized in Volume 3, Setting Performance Standards; evidence for validity and reliability/precision was collected and is reported in Volume 4, Evidence of Reliability and Validity; information on testing windows, test options, accommodations, training of test coordinators and administrators, and test security are provided in Volume 5, Test Administration; supporting documentation for tests, score uses and interpretation are included in Volume 6, Score Interpretation.

1.3.1 South Dakota Department of Education

The Office of Assessment manages test development, administration, scoring, and results reporting for the SDSA, including coordinating with other SDDOE offices, South Dakota public schools, and vendors.

1.3.2 South Dakota Educators

South Dakota educators participated in most aspects of the conceptualization and development of the SDSA. Science panels established by SDDOE assisted in the various activities that surround the academic content standards, such as developing the academic standards, clarifying how the standards are assessed, designing tests, and reviewing test questions and passages.

1.3.3 Technical Advisory Committee

The SDDOE convenes an advisory committee panel twice each year to discuss psychometrics, test development, and administrative and policy issues relevant to the current and future South Dakota assessments. This committee is comprised of highly experienced practitioners from several South Dakota school districts.

1.3.4 Cambium Assessment, Inc.

Cambium Assessment, Inc. (CAI; formerly the American Institutes for Research [AIR]) is the vendor that was selected through the state-mandated competitive procurement process. CAI is responsible for developing test content, building test forms, conducting psychometric analyses, administering and scoring test forms, and reporting test results for the SDSA described in this report. Additionally, CAI is responsible for developing and maintaining the ICCR item bank.

1.3.5 Caveon Test Security

Caveon Test Security monitored web pages and social media during the spring 2023 test administration to ensure that secure testing materials, such as items and prompts, were not leaked. Details of Caveon Test Security are described in Appendix 1-A, Caveon Test Security Overview.

1.4 AVAILABLE TEST FORMATS AND SPECIAL VERSIONS

The SDSA is administered online using an adaptive test design (see details of the adaptive test design in Section 3.3, Test Design. Science items are centered on a scientific phenomenon. They can consist of shorter (stand-alone) items or items with several parts (item clusters), requiring the student to interact with them in various ways.

Students unable to participate in the online test administration have the option to use print-on-demand—a feature that provides the same items administered to students online in a paper format. Spanish versions of the SDSA (developed to meet the same content standards as the English versions) are available for all tested grades. Students participating in the computer-based SDSA can use standard online testing features in the Test Delivery System (TDS), including a selection of font colors and sizes and the ability to zoom in and out or highlight text. In addition to the resources available to all students, options are available to accommodate students with an Individualized Education Program (IEP) or Section 504 Plan. These include braille, closed captioning, and large print. Students with disabilities have the option to take the SDSA with or without accommodations or to take an alternate assessment. For additional information about testing features and accommodations, refer to Volume 5, Test Administration, of this technical report.

1.5 STUDENT PARTICIPATION

The SDSA is administered in the spring. Table 2 shows the number of students who were tested (Number Tested) and the number of students whose scores were included for the analyses in this technical report (Number Reported). The number of students reported excludes students who tested

but did not have a valid score (e.g., the student opened the test and viewed the first item but abandoned the test without responding to any item). Table 3 shows the demographic characteristics of the student population, by counts and in percentages, in the spring administration of the 2022–2023 SDSA. The subgroups reported here are gender, ethnicity, students with limited English proficiency (LEP), and students in special education programs. In 2022–2023, all students in South Dakota were provided with free or reduced-priced lunch, and thus, no additional reporting was included for economically disadvantaged groups.

Table 2. Number of Students Participating in the SDSA in Spring 2023

Grade	Number Tested	Number Reported
5	10,605	10,605
8	10,748	10,744
11	9,475	9,474

Table 3. Distribution of Demographic Characteristics of Student Population

Group	Grade 5		Grade 8		Grade 11	
	N	%	N	%	N	%
All Students	10,605	100.00	10,744	100.00	9,474	100.00
Female	5,202	49.05	5,281	49.15	4,698	49.59
Male	5,403	50.95	5,463	50.85	4,776	50.41
African American	338	3.19	329	3.06	304	3.21
American Indian/Native Alaskan	1,230	11.60	1,145	10.66	746	7.87
Asian	200	1.89	142	1.32	163	1.72
Hispanic	856	8.07	825	7.68	627	6.62
Multi-Racial	641	6.04	559	5.20	437	4.61
Pacific Islander	15	0.14	16	0.15	15	0.16
White	7,325	69.07	7,728	71.93	7,182	75.81
Limited English Proficiency	502	4.73	400	3.72	247	2.61
Special-Education	1,747	16.47	1,259	11.72	869	9.17
Economically Disadvantaged	3,877	36.56	3,533	32.88	2,293	24.20

2. OPERATIONAL PRACTICES AND PROCEDURES

2.1 TEST ADMINISTRATION

Table 4 shows the testing windows for the 2022–2023 SDSA.

Table 4. SDSA Testing Windows

Assessment	Grades	Start Date	End Date	Test Options	Mode
Summative SDSA	5, 8, and 11	03/01/23	05/12/23	English Braille Spanish	Online
Summative SDSA	5, 8, and 11	03/15/23	04/21/23	English DEI Braille DEI	Paper*
Interim SDSA	5, 8, and 11	08/19/22	05/05/23	English	Online

*For the paper-pencil fixed-form tests, all student responses on the paper-pencil tests were entered in the Data Entry Interface (DEI) by test administrators. They included the paper-pencil Large Print tests in English and the paper-pencil Braille tests.

2.2 TEST ADMINISTRATORS

The key personnel involved with the SDSA test administration for the SDDOE are district administrators (DAs), district test coordinators (DCs), school test coordinators (SCs), teachers (TEs), and test administrators (TAs). *Test Administration Manuals* (available at [https://sd.portal.cambiumast.com/resources/educators/test-administration-manual%C2%A0\(tam\)---online,-summative/online-summative-ela,-math,-and-science-test-administration-manual-\(tam\)](https://sd.portal.cambiumast.com/resources/educators/test-administration-manual%C2%A0(tam)---online,-summative/online-summative-ela,-math,-and-science-test-administration-manual-(tam))) were provided so that personnel involved with the statewide assessment administrations could maintain both standardized administration conditions and test security.

2.3 TESTING ENVIRONMENT

The Cambium Assessment, Inc. (CAI) Secure Browser was required to access the online SDSA. The online browser provided a secure environment for student testing by disabling the hot keys, copy, and screen capture capabilities, and preventing access to the desktop (Internet, email, and other files or programs installed on school machines). During the online assessment, students could pause a test, review previously answered questions, and modify their responses if the test had not been paused for more than 20 minutes. Students do not have a required time limit for each test session, but schools were given approximate time estimates for how long each test may take for most students for test administration planning purposes. For additional information about the test administration, refer to Volume 5, Test Administration, of this technical report.

2.4 SIMULATIONS

CAI employs a simulation approach to all SDSA tests before the operational testing window begins. For adaptive tests, simulations were conducted to configure the item selection algorithm settings, to evaluate whether individual tests adhered to the test blueprint and correlated highly with student ability, to monitor item exposure rates, and to verify the scores produced by CAI’s scoring engine. Simulations were also conducted on fixed-form tests to quality check the scores. The simulation approaches and results are discussed in Volume 2, Test Development.

2.5 UNIVERSAL FEATURES, DESIGNATED SUPPORTS, AND ACCOMMODATIONS

The SDSA provides embedded (digitally provided) and non-embedded (non-digitally or locally provided) universal features for all students as they access instructional or assessment content. In addition to those universal features, designated supports and accommodations are provided for students who have special needs.

Accessibility supports are available to students when needed to remove barriers during testing while maintaining the constructs that are measured by the SDSA. The accessibility supports discussed in this technical report include embedded (digitally provided) and non-embedded (non-digitally or locally provided) universal features that are available to all students as they access instructional or assessment content; designated supports that are available to those students for whom the need has been identified by an informed educator or team of educators; and accommodations that are generally available for students for whom there is documentation on an Individualized Education Program (IEP) or Section 504 Plan. For English learners (ELs), Spanish language versions of the SDSA were available.

Scores achieved by students using designated supports are included for federal accountability purposes. All educators making these decisions were trained on the process and understand the range of designated supports available.

Accommodations are changes in procedures or materials that ensure equitable access to instructional and assessment content and generate valid assessment results for students who need them. Embedded accommodations (e.g., text-to-speech [TTS]) are provided digitally through instructional or assessment technology, and non-embedded designated features (e.g., scribe) are non-digital. State-approved accommodations do not compromise the learning expectations, constructs, or grade-level standards. These accommodations help students with a documented need to generate valid assessment outcomes that fully demonstrate what they know and are able to do. From the psychometric point of view, the purpose of providing accommodations is to “increase the validity of inferences about students with special needs by offsetting specific disability-related, construct-irrelevant impediments to performance” (Koretz & Hamilton, 2006, p. 562).

The school coordinators (SCs) and proctors (PRs) in South Dakota were responsible for ensuring that arrangements for accommodations were made before the test administration dates. The available accommodation options for eligible students include braille, streamline, abacus, assistive technology (e.g., adaptive keyboards, touch screen, switches), calculator, print-on-demand, and multiplication table.

Additional information about universal features, designated supports, and accommodations can be found in Volume 5, Test Administration, of this technical report. For more details of the accessibility support, refer to the *Tools, Supports, and Accommodations (TSA) Guidelines* at [https://sd.portal.cambiumast.com/resources/educators/tools,-supports,-and-accommodations-\(tsa\)-guidelines](https://sd.portal.cambiumast.com/resources/educators/tools,-supports,-and-accommodations-(tsa)-guidelines). Table 5 and Table 6 list the number of testing sessions in which a student was provided with each designated support or accommodation during the 2023 administration.

Table 5. Number of Testing Sessions with Allowed Designated Supports

Designated Supports	Grade		
	5	8	11
Embedded			
Color Contrast	9	8	9
Language—Spanish	10	21	12
Masking	5	1	
Mouse Pointer	-		
Streamlined Mode	6	4	2
Text-to-Speech: Items	83	39	23
Text-to-Speech: Stimuli	99	88	25
Text-to-Speech: Stimuli and Items	1,456	942	538
Zoom (5X-20X)	-		
Non-Embedded			
Amplification	3	2	
Color Contrast	-		
Color Overlay	2	-	
Magnification	7	3	2
Medical Support	3	9	3
Noise Buffer	9	6	2
Periodic Table	N/A	-	1
Read Aloud: Items (English)	-		
Read Aloud: Stimuli (English)	-		
Separate Setting	1,053	675	545
Simplified Test Directions	393	345	238

- indicates no data is available; N/A denotes “not applicable”.

Table 6. Number of Testing Sessions with Allowed Accommodations

Accommodations	Grade		
	5	8	11
Embedded			
Braille	-		
Permissive Mode	-	2	
Non-Embedded			
Alternate Response Options	-		
Braille Paper	-		
Calculator	-		
Large Print Paper	-		
Print on Demand	-		
Scribe	20	5	6
Special Considerations (Example: ASL)	1	-	1

- indicates no data is available.

3. ITEM BANK AND TEST DESIGN

3.1 SHARED SCIENCE ASSESSMENT ITEM BANK

CAI works with a group of states and one U.S. territory to develop science assessments to assess three-dimensional conceptualization of science understanding and other standards influenced by the same science framework. Many of these states and one U.S. territory have signed a Memorandum of Understanding (MOU) to share item specifications and items. CAI has coordinated this group of states and one U.S. territory and holds contracts to develop and deliver the items for most of them.

CAI also built the Independent College and Career Readiness (ICCR) science item pool in partnership with these states and one U.S. territory. These CAI-owned items make up a substantial part of the item bank and are shared with partner states and one U.S. territory.

In 2023, the Shared Science Assessment Item Bank was used for operational tests in 12 states and one U.S. territory, including South Dakota. The goals, uses, and claims that the Shared Science Assessment Item Bank and resulting tests are designed to support were identified in a collaborative meeting on August 22–23, 2016, in an attempt to facilitate the transition from a framework for three-dimensional science standards, specifically the NGSS, to statewide summative assessments for science. CAI invited content and assessment leaders from 10 states and four nationally recognized experts who helped co-author the NGSS. Two nationally recognized psychometricians also participated.

The SDSA uses only ICCR items, but because the ICCR science items are part of the Shared Science Assessment Item Bank, the latter will be described in this technical report. The Shared

Science Assessment Item Bank was used for operational accountability tests in 12 states and one U.S. territory in 2023. Tests for one additional state will become operational in 2024.

In 2017, cognitive lab studies were conducted to evaluate and refine the process of developing item clusters aligned to the three-dimensional science standards. The results of the cognitive lab studies confirmed the feasibility of the approach (refer to Volume 4, Appendix 4-F, Science Clusters Cognitive Lab Report, of this technical report).

A second set of cognitive lab studies was conducted in 2018 and 2019 to determine whether students using braille could understand the task demands of selected accommodated three-dimensional science-aligned item clusters. They also evaluated whether these students could navigate the interactive features of these item clusters in a manner that allowed them to fully display their knowledge and skills relative to the constructs of interest. In general, both the students who relied entirely on braille and/or Job Access With Speech (JAWS) and those who had some vision and were able to read the screen with magnification were able to find the information they needed to respond to the questions, navigate the various response formats, and finish within a reasonable amount of time (refer to Volume 4, Appendix 4-E, Braille Cognitive Lab Report, of this technical report).

In 2018, CAI field tested more than 540 item clusters and stand-alone items, of which 451 (including items from all sources) were accepted and made available as operational items in 2019 and future administrations. In 2019, 2021, and 2022, the numbers of item clusters and stand-alone items that were field tested were 347, 545, and 471, while the numbers of items that were accepted and made available for future operational use were 268, 458, and 403, respectively. In 2023, 348 item clusters and stand-alone items were field tested, of which 288 were accepted and made available for operational use in future administrations. All these items follow the same specifications, test development processes, and review processes, summarized below:

- CAI staff and participating states collaborated to develop item specifications, which are documents designed to guide item writers as they craft test questions and stakeholders while they review items. The item specifications were generally accompanied by sample items meeting those specifications. All specifications and sample items were reviewed by state content experts and committees of educators in at least one state.
- The specifications helped test developers create item clusters and stand-alone items that covered a range of difficulty, furthering the goal of measuring the full range of performance found in the population, but remaining at grade level. All item writers were trained in the principles of universal design, the appropriate use of item interactions, and the science item specifications.
- Items were reviewed by science experts in at least one state.
- Every item was reviewed by a content advisory committee (comprised of state educators) in at least one state or in a cross-state educator review process.
- Every item was reviewed by a committee of educators charged with evaluating language accessibility, bias, and sensitivity in at least one state or a cross-state educator review.

- Every item was field tested, all scoring protocols (i.e., rubrics) were validated using the field-test data, and items with questionable data were reviewed again by committees of educators.

A detailed description of the Shared Science Assessment Item Bank development process is included in Volume 2, Test Development.

3.2 FIELD TESTING

All items that were part of the operational pool of the Shared Science Assessment Item Bank were field tested in prior years, which is documented in Appendix 1-B, Shared Science Assessment Item Bank: Field-Testing. Field-testing for the current administration is described in this section.

3.2.1 2023 Field Tests

In 2023, items were field tested in 12 states and one U.S. territory. In all 12 states and one U.S. territory (Connecticut, Hawaii, Idaho, Montana, New Hampshire, North Dakota, Oregon, Rhode Island, South Dakota, U.S. Virgin Islands, Utah, West Virginia, and Wyoming), field-test items were administered as unscored items embedded among operational items. In total, 159 item clusters and 189 stand-alone items were administered as field-test items in the elementary, middle, and high school grade bands. Table 7 presents the number of field-test item clusters and stand-alone items administered in each grade band for each state. The numbers in parentheses in the column representing Total show the number of field-test items owned by ICCR.

Table 7. Number of Field-Test Items Administered, Spring 2023

Grade Band and Item Type	CT	HI	ID	MT	NH	ND	OR	RI	SD	USVI	UT	WV	WY	Total*
Elementary School	41	19	20	12	11	12	28	35	10 (10)	1	32	25	7	126 (35)
Cluster	16	7	12	4	7	4	12	12	4 (4)	1	32	7	3	60 (17)
Stand-Alone	25	12	8	8	4	8	16	23	6 (6)	0	0	18	4	66 (18)
Middle School	36	24	21	9	11	11	41	29	7 (7)	1	49	28	7	136 (41)
Cluster	6	8	5	5	7	4	10	9	5 (5)	1	49	8	3	59 (21)
Stand-Alone	30	16	16	4	4	7	31	20	2 (2)	0	0	20	4	77 (20)
High School	37	8	20	0	12	12	31	36	6 (6)	1	0	0	10	86 (23)
Cluster	21	4	8	0	1	3	25	12	4 (4)	1	0	0	2	40 (6)
Stand-Alone	16	4	12	0	11	9	6	24	2 (2)	0	0	0	8	46 (17)
Total	114	51	61	21	34	35	100	100	23 (23)	3	81	53	24	348 (99)

Note. The numbers in parentheses indicate ICCR items.

*The total count excludes 80 Oregon and ten South Dakota legacy field-tested items, nine HI enemy pilot study field-tested items, and three field-tested items intended for the interim pool while including several field-tested items which were moved to the comprehensive interim pool after rubric validation.

Two of the states (New Hampshire and Rhode Island) opted for a test in which operational items were grouped by science discipline. For these two states, the field-test items were presented together in a fourth group of items. The sequence of the four sets of items (corresponding to the three disciplines and a set of field-test items) was randomized across students. Ten other states and one U.S. territory (Connecticut, Hawaii, Idaho, Montana, North Dakota, Oregon, South Dakota, Utah, U.S. Virgin Islands, West Virginia, and Wyoming) opted for a test design in which the items were not grouped by discipline. In these 10 states and one U.S. territory, field-test items were administered at random positions throughout the test. A student received either a field-test item cluster or a set of four field-test stand-alone items. The test design for the SDSA is discussed in Section 3.3, Test Design.

A minimum sample size of 1,500 students per field-test item was targeted for any given state or territory. Most items were administered in two states or territory. Approximately 86% of the items met or exceeded the target sample size of 1,500 in at least one state, and 100% of the items had a sample size of at least 1,350 (90% of the target) in at least one state.

Table 8 to Table 10 present the number of item clusters and stand-alone items that were shared between the field-test pools of any two states or territory. The numbers below the shaded cells represent the number of common field-test items between any two states, and the numbers above the shaded cells represent the number of common field-test items that survived rubric validation and were included in the calibration. In each of the shaded cells, the number outside the parentheses represents the number of unique field-test items administered only in the given state or territory, and the number in the parentheses represents the number of unique and/or common items that were calibrated with only the data from that state. Table 8 presents the results for elementary schools, Table 9 presents the results for middle schools, and Table 10 presents the results for high schools. The numbers of field-test items administered are slightly different from the numbers of field-test items at calibration because some items were rejected during rubric validation.

Table 8. Number of Common Elementary School Field-Test Items Administered and Calibrated, Spring 2023

	State	CT	HI	ID	MT	NH	ND	OR	RI	SD	USVI	UT	WV	WY	
Item Clusters	CT	0 (0)	0	2	2	2	0	5	4	0	0	1	0	0	
	HI	0	0 (0)	1	0	0	0	0	0	0	0	6	0	0	
	ID	2	1	0 (0)	0	0	0	0	2	0	0	6	0	0	
	MT	2	0	0	0 (0)	0	0	0	0	0	0	2	0	0	
	NH	2	0	0	0	0 (0)	0	0	0	0	0	2	3	0	
	ND	0	0	0	0	0	0 (0)	3	0	0	0	1	0	0	
	OR	5	0	0	0	0	0	3	0 (0)	0	0	0	4	0	0
	RI	4	0	2	0	0	0	0	0	0 (0)	0	0	4	2	0
	SD	0	0	0	0	0	0	0	0	0	0 (0)	1	2	2	0
	USVI	0	0	0	0	0	0	0	0	0	1	0 (0)	1	0	0
	UT	1	6	7	2	2	1	4	4	4	2	1	0 (0)	0	3
	WV	0	0	0	0	0	3	0	0	2	2	0	0	0 (0)	0
WY	0	0	0	0	0	0	0	0	0	0	0	3	0	0 (0)	
Stand-Alone Items	CT	0 (0)	0	6	3	0	0	2	9	0	0	0	3	0	
	HI	0	0 (0)	0	0	0	0	8	0	0	0	0	0	4	
	ID	6	0	0 (0)	0	0	0	0	2	0	0	0	0	0	
	MT	4	0	0	0 (0)	0	0	0	4	0	0	0	0	0	
	NH	0	0	0	0	0 (0)	0	0	0	0	0	0	4	0	
	ND	0	0	0	0	0	0 (0)	6	2	0	0	0	0	0	
	OR	2	8	0	0	0	0	6	0 (0)	0	0	0	0	0	
	RI	10	0	2	4	0	2	0	0 (0)	0	0	0	5	0	
	SD	0	0	0	0	0	0	0	0	0	0 (0)	0	0	6	0
	USVI	0	0	0	0	0	0	0	0	0	0	0 (0)	0	0	0
	UT	0	0	0	0	0	0	0	0	0	0	0	0 (0)	0	0
	WV	3	0	0	0	0	4	0	0	5	6	0	0	0 (0)	0

	State	CT	HI	ID	MT	NH	ND	OR	RI	SD	USVI	UT	WV	WY	
	WY	0	4	0	0	0	0	0	0	0	0	0	0	0 (0)	
Total	CT	0 (0)	0	8	5	2	0	7	13	0	0	1	3	0	
	HI	0	0 (0)	1	0	0	0	8	0	0	0	6	0	4	
	ID	8	1	0 (0)	0	0	0	0	4	0	0	6	0	0	
	MT	6	0	0	0 (0)	0	0	0	4	0	0	2	0	0	
	NH	2	0	0	0	0 (0)	0	0	0	0	0	2	7	0	
	ND	0	0	0	0	0	0 (0)	9	2	0	0	1	0	0	
	OR	7	8	0	0	0	0	9	0 (0)	0	0	0	4	0	0
	RI	14	0	4	4	0	2	0	0	0 (0)	0	0	4	7	0
	SD	0	0	0	0	0	0	0	0	0	0 (0)	1	2	8	0
	USVI	0	0	0	0	0	0	0	0	0	1	0 (0)	1	0	0
	UT	1	6	7	2	2	2	1	4	4	2	1	0 (0)	0	3
	WV	3	0	0	0	0	7	0	0	7	8	0	0	0 (0)	0
WY	0	4	0	0	0	0	0	0	0	0	0	3	0	0 (0)	

Table 9. Number of Common Middle School Field-Test Items Administered and Calibrated, Spring 2023

	State	CT	HI	ID	MT	NH	ND	OR	RI	SD	USVI	UT	WV	WY	
Item Clusters	CT	0 (0)	0	0	0	0	0	0	1	0	1	4	1	0	
	HI	0	0 (0)	0	0	0	0	0	0	0	0	8	0	0	
	ID	0	0	0 (0)	0	0	0	0	1	0	0	3	2	0	
	MT	0	0	0	0 (0)	0	0	0	0	0	0	5	0	0	
	NH	0	0	0	0	0 (0)	0	0	0	0	0	7	0	0	
	ND	0	0	0	0	0	0 (0)	0	1	0	0	3	0	0	
	OR	0	0	0	0	0	0	0 (0)	1	0	0	8	1	0	
	RI	1	0	1	0	0	0	1	1	0 (0)	2	0	2	2	0
	SD	0	0	0	0	0	0	0	0	2	0 (0)	0	3	0	0
	USVI	1	0	0	0	0	0	0	0	0	0	0 (0)	1	0	0
	UT	4	8	3	5	7	3	8	2	3	1	0 (0)	3	2	
	WV	1	0	2	0	0	0	0	1	2	0	0	3	0 (0)	0
WY	0	0	0	0	0	0	0	0	0	0	0	3	0	0 (0)	
Stand-Alone Items	CT	0 (0)	4	10	0	4	0	7	0	0	0	0	4	0	
	HI	4	0 (0)	5	0	0	0	3	0	0	0	0	4	0	
	ID	10	5	0 (0)	0	0	0	0	1	0	0	0	0	0	
	MT	0	0	0	0 (0)	0	0	4	0	0	0	0	0	0	
	NH	4	0	0	0	0 (0)	0	0	0	0	0	0	0	0	
	ND	0	0	0	0	0	0 (0)	0	7	0	0	0	0	0	
	OR	8	3	0	4	0	0	0 (0)	6	0	0	0	7	3	
	RI	0	0	1	0	0	0	7	6	0 (0)	2	0	4	0	
	SD	0	0	0	0	0	0	0	0	2	0 (0)	0	0	0	0
	USVI	0	0	0	0	0	0	0	0	0	0	0 (0)	0	0	0
	UT	0	0	0	0	0	0	0	0	0	0	0	0 (0)	0	0
	WV	4	4	0	0	0	0	0	7	4	0	0	0	0 (0)	1

	State	CT	HI	ID	MT	NH	ND	OR	RI	SD	USVI	UT	WV	WY	
	WY	0	0	0	0	0	0	3	0	0	0	0	1	0 (0)	
Total	CT	0 (0)	4	10	0	4	0	7	1	0	1	4	5	0	
	HI	4	0 (0)	5	0	0	0	3	0	0	0	8	4	0	
	ID	10	5	0 (0)	0	0	0	0	2	0	0	3	2	0	
	MT	0	0	0	0 (0)	0	0	4	0	0	0	5	0	0	
	NH	4	0	0	0	0 (0)	0	0	0	0	0	7	0	0	
	ND	0	0	0	0	0	0 (0)	0	8	0	0	3	0	0	
	OR	8	3	0	4	0	0	0 (0)	7	0	0	8	8	3	
	RI	1	0	2	0	0	0	8	7	0 (0)	4	0	2	6	0
	SD	0	0	0	0	0	0	0	0	4	0 (0)	0	3	0	0
	USVI	1	0	0	0	0	0	0	0	0	0	0 (0)	1	0	0
	UT	4	8	3	5	7	3	8	2	3	1	0 (0)	3	2	
	WV	5	4	2	0	0	0	0	8	6	0	0	3	0 (0)	1
	WY	0	0	0	0	0	0	3	0	0	0	3	1	0 (0)	

Table 10. Number of Common High School Field-Test Items Administered and Calibrated, Spring 2023

	State	CT	HI	ID	MT	NH	ND	OR	RI	SD	USVI	UT	WV	WY
Item Clusters	CT	0 (0)	1	1	-	0	1	10	6	0	0	-	-	1
	HI	1	0 (0)	1	-	0	0	2	0	0	0	-	-	0
	ID	1	1	0 (0)	-	0	0	5	0	0	0	-	-	0
	MT	-	-	-	-	-	-	-	-	-	-	-	-	-
	NH	0	0	0	-	0 (0)	0	0	0	1	0	-	-	0
	ND	1	0	0	-	0	0 (0)	0	0	1	1	-	-	1
	OR	10	2	6	-	0	0	0 (0)	3	2	0	-	-	0
	RI	7	0	0	-	0	0	5	0 (0)	0	0	-	-	0
	SD	0	0	0	-	1	1	2	0	0 (0)	0	-	-	0
	USVI	0	0	0	-	0	1	0	0	0	0 (0)	-	-	1
	UT	-	-	-	-	-	-	-	-	-	-	-	-	-
	WV	-	-	-	-	-	-	-	-	-	-	-	-	-
WY	1	0	0	-	0	1	0	0	0	1	-	-	0 (0)	
Stand-Alone Items	CT	0 (0)	0	5	-	0	0	0	10	0	0	-	-	0
	HI	0	0 (0)	0	-	0	2	2	0	0	0	-	-	0
	ID	5	0	0 (0)	-	3	0	4	0	0	0	-	-	0
	MT	-	-	-	-	-	-	-	-	-	-	-	-	-
	NH	0	0	3	-	0 (0)	5	0	3	0	0	-	-	0
	ND	0	2	0	-	5	0 (0)	0	1	0	0	-	-	1
	OR	0	2	4	-	0	0	0 (0)	0	0	0	-	-	0
	RI	11	0	0	-	3	1	0	0 (0)	2	0	-	-	7
	SD	0	0	0	-	0	0	0	2	0 (0)	0	-	-	0
	USVI	0	0	0	-	0	0	0	0	0	0 (0)	-	-	0
	UT	-	-	-	-	-	-	-	-	-	-	-	-	-
	WV	-	-	-	-	-	-	-	-	-	-	-	-	-

	State	CT	HI	ID	MT	NH	ND	OR	RI	SD	USVI	UT	WV	WY
	WY	0	0	0	-	0	1	0	7	0	0	-	-	0 (0)
Total	CT	0 (0)	1	6	-	0	1	10	16	0	0	-	-	1
	HI	1	0 (0)	1	-	0	2	4	0	0	0	-	-	0
	ID	6	1	0 (0)	-	3	0	9	0	0	0	-	-	0
	MT	-	-	-	-	-	-	-	-	-	-	-	-	-
	NH	0	0	3	-	0 (0)	5	0	3	1	0	-	-	0
	ND	1	2	0	-	5	0 (0)	0	1	1	1	-	-	2
	OR	10	4	10	-	0	0	0 (0)	3	2	0	-	-	0
	RI	18	0	0	-	3	1	5	0 (0)	2	0	-	-	7
	SD	0	0	0	-	1	1	2	2	0 (0)	0	-	-	0
	USVI	0	0	0	-	0	1	0	0	0	0 (0)	-	-	1
	UT	-	-	-	-	-	-	-	-	-	-	-	-	-
	WV	-	-	-	-	-	-	-	-	-	-	-	-	-
	WY	1	0	0	-	0	2	0	7	0	1	-	-	0 (0)

Table 11. Number of Field-Test Item Administration, Rubric Validation, and Item Data Review, Spring 2023

Grade Band and Item Type	Number of Field-Test Items Administered	Number of Items Rejected Before/During Rubric Validation	Number of Items Sent to Data Review	Number of Items Rejected at Data Review	Number of Items Remaining
Elementary School	126 (10)	3 (0)	71 (7)	13 (0)	110 (10)
Cluster	60 (4)	1 (0)	17 (1)	1 (0)	58 (4)
Stand-Alone	66 (6)	2 (0)	54 (6)	12 (0)	52 (6)
Middle School	136 (7)	2 (0)	80 (4)	20 (0)	114 (7)
Cluster	59 (5)	1 (0)	21 (2)	5 (0)	53 (5)
Stand-Alone	77 (2)	1 (0)	59 (2)	15 (0)	61 (2)
High School	86 (6)	5 (0)	44 (1)	17 (0)	64 (6)
Cluster	40 (4)	4 (0)	19 (1)	6 (0)	30 (4)
Stand-Alone	46 (2)	1 (0)	25 (0)	11 (0)	34 (2)
Total	348 (23)	10 (0)	195 (12)	50 (0)	288 (23)

Note. The numbers in parentheses indicate ICCR items administered in South Dakota.

Following the administration, field-test items went through a substantial validation process. The process began with rubric validation. Rubric validation is a process in which a committee of state educators reviews student responses and the proposed scoring of those responses. The process is described in Volume 2, Section 2.7.1, Rubric Validation, of this technical report.

After rubric validation, classical item statistics were computed for the scoring assertions, including item difficulty and item discrimination statistics, testing time, and differential item functioning (DIF) statistics. The MOU established common standards for the statistics. Any items violating these standards were flagged for a second educator review. Even though the scoring assertions were the basic units of analysis used to compute classical item statistics, the business rules to flag items for another educator review were established at the item level because assertions cannot be reviewed in isolation. The statistics and business rules for flagging items are described in Section 4, Field-Test Classical Analysis. For each state, a data review committee consisting of educators (i.e., science teachers) supported by CAI content experts reviewed the items that were owned by the state and flagged for data review according to the established business rules. For ICCR, cross-state review committees were established.

Table 11 presents the number of field-test items administered in South Dakota, or another state or territory, the number of items rejected before or during rubric validation, the number of items sent for data review, and the number of items rejected during data review. The numbers in parentheses present the number of field-test items owned by ICCR.

Table 12 summarizes the Shared Science Assessment Item Bank after adding the field-test items that were administered in 2023 and passed rubric validation and item data review. The numbers in parentheses present the number of items owned by ICCR.

Table 12. Summary of Shared Science Assessment Item Bank, Spring 2023

Grade Band and Item Type	Science Discipline			Item Bank Total ^a
	Earth and Space Sciences	Life Sciences	Physical Sciences	
Elementary School	205 (47)	202 (44)	254 (55)	661 (146)
Cluster	112 (20)	102 (14)	131 (18)	345 (52)
Stand-Alone	93 (27)	100 (30)	123 (37)	316 (94)
Middle School	185 (37)	262 (60)	215 (45)	662 (142)
Cluster	88 (15)	129 (20)	106 (17)	323 (52)
Stand-Alone	97 (22)	133 40()	109 (28)	339 (90)
High School	110 (34)	207 (62)	151 (35)	468 (131)
Cluster	45 (13)	89 (20)	62 (18)	196 (51)
Stand-Alone	65 (21)	118 (42)	89 (17)	272 (80)
Total	500 (117)	671 (166)	620 (135)	1791 (419)

Note. The numbers in parentheses indicate ICCR items administered in South Dakota.

^aCount excludes nine MOU items that do not align to the NGSS.

3.3 TEST DESIGN

The science tests were assembled under an adaptive test design, with the exception of the braille, and paper-pencil forms. Tests were assembled using CAI’s adaptive testing algorithm. The adaptive item selection algorithm selects items based on their content value and information value. At any given point during the test, the content value of an item is determined by its contribution to meeting the blueprint, given the content characteristics of the items that have already been administered. During the test, the content value increases for items that exhibit features that have not met their designated minimum as the end of the test approaches. Vice versa, the content value decreases for items with content features for which the minimum has been met. The information value of an item is based on the item information function evaluated at the estimated proficiency. The proficiency estimate is updated throughout the test.

Under an adaptive test design, operational items are selected on the fly based on the performance of a student on past items while ensuring the test blueprint is followed for each individual student. The SDSA NGSS Assessment blueprints are presented in this technical report in Volume 2, Section 4.2, Test Blueprints. Details of CAI’s item selection algorithm are described in Volume 2, Appendix 2-L, Adaptive Algorithm Design.

The braille and paper-pencil tests were accommodated fixed forms. Form construction of the accommodated forms is discussed in Volume 2, Section 4.4, Paper-Pencil Accommodation Form Construction.

The main characteristics of the blueprint were that any performance expectation (PE) could be tested only once (indicated by the values of 0 and 1 for the minimum and maximum values of the individual PEs in the test blueprints; see Section 4.2, Test Blueprints of Volume 2). In general, no more than one item cluster or two stand-alone items could be sampled from the same disciplinary core idea (DCI), and no more than three total items could be sampled from the same DCI (as indicated by the minimum and maximum values in the rows representing DCIs).

Since 2021, a non-segmented test design was used for the operational test administration. Students received items from different disciplines in random order. Embedded field-test items were randomly positioned in the test and randomly distributed across students. Every student received either one item cluster or four stand-alone items as field-test items throughout the test.

4. FIELD-TEST CLASSICAL ANALYSIS

As explained in Section 3, Item Bank and Test Design, science items administered as field-test items in South Dakota or any of the states or territory that signed the Memorandum of Understanding (MOU) for item sharing underwent rubric validation and data review. Items were flagged for data review based on business rules defined on classical item statistics. Except for response times, the classical item statistics were computed for individual assertions, whereas the business rules for flagging were defined at the item level. In general, item statistics used to flag items for data review were computed using the student responses of the state or territory that owned the items. However, for Independent College and Career Readings (ICCR) item bank items, the flagging rules were defined on the item statistics computed from the combined data of states or territory that used ICCR items and that administered either an independent or operational test. In 2023, those states were Connecticut, Idaho, New Hampshire, North Dakota, Oregon, South Dakota, Rhode Island, Utah, and West Virginia. Furthermore, to compute the differential item functioning (DIF) statistics for the field-test items, the data from all states and territories with an operational or independent field test were combined to obtain a sufficient number of students for each demographic group. The criteria for flagging and reviewing items are provided in Table 13, and the statistics are described in Section 4.1, Item Discrimination, through Section 4.4, Differential Item Functioning Analysis. Items flagged for data review were reviewed by a committee, as explained in Section 3, Item Bank and Test Design.

Table 13. Thresholds for Flagging in Classical Item Analysis

Analysis Type	Flagging Criteria
Item Discrimination	Average biserial correlation < 0.25 (across the assertions within an item)
	One or more assertions with a biserial correlation < 0.05
Item Difficulty (Item Clusters)	Average <i>p</i> -value < 0.30 or > 0.85 (across the assertions within an item cluster)

Analysis Type	Flagging Criteria
Item Difficulty (Stand-Alone Items)	Average p -value < 0.15 or > 0.95 (across the assertions within a stand-alone item)
Timing (Item Clusters)	Percentile 80* > 15 minutes
Timing (Stand-Alone Items)	Percentile 80* > 3 minutes
Timing	Assertions per minute < 0.5
DIF (Item Clusters)	Two or more assertions show “C” DIF in the same direction
DIF (Stand-Alone Items)	One or more assertions show “C” DIF in the same direction

*A percentile 80 of x minutes: 80% of the students spent x minutes or less on the item.

4.1 ITEM DISCRIMINATION

The item discrimination index indicates the extent to which each item differentiated between those test takers who possessed the skills being measured and those who did not. Generally, the higher the value, the better the item is able to differentiate between high- and low-achieving students.

For each assertion within an item, the discrimination index was calculated as the biserial correlation between the assertion score and the ability estimate for students. The average biserial correlation was then calculated across the assertions within an item.

4.2 ITEM DIFFICULTY

Items that are either very difficult or very easy are flagged for review but are not necessarily removed if they are grade-level appropriate and aligned with the test specifications. Both the p -value for individual assertions and the average across all assertions of an item are calculated. Acceptable item p -values are summarized in Table 13.

4.3 RESPONSE TIME

Given that the science item clusters consist of multiple student interactions, they require more time for students to complete. To ensure a good balance between the amount of information an item provided, and the time students spend on the item, item response time was recorded and analyzed. Specifically, the statistic “percentile 80” was computed for each item. A percentile 80 of x minutes means that 80% of the students spent x minutes or fewer on the item. An item was flagged for review when the

- percentile 80 > 15 minutes, if the item is an item cluster;
- percentile 80 > 3 minutes, if the item is a stand-alone item; or
- assertions per (percentile 80) minute < 0.5 .

4.4 DIFFERENTIAL ITEM FUNCTIONING

DIF refers to items that appear to function differently across identifiable groups, typically across different demographic groups. Identifying DIF is important because it provides a statistical indicator that an item may contain cultural or other bias. DIF-flagged items are further examined by content experts who are asked to re-examine each flagged item to decide whether the item should be excluded from the pool due to bias. Not all items that exhibit DIF are biased, and various characteristics of the educational system may also lead to DIF.

CAI uses a generalized Mantel-Haenszel (MH) procedure to calculate DIF. The generalizations include adaptation to polytomous items and improved variance estimators to render the test statistics valid under complex sample designs. With this procedure, each student’s estimated theta score on the operational items on a given test is used as the ability-matching variable. That score is divided into 10 intervals to compute the MH chi-square ($MH\chi^2$) DIF statistics for balancing the stability and sensitivity of the DIF scoring category selection. For dichotomous items, the following statistics were computed: $MH\chi^2$ value, the conditional odds ratio, and the MH-delta. For polytomous items, the $GMH\chi^2$ and the standardized mean difference (SMD [Dorans & Schmitt, 1991]) were computed.

The MH chi-square statistic (Holland & Thayer, 1988) is calculated as:

$$MH\chi^2 = \frac{(|\sum_k n_{R1k} - \sum_k E(n_{R1k})| - 0.5)^2}{\sum_k var(n_{R1k})},$$

where $k = \{1, 2, \dots, K\}$ for the strata, n_{R1k} is the number of students with correct responses for the reference group in stratum k , and 0.5 is a continuity correction. The expected value is calculated as

$$E(n_{R1k}) = \frac{n_{+1k}n_{R+k}}{n_{++k}},$$

where n_{+1k} is the number of students with correct responses, n_{R+k} is the number of students in the reference group, and n_{++k} is the number of students in stratum k . The variance is calculated as

$$var(n_{R1k}) = \frac{n_{R+k}n_{F+k}n_{+1k}n_{+0k}}{n_{++k}^2(n_{++k}-1)},$$

where n_{F+k} is the number of students in the focal group, n_{+1k} is the number of students with correct responses, and n_{+0k} is the number of students with incorrect responses in stratum k .

The MH conditional odds ratio is calculated as

$$\alpha_{MH} = \frac{\sum_k n_{R1k}n_{F0k}/n_{++k}}{\sum_k n_{R0k}n_{F1k}/n_{++k}}.$$

The MH-delta (Δ_{MH} [Holland & Thayer, 1988]) is then defined as

$$\Delta_{MH} = -2.35\ln(\alpha_{MH}).$$

The generalized MH statistic generalizes the MH statistic to polytomous items (Somes, 1986), and is defined as

$$GMH\chi^2 = (\sum_k \mathbf{a}_k - \sum_k E(\mathbf{a}_k))' (\sum_k var(\mathbf{a}_k))^{-1} (\sum_k \mathbf{a}_k - \sum_k E(\mathbf{a}_k)),$$

where \mathbf{a}_k is a $(T - 1) \times 1$ vector of item response scores and $E(\mathbf{a}_k)$ is a $(T - 1) \times 1$ mean vector, both corresponding to the T response categories of a polytomous item (excluding one response); $var(\mathbf{a}_k)$ is a $(T - 1) \times (T - 1)$ covariance matrix calculated analogously to the corresponding elements in $MH\chi^2$ in stratum k .

The SMD (Dorans & Schmitt, 1991) is defined as

$$SMD = \sum_k p_{Fk} m_{Fk} - \sum_k p_{Rk} m_{Rk},$$

where

$$p_{Fk} = \frac{n_{F+k}}{n_{F++}}$$

is the proportion of the focal group students in stratum k ,

$$m_{Fk} = \frac{1}{n_{F+k}} \left(\sum_t a_t n_{Ftk} \right)$$

is the mean item score for the focal group in stratum k , and

$$m_{Rk} = \frac{1}{n_{R+k}} \left(\sum_t a_t n_{Rtk} \right)$$

is the mean item score for the reference group in stratum k .

DIF analysis was conducted for all field-test items with at least 200 responses per item in each subgroup (Zwick, 2012) to detect potential item bias for major demographic groups. Student responses from multiple states were combined to minimize the number of items with insufficient sample sizes for one or more demographic groups.

DIF statistics were calculated at the assertion level and were performed for the following groups (some items had insufficient sample sizes for DIF analyses in some groups):

- Female vs. male
- American Indian/Alaskan Native vs. White
- Hawaiian/Pacific Islander vs. White
- Asian vs. White
- African American vs. White
- Hispanic vs. White

- Multi-Racial vs. White
- English Learner (EL) vs. Non-EL
- Special Education (SPED) vs. Non-SPED
- Economically Disadvantaged vs. Non-Economically Disadvantaged

Similar to how the general MH statistic is used to classify items on traditional tests, assertions were classified into three categories (i.e., A, B, or C) for DIF, ranging from “no evidence of DIF” to “severe DIF”. The classification rules are shown in Table 13. Furthermore, assertions were categorized positively (i.e., +A, +B, or +C), signifying that an item favored the focal group (e.g., African American, Hispanic, female), or negatively (i.e., –A, –B, or –C), signifying that an item favored the reference group (e.g., White, male).

An item was flagged for data review according to the following criteria:

- **Item Clusters.** Two or more assertions showed “C” DIF in the same direction.
- **Stand-Alone Items.** One or more assertions showed “C” DIF in the same direction.

Table 14. DIF Classification Rules

Assertions	
Category	Rule
C	MH_{X^2} is significant and $ SMD / SD \geq 0.25$.
B	MH_{X^2} is significant and $ SMD / SD < 0.25$.
A	MH_{X^2} is not significant.

Note that for the 2018 field test, a slightly less strict criterion was used for item clusters with 10 or more assertions (i.e., three or more assertions with “C” DIF in the same direction). The change was made taking into consideration the feedback received from several Technical Advisory Committees (TACs) and modified such that the rate of flagging items for DIF was similar for item clusters and stand-alone items (based on the flagging rates computed on items field-tested in 2018).

4.5 CLASSICAL ANALYSIS RESULTS

This section presents a summary of results from classical item analysis of the field-test items administered in 2023. A total of 33 field-test items were administered in South Dakota, with ten items from the legacy science test (two, two, and six for elementary, middle, and high schools, respectively). Thirty-one field-test items passed rubric validation, with 23 ICCR field-test items and eight legacy items. For the 23 field-test ICCR items that passed rubric validation, one was flagged for item discrimination, four items were flagged for p -value, eight items were flagged for response time, and no item was flagged for DIF according to the criteria used in 2023 (as described in Section 4.1, Item Discrimination, through Section 4.4, Differential Item Functioning Analysis). Some items were flagged for multiple reasons. Flagged field-test items were reviewed by educators

during data review. For the ICCR items, the total number of field-test items flagged and the total number of field-test items that passed rubric validation in 2023 were summarized in Table 11.

For the ten legacy items that passed rubric validation, one item was flagged for *p*-value and item discrimination, three items were flagged for item discrimination, and two items were flagged for response time according to the criteria used in 2023 (as described in Section 4.1, Item Discrimination, through Section 4.4, Differential Item Functioning Analysis Analysis). The flagged field-test items were reviewed by CAI content specialists and SDDOE educators during item data review. The total number of field-test items flagged and the total number of field-test items that passed item data review in 2023 are summarized in Table 12.

Table 15 and Table 16 provide the summary of the *p*-values and biserial correlations for these field-test items. Table 17 presents the summary of the response times by item type (item cluster or stand-alone item). The statistics were computed using South Dakota data only. The average values across the assertions within an item were used to compute percentiles and ranges.

Table 15. Distribution of p-Values for Field-Test Items in Spring 2023

Grade	Total FT Items	Min	5th Percentile	25th Percentile	50th Percentile	75th Percentile	95th Percentile	Max
5	10	0.25	0.30	0.47	0.49	0.58	0.82	0.88
8	7	0.22	0.22	0.27	0.32	0.37	0.53	0.59
11	6	0.24	0.28	0.42	0.50	0.55	0.60	0.61

Table 16. Distribution of Item Biserial Correlations for Field-Test Items in Spring 2023

Grade	Total FT Items	Min	5th Percentile	25th Percentile	50th Percentile	75th Percentile	95th Percentile	Max
5	10	0.27	0.30	0.39	0.43	0.49	0.61	0.69
8	7	0.28	0.29	0.35	0.46	0.54	0.57	0.57
11	6	0.19	0.22	0.32	0.36	0.45	0.49	0.50

Table 17. Summary of Response Times for Field-Test Items Administered in Spring 2023

Grade	Item Type	Total FT Items	Min	25th Percentile	50th Percentile	75th Percentile	Max
5	Cluster	4	8.3	8.4	9.0	9.6	10.1
	Stand-Alone	6	2.1	2.1	2.1	2.2	2.4
8	Cluster	5	8.5	8.6	8.8	9.3	13.1
	Stand-Alone	2	2.0	2.1	2.3	2.5	2.8
11	Cluster	4	5.2	5.2	5.4	5.8	7.3
	Stand-Alone	2	1.9	1.9	1.9	2.0	2.0

Table 18 presents the number of items evaluated and the number of flagged items for DIF for each item type and demographic group included in the DIF analysis.

Table 18. Differential Item Functioning Classifications for Field-Test Items Administered in Spring 2023

DIF Flag	Item Type	Female/ Male	American Indian ^a / White	Asian/ White	African American /White	Hawaiian ^b / White	Hispanic/ White	Multi- Racial/ White	EL/Non- EL	SPED/ Non- SPED	Low Income/ Non-Low Income ^c
Grade 5											
Items Evaluated	Cluster	4	4	0	1	0	2	0	2	4	2
	Stand-Alone	6	2	0	0	0	0	1	0	6	0
Items Flagged C	Cluster	0	0	0	0	0	0	0	0	0	0
	Stand-Alone	0	0	0	0	0	0	0	0	0	0
% Items Flagged C	Cluster	0	0	-	0	-	0	-	0	0	0
	Stand-Alone	0	0	-	-	-	-	0	-	0	-
Grade 8											
Items Evaluated	Cluster	5	3	0	2	0	5	5	5	5	5
	Stand-Alone	2	2	0	2	0	2	0	2	2	2
Items Flagged C	Cluster	0	0	0	0	0	0	0	0	0	0
	Stand-Alone	0	0	0	0	0	0	0	0	0	0
% Items Flagged C	Cluster	0	0	-	0	-	0	0	0	0	0
	Stand-Alone	0	0	-	0	-	0	-	0	0	0
Grade 11											
Items Evaluated	Cluster	4	1	0	0	0	2	0	0	3	3
	Stand-Alone	2	0	0	0	0	2	0	0	2	2
Items Flagged C	Cluster	0	0	0	0	0	0	0	0	0	0
	Stand-Alone	0	0	0	0	0	0	0	0	0	0
% Items Flagged C	Cluster	0	0	-	-	-	0	-	-	0	0
	Stand-Alone	0	-	-	-	-	0	-	-	0	0

Note. Full DIF group names: ^aAmerican Indian/Alaskan Native; ^bHawaiian/Pacific Islander; and ^cEconomically Disadvantaged vs. Non-Economically Disadvantaged

5. ITEM CALIBRATION

5.1 MODEL DESCRIPTION

In discussing item response theory (IRT) models for South Dakota, we distinguish between the underlying latent structure of a model and the parameterization of the item response function conditional on that assumed latent structure. Subsequently, we discuss how group effects are considered.

5.1.1 Latent Structure

Most operational assessment programs rely on a unidimensional IRT model for item calibration and computing scores for students. These models assume a single underlying trait and that items are independent given that underlying trait. In other words, the models assume that given the value of the underlying trait, knowing the response to one item provides no information about responses to other items. This assumption of conditional independence implies that the conditional probability of a pattern of I item responses takes the relatively simple form of a product over items for a single student, as shown below:

$$P(\mathbf{z}_j|\theta_j) = \prod_{i=1}^I P(z_{ij}|\theta_j), \quad (1)$$

where z_{ij} represents the scored response of student j ($j = 1, \dots, N$) to item i ($i = 1, \dots, I$), \mathbf{z}_j represents the pattern of scored item responses for student j , and θ_j represents student j 's proficiency. Unidimensional IRT models differ with respect to the functional relation between the proficiency θ_j and the probability of obtaining a score z_{ij} on item i .

The SDSA items are more complex than traditional item types. A single item may contain multiple parts, and each part may contain multiple student interactions. For example, a student may be asked to select a term from a set of terms at several places in a single item. Instead of receiving a single score for each item, multiple inferences are made about the knowledge and skills that a student has demonstrated based on specific features of the student's responses to the item. These scoring units are called *assertions* and are the basic unit of analysis in our IRT analysis. That is, they fulfill the role of items in traditional assessments; however, for the SDSA items, multiple assertions are typically developed around a single item so that assertions are clustered within items.

One approach is to apply one of the traditional IRT models to the scored assertions; however, a substantial complexity that arises from using this new item type is that local dependencies exist between assertions pertaining to the same stimulus (i.e., item or item cluster). The local dependencies between the assertions pertaining to the same stimulus constitute a violation of the assumption that a single latent trait can explain all dependencies between assertions. Fitting a unidimensional model in the presence of local dependencies may result in biased item parameters and standard errors of measurement (SEMs). In particular, it is well documented that ignoring local

item dependencies leads to an overestimation of the amount of information conveyed by a set of responses and an underestimation of the SEM (e.g., Sireci, Thissen, & Wainer, 1991; Yen, 1993).

The effects of groups of assertions developed around a common stimulus can be accounted for by including additional dimensions corresponding to those groupings in the IRT model. These dimensions are considered to be nuisance dimensions¹. Whereas traditional unidimensional IRT models assume that all assertions (the basic units of analysis) are independent given a single underlying trait θ , we now assume the conditional independence of assertions, given the underlying latent trait θ and all nuisance dimensions:

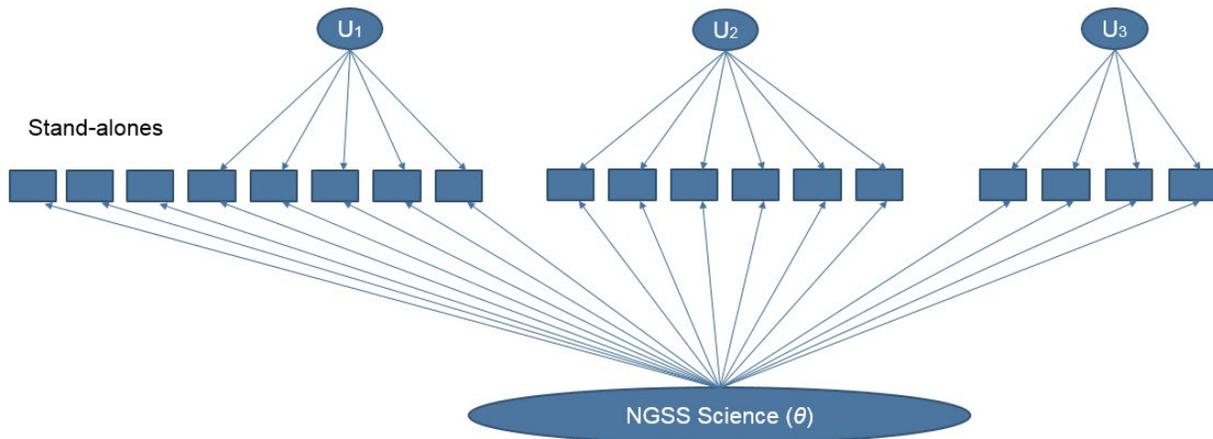
$$P(\mathbf{z}_j|\theta_j, \mathbf{u}_j) = \prod_{i \in SA} P(z_{ij}|\theta_j) \prod_{g=1}^G \prod_{i \in g} P(z_{ij}|\theta_j, u_{jg}), \quad (2)$$

where SA indicates stand-alone item assertions, u_g indicates the nuisance dimension for assertion group g (with the position of student j on that dimension denoted as u_{jg}), and \mathbf{u} is the vector of all G nuisance dimensions. It can be seen that the conditional probability $P(z_{ij}|\theta_j, u_{jg})$ becomes a function of two latent variables: the latent trait θ , representing a student’s proficiency in science (the underlying trait of interest), and the nuisance dimension u_g , accounting for the conditional dependencies between assertions of the same group. Furthermore, we assume that the nuisance dimensions are all uncorrelated with one another and with the general dimension. It is important to point out that even though every group of assertions introduces an additional dimension, models with this latent structure do not suffer from the complications of dimensionality like other multidimensional IRT models because one can take advantage of this special structure during model calibration (Gibbons & Hedeker, 1992). In this regard, Rijmen (2010) showed that it is unnecessary to assume all nuisance dimensions are uncorrelated; instead, it is sufficient that they are independent, given the general dimension θ .

The model structure of the IRT model for science is illustrated in Figure 1. Note that stand-alone items can be scored with more than one assertion. The assertions of stand-alone items with more than one assertion, but fewer than four assertions, were also modeled as stand-alone item assertions. Even though these assertions are likely to exhibit conditional dependencies, the variance of the nuisance dimension cannot be reliably estimated if it is based on a very small number of assertions. The few stand-alone items with four or more assertions were treated as item clusters to take into account the conditional dependencies.

¹ The term *nuisance dimension* here pertains to within-item local dependencies among scoring assertions and should not be confused with the three dimensions of *A Framework for K–12 Science Education*.

Figure 1. Directed Graph of the Science IRT Model



5.1.2 Item Response Function

The item response functions of the stand-alone item assertions are modeled with a unidimensional model. For the grouped assertions, like in unidimensional models, different parametric forms can be assumed for the conditional probability of obtaining a score of z_{ij} . The Rasch testlet model (Wang & Wilson, 2005) is adopted as the IRT model for the SDSA. For binary data, the Rasch testlet model is defined as:

$$P(z_{ij}|\theta_j, u_{jg}; b_i) = \frac{\exp(\theta_j + u_{jg} - b_i)}{1 + \exp(\theta_j + u_{jg} - b_i)} \quad (3)$$

The item response function of the Rasch testlet model is the probability of a correct answer (i.e., a true assertion), as a function of the overall proficiency θ , the overall nuisance dimension u_g , and the item (i.e., assertion) difficulty b_i . The Rasch testlet model does not include item discrimination parameters; however, the same model structure as presented in Figure 1 could be employed with discrimination parameters included in Equations (2) and (3). Furthermore, only models for binary data are considered. Assertions are always binary because they are either true or false. Nevertheless, the model could easily accommodate polytomous responses by using the same response function incorporated in unidimensional models for polytomous data.

5.1.3 Multigroup Model

The Shared Science Assessment Item Bank is calibrated concurrently using all the items administered in any state or territory that collaborates with CAI on their new science assessments. In the calibration, each state or territory is treated as a population of students or a group. Overall group differences are taken into account by allowing a group-specific distribution of the overall

proficiency variable θ . Specifically, for every student j belonging to group k , $k = 1, \dots, K$, a normal distribution is assumed,

$$\theta_j \sim N(\mu_k, \sigma_k^2),$$

where μ_k and σ_k^2 are the mean and variance of a normal distribution. The mean of the reference distribution ($k = 1$) is set to 0 to identify the model (for free item calibrations, where there are no anchor items with their location parameters set to specific values). For each of the nuisance variables u_g , a common variance parameter across groups is assumed, and the means are set to 0 in order to identify the model,

$$u_{jg} \sim N(0, \sigma_{u_g}^2).$$

5.2 ESTIMATION

5.2.1 Estimation

A separate IRT model is fit for each grade band. The parameters of the IRT model are estimated using the marginal maximum likelihood (MML) method. In the MML method, the latent proficiency variable θ_j and the vector of nuisance parameters \mathbf{u}_j for each student j are treated as random effects and integrated out to obtain the marginal log likelihood corresponding to the observed response pattern \mathbf{z}_j for student j ,

$$\ell_j = \log \int \int P(\mathbf{z}_j | \theta_j, \mathbf{u}_j) N(\theta_j | \mu_k, \sigma_k^2) N(\mathbf{u}_j | \mathbf{0}, \mathbf{\Sigma}) d\mathbf{u}_j d\theta_j,$$

where $\mathbf{\Sigma}$ is a diagonal matrix with diagonal elements $\sigma_{u_k}^2$, denoting nuisance variance for group k . Across all students and groups, the overall log likelihood to be maximized with respect to the vector $\boldsymbol{\gamma}$ of all model parameters (i.e., item difficulty parameters, and the mean and variance parameters of the latent variables) is

$$\ell(\boldsymbol{\gamma}) = \sum_k \sum_{j \in k} \ell_j.$$

Even though the number of latent variables in the overall log likelihood equation is very high, issues with dimensionality can be avoided because the integration over the high-dimensional latent (θ, \mathbf{u}) space can be carried out as a sequence of computations in two-dimensional space (θ, \mathbf{u}_g) (Gibbons & Hedeker, 1992; Rijmen, 2010).

The Shared Science Assessment Item Bank was calibrated freely in 2018 after the 2018 science test administrations concluded, and it was recalibrated in 2019 after the 2019 test administrations. Following 2019, field-test items are calibrated onto the scale of the Shared Science Assessment Item Bank by anchoring the operational items to their bank. In the anchored calibrations, the mean and variance of the overall science dimension are also estimated for each group.

Appendix 1-C, Calibration of the Shared Science Assessment Item Bank, contains a detailed description of the 2018 and 2019 calibration processes as well as a description of how the 2018 and 2019 scales were linked.

Starting in 2021, CAIRT (Cambium Assessment IRT) is used to calibrate item parameters. CAIRT was specifically developed by CAI to calibrate the multigroup Rasch model on very large data sets because estimation times in commercially available software (i.e., flexMIRT) became prohibitive. CAIRT relies on the same estimation methods as the Bayesian networks with the logistic regression (BNL; Rijmen, 2006), a suite of Matlab functions for estimating a wide variety of latent variable models. BNL uses an efficient expectation-maximization (EM) algorithm based on graphical model theory (Rijmen, 2010). CAI has cross-validated parameter estimates from CAIRT with BNL and flexMIRT under various scenarios (Rijmen, Liao, & Lin, 2021). CAIRT is a web application that is available at no cost to members of the MOU.

Table 19 provides an overview of the groups per grade band for calibration of the 2023 field-test items. All but 81 items were calibrated on at least 1,500 student responses, and all but five items had a sample size larger than 1,200. The five items with fewer than 1,200 responses were either Oregon legacy items or Hawaii items (for research purposes).

Table 19. Groups Per Grade Band for the Spring 2023 Calibration of Field-Test Items

Group	Elementary School	Middle School	High School
Connecticut	X	X	X
Hawaii	X	X	X
Idaho	X	X	X
Montana	X	X	
New Hampshire	X	X	X
North Dakota	X	X	X
Oregon	X	X	X
Rhode Island	X	X	X
South Dakota	X	X	X
U.S. Virgin Islands	X	X	X
Utah	X	X	
West Virginia	X	X	
Wyoming	X	X	X

5.2.2 Overview of the Operational Item Bank

Figure 2, Figure 3, and Figure 4 display the histogram of the difficulty parameters for grades 5, 8, and 11, respectively, for all items that are part of the SDSA operational pool. The figures also display the student proficiency distributions. The grade 5 items are slightly easier compared to the student proficiency level. For grades 8 and 11, items are slightly more difficult than the student proficiency in general.

Figure 2. SDSA Assertion Difficulty and Student Proficiency Distributions, Grade 5

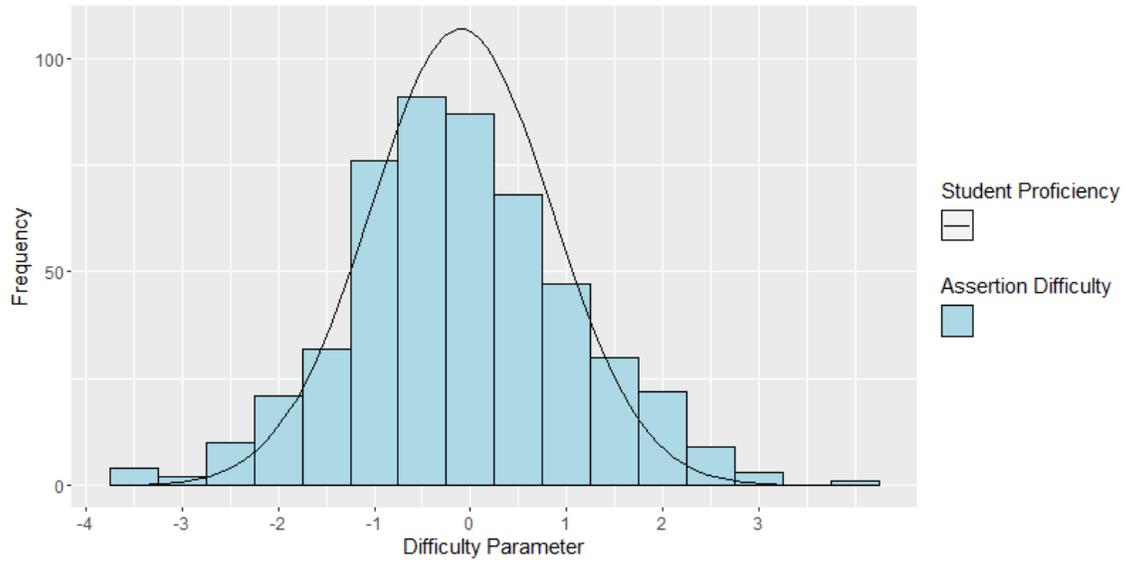


Figure 3. SDSA Assertion Difficulty and Student Proficiency Distributions, Grade 8

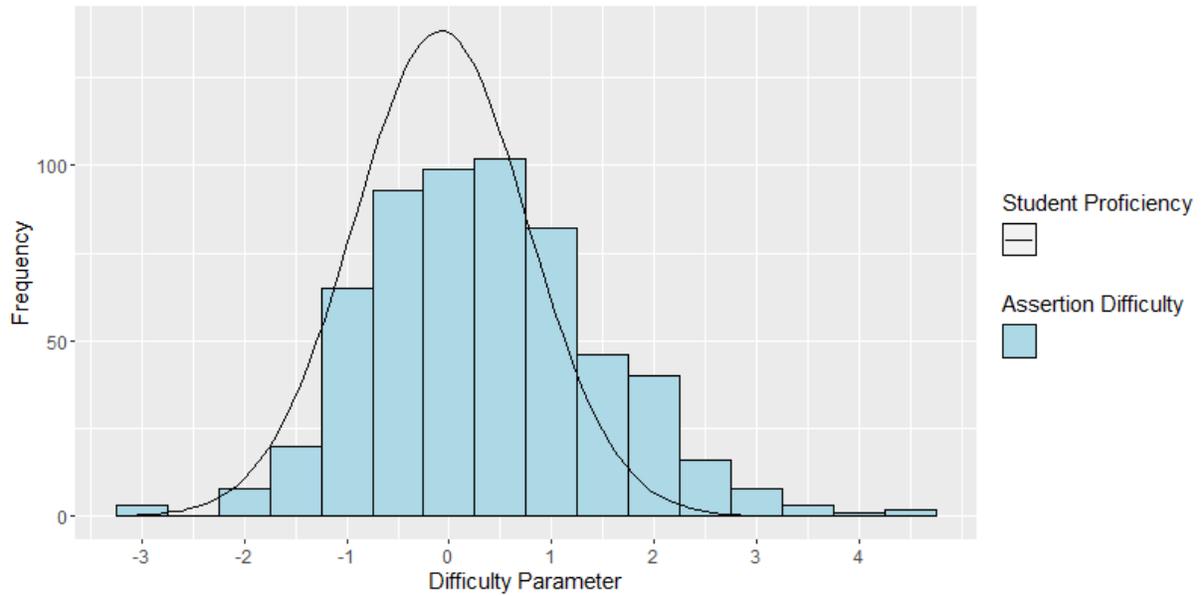
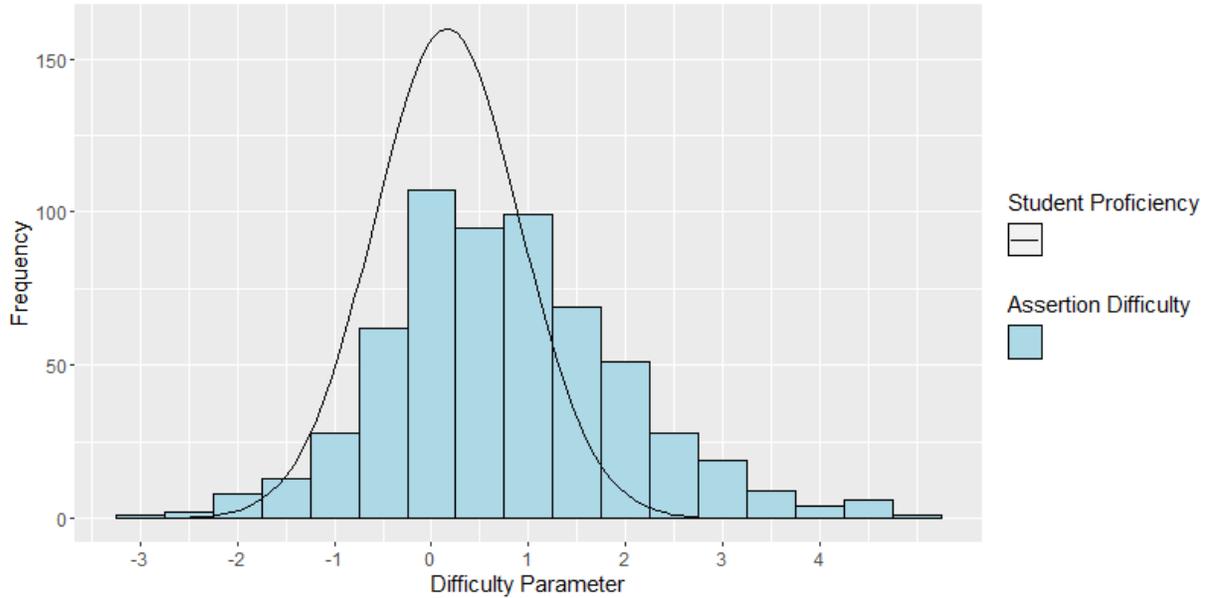


Figure 4. SDSA Assertion Difficulty and Student Proficiency Distributions, Grade 11



6. SCORING

6.1 MARGINAL MAXIMUM LIKELIHOOD FUNCTION

Student scores are obtained by marginalizing out the nuisance dimensions \mathbf{u}_j from the likelihood of the observed response pattern \mathbf{z}_j for student j ,

$$\ell_i(\theta_j) = \log \int_{\mathbf{u}_j} P(\mathbf{z}_j | \theta_j, \mathbf{u}_j) N(\mathbf{u}_j | \mathbf{0}, \boldsymbol{\Sigma}) d\mathbf{u}_j,$$

and maximizing this marginalized likelihood function for θ_j . The marginal maximum likelihood estimation (MMLE) estimator is a hybrid between the expected a posteriori (EAP) estimator (by marginalizing out the nuisance dimensions) and the maximum likelihood estimation MLE estimator (by maximizing the resulting marginal likelihood for θ). The marginal likelihood is maximized with respect to θ using the Newton Raphson method. See Rijmen, Jiang, and Turhan (2018) for more details of the MMLE estimator and the validation study by Connecticut State Department of Education (2019) for the use of this estimator.

The proposed model reduces to the unidimensional Rasch model when the nuisance variances are zero for all g . Likewise, the proposed MMLE is equivalent to the MLE of the unidimensional Rasch model when all the nuisance variances are zero. This can be shown by using the variable transformation $\mathbf{v} = \boldsymbol{\Sigma}^{-\frac{1}{2}}\mathbf{u}$. Then we have

$$\int_{\mathbf{u}_j} P(\mathbf{z}_j | \theta_j, \mathbf{u}_j) N(\mathbf{u}_j | \mathbf{0}, \boldsymbol{\Sigma}) d\mathbf{u}_j = \int_{\mathbf{v}_j} P(\mathbf{z}_j | \theta_j, \boldsymbol{\Sigma}^{\frac{1}{2}}\mathbf{v}_j) N(\mathbf{v}_j | \mathbf{0}, \mathbf{I}) d\mathbf{v}_j.$$

If $\sigma_{u_g}^2 = 0$ for all g , then

$$\int_{\mathbf{u}_j} P(\mathbf{z}_j|\theta_j, \mathbf{u}_j) N(\mathbf{u}_j|\mathbf{0}, \boldsymbol{\Sigma}) d\mathbf{u}_j = P(\mathbf{z}_j|\theta_j),$$

which is the likelihood under the unidimensional Rasch model.

6.2 DERIVATIVE

The marginal log likelihood function based on the item response theory (IRT) model with one overall dimension and one nuisance dimension for each grouping of assertions can be written as

$$l(\theta) = \sum_{i \in \text{SA}} \log(P(z_i|\theta)) + \sum_{g=1}^G \log \left\{ \int \text{Exp} \left[\sum_{i \in g} \log(P(z_{ig}|\theta, u_g)) \right] N(u_g|0, \sigma_{u_g}^2) du_g \right\}.$$

The first derivative of the marginal log likelihood function with respect to θ is

$$\begin{aligned} & \frac{dl(\theta)}{d\theta} \\ &= \sum_{i \in \text{SA}} \frac{\frac{dP(z_i|\theta)}{d\theta}}{P(z_i|\theta)} \\ & \quad + \sum_{g=1}^G \frac{\int \left\{ \text{Exp} \left[\sum_{i \in g} \log(P(z_{ig}|\theta, u_g)) \right] \left(\sum_{i \in g} \frac{\frac{dP(z_{ig}|\theta, u_g)}{d\theta}}{P(z_{ig}|\theta, u_g)} \right) N(u_g|0, \sigma_{u_g}^2) \right\} du_g}{\int \left\{ \text{Exp} \left[\sum_{i \in g} \log(P(z_{ig}|\theta, u_g)) \right] N(u_g|0, \sigma_{u_g}^2) \right\} du_g} \end{aligned}$$

and the second derivative of the marginal log likelihood function with respect to θ is

$$\begin{aligned}
 & \frac{d^2 l(\theta)}{d\theta^2} \\
 &= \sum_{i \in \text{SA}} \left[\frac{d^2 P(z_i|\theta)}{d\theta^2} - \left(\frac{d P(z_i|\theta)}{d\theta} \right)^2 \right] \\
 &+ \sum_{g=1}^G \frac{\int \text{Exp} \left[\sum_{i \in g} \log \left(P(z_{ig}|\theta, u_g) \right) \right] \left(\sum_{i \in g} \frac{d P(z_{ig}|\theta, u_g)}{d\theta} \right)^2 N(u_g|0, \sigma_{u_g}^2) du_g}{\int \left\{ \text{Exp} \left[\sum_{i \in g} \log \left(P(z_{ig}|\theta, u_g) \right) \right] N(u_g|0, \sigma_{u_g}^2) \right\} du_g} \\
 &+ \sum_{g=1}^G \frac{\int \text{Exp} \left[\sum_{i \in g} \log \left(P(z_{ig}|\theta, u_g) \right) \right] \left(\sum_{i \in g} \left[\frac{d^2 P(z_{ig}|\theta, u_g)}{d\theta^2} - \left(\frac{d P(z_{ig}|\theta, u_g)}{d\theta} \right)^2 \right] \right) N(u_g|0, \sigma_{u_g}^2) du_g}{\int \left\{ \text{Exp} \left[\sum_{i \in g} \log \left(P(z_{ig}|\theta, u_g) \right) \right] N(u_g|0, \sigma_{u_g}^2) \right\} du_g} \\
 &- \sum_{g=1}^G \left\{ \frac{\int \text{Exp} \left[\sum_{i \in g} \log \left(P(z_{ig}|\theta, u_g) \right) \right] \left(\sum_{i \in g} \frac{d P(z_{ig}|\theta, u_g)}{d\theta} \right)^2 N(u_g|0, \sigma_{u_g}^2) du_g}{\int \left\{ \text{Exp} \left[\sum_{i \in g} \log \left(P(z_{ig}|\theta, u_g) \right) \right] N(u_g|0, \sigma_{u_g}^2) \right\} du_g} \right\}^2.
 \end{aligned}$$

Based on the above equations, we need to define only the ratios of the first and second derivatives of the item response probabilities with respect to θ to the response probabilities. For the Rasch testlet model, these are obtained as

$$p_i = P(z_i = 1|\theta) = \frac{\text{Exp}(\theta - b_i)}{1 + \text{Exp}(\theta - b_i)}, q_i = P(z_i = 0|\theta) = 1 - p_i,$$

and

$$p_{ig} = P(z_{ig} = 1|\theta, u_g) = \frac{\text{Exp}(\theta + u_g - b_i)}{1 + \text{Exp}(\theta + u_g - b_i)}, q_{ig} = P(z_{ig} = 0|\theta, u_g) = 1 - p_{ig}.$$

Therefore, we have,

$$\begin{aligned}
 \frac{\frac{dp_i}{d\theta}}{p_i} &= q_i, & \frac{\frac{dq_i}{d\theta}}{q_i} &= -p_i, \\
 \frac{\frac{dp_{ig}}{d\theta}}{p_{ig}} &= q_{ig}, & \frac{\frac{dq_{ig}}{d\theta}}{q_{ig}} &= -p_{ig},
 \end{aligned}$$

$$\frac{\frac{d^2 p_i}{d\theta^2}}{p_i} - \left(\frac{\frac{d p_i}{d\theta}}{p_i}\right)^2 = -p_i q_i,$$

$$\frac{\frac{d^2 q_i}{d\theta^2}}{q_i} - \left(\frac{\frac{d q_i}{d\theta}}{q_i}\right)^2 = -p_i q_i,$$

$$\frac{\frac{d^2 p_{ig}}{d\theta^2}}{p_{ig}} - \left(\frac{\frac{d p_{ig}}{d\theta}}{p_{ig}}\right)^2 = -p_{ig} q_{ig}, \text{ and}$$

$$\frac{\frac{d^2 q_{ig}}{d\theta^2}}{q_{ig}} - \left(\frac{\frac{d q_{ig}}{d\theta}}{q_{ig}}\right)^2 = -p_{ig} q_{ig}.$$

6.3 EXTREME CASE HANDLING

As with the MLE, the MMLE is not defined for zero and perfect scores. These cases are handled by assigning the lowest obtainable theta (LOT) scores and highest obtainable theta (HOT) scores, respectively. Table 19 contains the LOT and HOT values for each grade.

6.4 STANDARD ERROR OF MEASUREMENT

The standard error of measurement (SEM) of the MMLE score estimate is:

$$SEM(\hat{\theta}_{MMLE}) = \frac{1}{\sqrt{I(\hat{\theta}_{MMLE})}}$$

where $I(\hat{\theta}_{MMLE})$ is the observed information evaluated at $\hat{\theta}_{MMLE}$. The observed information is calculated as $I(\theta^2) = -\frac{d^2 l(\theta)}{d\theta^2}$, where $\frac{d^2 l(\theta)}{d\theta^2}$ is defined in Section 6.2, Derivative. Note that the calculation of the SEM depends on the unique set of items that each student answers and their estimate of θ . Different students have different SEM values, even if they have the same raw score and/or theta estimate. Standard errors are truncated at 1 for the overall science scores and truncated at 1.4 for the discipline scores.

Standard errors for MMLE estimates truncated at the LOT and HOT are computed by evaluating the observed information at the MMLE before truncation. For all incorrect or all correct answers, the reported SEM is set at the truncation value for the standard error.

6.5 SCORING INCOMPLETE TESTS

The SDSA is assembled on-the-fly using an adaptive testing design. For science, a test is considered “attempted” if a student responded to at least one item (cluster or stand-alone). An attempted test is considered complete if the student responds to all the operational items. Otherwise, the test is “incomplete.”

Tests that are attempted but incomplete receive overall science scores. In order to receive a discipline score (e.g., Life Sciences, Physical Sciences, Earth and Space Sciences), a student must have attempted at least one item in the corresponding discipline of the test. The MMLE is used to score the attempted incomplete tests, counting unanswered items as incorrect. If the identities of the unanswered items are unknown due to the test being assembled on-the-fly, the item parameters for a “typical” item are used. If a missing item is an item cluster, the simulated item parameters of the missing item are the item parameters of item cluster 442 for grade 5, item cluster 324 for grade 8, and item cluster 390 for grade 11, which are operational item clusters that are typical for the Independent College and Career Readiness (ICCR) item bank in terms of the number of assertions and estimated parameters. Likewise, if a missing item is a stand-alone item, the simulated item parameters of the missing item are the item parameters of stand-alone item 371 for grade 5, item 354 for grade 8, and item 401 for grade 11, which are operational stand-alone items that are typical for the ICCR bank in terms of the number of assertions and estimated parameters.

If the identities of items that have not been answered are known because they have already been lined up through the pre-fetch process, the item parameters of the lined-up items are used. Similarly, for the accommodated forms that are fixed forms, the item parameters of the unanswered items on the form are used.

6.6 STUDENT-LEVEL SCALE SCORE

At the student level, scale scores are computed for

1. Overall Science
2. Life Sciences
3. Physical Sciences
4. Earth and Space Sciences

Scores are computed using the MMLE method outlined in this report, with all items from overall science or only items within the given discipline. Scores are truncated on the “theta” scale at the LOT and HOT values specified in Table 20, which correspond to values of the estimated mean minus/plus four times the estimated standard deviation of θ .

The reporting scales are a linear transformation of the theta scales

$$SS = a * \hat{\theta}_{MMLE} + b,$$

where a and b are the slope and intercept of the linear transformation that transforms $\hat{\theta}_{MMLE}$ to the reporting scale (refer to Table 20). The SEM for the estimated scale score is obtained as

$$SEM_{SS} = a * SEM_{\hat{\theta}_{MMLE}}.$$

In 2021, the slope a and intercept b were chosen so that the center of the reporting scale of each grade (500, 800, and 1100, respectively) is at the grade mean of the 2021 base-year and has a standard deviation of 25. Furthermore, for each grade, the reporting scale ranges from the base-year mean minus 4 times the standard deviation to the base-year mean plus 4 times the standard deviation. Specifically, for grade 5, the slope and intercept were obtained as

$$\begin{aligned}
 SS &= 25\theta^* + 500 \\
 &= 25 \frac{\theta - \hat{\mu}_\theta}{\hat{\sigma}_\theta} + 500 \\
 &= \frac{25}{\hat{\sigma}_\theta} \theta + \left(500 - \frac{25\hat{\mu}_\theta}{\hat{\sigma}_\theta} \right),
 \end{aligned}$$

where the second line stems from standardizing theta, $\theta^* = \frac{\theta - \hat{\mu}_\theta}{\hat{\sigma}_\theta}$. For grades 8 and 11, the slope and intercept can also be derived similarly.

Per grade, Table 20 presents the intercept, slope, LOT, HOT, lowest obtainable scale score (LOSS), and highest obtainable scale score (HOSS) values used for the 2021 reporting scale. The scale score distribution is reported for overall science in Appendix 1-D, Distribution of Scale Scores and Achievement Levels. The scale score distribution is reported for the science disciplines in Appendix 1-E, Distribution of Scale Scores by Science Discipline.

Table 20. Science Reporting Scale Linear Transformation Constants, Theta, and Corresponding Scaled-Score Limits for Extreme Ability Estimates (for 2021 θ Scale)

Grade	Slope (a)	Intercept (b)	Lowest Obtainable Theta (LOT)	Highest Obtainable Theta (HOT)	Lowest Obtainable Scale Score (LOSS)	Highest Obtainable Scale Score (HOSS)
5	28.347	502.475	-3.61	3.44	400	600
8	36.434	801.964	-2.79	2.69	700	900
11	35.552	1095.576	-2.68	2.93	1000	1200

6.7 RULES FOR CALCULATING ACHIEVEMENT LEVELS

Achievement levels and corresponding cut scores were set during standard setting in September of 2021 and were evaluated during a standards confirmation workshop in July of 2022. Students are classified into one of four achievement levels, based on their total score. The distribution of achievement levels is summarized in Appendix 1-D, Distribution of Scale Scores and Achievement Levels. Further, the distribution of scale scores and achievement levels for subgroups described in Section 4.4, Differential Item Functioning Analysis, are presented in Appendix 1-F, Distribution of Scale Scores and Achievement Levels by Subgroup. Following the guidance for statewide longitudinal data systems (National Center for Education Statistics, 2010), results for groups with fewer than 10 students are not provided.

Table 21 lists the cut scores on the reporting scale metrics for each grade.

Table 21. Achievement-Level Cut Scores

Grade	Cut 1	Cut 2	Cut 3
5	477	508	527
8	773	810	836
11	1073	1102	1134

6.7.1 Strengths and Weaknesses for Disciplines Relative to Proficiency Cut Score

Discipline-level classifications are computed to classify student achievement levels for each of the science disciplines. The following are the classification rules:

- If $(\hat{\theta}_{discipline} < \theta_{proficient} - 1.5 * SEM(\hat{\theta}_{discipline}))$, then achievement is classified as *Below Standard*.
- If $(\theta_{proficient} - 1.5 * SEM(\hat{\theta}_{discipline}) \leq \hat{\theta}_{discipline} < \theta_{proficient} + 1.5 * SEM(\hat{\theta}_{discipline}))$, then achievement is classified as *At/Near Standard*.
- If $(\hat{\theta}_{discipline} \geq \theta_{proficient} + 1.5 * SEM(\hat{\theta}_{discipline}))$, then achievement is classified as *Above Standard*.

Note that $\theta_{proficient}$ is the proficiency cut score of the overall test. Standard errors are truncated at 1.4. The LOT is always classified as *Below Standard*, and the HOT is always classified as *Above Standard*.

6.8 DISCIPLINARY CORE IDEA-LEVEL REPORTING

6.8.1 Relative to Overall Achievement

For aggregated units (i.e., classrooms, schools, and districts), there is reporting at levels below the science discipline level. In 2022–2023, reports were provided at the level of Disciplinary Core Ideas (DCI). The method for reporting at levels below the science discipline level is based on the use of residuals. The equations are presented for DCIs.

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

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

The expected score is computed for a student’s estimated overall ability. For the assertions clustered within an item, the expected score is marginalized over the nuisance dimensions for the assertions clustered within an item,

$$E(z_{ijg} = 1; \theta_{j,overall}, \boldsymbol{\tau}_i) = \int P(z_{ijg} = 1|u_{jg}; \theta_{j,overall}, \boldsymbol{\tau}_i)N(u_{jg})du_{jg},$$

where $\boldsymbol{\tau}_i$ is the vector of parameters for assertion i (e.g., for the Rasch testlet model, $\boldsymbol{\tau}_i = b_i$), and $P(z_{ijg} = 1|u_{jg}; \theta_{j,overall}, \boldsymbol{\tau}_i)$ is defined in Section 6.2, Derivative. Next, residuals are aggregated over assertions within each students,

$$\delta_{jDCI} = \frac{\sum_{i \in DCI} \delta_{ij}}{n_{jDCI}},$$

and over students of the group on which is reported,

$$\bar{\delta}_{DCIm} = \frac{1}{n_m} \sum_{j \in m} \delta_{jDCI},$$

where n_{jDCI} is the number of assertions related to the DCI for student j , and n_m is the number of students in a group assessed on the DCI. If a student did not see any items on a DCI, the student is not included in the n_m count for the aggregate. The standard error of the average residual is computed as

$$SEM(\bar{\delta}_{DCIm}) = \sqrt{\frac{1}{n_m(n_m-1)} \sum_{j \in m} (\delta_{jDCI} - \bar{\delta}_{DCIm})^2}.$$

A statistically significant difference from zero in these aggregates is evidence that a class, teacher, school, or district is more effective (if $\bar{\delta}_{DCIm}$ is positive) or less effective (negative $\bar{\delta}_{DCIm}$) in teaching a given DCI.

We do not suggest direct reporting of the statistic $\bar{\delta}_{DCIm}$; instead, we recommend reporting in the aggregate whether a group of students perform better, worse, or as expected on this DCI. It will also be indicated that, in some cases, sufficient information is not available.

For target-level strengths/weakness, the following is reported:

- If $\bar{\delta}_{DCIm} \leq -1.5 * SEM(\bar{\delta}_{DCIm})$, then achievement is *worse than* on the overall test.
- If $\bar{\delta}_{DCIm} \geq 1.5 * SEM(\bar{\delta}_{DCIm})$, then achievement is *better than* on the overall test.
- Otherwise, achievement is *similar to* on the overall test.
- If $SEM(\bar{\delta}_{DCIm}) > 0.2$, data are insufficient.

6.8.2 Relative to Proficiency Cut Score

DCI-level scores for aggregated units can be computed using the same method as outlined in Section 6.8.1, Relative to Overall Achievement, but with the expected score computed at the theta value corresponding to the proficiency cut score:

$$E(z_{ijg} = 1; \theta_{proficiency}, \boldsymbol{\tau}_i) = \int P(z_{ijg} = 1|u_{jg}; \theta_{proficiency}, \boldsymbol{\tau}_i)N(u_{jg})du_{jg}.$$

The following is reported for DCIs for aggregate units:

- If $\bar{\delta}_{DCIm} \leq -1.5 * SEM(\bar{\delta}_{DCIm})$, then achievement is *below* the proficiency cut score.

- If $\bar{\delta}_{DCIm} \geq 1.5 * SEM(\bar{\delta}_{DCIm})$, then achievement is *above* the proficiency cut score.
- Otherwise, achievement is *near* the proficiency cut score.
- If $SEM(\bar{\delta}_{DCIm}) > 0.2$, data are insufficient.

7. QUALITY CONTROL PROCEDURES

CAI’s quality assurance (QA) procedures are built on two key principles: (1) automation and (2) replication. Certain procedures can be automated, which removes the potential for human error. Procedures that cannot be reasonably automated are replicated by two independent analysts at CAI.

Although the quality of any test is monitored as an ongoing activity, several sources of CAI’s quality control system are described here. First, QA reports are routinely generated and evaluated throughout the testing window to ensure that each test performs as anticipated. Second, the score quality is ensured by employing a second independent scoring verification system.

7.1 QUALITY ASSURANCE REPORTS

Test monitoring occurs while tests are administered in a live environment to ensure that item behavior is consistent with expectations. This is accomplished using CAI’s Quality Monitoring System that yields item statistics, blueprint match rates, and item exposure rate reports.

7.1.1 Item Analysis

The item analysis report is 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. To examine the performance of test items, this report generates classical item analysis indicators of difficulty and discrimination, including proportion correct, biserial/polyserial correlation, and item fit statistics based on the item response theory (IRT). The report is configurable and can be produced to flag only items with statistics falling outside a specified range or to generate reports based on all items in the pool. For science, statistics reports at the assertion level (which are the units of analysis for science) are not yet available; however, as a routine and continuing practice, CAI psychometricians compute and monitor classical item statistics at the end of each testing window.

7.1.2 Blueprint Match

The QA system generates Blueprint Match reports at the content-standards level and for other content requirements, such as strand and affinity group for science. For each blueprint element, the report indicates the minimum and maximum number of items specified in the blueprint, the number of test administrations in which those specifications were met, the number of administrations in which the blueprint requirements were not met, and, for administrations in which specifications were not met, the number of items by which the requirement was not met.

For all three grades, every test met the blueprint specifications at all levels. Blueprint match is discussed in detail in this technical report in Volume 2, Test Development, for both simulated and operational test administrations.

7.1.3 Item Exposure Rates

The QA system also generates item exposure reports that allow test items to be monitored for unexpectedly large exposure rates or unusually low item-pool usage throughout the testing window. As with other reports, it is possible to examine the exposure rate for all items or flagged items with exposure rates that exceed an acceptable range. Often, item overexposure indicates a blueprint element or combination of blueprint elements that are underrepresented in the item pool and should be targeted for future item development. Such item overexposure is also usually anticipated in the simulation studies used to configure the adaptive algorithm. A total of 13.99% of the items in grade 5, 19.86% of the items in grade 8, and 25.93% of the items in grade 11 were administered to 20% or more online English version test takers at that grade. More details are discussed in Volume 2, Test Development, of this technical report, of this technical report.

7.1.4 Cheating Detection Analysis

The CAI QA system can also provide a forensics report to identify possible irregularities in test administration for further investigation. Unusual patterns of responding at the student level can be aggregated to the test session, test administrator, and school levels to identify possible group-level testing anomalies. CAI psychometricians can monitor testing anomalies throughout the testing window. Evidence can be evaluated, including changes in test scores across administrations, item response times, and item response patterns using the person-fit index. The flagging criteria used for these analyses are configurable and can be changed by the user. The analyses used to detect the testing anomalies can be run anytime within the testing window.

7.2 SCORING QUALITY CHECK

All student test scores are produced using CAI’s scoring engine. Before releasing any scores, a second score verification system is used to verify that all test scores match with 100% agreement in all tested grades. The second system is independently constructed and maintained from the main scoring engine and estimates scores separately using the procedures described within this report.

8. REFERENCES

- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for educational and psychological testing*. Washington, DC: Author.
- Cai, L. (2017). flexMIRT®: Flexible multilevel multidimensional item analysis and test scoring (version 3.51) [Computer software]. Chapel Hill, NC: Vector Psychometric Group.
- Connecticut State Department of Education. (2019). *Validating American Institutes for Research's calibration and scoring processes for science assessments* (Research Report). Hartford, CT: Author.
- Dorans, N. J., & Schmitt, A. P. (1991). Constructed response and differential item functioning: A pragmatic approach (ETS Research Report No. 91–47). Princeton, NJ: Educational Testing Service.
- Gibbons, R. D., & Hedeker, D. R. (1992). Full-information item bi-factor analysis. *Psychometrika*, 57(3), 423–436.
- Holland, P. W., & Thayer, D. T. (1988). Differential item functioning and the Mantel-Haenszel procedure. In H. Wainer & H. I. Braun (eds.), *Test validity* (pp. 129–145). Hillsdale, NJ: Erlbaum.
- Koretz, D., & Hamilton, L. S. (2006). Testing for accountability in K–12. In R. L. Brennan (Ed.), *Educational measurement* (4th ed., pp. 531–578). Westport, CT: American Council on Education/Praeger.
- National Center for Education Statistics. (2010). *Statistical methods for protecting personally identifiable information in aggregate reporting* (Statewide Longitudinal Data System Technical Brief, Brief 3). Retrieved from <https://nces.ed.gov/pubs2011/2011603.pdf>.
- National Research Council. (2012). *A framework for K–12 science education: Practices, crosscutting concepts, and core ideas*. Washington, DC: The National Academies Press.
- Rijmen, F. (2006). *BNL: A Matlab toolbox for Bayesian networks with logistic regression nodes* (Technical Report). Amsterdam: VU University Medical Center.
- Rijmen, F. (2010). Formal relations and empirical comparison among the bi-factor, the testlet, and a second-order multidimensional IRT model. *Journal of Educational Measurement*, 47(3), 361–372.
- Rijmen, F., Jiang, T., & Turhan, A. (2018, April). *An item response theory model for new science assessments*. Paper presented at the National Council on Measurement in Education, New York, NY.
- Rijmen, F., Liao, D., & Lin, Z. (2021). *The Rasch testlet model for the calibration of three-dimensional science assessments: A software comparison* [White paper]. Washington, DC: Cambium Assessment, Inc.

- Sireci, S. G., Wainer, H., & Thissen, D. (1991). On the reliability of testlet-based tests. *Journal of Educational Measurement, 28*, 237–247.
- Somes, G. W. (1986). The generalized Mantel Haenszel statistic. *The American Statistician, 40*, 106–108.
- Wang, W. C., & Wilson, M. (2005). The Rasch testlet model. *Applied Psychological Measurement, 29*, 126–149.
- Yen, W. (1993). Scaling performance assessments: Strategies for managing local item dependence. *Journal of Educational Measurement, 30*, 187–213.
- Zwick, R. (2012). *A review of ETS differential item functioning assessment procedures: Flagging rules, minimum sample size requirements, and criterion refinement* (ETS Research Report No. 12–08). Princeton, NJ: Educational Testing Service.

Appendix 1-A

Caveon Test Security Overview

Caveon Test Security Overview

TEST ADMINISTRATION SECURITY – CAVEON

The Cambium Assessment, Inc (CAI) utilizes the Caveon Web Patrol™ service to support test security compliance. Caveon is recognized as the only full-service test security organization that has national experience and expertise in this area. Caveon has been successfully providing Web Patrol monitoring services to influential clients since 2003 and has been delivering Web Patrol services on behalf of State Education Agencies since 2005. Caveon currently provides full-scope Web Patrol services in twenty-nine (29) states plus the WIDA consortium, the Smarter Balanced Assessment Consortium, and nearly fifty (50) certification and licensure programs.

By scouring the Internet and public-facing social media sites for breaches in test security, Caveon can systematically find and track threats to the testing program.

Web Patrol leverages the best of both automated technologies and the human capacity to judge and analyze. The result of this unique combination is a service that continually and systematically finds and tracks threats to the testing program.

DESCRIPTION

Caveon Web Patrol leverages technology tools and human expertise to identify, prioritize, and monitor websites, discussion forums, public social media platforms, etc., where sensitive test information may be disclosed or at risk of disclosure.

Patrolling efforts routinely find and evaluate “brain-dumps” (websites where test questions have been posted, supposedly by individuals who memorized them and/or where disclosed test content may be resold), test preparation training/education sites that may use actual (operational) test questions in the training, online auctions and classifieds such as eBay and Craigslist, and social media channels, forums and groups in which actual test items may be revealed or proxy test-takers offer their services. Real-time updates are generated in Caveon’s incident reporting platform, Caveon Core, that categorize identified incidents by level of actual or potential risk to the testing program based on the representations made on the websites, or actual analysis of the proffered content. Websites and Internet extracts are ranked from CLEARED (Lowest risk but should be monitored) to SEVERE (Highest risk). The reports contain specific URLs and other content extractions that represent and depict the categorized threat. Additionally, Caveon Core includes overall and specific threat analytics and actionable recommendations as well as any anticipated mitigation strategies from Caveon Web Patrol’s highly experienced team of Web Patrol analysts for the South Dakota Department of Education (SDDOE) to follow in minimizing and removing the dangers.

COMPREHENSIVE, CONSISTENT MONITORING

In conducting web patrol operations, Caveon utilizes a team of specialists who spend days and evenings continually trolling the Internet for intellectual property, the team leverages numerous

search technologies, some licensed and some publicly accessible (e.g., “Open Source”), to ensure comprehensive, consistent, and continual monitoring of the web.

VERIFYING AND MANAGING THREATS

Casting such a broad net across the web means the team must cull through thousands of search results (each is a possible threat) and dig deeper to explore whether a result is benign or a legitimate worry. Team members have, after years of service, become experts at quickly reviewing a search hit and discerning a level of risk. Despite technology innovations in other aspects of the service, this work requires human judgment and is vitally necessary to take action against real threats to test security.

Once a threat is verified, CAI and Caveon coordinate with SDDOE to systematically work through the steps necessary to have infringing content removed.

An escalation path of legal remedies is available. That path begins with formal takedown request letters leveraging the Digital Millennium Copyright Act (DMCA). The path ends when the website operators remove copyrighted material and/or cease operations, either voluntarily or by compulsion. CAI endeavors to complement existing activities of SDDOE, including issuing formal notices under existing U.S. copyright laws to offending website owners, ISPs, search engines, etc. Keys to successful threat removal include the following:

Timeliness of Notification

By continually, systematically patrolling for new threats and monitoring existing ones, Caveon Web Patrol quickly ascertains when a breach has or may occur. When a breach has been discovered, CAI will immediately notify the SDDOE.

Assistance Taking Down Material

Immediate notification of dangerous threats to the testing program is only half the solution. With direction and support from the SDDOE, CAI provides quick front line support through various means to take the next step, neutralizing the hazard. There are multiple options at SDDOE’s disposal to help protect its IP. CAI has experience with:

- DMCA Takedown Request Letters

DMCA Takedown Request Letters can be sent immediately to website operators upon threat detection by Caveon. In most cases, simply alerting operators that copyrighted materials may be published on their websites is enough to get it removed.

Caveon Web Patrol service begins one week prior to the opening of the administration window and continues for one week after the test administration window closes.

Appendix 1-B

Shared Science Assessment Item Bank: Field-Testing

TABLE OF CONTENTS

1.	2018 FIELD TESTS.....	1
2.	2019 FIELD TESTS.....	7
3.	2021 FIELD TESTS.....	14
4.	2022 FIELD TESTS.....	25

LIST OF TABLES

Table 1. Number of Field-Test Items Administered in Spring 2018	1
Table 2. Number of Common Elementary School Field-Test Items Administered and Calibrated in Spring 2018.....	3
Table 3. Number of Common Middle School Field-Test Items Administered and Calibrated in Spring 2018.....	4
Table 4. Number of Common High School Field-Test Items Administered and Calibrated in Spring 2018.....	5
Table 5. Overview of Field-Test Item Administration, Rubric Validation, and Item Data Review in Spring 2018.....	6
Table 6. Summary of Shared Science Assessment Item Bank in Spring 2018.....	7
Table 7. Number of Field-Test Items Administered in Spring 2019	8
Table 8. Number of Common Elementary School Field-Test Items Administered and Calibrated in Spring 2019.....	10
Table 9. Number of Common Middle School Field-Test Items Administered and Calibrated in Spring 2019.....	11
Table 10. Number of Common High School Field-Test Items Administered and Calibrated in Spring 2019.....	12
Table 11. Number of Field-Test Item Administration, Rubric Validation, and Item Data Review in Spring 2019.....	13
Table 12. Summary of Shared Science Assessment Item Bank in Spring 2019.....	14
Table 13. Number of Field-Test Items Administered in Spring 2021	15
Table 14. Number of Common Elementary School Field-Test Items Administered and Calibrated in Spring 2021.....	17
Table 15. Number of Common Middle School Field-Test Items Administered and Calibrated in Spring 2021	19
Table 16. Number of Common High School Field-Test Items Administered and Calibrated in Spring 2021	21
Table 17. Number of Field-Test Item Administration, Rubric Validation, and Item Data Review in Spring 2021.....	24
Table 18. Summary of Shared Science Assessment Item Bank in Spring 2021.....	25
Table 19. Number of Field-Test Items Administered in Spring 2022	26
Table 20. Number of Common Elementary School Field-Test Items Administered and Calibrated in Spring 2022.....	28
Table 21. Number of Common Middle School Field-Test Items Administered and Calibrated in Spring 2022.....	30
Table 22. Number of Common High School Field-Test Items Administered and Calibrated in Spring 2022.....	32
Table 23. Overview of Field-Test Item Administration, Rubric Validation, and Item Data Review in Spring 2022.....	35
Table 24. Summary of Shared Science Assessment Item Bank in Spring 2022.....	36

The Shared Science Assessment Item Bank is the product of a collaboration between Cambium Assessment, Inc. (CAI), multiple states and one US territory that share a Memorandum of Understanding (MOU). Every participant of the MOU contributes items to the Shared Science Assessment Item Bank that underwent the same development process. The portion of the bank contributed by CAI is part of the Independent College and Career Readiness (ICCR) bank. Starting in the school year of 2017–2018, items are field-tested every year. This appendix describes how field tests were conducted from 2017–2018 until 2021–2022 across the states relying on the Shared Science Assessment Item Bank, or the ICCR portion thereof, for their three-dimensional science assessments.

1. 2018 FIELD TESTS

In 2018, a large pool of items was field tested in nine states. For three states (Hawaii, Oregon, and Wyoming), unscored field-test items were added as an additional segment to the operational (scored) legacy science test. Two other states (Connecticut and Rhode Island) conducted an independent field test in which all students participated and were administered a full set of items, but no scores were reported. In the remaining four states that field-tested items from the Shared Science Assessment Item Bank (New Hampshire, Utah, Vermont, and West Virginia), an operational field test was administered, meaning tests consisted of scored field-test items. Items became operational and were scored after the test administration if they were not rejected during rubric validation or item data review, as described later in this section. In total, 340 item clusters and 205 stand-alone items were administered in the elementary, middle, and high school grade bands. Table 1 presents the number of item clusters and stand-alone items administered in each grade band for each state.

Table 1. Number of Field-Test Items Administered in Spring 2018

Grade Band and Item Type	CT	HI	MSSA ^a	NH	OR	UT	WV	WY	Total
Elementary School	135	24	69	58	26	–	91	14	153 (65)
Cluster	78	13	40	34	20	–	56	6	86 (34)
Stand-Alone	57	11	29	24	6	–	35	8	67 (31)
Middle School	174	27	56	55	28	98	123	17	241 (59)
Cluster	115	13	26	30	22	98	90	5	171 (31)
Stand-Alone	59	14	30	25	6	–	33	12	70 (28)
High School	149	23	75	60	38	–	–	14	151 (63)
Cluster	81	14	34	33	30	–	–	6	83 (34)
Stand-Alone	68	9	41	27	8	–	–	8	68 (29)
Total	458	74	200	173	92	98	214	45	545 (187)

Note. The numbers in parentheses indicate ICCR items.

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment.

For the states with a separate field-test segment (states with a legacy science test), and one of the states with an operational field test (Utah), fixed field-test forms were constructed (using a balanced incomplete design, except for Utah) and randomly assigned so that the group of students administered one form was comparable to the groups of students that were assigned other forms.

For the independent and operational field tests (except for Utah), items were administered using a linear on-the-fly (LOFT) test design in which items are selected on the fly, resulting in a unique test form for each student. The difference between the test design for the independent field tests and operational field tests depended on the test blueprint. The only blueprint constraint imposed on the independent field tests was that students received four stand-alone items and two item clusters for each of the three science disciplines. In contrast, a full blueprint was implemented for the states with an operational field test.

For any given state, there was a target of a minimum sample size of 1,500 students per item for any given state. Most items were administered in two or more states so that the item pools for all individual states were linked through common items. Approximately 98.3% of the items met or exceeded the target sample size of 1,500 in at least one state, while 98.8% of the items had a sample size of at least 1,350 (10% of the target) in at least one state. The common item design was used to calibrate all the items on a common science scale.

Table 2, Table 3, and Table 4 present the number of item clusters and stand-alone items that were in common between the item pools of any two states. The numbers below the shaded cells represent the number of common field-test items between any two states, and the numbers above the shaded cells represent the number of common items that survived rubric validation and were included in the 2018 calibration. In each of the shaded cells, the number outside the parentheses represents the number of unique items administered only in the given state, and the number provided in parentheses represents the number of unique and/or common items that were calibrated with only the data from that state. Table 2 presents the results for elementary school, Table 3 presents the results for middle school, and Table 4 presents the results for high school. The numbers at field-testing differ slightly from the numbers at calibration for a variety of reasons, such as items not passing rubric validation and versioning issues for some items in some states.

Table 2. Number of Common Elementary School Field-Test Items Administered and Calibrated in Spring 2018

	State	CT	HI	MSSA ^a	NH	OR	UT	WV	WY
Cluster	CT	3 (3)	9	36	28	16	–	49	6
	HI	10	0 (0)	7	8	5	–	12	1
	MSSA	36	8	0 (2)	15	12	–	26	2
	NH	30	8	17	1 (3)	5	–	22	2
	OR	17	5	13	5	1 (1)	–	5	1
	UT	–	–	–	–	–	–	–	–
	WV	49	12	27	25	5	–	0 (4)	2
	WY	6	1	2	2	1	–	2	0 (0)
Stand-Alone	CT	1 (3)	5	25	22	2	–	33	7
	HI	5	6 (6)	0	0	0	–	4	0
	MSSA	26	0	0 (1)	10	4	–	13	3
	NH	24	0	11	0 (2)	0	–	15	2
	OR	2	0	4	0	1 (1)	–	0	0
	UT	–	–	–	–	–	–	–	–
	WV	35	4	14	17	0	–	0 (2)	1
	WY	8	0	3	3	0	–	2	0 (1)
Total	CT	4 (6)	14	61	50	18	–	82	13
	HI	15	6 (6)	7	8	5	–	16	1
	MSSA	62	8	0 (3)	25	16	–	39	5
	NH	54	8	28	1 (5)	5	–	37	4
	OR	19	5	17	5	2 (2)	–	5	1
	UT	–	–	–	–	–	–	–	–
	WV	84	16	41	42	5	–	0 (6)	3
	WY	14	1	5	5	1	–	4	0 (1)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 3. Number of Common Middle School Field-Test Items Administered and Calibrated in Spring 2018

	State	CT	HI	MSSA ^a	NH	OR	UT	WV	WY
Cluster	CT	2 (6)	12	22	26	19	44	77	5
	HI	11	1 (0)	3	6	6	0	9	1
	MSSA	23	3	0 (1)	9	1	7	22	2
	NH	26	6	10	1 (2)	7	0	17	3
	OR	19	6	1	7	2 (2)	0	5	1
	UT	48	0	7	0	0	48 (52)	43	0
	WV	83	10	21	18	6	48	1 (9)	2
	WY	5	1	2	3	1	0	2	0 (0)
Stand-Alone	CT	2 (3)	6	27	25	3	0	33	12
	HI	6	8 (8)	2	0	0	0	2	0
	MSSA	27	2	0 (0)	18	3	0	20	2
	NH	25	0	18	0 (0)	0	0	21	3
	OR	3	0	3	0	0 (0)	0	0	0
	UT	0	0	0	0	0	0 (0)	0	0
	WV	33	2	20	21	0	0	0 (0)	2
	WY	12	0	2	3	0	0	2	0 (0)
Total	CT	4 (9)	18	49	51	22	44	110	17
	HI	17	9 (8)	5	6	6	0	11	1
	MSSA	50	5	0 (1)	27	4	7	42	4
	NH	51	6	28	1 (2)	7	0	38	6
	OR	22	6	4	7	2 (2)	0	5	1
	UT	48	0	7	0	0	48 (52)	43	0
	WV	116	12	41	39	6	48	1 (9)	4
	WY	17	1	4	6	1	0	4	0 (0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 4. Number of Common High School Field-Test Items Administered and Calibrated in Spring 2018

	State	CT	HI	MSSA ^a	NH	OR	UT	WV	WY
Cluster	CT	10 (16)	13	30	29	30	–	–	5
	HI	13	0 (0)	7	7	8	–	–	1
	MSSA	32	7	0 (2)	13	12	–	–	1
	NH	32	7	14	0 (3)	12	–	–	3
	OR	30	8	12	12	0 (0)	–	–	1
	UT	–	–	–	–	–	–	–	–
	WV	–	–	–	–	–	–	–	–
	WY	6	1	1	3	1	–	–	0 (1)
Stand-Alone	CT	4 (4)	9	40	27	8	–	–	8
	HI	9	0 (0)	4	0	0	–	–	0
	MSSA	39	4	0 (1)	20	3	–	–	1
	NH	25	0	20	0 (0)	0	–	–	1
	OR	8	0	3	0	0 (0)	–	–	0
	UT	–	–	–	–	–	–	–	–
	WV	–	–	–	–	–	–	–	–
	WY	7	0	1	1	0	–	–	0 (0)
Total	CT	14 (20)	22	70	56	38	–	–	13
	HI	22	0 (0)	11	7	8	–	–	1
	MSSA	71	11	0 (3)	33	15	–	–	2
	NH	57	7	34	0 (3)	12	–	–	4
	OR	38	8	15	12	0 (0)	–	–	1
	UT	–	–	–	–	–	–	–	–
	WV	–	–	–	–	–	–	–	–
	WY	13	1	2	4	1	–	–	0 (1)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Following the (operational) field-test administration, items went through rubric validation and item data review, as described in Volume 2, Section, 2.7.1, Rubric Validation, and Section 2.7.2, Data Review.

Table 5 presents the number of field-test items administered, the number of items rejected before or during rubric validation, the number of items submitted for data review, and the number of items rejected during data review. The numbers in parentheses indicate the ICCR owned items.

Table 5. Overview of Field-Test Item Administration, Rubric Validation, and Item Data Review in Spring 2018

Grade Band and Item Type	Number of Field-Test Items Administered	Number of Items Rejected Before/During Rubric Validation	Number of Items Submitted for Data Review	Number of Items Rejected at Data Review ^a	Number of Items Remaining
Elementary School	153 (65)	3 (0)	65 (26)	13 (3)	137 (62)
Cluster	86 (34)	3 (0)	24 (7)	5 (1)	78 (33)
Stand-Alone	67 (31)	0 (0)	41 (19)	8 (2)	59 (29)
Middle School	241 (59)	16 (0)	102 (26)	24 (3)	201 (56)
Cluster	171 (31)	12 (0)	65 (11)	15 (1)	144 (30)
Stand-Alone	70 (28)	4 (0)	37 (15)	9 (2)	57 (26)
High School	151 (63)	10 (2)	80 (31)	13 (2)	128 (59)
Cluster	83 (34)	8 (2)	35 (14)	4 (0)	71 (32)
Stand-Alone	68 (29)	2 (0)	45 (17)	9 (2)	57 (27)
Total	545 (187)	29 (2)	247 (83)	50 (8)	466 (177)

Note. The numbers in parentheses indicate ICCR items.

^aFigures in this column include three middle school clusters rejected after item data review.

Table 6 summarizes the operational Shared Science Assessment Item Bank for each of the three science disciplines after adding the 2018 field-test items that passed rubric validation and item data review. The numbers in parentheses indicate the ICCR owned items.

Table 6. Summary of Shared Science Assessment Item Bank in Spring 2018

Grade Band and Item Type	Science Discipline			Total ^a
	Earth and Space Sciences	Life Sciences	Physical Sciences	
Elementary School	41 (19)	47 (21)	49 (22)	137 (62)
Cluster	23 (11)	29 (11)	26 (11)	78 (33)
Stand-Alone	18 (8)	18 (10)	23 (11)	59 (29)
Middle School	56 (16)	72 (19)	70 (21)	198 (56)
Cluster	41 (9)	49 (7)	51 (14)	141 (30)
Stand-Alone	15 (7)	23 (12)	19 (7)	57 (26)
High School	37 (19)	53 (23)	38 (17)	128 (59)
Cluster	19 (8)	32 (15)	20 (9)	71 (32)
Stand-Alone	18 (11)	21 (8)	18 (8)	57 (27)
Total	134 (54)	172 (63)	157 (60)	463 (177)

Note. The numbers in parentheses indicate ICCR items.

^aTotals exclude three Utah-owned middle school item clusters aligned to Utah-specific standards.

2. 2019 FIELD TESTS

In 2019, a second wave of items was field tested in nine states. For three states (Hawaii, Idaho [elementary school only], and Wyoming), unscored field-test items were added as a separate segment to the operational scored legacy science test. An independent field test, in which students were administered a full set of items, was conducted for a sample of Idaho middle schools. In the remaining six states (Connecticut, New Hampshire, Oregon, Rhode Island, Vermont, and West Virginia), field-test items were administered as unscored items embedded among the operational items. In total, 123 item clusters and 224 stand-alone items were administered as field-test items in the elementary, middle, and high school grade bands. Table 7 presents the number of field-test item clusters and stand-alone items administered in each grade band for each state. The numbers in parentheses in the column Total indicate the ICCR owned items.

Table 7. Number of Field-Test Items Administered in Spring 2019

Grade Band and Item Type	CT	HI	ID	MSSA ^a	NH	OR	WV	WY	Total
Elementary School	47	31	53	42	18	27	18	16	117 (18)
Cluster	18	19	20	17	0	16	10	5	50 (0)
Stand-Alone	29	12	33	25	18	11	8	11	67 (18)
Middle School	56	23	53	46	28	26	26	15	127 (28)
Cluster	14	9	17	10	4	9	8	5	38 (4)
Stand-Alone	42	14	36	36	24	17	18	10	89 (24)
High School	69	21	–	37	29	28	–	25	103 (29)
Cluster	25	14	–	18	2	13	–	2	35 (2)
Stand-Alone	44	7	–	19	27	15	–	23	68 (27)
Total	172	75	106	125	75	81	44	56	347 (75)

Note. The numbers in parentheses indicate ICCR items.

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

For the three states with a separate field-test segment (i.e., states with a legacy science test), field-test forms were constructed using a balanced incomplete design and randomly assigned so that the group of students administered one form was comparable to the groups of students that were assigned other forms. For the independent field test, items were administered under a LOFT design, where the only blueprint constraint imposed was that students received four stand-alone items and two item clusters for each of the three science disciplines.

In the states with an operational test, field-test items were embedded within the test. Three states with an operational test—New Hampshire, Rhode Island, and Vermont—opted for a test in which operational items were grouped by science discipline. For these states, the field-test items were presented together in a fourth group of items. The sequence of the four sets of items (corresponding to the three science disciplines and a set of field-test items) was randomized across students. Three other states—Connecticut, Oregon, and West Virginia—opted for a test design in which the items were not grouped by discipline. In these three states, field-test items were administered at random positions throughout the test. A student received either a field-test item cluster or a set of five field-test stand-alone items. The test design for the SDSA is discussed in Section **Error! Reference source not found.** of this volume, Test Design.

A minimum sample size of 1,500 students per field-test item was targeted for any given state. Most items were administered in two or more states. Approximately 88.8% of the items met or exceeded the target sample size of 1,500 in at least one state, while 96.4% of the items had a sample size of at least 1,350 (10% of the target) in at least one state.

Table 8 to Table 10 present the number of item clusters and stand-alone items that were shared between the field-test pools of any two states. The numbers below the shaded cells represent the number of common field-test items between any two states, while the numbers above the shaded cells represent the number of common field-test items that survived rubric validation and were included in the calibration. In each of the shaded cells, the number outside the parentheses represents the number of unique field-test items administered only in the given state, and the number provided in parentheses represents the number of unique and/or common items that were calibrated with the data from that state only.

Table 8 presents the results for elementary schools, Table 9 presents the results for middle schools, and Table 10 presents the results for high schools. The numbers of field-test items administered differ slightly from the numbers of field-test items at calibration because some items were rejected during rubric validation.

Table 8. Number of Common Elementary School Field-Test Items Administered and Calibrated in Spring 2019

	State	CT	HI	ID	MSSA ^a	NH	OR	WV	WY
Cluster	CT	2 (2)	2	10	3	0	2	1	4
	HI	2	0 (0)	3	8	0	14	2	0
	ID	10	3	4 (4)	0	0	1	3	3
	MSSA	3	8	0	3 (3)	0	9	4	1
	NH	0	0	0	0	0 (0)	0	0	0
	OR	2	14	1	9	0	1 (1)	0	0
	WV	1	2	3	4	0	0	1 (0)	1
	WY	4	0	3	1	0	0	1	0 (0)
Stand-Alone	CT	5 (5)	1	13	1	9	0	0	2
	HI	1	0 (0)	10	6	0	6	0	0
	ID	13	11	1 (1)	12	1	9	2	4
	MSSA	1	7	13	3 (3)	5	8	5	6
	NH	9	0	1	5	2 (3)	0	0	6
	OR	0	7	10	9	0	1 (1)	0	0
	WV	0	0	2	5	0	0	1 (1)	0
	WY	2	0	4	6	7	0	0	0 (0)
Total	CT	7 (7)	3	23	4	9	2	1	6
	HI	3	0 (0)	13	14	0	20	2	0
	ID	23	14	5 (5)	12	1	10	5	7
	MSSA	4	15	13	6 (6)	5	17	9	7
	NH	9	0	1	5	2 (3)	0	0	6
	OR	2	21	11	18	0	2 (2)	0	0
	WV	1	2	5	9	0	0	2 (1)	1
	WY	6	0	7	7	7	0	1	0 (0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 9. Number of Common Middle School Field-Test Items Administered and Calibrated in Spring 2019

	State	CT	HI	ID	MSSA ^a	NH	OR	WV	WY
Cluster	CT	5 (5)	3	4	2	0	2	1	0
	HI	3	0 (0)	4	4	0	5	1	0
	ID	4	4	2 (2)	4	0	4	3	3
	MSSA	2	4	4	1 (1)	0	2	3	1
	NH	0	0	1	0	3 (0)	0	0	0
	OR	2	5	4	2	0	1 (1)	1	2
	WV	1	1	3	3	0	1	0 (0)	2
	WY	0	0	3	1	0	2	2	0 (0)
Stand-Alone	CT	10 (9)	2	13	9	10	3	6	0
	HI	2	0 (0)	9	9	0	6	3	0
	ID	13	9	2 (2)	11	1	12	6	5
	MSSA	9	9	11	1 (1)	6	11	9	7
	NH	10	0	2	6	3 (1)	0	0	2
	OR	3	6	12	11	0	0 (0)	2	7
	WV	6	3	6	9	1	2	0 (0)	0
	WY	0	0	5	7	2	7	0	0 (0)
Total	CT	15 (14)	5	17	11	10	5	7	0
	HI	5	0 (0)	13	13	0	11	4	0
	ID	17	13	4 (4)	15	1	16	9	8
	MSSA	11	13	15	2 (2)	6	13	12	8
	NH	10	0	3	6	6 (1)	0	0	2
	OR	5	11	16	13	0	1 (1)	3	9
	WV	7	4	9	12	1	3	0 (0)	2
	WY	0	0	8	8	2	9	2	0 (0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 10. Number of Common High School Field-Test Items Administered and Calibrated in Spring 2019

	State	CT	HI	ID	MSSA ^a	NH	OR	WV	WY
Cluster	CT	9 (9)	10	–	11	0	8	–	1
	HI	11	0 (0)	–	8	0	11	–	0
	ID	–	–	–	–	–	–	–	–
	MSSA	12	9	–	3 (2)	0	7	–	2
	NH	0	0	–	0	1 (0)	1	–	0
	OR	8	11	–	7	1	1 (1)	–	0
	WV	–	–	–	–	–	–	–	–
	WY	1	0	–	2	0	0	–	0 (0)
Stand-Alone	CT	14 (13)	7	–	7	6	13	–	13
	HI	7	0 (0)	–	0	0	6	–	0
	ID	–	–	–	–	–	–	–	–
	MSSA	8	0	–	3 (3)	6	5	–	12
	NH	8	0	–	6	10 (10)	0	–	7
	OR	14	6	–	6	0	0 (1)	–	8
	WV	–	–	–	–	–	–	–	–
	WY	14	0	–	13	7	9	–	0 (0)
Total	CT	23 (22)	17	–	18	6	21	–	14
	HI	18	0 (0)	–	8	0	17	–	0
	ID	–	–	–	–	–	–	–	–
	MSSA	20	9	–	6 (5)	6	12	–	14
	NH	8	0	–	6	11 (10)	1	–	7
	OR	22	17	–	13	1	1 (1)	–	8
	WV	–	–	–	–	–	–	–	–
	WY	15	0	–	15	7	9	–	0 (0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Following the (operational) administration, field-test items underwent rubric validation and item data review, as described in Volume 2, Section, 2.7.1, Rubric Validation, and Section 2.7.2, Data Review . Table 11 presents the number of field-test items administered in South Dakota, or another state, the number of items rejected before or during rubric validation, the number of items submitted for data review, and the number of items rejected during data review. The numbers in parentheses indicate the items owned by ICCR.

Table 11. Number of Field-Test Item Administration, Rubric Validation, and Item Data Review in Spring 2019

Grade Band and Item Type	Number of Field-Test Items Administered	Number of Items Rejected Before/During Rubric Validation	Number of Items Submitted for Data Review	Number of Items Rejected at Data Review	Number of Items Remaining ^a
Elementary School	117 (18)	2 (0)	72 (16)	24 (0)	91 (18)
Cluster	50 (0)	1 (0)	16 (0)	10 (0)	39 (0)
Stand-Alone	67 (18)	1 (0)	56 (16)	14 (0)	52 (18)
Middle School	127 (28)	6 (8)	66 (15)	21 (1)	97 (19)
Cluster	38 (4)	1 (4)	12 (0)	5 (0)	29 (0)
Stand-Alone	89 (24)	5 (4)	54 (15)	16 (1)	68 (19)
High School	103 (29)	6 (4)	52 (12)	15 (1)	80 (24)
Cluster	35 (2)	2 (2)	15 (0)	5 (0)	26 (0)
Stand-Alone	68 (27)	4 (2)	37 (12)	10 (1)	54 (24)
Total	347 (75)	14 (12)	190 (43)	60 (2)	268 (61)

Note: The numbers in parentheses indicate ICCR items.

^aNumber of items remaining excludes five AI-scored items (four ICCR and one MSSA-owned) field tested in spring 2019 that were not brought to item data review.

Table 12 summarizes the Shared Science Assessment Item Bank after adding the field-test items that were administered in 2019 and passed rubric validation and item data review. The numbers in parentheses indicate the ICCR owned items.

Table 12. Summary of Shared Science Assessment Item Bank in Spring 2019

Grade Band and Item Type	Science Discipline			Total
	<i>Earth and Space Sciences</i>	<i>Life Sciences</i>	<i>Physical Sciences</i>	
Elementary School	68 (23)	77 (28)	80 (28)	225 (79)
Cluster	33 (11)	40 (11)	40 (10)	113 (32)
Stand-Alone	35 (12)	37 (17)	40 (18)	112 (47)
Middle School	83 (23)	108 (26)	92 (20)	287 (69)
Cluster	44 (8)	62 (6)	53 (11)	163 (25)
Stand-Alone	39 (15)	46 (20)	39 (9)	124 (44)
High School	39 (17)	108 (46)	53 (16)	200 (79)
Cluster	18 (6)	48 (14)	24 (8)	90 (28)
Stand-Alone	21 (11)	60 (32)	29 (8)	110 (51)
Total	190 (63)	293 (100)	225 (64)	712 (227)

Note. The numbers in parentheses indicate ICCR items.

3. 2021 FIELD TESTS

In 2021, a third wave of items was field-tested in 12 states. For one state (Wyoming), unscored field-test items were added as a separate segment to the operational scored legacy science test. An independent field test, in which students were administered a full set of items, was conducted in Idaho and Montana. In the remaining nine states (Connecticut, Hawaii, New Hampshire, North Dakota, Rhode Island, South Dakota, Utah, Vermont, and West Virginia), field-test items were administered as unscored items embedded among the operational items. In total, 223 item clusters and 322 stand-alone items were administered as field-test items in the elementary, middle, and high school grade bands. Table 13 presents the number of field-test item clusters and stand-alone items administered in each grade band for each state. The numbers in parentheses indicate the ICCR owned field-test items.

Table 13. Number of Field-Test Items Administered in Spring 2021

Grade Band and Item Type	CT	HI	ID	MSSA ^a	MT	ND	NH	SD	UT	WV	WY	Total
Elementary School	36	22	140	55	21	11	19	8	54	19	17	214 (51)
Cluster	16	6	58	18	7	3	3	3	54	7	5	106 (9)
Stand-Alone	20	16	82	37	14	8	16	5	0	12	12	108 (42)
Middle School	33	19	129	54	20	11	18	11	45	19	20	159 (47)
Cluster	17	6	44	18	7	3	2	2	45	7	4	60 (8)
Stand-Alone	16	13	85	36	13	8	16	9	0	12	16	99 (39)
High School	49	17	156	49	–	11	12	8	–	–	20	172 (43)
Cluster	11	5	54	16	–	3	4	3	–	–	3	57 (15)
Stand-Alone	38	12	102	33	–	8	8	5	–	–	17	115 (28)
Total	118	58	425	158	41	33	49	27	99	38	57	545 (141)

Note. The numbers in parentheses indicate ICCR items.

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

For the state with a separate field-test segment (i.e., Wyoming), field-test forms were constructed using a balanced incomplete design and randomly assigned so that the group of students administered one form was comparable to the groups of students that were assigned other forms. For the independent field test, items were administered under a LOFT design, where the only blueprint constraint imposed was the discipline-level constraint for stand-alone items and item clusters.

For the states with an operational test, field-test items were embedded within the test. Three states with an operational test—New Hampshire, Rhode Island, and Vermont—chose to administer a test in which operational items were grouped by science discipline. For these states, the field-test items were presented together as a fourth group of items. The sequence of the four sets of items (corresponding to the three disciplines and a set of field-test items) was randomized across students. Six other states—Connecticut, Hawaii, North Dakota, South Dakota, Utah, and West Virginia—opted for a test design in which the items were not grouped by discipline. In these states, field-test items were administered at random positions throughout the test. A student received either a field-test item cluster or a set of four field-test stand-alone items. The test design for the SDSA is discussed in Section **Error! Reference source not found.** in this volume, **Error! Not a valid bookmark self-reference..**

A minimum sample size of 1,500 students per field-test item was targeted for any given state. Most items were administered in two or more states. Approximately 96.7% of the items met or exceeded the target sample size of 1,500 in at least one state, while 99.1% of the items had a sample size of at least 1,350 (10% of the target) in at least one state.

Table 14 to Table 16 present the number of item clusters and stand-alone items that were shared between the field-test pools of any two states. The numbers below the shaded cells represent the number of common field-test items between any two states, and the numbers above the shaded diagonal elements represent the number of common field-test items that survived rubric validation and were included in the calibration. In each of the shaded diagonal elements, the number outside the parentheses represents the number of unique field-test items administered only in the given state, and the number in parentheses represents the number of unique and/or common items that were calibrated with the data from that state only. Table 14 presents the results for elementary schools, Table 15 presents the results for middle schools, and Table 16 presents the results for high schools. The numbers of field-test items administered differ slightly from the numbers of field-test items at calibration because some items were rejected during rubric validation.

Table 14. Number of Common Elementary School Field-Test Items Administered and Calibrated in Spring 2021

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	SD	UT	WV	WY
Cluster	CT	3 (3)	0	13	0	0	0	0	0	0	0	0
	HI	0	1 (1)	3	0	0	0	0	0	0	1	0
	ID	13	4	3 (2)	5	5	2	0	2	20	1	4
	MSSA	0	0	6	2 (2)	2	0	0	0	7	0	0
	MT	0	0	5	2	0 (0)	0	0	0	0	0	0
	ND	0	0	2	0	0	0 (0)	0	1	0	1	0
	NH	0	0	0	0	0	0	0 (0)	0	0	3	0
	SD	0	0	2	0	0	1	0	0 (0)	0	2	0
	UT	0	0	20	8	0	0	0	0	25 (24)	0	2
	WV	0	1	1	0	0	1	3	2	0	1 (1)	0
WY	0	0	4	0	0	0	0	0	2	0	0 (0)	
Stand-Alone	CT	3 (3)	0	14	2	0	0	0	0	0	0	1
	HI	0	0 (0)	12	1	0	0	2	3	0	1	0
	ID	14	12	3 (3)	30	13	4	3	3	0	4	9
	MSSA	2	1	30	0 (0)	12	0	3	1	0	0	0
	MT	0	0	13	12	0 (0)	0	0	0	0	0	0
	ND	0	0	4	0	0	0 (0)	2	0	0	0	1
	NH	0	2	4	3	0	2	0 (0)	2	0	3	1
	SD	0	3	3	1	0	0	2	0 (0)	0	0	0
	UT	0	0	0	0	0	0	0	0	0 (0)	0	0
	WV	0	1	4	0	0	1	3	0	0	3 (3)	0
WY	1	0	9	0	0	1	1	0	0	0	0 (0)	
Total	CT	6 (6)	0	27	2	0	0	0	0	0	0	1
	HI	0	1 (1)	15	1	0	0	2	3	0	2	0

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	SD	UT	WV	WY
	ID	27	16	6 (5)	35	18	6	3	5	20	5	13
	MSSA	2	1	36	2 (2)	14	0	3	1	7	0	0
	MT	0	0	18	14	0 (0)	0	0	0	0	0	0
	ND	0	0	6	0	0	0 (0)	2	1	0	1	1
	NH	0	2	4	3	0	2	0 (0)	2	0	6	1
	SD	0	3	5	1	0	1	2	0 (0)	0	2	0
	UT	0	0	20	8	0	0	0	0	25 (24)	0	2
	WV	0	2	5	0	0	2	6	2	0	4 (4)	0
	WY	1	0	13	0	0	1	1	0	2	0	0 (0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 15. Number of Common Middle School Field-Test Items Administered and Calibrated in Spring 2021

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	SD	UT	WV	WY
Cluster	CT	0 (0)	0	9	2	0	0	0	0	10	0	0
	HI	0	0 (0)	2	3	0	0	0	0	3	1	0
	ID	11	2	1 (1)	10	6	2	1	1	31	0	4
	MSSA	4	3	11	0 (0)	0	2	0	0	9	1	1
	MT	0	0	6	0	1 (1)	0	1	1	4	0	0
	ND	0	0	3	2	0	0 (0)	0	0	2	0	0
	NH	0	0	1	0	1	0	0 (0)	1	0	1	0
	SD	0	0	1	0	1	0	1	0 (0)	0	0	0
	UT	14	3	36	11	4	3	0	1	0 (0)	2	2
	WV	0	1	1	1	0	0	1	1	5	0 (0)	0
WY	0	0	4	1	0	0	0	0	2	0	0 (0)	
Stand-Alone	CT	2 (2)	0	12	2	0	0	0	3	0	0	2
	HI	0	0 (0)	10	1	0	0	0	0	0	2	0
	ID	13	10	2 (2)	29	10	6	12	7	0	5	15
	MSSA	2	1	29	0 (0)	10	2	1	1	0	2	4
	MT	0	0	12	10	0 (0)	0	0	0	0	0	0
	ND	0	0	7	2	0	0 (0)	1	0	0	0	0
	NH	0	0	12	1	0	1	0 (0)	2	0	1	3
	SD	3	0	7	1	0	0	2	0 (0)	0	3	4
	UT	0	0	0	0	0	0	0	0	0 (0)	0	0
	WV	0	2	6	3	0	1	1	3	0	0 (0)	0
WY	2	0	15	4	0	0	3	4	0	0	0 (0)	
Total	CT	2 (2)	0	21	4	0	0	0	3	10	0	2
	HI	0	0 (0)	12	4	0	0	0	0	3	3	0

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	SD	UT	WV	WY
	ID	24	12	3 (3)	39	16	8	13	8	31	5	19
	MSSA	6	4	40	0 (0)	10	4	1	1	9	3	5
	MT	0	0	18	10	1 (1)	0	1	1	4	0	0
	ND	0	0	10	4	0	0 (0)	1	0	2	0	0
	NH	0	0	13	1	1	1	0 (0)	3	0	2	3
	SD	3	0	8	1	1	0	3	0 (0)	0	3	4
	UT	14	3	36	11	4	3	0	1	0 (0)	2	2
	WV	0	3	7	4	0	1	2	4	5	0 (0)	0
	WY	2	0	19	5	0	0	3	4	2	0	0 (0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 16. Number of Common High School Field-Test Items Administered and Calibrated in Spring 2021

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	SD	UT	WV	WY
Cluster	CT	1 (1)	0	8	0	0	0	0	0	0	0	0
	HI	0	0 (0)	5	0	0	0	0	0	0	0	0
	ID	10	5	16 (15)	12	0	2	2	3	0	0	3
	MSSA	0	0	15	0 (0)	0	0	1	2	0	0	0
	MT	0	0	0	0	0 (0)	0	0	0	0	0	0
	ND	0	0	2	0	0	0 (0)	1	0	0	0	0
	NH	0	0	2	1	0	1	0 (0)	0	0	0	0
	SD	0	0	3	2	0	0	0	0 (0)	0	0	0
	UT	0	0	0	0	0	0	0	0	0 (0)	0	0
	WV	0	0	0	0	0	0	0	0	0	0 (0)	0
WY	0	0	3	0	0	0	0	0	0	0	0 (0)	
Stand-Alone	CT	3 (3)	0	31	3	0	0	0	0	0	0	1
	HI	0	0 (0)	11	1	0	0	0	0	0	0	0
	ID	31	11	9 (8)	24	0	7	4	5	0	0	14
	MSSA	3	1	25	0 (0)	0	0	3	4	0	0	1
	MT	0	0	0	0	0 (0)	0	0	0	0	0	0
	ND	0	0	7	0	0	0 (0)	1	0	0	0	0
	NH	0	0	4	3	0	1	0 (0)	0	0	0	0
	SD	0	0	5	4	0	0	0	0 (0)	0	0	1
	UT	0	0	0	0	0	0	0	0	0 (0)	0	0
	WV	0	0	0	0	0	0	0	0	0	0 (0)	0
WY	1	0	15	1	0	0	0	0	1	0	0 (0)	
Total	CT	4 (4)	0	39	3	0	0	0	0	0	0	1
	HI	0	0 (0)	16	1	0	0	0	0	0	0	0

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	SD	UT	WV	WY
	ID	41	16	25 (23)	36	0	9	6	8	0	0	17
	MSSA	3	1	40	0 (0)	0	0	4	6	0	0	1
	MT	0	0	0	0	0 (0)	0	0	0	0	0	0
	ND	0	0	9	0	0	0 (0)	2	0	0	0	0
	NH	0	0	6	4	0	2	0 (0)	0	0	0	0
	SD	0	0	8	6	0	0	0	0 (0)	0	0	1
	UT	0	0	0	0	0	0	0	0	0 (0)	0	0
	WV	0	0	0	0	0	0	0	0	0	0 (0)	0
	WY	1	0	18	1	0	0	0	1	0	0	0 (0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Following the administration, field-test items went through rubric validation and item data review, as described in Volume 2, Section, 2.7.1, Rubric Validation, and Section 2.7.2, Data Review.

Table 17 presents the number of field-test items administered in South Dakota, or another state, the number of items rejected before or during rubric validation, the number of items sent for data review, and the number of items rejected during data review. The numbers in parentheses indicate the ICCR owned field-test items.

Table 17. Number of Field-Test Item Administration, Rubric Validation, and Item Data Review in Spring 2021

Grade Band and Item Type	Number of Field-Test Items Administered	Number of Items Rejected Before/During Rubric Validation	Number of Items Submitted for Data Review	Number of Items Rejected at Data Review	Number of Items Remaining ^a
Elementary School	214 (51)	7 (2)	100 (32)	19 (2)	188 (47)
Cluster	106 (9)	5 (0)	24 (3)	7 (0)	94 (9)
Stand-Alone	108 (42)	2 (2)	76 (29)	12 (2)	94 (38)
Middle School	159 (47)	15 (4)	87 (27)	13 (0)	129 (43)
Cluster	60 (8)	10 (3)	22 (2)	5 (0)	43 (5)
Stand-Alone	99 (39)	5 (1)	65 (25)	8 (0)	86 (38)
High School	172 (43)	9 (3)	94 (16)	22 (4)	141 (36)
Cluster	57 (15)	6 (2)	27 (6)	4 (2)	47 (11)
Stand-Alone	115 (28)	3 (1)	67 (10)	18 (2)	94 (25)
Total	545 (141)	31 (9)	281 (75)	54 (6)	458 (126)

Note. The numbers in parentheses indicate ICCR items.

^aTwo Hawaii-owned items were not shared to the Shared Science Assessment Item bank.

Table 18 summarizes the Shared Science Assessment Item Bank after addition of the field-test items that were administered in 2021 and passed rubric validation and item data review. The numbers in parentheses indicate the ICCR owned items.

Table 18. Summary of Shared Science Assessment Item Bank in Spring 2021

Grade Band and Item Type	Science Discipline			Total ^a
	<i>Earth and Space Sciences</i>	<i>Life Sciences</i>	<i>Physical Sciences</i>	
Elementary School	136 (42)	128 (41)	149 (43)	413 (126)
Cluster	65 (14)	66 (14)	76 (13)	207 (41)
Stand-Alone	71 (28)	62 (27)	73 (30)	206 (85)
Middle School	114 (32)	156 (46)	137 (34)	407 (112)
Cluster	55 (9)	76 (9)	67 (12)	198 (30)
Stand-Alone	59 (23)	80 (37)	70 (22)	209 (82)
High School	68 (21)	163 (66)	106 (27)	337 (114)
Cluster	27 (9)	64 (18)	42 (11)	133 (38)
Stand-Alone	41 (12)	99 (48)	64 (16)	204 (76)
Total	318 (95)	447 (153)	392 (104)	1157 (352)

Note. The numbers in parentheses indicate ICCR items.

^aTwo Hawaii -owned items were not shared to the Shared Science Assessment Item bank.

4. 2022 FIELD TESTS

In 2022, a fourth wave of items was field tested in 13 states and one U.S. territory. In all of these locations—Connecticut, Hawaii, Idaho, Montana, New Hampshire, North Dakota, Oregon, Rhode Island, South Dakota, the U.S. Virgin Islands, Utah, Vermont, West Virginia, and Wyoming—the field-test items were administered as unscored items embedded among the operational items. In total, 217 item clusters and 254 stand-alone items were administered as field-test items in the elementary, middle, and high school grade bands. Table 19 presents the number of field-test item clusters and stand-alone items administered in each grade band for each state and territory. The numbers in parentheses indicate the ICCR owned field-test items.

Table 19. Number of Field-Test Items Administered in Spring 2022

Grade Band and Item Type	CT	HI	ID	MSSA ^a	MT	ND	NH	OR	SD	VI	UT	WV	WY	Total ^b
Elementary School	34	28	22	66	12	12	17	41	10 (10)	1	62	19	10	170
Cluster	22	8	11	22	4	4	5	15	4 (4)	1	62	11	2	88
Stand-Alone	12	20	11	44	8	8	12	26	6 (6)	0	-	8	8	82
Middle School	40	30	35	64	12	12	17	39	10 (10)	1	76	33	10	190
Cluster	20	10	7	21	4	4	5	16	4 (4)	1	76	5	2	88
Stand-Alone	20	20	28	43	8	8	12	23	6 (6)	0	-	28	8	102
High School	46	14	14	58	-	12	16	43	9 (9)	1	-	-	10	111
Cluster	18	6	10	19	-	4	4	16	3 (3)	1	-	-	2	41
Stand-Alone	28	8	4	39	-	8	12	27	6 (6)	0	-	-	8	70
Total	120	72	71	188	24	36	50	123	29 (29)	3	138	52	30	471

Note. The numbers in parentheses indicate ICCR items.

^aMSSA = Rhode Island and Vermont's Multi-State Science Assessment

^bCount excludes seventy (sixty Oregon and ten South Dakota) legacy field-test items.

In the states where an operational test was administered in spring 2022, field-test items were embedded within the test. Three states with an operational test—New Hampshire, Rhode Island, and Vermont—chose to administer a test in which operational items were grouped by science discipline. For these states, the field-test items were presented together as a fourth group of items. The sequence of the four sets of items (corresponding to the three disciplines and a set of field-test items) was randomized across students. Ten other states and one U.S. territory—Connecticut, Hawaii, Idaho, Montana, North Dakota, Oregon, South Dakota, Utah, the U.S. Virgin Islands, West Virginia, and Wyoming—opted for a test design in which the items were not grouped by discipline. In these 10 states and one U.S. territory, field-test items were administered at random positions throughout the test. A student received either a field-test item cluster or a set of four field-test stand-alone items. The test design for the SDSA is discussed in Section **Error! Reference source not found.** of this volume, **Error! Not a valid bookmark self-reference..**

A minimum sample size of 1,500 students per field-test item was targeted for any given state or territory. Most items were administered in two states or territory. Approximately 61.6% of the items met or exceeded the target sample size of 1,500 in at least one state, while 88.0% of the items had a sample size of at least 1,350 (10% of the target) in at least one state. In addition, 98.3% of the items had a sample size of at least 1,200 in at least one state.

Table 20 to Table 22 present the number of item clusters and stand-alone items that were shared between the field-test pools of any two states or one state and one territory. The numbers below the shaded diagonal elements represent the number of common field-test items between any two states, and the numbers above the shaded cells represent the number of common field-test items that survived rubric validation and were included in the calibration. In each of the shaded cells, the number outside the parentheses represents the number of unique field-test items administered only in the given state or territory, and the number in parentheses represents the number of unique and/or common items that were calibrated with the data from that state only. Table 20 presents the results for elementary schools, Table 21 presents the results for middle schools, and Table 22 presents the results for high schools. The numbers of field-test items administered differ slightly from the numbers of field-test items at calibration because some items were rejected during rubric validation.

Table 20. Number of Common Elementary School Field-Test Items Administered and Calibrated in Spring 2022

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	OR	SD	UT	WV	WY	USVI
Cluster	CT	0 (0)	0	3	1	0	0	0	3	0	15	0	0	0
	HI	0	0(0)	0	0	0	0	0	5	0	2	0	0	0
	ID	3	0	0(0)	3	0	0	0	0	0	5	0	0	0
	MSSA	1	0	3	0(0)	0	0	0	5	1	12	0	0	0
	MT	0	0	0	0	0(0)	0	0	0	0	4	0	0	0
	ND	0	0	0	0	0	0(0)	4	0	0	0	0	0	0
	NH	0	0	0	0	0	4	0(0)	0	0	1	0	0	1
	OR	3	6	0	5	0	0	0	0(0)	0	1	0	0	0
	SD	0	0	0	1	0	0	0	0	0(0)	3	0	0	0
	UT	15	2	5	12	4	0	1	1	3	6 (6)	11	2	1
	WV	0	0	0	0	0	0	0	0	0	11	0(0)	0	0
	WY	0	0	0	0	0	0	0	0	0	2	0	0(0)	0
	USVI	0	0	0	0	0	0	0	1	0	0	0	0	0(0)
Stand-Alone	CT	0(0)	2	0	4	4	0	0	0	2	0	0	0	0
	HI	2	0(0)	3	7	0	0	0	8	0	0	0	0	0
	ID	0	3	0(0)	1	1	4	0	2	0	0	0	0	0
	MSSA	4	7	1	0(0)	3	0	1	7	4	0	8	8	0
	MT	4	0	1	3	0(0)	0	0	0	0	0	0	0	0
	ND	0	0	4	0	0	0(0)	3	1	0	0	0	0	0
	NH	0	0	0	1	0	3	1(0)	7	0	0	0	0	0
	OR	0	8	2	8	0	1	7	0(0)	0	0	0	0	0
	SD	2	0	0	4	0	0	0	0	0(0)	0	0	0	0
	UT	0	0	0	0	0	0	0	0	0	0(0)	0	0	0
	WV	0	0	0	8	0	0	0	0	0	0	0(0)	0	0
	WY	0	0	0	8	0	0	0	0	0	0	0	0(0)	0

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	OR	SD	UT	WV	WY	USVI
	USVI	0	0	0	0	0	0	0	0	0	0	0	0	0 (0)
Total	CT	0(0)	2	3	5	4	0	0	3	2	15	0	0	0
	HI	2	0(0)	3	7	0	0	0	13	0	2	0	0	0
	ID	3	3	0(0)	4	1	4	0	2	0	5	0	0	0
	MSSA	5	7	4	0(0)	3	0	1	12	5	12	8	8	0
	MT	4	0	1	3	0(0)	0	0	0	0	4	0	0	0
	ND	0	0	4	0	0	0(0)	7	1	0	0	0	0	0
	NH	0	0	0	1	0	7	1(0)	7	0	1	0	0	1
	OR	3	14	2	13	0	1	7	0(0)	0	1	0	0	0
	SD	2	0	0	5	0	0	0	0	0(0)	3	0	0	0
	UT	15	2	5	12	4	0	1	1	3	6(6)	11	2	1
	WV	0	0	0	8	0	0	0	0	0	11	0(0)	0	0
	WY	0	0	0	8	0	0	0	0	0	2	0	0(0)	0
	USVI	0	0	0	0	0	0	0	1	0	0	1	0	0 (0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 21. Number of Common Middle School Field-Test Items Administered and Calibrated in Spring 2022

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	OR	SD	UT	WV	WY	USVI
Cluster	CT	0(0)	1	1	0	0	0	0	1	0	17	0	0	0
	HI	1	0(0)	0	1	0	0	0	1	0	5	0	0	0
	ID	1	0	0(0)	0	0	0	0	0	1	5	0	0	0
	MSS A	0	1	0	0(0)	0	0	0	2	0	18	0	0	0
	MT	0	0	0	0	0(0)	0	0	0	0	4	0	0	0
	ND	0	0	0	0	0	0(0)	0	0	0	4	0	0	0
	NH	0	0	0	0	0	0	0(0)	3	0	2	0	0	0
	OR	1	2	0	2	0	0	3	0(0)	0	8	0	0	0
	SD	0	0	1	0	0	0	0	0	0(0)	2	0	0	1
	UT	17	6	5	18	4	4	2	8	3	2(2)	5	2	1
	WV	0	0	0	0	0	0	0	0	0	5	0(0)	0	0
	WY	0	0	0	0	0	0	0	0	0	2	0	0(0)	0
	USVI	0	0	0	0	0	0	0	0	0	1	1	0	0(0)
Stand-Alone	CT	0(0)	0	0	12	0	0	0	4	1	0	3	0	0
	HI	0	0(0)	8	5	0	0	0	6	0	0	1	0	0
	ID	0	8	0(0)	5	8	0	0	3	4	0	0	0	0
	MSS A	12	5	5	0(0)	0	0	0	4	0	0	9	8	0
	MT	0	0	8	0	0(0)	0	0	0	0	0	0	0	0
	ND	0	0	0	0	0	0(0)	0	0	0	0	8	0	0
	NH	0	0	0	0	0	0	0(0)	6	0	0	5	0	0
	OR	4	6	3	4	0	0	6	0(0)	0	0	0	0	0
	SD	1	0	4	0	0	0	0	0	0(0)	0	1	0	0
	UT	0	0	0	0	0	0	0	0	0	0(0)	0	0	0
	WV	3	1	0	9	0	8	6	0	1	0	0(0)	0	0

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	OR	SD	UT	WV	WY	USVI
	WY	0	0	0	8	0	0	0	0	0	0	0	0(0)	0
	USVI	0	0	0	0	0	0	0	0	0	0	0	0	0(0)
Total	CT	0(0)	1	1	12	0	0	0	5	1	17	3	0	0
	HI	1	0(0)	8	6	0	0	0	7	0	5	1	0	0
	ID	1	8	0(0)	5	8	0	0	3	5	5	0	0	0
	MSSA	12	6	5	0(0)	0	0	0	6	0	18	9	8	0
	MT	0	0	8	0	0(0)	0	0	0	0	4	0	0	0
	ND	0	0	0	0	0	0(0)	0	0	0	4	8	0	0
	NH	0	0	0	0	0	0	0(0)	9	0	2	5	0	0
	OR	5	8	3	6	0	0	9	0(0)	0	8	0	0	0
	SD	1	0	5	0	0	0	0	0	0(0)	2	1	0	1
	UT	17	6	5	18	4	4	2	8	3	2(2)	5	2	1
	WV	3	1	0	9	0	8	6	0	1	5	0(0)	0	0
	WY	0	0	0	8	0	0	0	0	0	2	0	0(0)	0
	USVI	0	0	0	0	0	0	0	0	0	1	1	0	0(0)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 22. Number of Common High School Field-Test Items Administered and Calibrated in Spring 2022

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	OR	SD	UT	WV	WY	USVI
Cluster	CT	0(0)	0	2	6	-	2	1	5	1	-	-	1	1
	HI	0	0(0)	3	0	-	0	0	2	0	-	-	0	0
	ID	2	3	0(0)	2	-	0	0	2	0	-	-	1	0
	MSSA	6	1	2	0(0)	-	2	1	4	2	-	-	0	0
	MT	-	-	-	-	-	-	-	-	-	-	-	-	-
	ND	2	0	0	2	-	0(0)	0	0	0	-	-	0	0
	NH	1	0	0	1	-	0	0(0)	2	0	-	-	0	0
	OR	5	2	2	5	-	0	2	0(0)	0	-	-	0	1
	SD	1	0	0	2	-	0	0	0	0(0)	-	-	0	0
	UT	-	-	-	-	-	-	-	-	-	-	-	-	-
	WV	-	-	-	-	-	-	-	-	-	-	-	-	-
	WY	1	0	1	0	-	0	0	0	0	-	-	0(0)	0
USVI	1	0	0	0	-	0	0	0	1	0	-	-	0	0(0)
Stand-Alone	CT	0(0)	0	1	19	-	6	0	1	1	-	-	0	0
	HI	0	0(0)	1	1	-	0	0	6	0	-	-	0	0
	ID	1	1	0(0)	1	-	0	0	1	0	-	-	0	0
	MSSA	19	1	1	0(0)	-	2	0	5	3	-	-	8	0
	MT	-	-	-	-	-	-	-	-	-	-	-	-	-
	ND	6	0	0	2	-	0(0)	0	0	0	-	-	0	0
	NH	0	0	0	0	-	0	0(0)	12	0	-	-	0	0
	OR	1	6	1	5	-	0	12	0(0)	2	-	-	0	0
	SD	1	0	0	3	-	0	0	2	0(0)	-	-	0	0
	UT	-	-	-	-	-	-	-	-	-	-	-	-	-
	WV	-	-	-	-	-	-	-	-	-	-	-	-	-
WY	0	0	0	8	-	0	0	0	0	-	-	0(0)	0	

	State	CT	HI	ID	MSSA ^a	MT	ND	NH	OR	SD	UT	WV	WY	USVI
	USVI	0	0	0	0	-	0	0	0	0	-	-	0	0(0)
Total	CT	0(0)	0	3	25	-	8	1	6	2	-	-	1	1
	HI	0	0(0)	4	1	-	0	0	8	0	-	-	0	0
	ID	3	4	0(0)	3	-	0	0	3	0	-	-	1	0
	MSSA	25	2	3	0(0)	-	4	1	9	5	-	-	8	0
	MT	-	-	-	-	-	-	-	-	-	-	-	-	-
	ND	8	0	0	4	-	0(0)	0	0	0	-	-	0	0
	NH	1	0	0	1	-	0	0(0)	14	0	-	-	0	0
	OR	6	8	3	10	-	0	14	0(0)	2	-	-	0	1
	SD	2	0	0	5	-	0	0	2	0(0)	-	-	0	0
	UT	-	-	-	-	-	-	-	-	-	-	-	-	-
	WV	-	-	-	-	-	-	-	-	-	-	-	-	-
	WY	1	0	1	8	-	0	0	0	0	-	-	0(0)	0
	USVI	1	0	0	0	-	0	0	0	1	0	-	-	0

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Following the (operational) administration, field-test items went through rubric validation and item data review, as described in Volume 2, Section, 2.7.1, Rubric Validation, and Section 2.7.2, Data Review.

Table 23 presents the number of field-test items administered in South Dakota, or another state or territory, the number of items rejected before or during rubric validation, the number of items submitted for data review, and the number of items rejected during data review. The numbers in parentheses indicate the ICCR owned field-test items.

Table 23. Overview of Field-Test Item Administration, Rubric Validation, and Item Data Review in Spring 2022

Grade Band and Item Type	Number of Field-Test Items Administered	Number of Items Rejected Before/During Rubric Validation	Number of Items Submitted for Data Review	Number of Items Rejected at Data Review	Number of Items Remaining ^a
Elementary School	170 (32)	3 (1)	82 (20)	14 (5)	153 (26)
Cluster	88 (9)	1 (0)	18 (1)	4 (0)	83 (9)
Stand-Alone	82 (23)	2 (1)	64 (19)	10 (5)	70 (17)
Middle School	190 (62)	4 (2)	94 (26)	26 (5)	160 (55)
Cluster	88 (26)	3 (1)	26 (4)	13 (1)	72 (24)
Stand-Alone	102 (36)	1 (1)	68 (22)	13 (4)	88 (31)
High School	111 (42)	2 (0)	63 (22)	19 (4)	90 (38)
Cluster	41 (13)	2 (0)	21 (7)	3 (0)	36 (13)
Stand-Alone	70 (29)	0 (0)	42 (15)	16 (4)	54 (25)
Total	471 (136)	9 (3)	239 (68)	59 (14)	403 (119)

Note. The numbers in parentheses indicate ICCR items.

^aCount excludes seventy (sixty Oregon and ten South Dakota) legacy field-test items.

Table 24 summarizes the Shared Science Assessment Item Bank after adding the field-test items that were administered in 2022 and passed rubric validation and item data review. The numbers in parentheses indicate the ICCR owned items.

Table 24. Summary of Shared Science Assessment Item Bank in Spring 2022

Grade Band and Item Type	Science Discipline			Total ^a
	Earth and Space Sciences	Life Sciences	Physical Sciences	
Elementary School	180 (46)	162 (44)	214 (52)	556 (142)
Cluster	96 (18)	82 (14)	111 (17)	289 (49)
Stand-Alone	84 (28)	80 (30)	103 (35)	267 (93)
Middle School	150 (43)	220 (65)	187 (48)	557 (156)
Cluster	70 (15)	110 (21)	90 (17)	270 (53)
Stand-Alone	80 (28)	110 (44)	97 (31)	287 (103)
High School	91 (34)	194 (68)	129 (39)	414 (141)
Cluster	35 (13)	78 (20)	53 (17)	166 (50)
Stand-Alone	56 (21)	116 (48)	76 (22)	248 (91)
Total	421 (123)	576 (177)	530 (139)	1527 (439)

Note. The numbers in parentheses indicate ICCR items.

^aCount excludes eighty-five (sixty Oregon and twenty-five South Dakota) legacy items and nine MOU items that do not align to the NGSS.

Appendix 1-C

Calibration of the Shared Science Assessment Item Bank



TABLE OF CONTENTS

1.	2018 CALIBRATION SEQUENCE.....	1
2.	2019 CALIBRATION	4
3.	LINKING THE 2018 SCALE TO THE 2019 SCALE.....	8
4.	CALIBRATION OF FIELD-TEST ITEMS IN 2021 AND BEYOND.....	9
5.	CALIBRATION SOFTWARE	11
6.	REFERENCES	12

LIST OF TABLES

Table 1. Groups Per Grade Band for the Spring 2018 Core Calibration	1
Table 2. Spring 2018 State-Sharing Matrix	3
Table 3. Groups Per Grade Band for the Spring 2019 Calibration of Operational Items	4
Table 4. Number of Common Elementary School Operational Items Administered and Calibrated in Spring 2019.....	5
Table 5. Number of Common Middle School Operational Items Administered and Calibrated in Spring 2019.....	6
Table 6. Number of Common High School Operational Items Administered and Calibrated in Spring 2019.....	7
Table 7. Groups Per Grade Band for the Spring 2019 Calibration of Field-Test Items	7
Table 8. Estimated Latent Means and Number of Students Per State	9
Table 9. Groups Per Grade Band for the Spring 2021 Calibration of Field-Test Items	10
Table 10. Groups Per Grade Band for the Spring 2022 Calibration of Field-Test Items	11

The Shared Science Assessment Item Bank was first calibrated in 2018 after the 2018 science test administrations concluded, and it was recalibrated in 2019 following the 2019 test administrations. The calibration sequences are documented in this appendix, which also includes details on scale linking and the creation of the anchor scale in 2019. The calibration of field-test items in 2021 and beyond as well as the calibration software are addressed.

1. 2018 Calibration Sequence

Table 1 provides an overview of the groups per grade band for the 2018 calibration.

Table 1. Groups Per Grade Band for the Spring 2018 Core Calibration

Group	Elementary School	Middle School	High School
Connecticut	X	X	X
Hawaii	X	X	X
New Hampshire	X	X	X
Rhode Island	X	X	X
Utah Grade 6		X	
Utah Grade 7		X	
Utah Grade 8		X	
Vermont	X	X	X
West Virginia	X	X	

Items were calibrated in three steps for two reasons. First, the rubric validations for some states took place at a later date, and the student responses for the items owned by those states could not be included in the first round of calibrations without jeopardizing the reporting schedule of the two states with operational field tests (i.e., those two states did not have any of the items with late rubric validation in their item pool). Second, to divide the large set of items and assertions into more manageable portions, a separate calibration was carried out for two states with many items administered in those states only. Specifically, the following sequence of calibrations was carried out:

1. **Core Calibration.** The core calibration was performed on the following items:
 - a. All item responses for New Hampshire and West Virginia. These states administered items from the following sources (as described in the state-sharing matrix in Table 2):
 - i. ICCR item bank
 - ii. Connecticut
 - iii. Hawaii
 - iv. Rhode Island
 - v. Vermont

vi. Utah

vii. West Virginia

A more detailed overlap of the common items at the time of the 2018 calibration was given in Section 1 of Appendix 1-B, Shared Science Assessment Item Bank: Field Testing (refer to Table 2 - Table 4).

- b. All item responses from Connecticut, Rhode Island, and Vermont, excluding responses to Wyoming and Oregon items. These states administered items from the following sources:
- i. ICCR item bank
 - ii. Connecticut
 - iii. Hawaii
 - iv. Rhode Island
 - v. Vermont
 - vi. Utah
 - vii. West Virginia
 - viii. Wyoming (items were treated as “not administered”; responses were replaced by missing code.)
 - ix. Oregon (items were treated as “not administered”; responses were replaced by missing code.)
- c. Item responses from Hawaii to items also administered in another state (Hawaii items were used in Hawaii, Connecticut, Rhode Island, Vermont, and West Virginia.).
- d. Item responses from Utah to items also administered in another state (Utah items were used in Utah, Connecticut, Rhode Island, Vermont, and West Virginia) underwent core calibration. Utah tested only middle-school students. One-third of students were selected at random to balance the large population size for Utah.

Table 2. Spring 2018 State-Sharing Matrix

Source Bank	CT	HI	MSSA ^a	NH	OR	UT	WV	WY
ICCR	X	X	X	X	X		X	X
Connecticut	X		X				X	
Hawaii	X	X	X				X	
MSSA ^a	X		X				X	
Oregon	X		X		X			
Utah	X		X			X	X	
West Virginia	X		X				X	
Wyoming	X		X					X

Note. The core calibration provided parameters for all items used in New Hampshire and West Virginia.

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

2. **Calibration of State-Specific Items.** In terms of the calibration of state-specific items, both Hawaii and Utah had a substantial proportion of items that were administered only in Hawaii and Utah, respectively. Hawaii had both Hawaii and ICCR items in common with the states involved in the core calibration (Hawaii administered Hawaii and ICCR items only); whereas Utah had only Utah items in common (Utah administered Utah items only). The parameters for the unique Hawaii items depended on responses from Hawaii students only, and the parameters for the unique Utah items depended on responses from Utah students only. For both states, the state-specific items were calibrated through a separate calibration based on the state data only, with the items in common with the core states mentioned in Step 1 anchored to the estimates from Step 1. These calibrations were performed separately for each group under a single-group item response theory (IRT) model. The mean and variance of the groups were fixed to the estimated mean and variance from the core calibration.
3. **Calibration of States with Late Rubric Validation.** Oregon and Wyoming items were administered in some of the states involved in the core calibration (Connecticut, Rhode Island, and Vermont) but could not be calibrated in Step 1 because of their late rubric validation dates. In a later stage, items from Oregon and Wyoming were calibrated by
 - a. adding Oregon and Wyoming student responses to the core calibration;
 - b. keeping the responses from Connecticut, Rhode Island, and Vermont to Wyoming and Oregon items (as opposed to treating them as missing in Step 1);
 - c. removing the responses from Hawaii, Utah, New Hampshire, and West Virginia, who did not administer Oregon or Wyoming items (as the item parameters for the Oregon and Wyoming items did not depend on the students from these states); and
 - d. fixing the parameters of all other items to the values obtained in Step 1 and the group means and standard deviations that were estimated in Step 1.

2. 2019 CALIBRATION

Calibration was performed in two steps. First, CAI calibrated all items in operational use in 2019, for which 1,000 or more student responses were available (among these, there were 1,500 or more student responses for all but three items). In this step, only the data of states with an operational test were included. Table 3 provides an overview of the groups per grade band for this first calibration. All students who attempted the test were included in the calibration. The assertions of skipped items were scored as incorrect. Note that only Rhode Island allowed students to skip items. Out of a total of 438 items, there were nine items administered as operational items in 2019, for which the sample size was smaller than 1,000.

Table 4 through Table 6 present the number of operational item clusters and stand-alone items that were shared between the item pools of any two states. The numbers below the shaded cells represent the number of all the operational items administered, and the numbers above the shaded cells represent the number of common operational items at the time of the 2019 calibration. The shaded cells represent the number of operational items administered only in the given state (the number of unique operational items at the time of calibration are provided in parentheses). Since the items that were administered but not calibrated were administered only in one state, the numbers above the diagonal are the same as the numbers below the diagonal.

Table 4 presents the results for elementary schools, Table 5 presents the results for middle schools, and Table 6 presents the results for high schools. The numbers at the operational administration are slightly different from the numbers at the calibration because items with sample sizes smaller than 1,000 were excluded from the calibration.

Table 3. Groups Per Grade Band for the Spring 2019 Calibration of Operational Items

Group	Elementary School	Middle School	High School
Connecticut	X	X	X
New Hampshire	X	X	X
Oregon	X	X	X
Rhode Island	X	X	X
Vermont	X	X	X
West Virginia	X	X	

Table 4. Number of Common Elementary School Operational Items Administered and Calibrated in Spring 2019

	State	CT	MSSA ^a	NH	OR	WV
Cluster	CT	1 (1)	44	24	42	55
	MSSA	44	0 (0)	17	37	41
	NH	24	17	0 (0)	14	27
	OR	42	37	14	0 (0)	41
	WV	55	41	27	41	1 (1)
Stand-Alone	CT	3 (3)	34	26	30	47
	MSSA	34	0 (0)	20	23	32
	NH	26	20	0 (0)	14	25
	OR	30	23	14	0 (0)	25
	WV	47	32	25	25	1 (1)
Total	CT	4 (4)	78	50	72	102
	MSSA	78	0 (0)	37	60	73
	NH	50	37	0 (0)	28	52
	OR	72	60	28	0 (0)	66
	WV	102	73	52	66	2 (2)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 5. Number of Common Middle School Operational Items Administered and Calibrated in Spring 2019

	State	CT	MSSA ^a	NH	OR	WV
Cluster	CT	3 (3)	26	24	54	92
	MSSA	26	0 (0)	11	14	21
	NH	24	11	1 (1)	9	18
	OR	54	14	9	2 (2)	56
	WV	92	21	18	56	12 (4)
Stand-Alone	CT	0 (0)	42	26	34	50
	MSSA	42	0 (0)	25	30	37
	NH	26	25	0 (0)	16	21
	OR	34	30	16	1 (0)	29
	WV	50	37	21	29	0 (0)
Total	CT	3 (3)	68	50	88	142
	MSSA	68	0 (0)	36	44	58
	NH	50	36	1 (1)	25	39
	OR	88	44	25	3 (2)	85
	WV	142	58	39	85	12 (4)

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

Table 6. Number of Common High School Operational Items Administered and Calibrated in Spring 2019

	State	CT	MSSA ^a	NH	OR	WV
Cluster	CT	5 (5)	33	22	30	–
	MSSA	33	0 (0)	20	31	–
	NH	22	20	2 (2)	15	–
	OR	30	31	15	1 (1)	–
	WV	–	–	–	–	–
Stand-Alone	CT	0 (0)	39	27	40	–
	MSSA	39	2 (2)	23	32	–
	NH	27	23	0 (0)	20	–
	OR	40	32	20	4 (4)	–
	WV	–	–	–	–	–
Total	CT	5 (5)	72	49	70	–
	MSSA	72	2 (2)	43	63	–
	NH	49	43	2 (2)	35	–
	OR	70	63	35	5 (5)	–
	WV	–	–	–	–	–

^aMSSA = Rhode Island and Vermont’s Multi-State Science Assessment

In Step 2, the field-test items were calibrated. The calibration included the operational items that were calibrated in Step 1 and the field-test items across all states in which they were administered. All students who attempted at least one field-test item were included in the calibration. Table 7 provides an overview of the groups per grade band for calibration of the field-test items.

Table 7. Groups Per Grade Band for the Spring 2019 Calibration of Field-Test Items

Group	Elementary School	Middle School	High School
Connecticut	X	X	X
Hawaii	X	X	X
Idaho	X	X	
New Hampshire	X	X	X
Oregon	X	X	X
Rhode Island	X	X	X
Vermont	X	X	X
West Virginia	X	X	
Wyoming	X	X	X

3. LINKING THE 2018 SCALE TO THE 2019 SCALE

The item parameter estimates obtained from the 2018 student responses were highly correlated with the item parameters obtained from the 2019 student responses. For item difficulties, the correlation between the 2018 and 2019 estimates was 0.993 for elementary school, 0.986 for middle school, and 0.994 for high school. For the standard deviations of the clusters, these correlations were 0.971 for elementary school, 0.972 for middle school, and 0.964 for high school. These high correlations indicate that items functioned similarly in 2018 and 2019. Nevertheless, item parameters from separate calibrations cannot be directly compared because the scale of an IRT model was not determined. In the multigroup Rasch testlet model, the only scale indeterminacy was the origin of the scale. The models can be identified by setting the mean of the overall proficiency variable θ to zero for the reference distribution. As a result, the 2018 and 2019 variable θ and item parameters were on the same scale except for an overall shift parameter B . Specifically, the 2018 scale can be linked to the 2019 scale as follows:

$$\begin{aligned} P(z_{ij} | \theta_{j\ 2018}, u_{jg}; b_{i\ 2018}) &= \frac{\exp(\theta_{j\ 2018} + u_{jg} - b_{i\ 2018})}{1 + \exp(\theta_{j\ 2018} + u_{jg} - b_{i\ 2018})} \\ &= \frac{\exp(\theta_{j\ 2018} + B + u_{jg} - b_{i\ 2018} - B)}{1 + \exp(\theta_{j\ 2018} + B + u_{jg} - b_{i\ 2018} - B)} \\ &= \frac{\exp(\theta_{j\ 2019} + u_{jg} - b_{i\ 2019})}{1 + \exp(\theta_{j\ 2019} + u_{jg} - b_{i\ 2019})}. \end{aligned}$$

Because $\theta_{j\ 2019} = \theta_{j\ 2018} + B$, the population means of θ must be transformed accordingly,

$$\theta_{j\ 2019} \sim N(\mu_{k\ 2018} + B, \sigma_k^2) \text{ and}$$

$$\theta_{j\ 2018} \sim N(\mu_{k\ 2018}, \sigma_k^2).$$

Item parameters based on 2018 student responses were expressed on the 2019 scale by adding the constant B to the 2018 item parameter. The 2018 parameters were expressed on the 2019 scale for items that were part of the pool in both 2018 and 2019 but were not administered in any states in 2019 (13 items), and for items that were administered in 2019 but the number of student responses from the 2019 assessments was lower than 1,000 (nine items). Therefore, the linking process was performed for 22 items only.

All items that were operational in 2019 were also administered in 2018. Therefore, the shift parameter B was estimated from a separate calibration of the 2019 operational items using the 2019 student responses (from the six operational states), but with the item parameters fixed to the estimates obtained from the 2018 calibrations. By fixing a subset of the item parameters, the model is identified so that the means and variances of θ can be estimated for all groups. Parameter B can be obtained by equating the overall mean of θ across all groups for the 2019 student response data from the free calibration (i.e., the 2019 overall mean expressed on the 2019 scale) to the overall mean of θ across all groups for the 2019 student response data from the calibration with items anchored to their 2018 parameters values (i.e., the 2019 overall mean expressed on the 2018 scale):

$$\frac{1}{K} \sum_{k=1}^K \mu_{k\ 2019} = \frac{1}{K} \sum_{k=1}^K (\mu_{k\ 2018} + B).$$

Therefore, an estimate of parameter B can be obtained as

$$\hat{B} = \frac{1}{K} \sum_{k=1}^K (\hat{\mu}_{k\ 2019} - \hat{\mu}_{k\ 2018}).$$

Table 8 presents the estimated means of θ under both the free and anchored calibrations and the number of students per state. Table 8 also presents the overall means and estimated shift in parameter B . Note that the parameters for three items were not anchored, but instead were freely estimated together with the means and variances in the anchored calibration. The reason for not treating these items as common items across the 2018 and 2019 administrations is that they had an omit rate of 4% or higher for the last item interaction in the 2018 administration in at least one state. In 2019, these interactions could no longer be omitted because all interactions of an item needed to be responded to in states where skipping was not allowed (i.e., all states excluding Rhode Island). Therefore, out of an abundance of caution, these three items are not anchored to their 2018 parameter values.

Table 8. Estimated Latent Means and Number of Students Per State

Group	Elementary School			Middle School			High School		
	$\hat{\mu}_{k\ 2019}$	$\hat{\mu}_{k\ 2018}$	N	$\hat{\mu}_{k\ 2019}$	$\hat{\mu}_{k\ 2018}$	N	$\hat{\mu}_{k\ 2019}$	$\hat{\mu}_{k\ 2018}$	N
Connecticut	0.0000	0.0518	38,549	0.0000	0.0234	39,347	0.0000	0.1443	37,616
New Hampshire	0.0631	0.1083	13,187	0.0940	0.1108	12,060	0.0798	0.2278	11,385
Oregon	-0.0101	0.0096	44,989	0.0028	0.0156	42,043	-0.0383	0.1030	41,630
Rhode Island	-0.0312	0.0142	10,751	-0.1044	-0.0692	10,306	-0.2261	-0.0879	9,612
Vermont	0.1069	0.1504	6,017	0.0781	0.1133	5,894	0.0179	0.1545	5,332
West Virginia	-0.1970	-0.1529	19,540	-0.3012	-0.2783	19,043	–	–	–
	$\frac{1}{K} \sum_{k=1}^K \hat{\mu}_{k\ 2019}$	$\frac{1}{K} \sum_{k=1}^K \hat{\mu}_{k\ 2018}$	\hat{B}	$\frac{1}{K} \sum_{k=1}^K \hat{\mu}_{k\ 2019}$	$\frac{1}{K} \sum_{k=1}^K \hat{\mu}_{k\ 2018}$	\hat{B}	$\frac{1}{K} \sum_{k=1}^K \hat{\mu}_{k\ 2019}$	$\frac{1}{K} \sum_{k=1}^K \hat{\mu}_{k\ 2018}$	\hat{B}
Overall	-0.0114	0.0303	-0.0416	-0.0385	-0.0141	-0.0244	-0.0333	0.1083	-0.1417

4. CALIBRATION OF FIELD-TEST ITEMS IN 2021 AND BEYOND

Starting in 2021, Starting in 2021, field-test items were calibrated with one multigroup calibration per grade band. In each calibration, the parameters of the operational items were fixed to their bank values (anchor items) and the item parameters of the field-test items, as well as the mean and variance of each group, were estimated using the marginal maximum likelihood (MML) method. The calibration included the field-test items across all states in which the items were administered. All students who attempted at least one field-test item were included in the calibration. In 2021 and 2022, the same groups were included for each grade band for the field-test calibration. Refer

to Table 9 and Table 10 for an overview of the groups per grade band for calibration of the field-test items in 2021 and 2022, respectively.

In 2021, all but 12 items were calibrated on at least 1,500 student responses, and all but one item had a sample size larger than 1,200. The item with fewer than 1,200 responses had a sample size of 981 and was an interim item. In 2022, all but 64 items were calibrated on at least 1,500 student responses, and all but nine items were calibrated on at least 1,200 responses. The nine items with fewer than 1,200 responses were either Oregon legacy items or interim items.

Table 9. Groups Per Grade Band for the Spring 2021 Calibration of Field-Test Items

Group	Elementary School	Middle School	High School
Connecticut	X	X	X
Hawaii	X	X	X
Idaho	X	X	X
Montana	X	X	
North Dakota	X	X	X
New Hampshire	X	X	X
Rhode Island	X	X	X
South Dakota	X	X	X
Utah	X	X	
Vermont	X	X	X
West Virginia	X	X	
Wyoming	X	X	X

Table 10. Groups Per Grade Band for the Spring 2022 Calibration of Field-Test Items

Group	Elementary School	Middle School	High School
Connecticut	X	X	X
Hawaii	X	X	X
Idaho	X	X	X
Montana	X	X	
North Dakota	X	X	X
New Hampshire	X	X	X
Oregon	X	X	X
Rhode Island	X	X	X
South Dakota	X	X	X
Utah	X	X	
Vermont	X	X	X
West Virginia	X	X	
Wyoming	X	X	X

5. CALIBRATION SOFTWARE

In 2018 and 2019, the IRT models were fitted using the Bayesian networks with the logistic regression (BNL) suite of Matlab functions (Rijmen, 2006) and flexMIRT (Cai, 2017). The resulting parameters from BNL were used as starting values for flexMIRT to reduce the estimation time for flexMIRT. The flexMIRT estimates were taken to be the operational parameters, except for the middle-school items calibrated in 2018 during the core calibration. For the 2018 core calibration of middle-school items, flexMIRT did not converge after several weeks, and the estimates obtained from BNL were used as operational parameters. Note that the parameters estimates were very similar across software packages.

Starting in 2021, Cambium Assessment IRT (CAIRT) was used for all calibrations because the estimation time in flexMIRT became prohibitive. CAIRT was specifically developed by CAI to calibrate the multigroup Rasch model on very large data sets. It relies on the same estimation methods as BNL. CAI has cross-validated parameter estimates from CAIRT with BNL and flexMIRT under various scenarios (Rijmen, Liao, & Lin, 2021). In 2023, field-test items were calibrated in CAIRT using the same procedure used in 2021.

6. REFERENCES

- Cai, L. (2017). flexMIRT[®]: Flexible multilevel multidimensional item analysis and test scoring (version 3.51) [Computer software]. Chapel Hill, NC: Vector Psychometric Group.
- Rijmen, F. (2006). *BNL: A Matlab toolbox for Bayesian networks with logistic regression nodes* (Technical Report). Amsterdam: VU University Medical Center.
- Rijmen, F., Liao, D., & Lin, Z. (2021). *The Rasch testlet model for the calibration of three-dimensional science assessments: A software comparison* [White paper]. Washington, DC: Cambium Assessment, Inc.

Appendix 1-D

Distribution of Scale Scores and Achievement Levels

Distribution of Scale Scores and Achievement Levels

Table 1-D-1. Scale Score Mean and Standard Deviation by Grade, Science

	Grade		
	5	8	11
<i>Number of Students</i>	10,605	10,744	9,474
<i>Mean Scale Score</i>	499.71	798.56	1,101.11
<i>SD of Scale Score</i>	29.01	33.35	29.05

Table 1-D-2. Proportion of Students in Each Achievement Level by Grade, Science

Achievement Level	Grade		
	5	8	11
<i>Number of Students</i>	10,605	10,744	9,474
Level 1	0.22	0.23	0.17
Level 2	0.38	0.39	0.32
Level 3	0.22	0.25	0.4
Level 4	0.19	0.13	0.11

Appendix 1-E

Distribution of Scale Scores by Science Discipline

Distribution of Scale Scores by Science Discipline

Table 1-E-1. Science Disciplines

Grade	Discipline
5, 8, and 11	Physical Sciences (PS) Life Sciences (LS) Earth and Space Sciences (ESS)

Table 1-E-2. Overall Discipline Score Mean and Standard Deviation, Grade 5 Science

N	Scale Score	Discipline		
		Physical Sciences	Life Sciences	Earth and Space Sciences
10,605	Mean	499.68	499.01	499.97
	SD	32.53	36.01	33.77

Table 1-E-3. Overall Discipline Score Mean and Standard Deviation, Grade 8 Science

N	Scale Score	Discipline		
		Physical Sciences	Life Sciences	Earth and Space Sciences
10,744	Mean	799.52	797.03	798.50
	SD	38.70	39.65	36.58

Table 1-E-4. Overall Discipline Score Mean and Standard Deviation, Grade 11

N	Scale Score	Discipline		
		Physical Sciences	Life Sciences	Earth and Space Sciences
9,474	Mean	1,102.92	1,100.16	1,099.77
	SD	33.73	35.78	31.54

Appendix 1-F

Distribution of Scale Scores and Achievement Levels by Subgroup

Distribution of Scale Scores and Achievement Levels by Subgroup

Table 1-F-1. Mean and Standard Deviation of Scale Scores by Subgroup

Group	Grade 5			Grade 8			Grade 11		
	<i>N</i>	<i>Mean</i>	<i>SD</i>	<i>N</i>	<i>Mean</i>	<i>SD</i>	<i>N</i>	<i>Mean</i>	<i>SD</i>
All Students	10,605	499.71	29.01	10,744	798.56	33.35	9,474	1,101.11	29.05
Female	5,202	498.97	27.81	5,281	798.74	31.41	4,698	1,100.19	26.42
Male	5,403	500.43	30.11	5,463	798.38	35.13	4,776	1,102.01	31.41
African American	338	485.57	25.46	329	781.04	32.26	304	1,086.77	27.57
American Indian/Alaskan Native	1,230	478.04	25.49	1,145	773.33	30.40	746	1,079.54	24.94
Asian	200	501.96	27.28	142	798.12	34.60	163	1,102.71	33.97
Hispanic	856	489.17	27.21	825	784.22	33.01	627	1,089.38	29.20
Multi-Racial	641	495.56	28.97	559	794.57	33.11	437	1,098.14	26.51
Pacific Islander	15	496.39	23.00	16	780.28	29.43	15	1,078.88	26.73
White	7,325	505.54	27.61	7,728	804.91	31.22	7,182	1,105.17	28.02
Limited English Proficiency	502	474.51	21.06	400	765.13	25.21	247	1,065.9	20.72
Special Education	1,747	480.02	29.17	1,259	769.88	30.33	869	1,073.16	24.95
Economically Disadvantaged	3,877	487.34	27.63	3,533	783.82	32.49	2,293	1,088.72	27.49

Table 1-F-2. Percentage of Achievement Level by Subgroup

Group	Grade 5					Grade 8					Grade 11				
	N	L1	L2	L3	L4	N	L1	L2	L3	L4	N	L1	L2	L3	L4
All Students	10,605	22%	38%	22%	19%	10,744	23%	39%	25%	13%	9,474	17%	32%	40%	11%
Female	5,202	22%	40%	21%	17%	5,281	21%	42%	26%	11%	4,698	15%	35%	41%	9%
Male	5,403	22%	36%	22%	20%	5,463	25%	37%	24%	14%	4,776	19%	29%	38%	14%
African American	338	36%	42%	16%	6%	329	44%	38%	13%	5%	304	32%	38%	24%	6%
American Indian/Alaskan Native	1,230	51%	35%	10%	4%	1,145	52%	36%	9%	2%	746	40%	40%	18%	2%
Asian	200	21%	34%	26%	18%	142	25%	36%	23%	17%	163	15%	36%	34%	15%
Hispanic	856	33%	42%	16%	9%	825	37%	40%	17%	6%	627	31%	33%	30%	6%
Multi-Racial	641	26%	40%	19%	15%	559	27%	39%	23%	11%	437	19%	36%	38%	7%
Pacific Islander	15	27%	53%	7%	13%	16	44%	38%	19%	-	15	40%	40%	20%	-
White	7,325	15%	38%	25%	23%	7,728	16%	40%	29%	16%	7,182	12%	31%	44%	13%
Limited English Proficiency	502	55%	37%	7%	1%	400	61%	36%	2%	1%	247	64%	31%	4%	-
Special Education	1,747	52%	32%	9%	8%	1,259	59%	31%	6%	4%	869	53%	34%	12%	1%
Economically Disadvantaged	3,877	37%	39%	16%	8%	3,533	38%	40%	16%	6%	2,293	29%	38%	27%	5%

- No data are available.