# Rhode Island and Vermont Multi-State Science Assessment

2020-2021

# Volume 2: Test Development





# TABLE OF CONTENTS

1.	Introduction	1
1.1	Claim Structure	2
1.2	Underlying Principles Guiding Development	2
1.3	Organization of this Volume	3
2.	ITEM DEVELOPMENT PROCESS THAT SUPPORTS VALIDITY OF CLAIMS	3
2.1	Overview	3
2.2	Item Specifications	
2.3	Selection and Training of Item Writers	7
2.4	Internal Review	
	2.4.1 Preliminary Review	8
	2.4.2 Scoring Entry and Review	
	2.4.3 Content Review One	9
	2.4.4 Edit Review	
	2.4.5 Senior Review	
2.5	Review by State Personnel and Stakeholder Committees	
	2.5.1 State Review	
	2.5.2 Content Advisory Committee Reviews	
	2.5.3 Language Accessibility, Bias, and Sensitivity Committee Reviews	
2.6	2.5.4 Markup for Translation and Accessibility Features	
2.6	Field Testing	
2.7	Post-Field-Test Review	
	2.7.1 Rubric Validation	
	2.7.2 Data Review	16
3.	SCIENCE ITEM BANK SUMMARY	20
3.1	Current Composition of the Shared Science Assessment Item Bank	20
3.2	Strategy for Bank Evaluation and Replenishment	27
4.	MULTI-STATE SCIENCE ASSESSMENT TEST CONSTRUCTION	27
4.1	Test Design	27
4.2		28
4.3	Online Test Construction.	
4.4	Paper-Pencil Accommodation Form Construction	
5.	SIMULATION SUMMARY REPORT	45
5.1	Factors Affecting Simulation Results	46
5.2	Results of Simulated Test Administrations: English	
	5.2.1 Summary of Blueprint Match	
	5.2.2 Item Exposure	
5.3	Results of Simulated Test Administrations: Spanish	47

	5.3.1 Summary of Blueprint Match	47
	5.3.2 Item Exposure	
6.	OPERATIONAL TEST ADMINISTRATION SUMMARY REPORT	48
6.1	Blueprint Match	48
6.2	Item Exposure	48
7.	References	50

# LIST OF TABLES

Table 1. Science Interaction Types and Descriptions	. 21
Table 2. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank	. 22
Гable 3. Spring 2021 Shared Science Assessment Operational Item Bank	. 23
Гable 4. Spring 2021 Shared Science Assessment Field-Test Item Bank	. 23
Гable 5. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank by	
Science Discipline	. 24
Гable 6. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank by	
Disciplinary Core Idea	. 25
Гable 7. Science Test Blueprint, Grade 5	. 29
Гable 8. Science Test Blueprint, Grade 8	. 32
Гable 9. Science Test Blueprint, Grade 11	. 35
Table 10. Combined Percentile 85 Testing Times by Grade	. 39
Гable 11. Rhode Island Percentile 85 Testing Times by Grade	. 39
Table 12. Vermont Percentile 85 Testing Times by Grade	. 39
Гable 13. MSSA Spring 2021 Operational and Field-Test Item Pool	. 40
Гable 14. MSSA Spring 2021 Operational Item Pool	. 40
Гable 15. MSSA Spring 2021 Field-Test Item Pool	. 41
Table 16. MSSA Spring 2021 Operational and Field-Test Item Pool by Science Discipline	. 41
Гable 17. MSSA Spring 2021 Operational and Field-Test Item Pool by Disciplinary Core Idea	a 43
Table 18. Item Exposure Rates by Grade: Percentage of Items by Exposure Rate, Across All	
English Online Simulation Sessions	. 47
Гable 19. Spring 2019 Spanish Operational Item Pool	. 47
Гable 20. Item Exposure Rates by Grade: Percentage of Items by Exposure Rate, Across All	
Spanish Simulation Sessions	. 48
Table 21. Item Exposure Rates by Grade: Percentage of Items by Exposure Rate, Across All	
Spring 2019 Test Administrations in Rhode Island	. 48
Γable 22. Item Exposure Rates by Grade: Percentage of Items by Exposure Rate, Across All	
Spring 2019 Test Administrations in Vermont	. 49

# **LIST OF EXHIBITS**

Exhibit A. Summary of How Each Step of Development Supports the Validity of Claims	4
Exhibit B. Sample Science Item Cluster Specifications for Middle School	6
Exhibit C. Summary of Content Advisory Committee Meetings	11
Exhibit D. Summary of Fairness Committee Meetings	13
Exhibit E. Features of the REVISE Software	16
Exhibit F. Summary of Data Review Committee Meetings	17

# LIST OF APPENDICES

- Appendix A. Item Writer Training Materials
- Appendix B. Item Review Checklist
- Appendix C. Content Advisory Committee Participant Details
- Appendix D. Fairness Committee Participant Details
- Appendix E. Sample Data Review Training Materials
- Appendix F. Data Review Committee Participant Details
- Appendix G. Example Item Interactions
- Appendix H. Science Item Bank
- Appendix I. Multi-State Science Assessment Item Pool
- Appendix J. Adaptive Algorithm Design
- Appendix K. Content Advisory Committee Review Training Slides
- Appendix L. Fairness Committee Review Training Slides

### 1. Introduction

In 2013, the Rhode Island Department of Education (RIDE) and Vermont Agency of Education (VT AOE) adopted the Next Generation Science Standards (NGSS). The RIDE and the VT AOE and their assessment vendor, Cambium Assessment, Inc. (CAI; formerly the American Institutes for Research [AIR]), developed and administered a new online assessment to measure the new standards. In 2017–2018, the Rhode Island Next Generation Science Assessment (RI NGSA) was administered as an independent field test in Rhode Island, and the Vermont Science Assessment (VTSA) was administered as an operational field test in Vermont. The RI NGSA and VTSA were administered operationally for the first time in 2018–2019. The RI NGSA and the VTSA measure the science knowledge and skills of Rhode Island and Vermont students in grades 5, 8, and 11 as an online assessment, constructed linearly on the fly, making use of several technology-enhanced item types. The content measures the three-dimensional science standards based on the National Research Council's *A Framework for K–12 Science Education* published in 2012.

In the remainder of this volume, the term *Multi-State Science Assessment* (MSSA) will refer to the RI NGSA and VTSA.

Additional details on the implementation of the assessments can be found in Volume 1, Annual Technical Report.

The interpretation, usage, and validity of test scores rely heavily upon the process of developing the test itself. This volume provides details on the test development process of the MSSA, which contributes to the validity of the test scores. Specifically, this volume provides evidence to support the following:

- The test item specifications provided detailed guidance for item writers and reviewers to ensure that science items were aligned to the performance expectations (PEs) they were intended to measure (Appendix A, Item Writer Training Materials, and Appendix B, Item Review Checklist).
- The item development procedures employed for MSSA tests were consistent with industry standards.
- The development and maintenance of the Shared Science Assessment Item Bank, in which test items cover the range of measured PEs, grade-level difficulties, and levels of cognitive engagement through the use of both item clusters and stand-alone items.
- The Test Design Summary/Blueprint stipulated the range of operational items from each item type and content category required on each test administration. This document was implemented in the item selection algorithm for science (Appendix J, Adaptive Algorithm Design).

Note that for the science assessments, as outlined in Volume 1, Annual Technical Report, CAI works with a group of states that share common item development processes. In addition to developing items for each of those states, CAI develops and maintains the Independent College and Career Readiness (ICCR) item bank, which consists of items that are developed according to

the same principles that are followed for the items owned by each of the states. Therefore, this volume focuses on the general test development activities.

For the MSSA test, items are drawn from an item bank that consists of ICCR items, items owned by Rhode Island and Vermont, and items owned by several other states that share a Memorandum of Understanding (MOU) to share content, leadership, and new ideas and methods. Specifically, all items developed under the MOU went through the same development process. For the remainder of this volume, the term *item bank* will refer to all items developed under the MOU unless stated explicitly otherwise.

# 1.1 CLAIM STRUCTURE

The goals, uses, and claims that the science item bank and subsequent tests would be designed to support were identified in a series of collaborative meetings held over August 22–23, 2016. The overarching goal was to support the development of statewide summative assessments using science content that measures the three-dimensional science standards based on *A Framework for K–12 Science Education* (National Research Council, 2012).

To this end, CAI invited content and assessment leaders from 10 states as well as four nationally recognized experts that helped author the NGSS. Two nationally recognized psychometricians also participated.

CAI staff and participating states collaborated to develop items and test specifications to measure the three-dimensional science standards. The item specifications were generally accompanied by sample item clusters meeting those specifications. All specifications and sample item clusters were reviewed by state content experts and committees of educators in at least one of the states.

#### 1.2 Underlying Principles Guiding Development

The Shared Science Assessment Item Bank for science was established using a highly structured, evidence-centered design. The process began with detailed item specifications. The specifications, discussed in Section 2.2, Item Specifications, described the interaction types that can be used, gave guidelines for targeting the appropriate cognitive engagement, offered suggestions for controlling item difficulty, and provided sample items.

Items were written with the goal that virtually every item would be accessible to all students, either by itself or in conjunction with accessibility tools, such as text-to-speech (TTS), translations, or assistive technologies. This goal is supported by the delivery of the items on CAI's Test Delivery System (TDS), which has received Web Content Accessibility Guidelines (WCAG) 2.0 AA certification. This platform offers a wide array of accessibility tools and is compatible with most assistive technologies.

Item development supported the goal of high-quality item clusters and stand-alone items through rigorous development processes managed and tracked by a content development platform. This system ensures that every item flows through the correct sequence of reviews and captures every comment and change to the item.

CAI sought to ensure that the items measured the PEs in a fair and meaningful way by engaging educators and other stakeholders at each step of the process. Educators evaluated the alignment of

items to the PEs and offered guidance and suggestions for improvement. They participated in the review of items for fairness and sensitivity. Following item field testing, educators engaged in rubric validation, a process that refines rule-based rubrics upon review of student responses.

Combined, these principles and the processes that support them have been incorporated into an item bank that measures the PEs with fidelity and does so in a way that minimizes construct-irrelevant variance and barriers to access. The details of these processes are described in this volume of the technical report.

# 1.3 ORGANIZATION OF THIS VOLUME

This volume is organized in three subsequent sections:

- 1. An overview of the science item development process that supports the validity of the claims that science tests are designed to support.
- 2. An overview of the science item bank, the types of assessments the bank is designed to support, and methods for refreshing the bank.
- 3. A description of the test construction process followed for the MSSA, including the blueprint, the test design, an evaluation of simulated test sessions, the operational blueprint match results, and the item exposure rates.

#### 2. ITEM DEVELOPMENT PROCESS THAT SUPPORTS VALIDITY OF CLAIMS

#### 2.1 OVERVIEW

Cambium Assessment, Inc. (CAI) developed the Shared Science Assessment Item Bank in collaboration with the states that were part of the Memorandum of Understanding (MOU) using a rigorous, structured process that engaged stakeholders at critical junctures. This process was managed by CAI's Item Tracking System (ITS), which is an auditable content development tool that enforces rigorous workflow and captures each item change and comment. Reviewers, including internal CAI reviewers or stakeholders in committee meetings, can review items in ITS as they will appear to the student, with all accessibility features and tools.

The process begins with the definition of item specifications, and continues with

- selection and training of item writers;
- writing and internal review of items;
- review by state personnel and stakeholder committees;
- markup for translation and accessibility features;
- field testing; and
- post-field-test reviews.

Each of these steps has a role in ensuring that the items can support the claims on which they will be based. Exhibit A describes how each step contributes to these goals and describes each step in the process in more detail.

Exhibit A. Summary of How Each Step of Development Supports the Validity of Claims

Developmental Steps	Supports Alignment to the Performance Expectations	Reduces Construct- Irrelevant Variance Through Universal Design	Expands Access Through Linguistic and Other Supports
Item specifications	Specifies item interactions, content limits, and guidelines for meeting task demands and levels of cognitive engagement requirements and adjusting difficulty.	Avoids the use of any item interactions with accessibility constraints and provides language guidelines. Allows for multiple response modes to accommodate different styles.	
Selection and training of item writers	Ensures that item writers have the background to understand the PEs and item specifications. Teaches item writers how to select item interactions for measurement and accessibility.	Training in language accessibility, bias, and sensitivity helps item writers avoid unnecessary barriers.	
Writing and internal review of items	Checks content alignment and evaluates and improves overall quality.	Eliminates editorial issues and flags and removes bias and accessibility issues.	
Markup for translation and accessibility features		Adds universal features, such as text-to-speech (TTS) for science, that reduce barriers.	Adds TTS, braille, American Sign Language (ASL), translations, and glossaries.
Review by state personnel and stakeholder committees	Checks content and cognitive complexity alignment; evaluates and improves overall quality.	Flags sensitivity issues.	
Field testing	Provides statistical checks on quality and flags issues.	Flags items that appear to function differently for subsequent review to identify issues.	May reveal usability or implementation issues with markup.
Post-field-test reviews	Provides final, more focused checks on flagged items. Rubric validation ensures that scoring reflects PEs.	Provides final, focused review on items flagged for differential item functioning (DIF).	

### 2.2 ITEM SPECIFICATIONS

CAI is working with a group of states, psychometricians, and science experts, including the authors of the Next Generation Science Standards (NGSS), to develop powerful innovative solutions to the challenges of measuring three-dimensional science standards based on the National Research Council's *A Framework for K–12 Science Education* published in 2012. Participating states included Connecticut, Hawaii, Idaho, Montana, Oregon, Rhode Island, Utah, Vermont, West Virginia, and Wyoming. New Hampshire, North Dakota, and South Dakota participate in some activities. This collaboration has yielded item specifications for PEs, sample item clusters for some specifications, and hundreds of science item clusters and stand-alone items in various stages of development. Under this collaboration, using guidelines for item specifications proposed by WestEd in collaboration with the Council of Chief State School Officers (CCSSO), state members, and content experts (CCSSO, 2015), states developed item specifications jointly.

Item specifications are documents designed to guide item writers as they craft test items and stakeholders as they review those items. These specifications are intended to serve as a roadmap for writers to facilitate the creation of items that are properly aligned to the three dimensions that comprise each science standard and that together form coherent item clusters and stand-alone items. Exhibit B provides a sample of the item specifications developed by content experts for a middle school Life Sciences PE. Item specifications in science include the following:

- **Performance Expectation.** This identifies the PE being assessed.
- **Dimensions.** This identifies the Science and Engineering Practices (SEPs), crosscutting concepts (CCCs), and Disciplinary Core Ideas (DCIs) that the PE assesses.
- Clarifications and Content Limits. This delineates the specific content that the PE measures and the parameters in which items must be developed to assess the PE accurately, including the lower and upper complexity limits of items. Specifically, content limits refine the intent of the PE and provide limits of what may be asked of test takers. For example, content limits may identify the specific formulae that students are expected to know or not know.
- Science Vocabulary. This section identifies the relevant technical words that students are expected to know, and related words that they are explicitly not expected to know. These categories should not be considered exhaustive, as the boundaries of relevance are ambiguous, and the list is limited by the imagination of the writers.
- Content/Phenomena. This section provides examples of the types of phenomena that would support the effective items related to the PE in question. In general, these are guideposts, and item writers seek comparable phenomena, rather than drawing on those within the documents.
- Task Demands. In this section, the PEs and associated evidence statements are broken down into specific task demands aligned to each PE. Task demands denote the specific ways in which students will provide evidence of their understanding of the concept or skill. Specifically, the task demands identify the types of interactions and activities that

item writers should employ. Each item should be clearly linked to one or more of the task demands, and the verbs guide the types of interactions writers might employ to elicit the student response.

Exhibit B. Sample Science Item Cluster Specifications for Middle School Life Sciences
Performance Expectation

Performance	MS-LS1-1 <sup>a</sup>					
Expectation	Conduct an investigation	Conduct an investigation to provide evidence that living things are made of cells; either				
	one cell or many different numbers and types of cells.					
Dimensions	Planning and Carrying	LS1.A: Structure and Function	Scale, Proportion, and			
	Out Investigations	All living things are made up	Quantity			
	<ul> <li>Conduct an</li> </ul>	of cells, which is the smallest	Phenomena that can			
	investigation to	unit that can be said to be	be observed at one			
	produce data to	alive. An organism may	scale may not be			
	serve as the basis	consist of one single cell	observable at			
	for evidence that	(unicellular) or many different	another scale.			
	meets the goals of	numbers and types of cells				
	an investigation.	(multicellular).				
Clarifications	Clarification Statements					
and Content	<ul> <li>Emphasis is on de</li> </ul>	eveloping evidence that living things	are made of cells,			
Limits	distinguishing bet	ween living and non-living things, and	d understanding that living			
	things may be ma	de of one cell or many varying cells.				
	Content Limits					
	Students do not need to know the following:					
		tures or functions of specific organell	es or different proteins			
	<ul> <li>Systems of specialized cells</li> </ul>					
	The mechanisms by which cells are alive					
	•	of DNA and proteins or of cell growth	n and division			
	<ul> <li>Endosymbiotic theory</li> </ul>					
	Histological procedures					
Science Vocabulary		ell, tissue, organ, system, organism hote, magnify, microscope, DNA, nucle	-			
Students are	membrane, algae, chlorop	plast(s), chromosome, cork				
Expected to						
Know						
Science		eiosis, genetics, cellular respiration, e	0,			
Vocabulary	protozoa, amoeba, histology, protista, archaea, nucleoid, plasmid, diatoms, cyanobacteria					
Students are						
Not Expected						
to Know						
		Phenomena				
Context/	Some example phenomer					
Phenomena	Plant leaves and roots have tiny box-like structures that can be seen under a					
	microscope.					
		an be seen swimming in samples of	pond water viewed through			
	a microscope.					
		a frog's body (e.g., muscles, skin, tor	ngue) are observed under a			
	microscope, and are seen to be composed of cells.					

- One-celled organisms (e.g., bacteria, protists) perform the eight necessary functions of life, but nothing smaller has been seen to do this.
- Swabs from the human cheek are observed under a microscope. Small cells can be seen.

This Performance Expectation and associated Evidence Statements support the following Task Demands.

Task Demands

- 1. Identify from a list, including distractors, the materials/tools needed for an investigation to find the smallest unit of life (cell).
- 2. Identify the outcome data that should be collected in an investigation of the smallest unit of living things.
- 3. Evaluate the sufficiency and limitations of data collected to explain that the smallest unit of living things is the cell.
- 4. Make and/or record observations about whether the sample contains cells.b
- 5. Interpret and/or communicate data from the investigation to determine if a specimen is alive.
- 6. Construct a statement to describe the overall trend suggested by the observed data.

Note. <sup>a</sup>MS-LS1-1 is the performance expectation code for Middle School Life Sciences 1-1.

<sup>b</sup>Denotes those task demands which are deemed appropriate for use in stand-alone item development.

The specifications help test developers create item clusters and stand-alone items that will support a range of difficulty, furthering the goal of measuring the full range of performance found in the population, but remaining on grade level.

#### 2.3 SELECTION AND TRAINING OF ITEM WRITERS

All item writers developing science items at CAI have at least a bachelor's degree, and many bring teaching experience. All item writers are trained in

- the principles of universal design;
- the appropriate use of item interactions; and
- the science item specifications.

Key materials are shown in Appendix A, Item Writer Training Materials. These include

- CAI's Language Accessibility, Bias, and Sensitivity Guidelines; and
- a training (presented using Microsoft PowerPoint) for the appropriate use of item interactions.

#### 2.4 INTERNAL REVIEW

CAI's test development structure uses highly effective units organized around each content area. Unit directors oversee team leaders who work with team members to ensure item quality and adherence to best practices. All team members, including item writers, are content-area experts.

Teams include senior content specialists who review items before client review and provide training and feedback for all content-area team members.

ICCR and MOU science items go through a rigorous, multiple-level internal review process before they are sent to external review. Staff members are trained to review items for both content and accessibility throughout the process. A sample item review checklist that our test developers use is included in Appendix B, Item Review Checklist. The ICCR and MOU science internal review cycle includes the following phases:

- Preliminary Review
- Scoring Entry and Review
- Content Review One
- Edit Review
- Senior Review

# 2.4.1 Preliminary Review

Team leads or senior content staff conduct Preliminary Review. Sometimes Preliminary Review is conducted in a group setting, led by a senior test developer. During the Preliminary Review process, team leads or senior content staff analyze items to ensure the following:

- The item aligns with the PE.
- The item matches the item specification for the skills being assessed.
- The item is based on a quality scientific phenomenon (i.e., it assesses something in a reasonable way and it is a discrete observation that grounds a scenario, which allows for the assessment of something worthwhile in a meaningful way).
- The item aligns appropriately with the task demands.
- The vocabulary used in the item is appropriate for the grade and subject matter.
- The item considers language accessibility, bias, and sensitivity.
- The content is accurate and straightforward.
- The graphic and stimulus materials are necessary to answer the question.
- The item follows the approved style guide.
- The stimulus is clear, concise, and succinct (i.e., it contains enough information to convey what is being asked, it is stated positively, and it does not rely on negatives—such as *no*, *not*, *none*, *never*—unless absolutely necessary).

For selected-response item interactions, test developers also check to ensure that the set of response options are

- as succinct and short as possible (without repeating text);
- parallel in structure, grammar, length, and content;
- sufficiently distinct from one another;
- all plausible (but with only one correct option); and
- free of obvious or subtle cuing.

# 2.4.2 Scoring Entry and Review

During Scoring Entry, the item writer inputs the machine scoring for review by the team lead or senior staff before the Content Review One Level. This step is separate from Preliminary Review to allow senior staff to suggest changes to the interaction at Preliminary Review without requiring the writer to overhaul scoring that they already created. This step also allows senior staff to ensure that the scoring suggested by the writer at Preliminary Review is appropriate. This process ensures that the scoring is entered once, streamlining the process. At this level, the scoring is analyzed to ensure the following criteria:

- The scoring works as intended (i.e., the student gets a point for ALL correct responses and no points for ALL incorrect responses).
- The student receives a point for every unique piece of information they reveal about their understanding through their responses.
- Dependent scoring between and within interactions is captured.
- The way in which the scoring is set up is unambiguous and matches the questions asked (i.e., if we ask students to round a number to a certain decimal place, we score accordingly).

The senior staff approves the intent of the scoring from the Preliminary Review. At the Scoring Entry level, the writer inputs this approved scoring, after which senior staff checks the functionality of the scoring. Once the scoring is determined to be working correctly, the senior staff signs off on it and moves it to Content Review One.

#### 2.4.3 Content Review One

Content Review One is conducted by a senior content specialist who was not part of the Preliminary Review. This reviewer carefully examines each item based on the same criteria identified for Preliminary Review. He or she also ensures that the revisions made during the Preliminary Review did not introduce errors or content inaccuracies. This reviewer approaches the item from the perspective of potential clients and his or her own experience in test development.

#### 2.4.4 Edit Review

During Edit Review, editors have four primary tasks:

1. Editors perform basic line editing for correct spelling, punctuation, grammar, and mathematical and scientific notation, ensuring consistency of style across the items.

- 2. Editors ensure that all items are accurate in content. Editors compare reading passages against the original publications to ensure that all information is internally consistent across stimulus materials and items, including names, facts, or cited lines of text that appear in the item. They ensure that the keys and all information in the item are correct. For items with mathematical tasks, editors perform all calculations to ensure accuracy.
- 3. Editors review all material for fairness and language accessibility issues.
- 4. Editors confirm that items reflect the accepted guidelines for good item construction. They examine all items for language that is simple, direct, and free of ambiguity with minimal verbal difficulty. Editors confirm that a problem or task and its stem are clearly defined and concisely worded with no unnecessary information. For multiple-choice interactions, editors check that options are parallel in structure and fit logically and grammatically with the stem. They also ensure that the key answers the question posed accurately and correctly, is not inappropriately obvious, and is the only correct answer to an item among the distractors. For constructed-response interactions, editors review the rubrics for appropriate style and grammar.

# 2.4.5 Senior Review

By the time a science item arrives at Senior Review, both content reviewers and editors have thoroughly vetted it. Senior reviewers (in particular, senior content specialists) look at the item's entire review history, ensuring that all the issues identified in that item have been adequately addressed. Senior reviewers verify the overall content of each item, confirming its accuracy, alignment to the PE, and consistency with expectations for the highest quality. They check whether the scoring is working as intended and scoring assertions adequately address the evidence the student provides with each type of response.

#### 2.5 REVIEW BY STATE PERSONNEL AND STAKEHOLDER COMMITTEES

All science items have been through an exhaustive external review process. Items in the Shared Science Assessment Item Bank were reviewed by content experts in one or several states and reviewed and approved by multiple stakeholder committees to evaluate both content and bias/sensitivity.

#### 2.5.1 State Review

After items have been developed for a state participating in the MOU, content experts from the state that owns the item review any eligible items before committee review. At this stage in the review process, clients can request edits, such as wording edits, scoring edits, alignment changes, or task demand updates. A CAI science content expert reviews all client-requested edits considering the science item specifications, other clients' requests, and existing items in the bank to determine whether the requested edits will be made. At this stage, clients have the option to present these items to the committee (based on the edits made) or withhold them from committee review.

ICCR items are reviewed by at least one or two states. The states provide feedback on the ICCR items, and CAI science leadership gathers suggestions and makes edits that improve the ICCR item. Not all suggestions are implemented, as these items are owned by CAI. Further,

most MOU states accept or reject ICCR and MOU items (as they appear at the time), to be presented to their committees. Some clients skip this step and allow CAI to review all items with their committees before reviewing them. These items can either be set for field testing in a future administration or become a part of the locked operational pool.

# 2.5.2 Content Advisory Committee Reviews

During the Content Advisory Committee (CAC) reviews, items are reviewed for content validity, grade-level appropriateness, and alignment to the PE. CAC members are typically grade-level and subject-matter experts. During this review, educators also ensure that the scoring assertions clearly identify what is being scored as correct and give credit where they should (see Section Error! Reference source not found.). Before the CAC review begins, CAI provides a presentation on the three-dimensional science standards, the item development process, the CAI systems that will be used in the review, and how to review the items for content. Appendix K, Content Advisory Committee Review Training Slides, provides the slides used during the CAC review training.

Items developed for each state under the MOU are reviewed by the state that owns the items. ICCR items are reviewed by the CAC of one or more states. In most cases, items are seen by multiple state committees prior to their field-test or operational use.

In 2021, MOU states were all involved in a single CAC process where participants from multiple states reviewed items. The items were edited and then returned to the owning state for final approval.

A summary of the committee meetings appears in Exhibit C, with further details about the participants in Appendix C, Content Advisory Committee Participant Details.

Exhibit C. Summary of Content Advisory Committee Meetings

State/Item Bank	Meeting	Number of Committee Members	Number of Items Reviewed
	February 2017	41	45
	May 2017	42	40
	October 2017	41	75
	November 2017	35	41
	January 2018	33	42
Connecticut	October 2018	45	84
	November 2018	49	235
	December 2018	32	56
	January 2019	44	65
	September 2019	50	60
	July 2021	b	24
Hawaii	July 2017	22	25

State/Item Bank	Meeting	Number of Committee Members	Number of Items Reviewed
	September 2017	20	65
	October 2018	29	85
	February 2019	21	44
ICCD	March 2018	26	152
ICCR	July 2021	b	164
ldaho	December 2018	21	111
Montana	July 2021	b	41
	January 2018	42	73
MOOA	March 2018	28	100
MSSA	January 2019	21	116
	July 2021	b	30
	August 2017	10	110
0	August 2018	20	257
Oregon	December 2018	16	62
	July 2021	b	22
	July 2017	23	55
Utah	December 2017	36	48
	July 2021	b	65
	January 2017	28ª	39
NATA A A Manada da	October 2018	10	191
West Virginia	July 2019	12	50
	July 2021	В	12
	December 2017	17	51
Wyoming	October 2018	14	37
	July 2021	b	32

*Note*. <sup>a</sup>Number of Committee Members includes total committee members for English language arts (ELA), mathematics, and science. The number for science-only committee members is not available.

# 2.5.3 Language Accessibility, Bias, and Sensitivity Committee Reviews

During the bias and sensitivity reviews, stakeholders review items to check for issues that might unfairly impact students based on their background. For example, some states include representatives from student populations such as Special Education, low vision, and the hearing impaired. Further, diverse members of this committee represent students of various ethnic and economic backgrounds to ensure that all items are free of bias and sensitivity concerns. Before the bias and sensitivity review begins, CAI provides a presentation on the

<sup>&</sup>lt;sup>b</sup>Multi-State review occurred over two weeks, with participants from multiple states involved. Items were reviewed by at least four participants.

three-dimensional science standards, the item development process, the CAI systems that will be used in the review, and how to review the items for fairness. Appendix L, Fairness Committee Review Training Slides, provides the slides used during the bias and sensitivity review training.

Due to the Covid-19 pandemic during 2020 and 2021, CAI reviewed items that contained references to virus, vaccine, bacteria, disease, infection, and related words and phrases. CAI content experts reviewed 65 items and rejected one item for sensitivity concerns.

In 2021, MOU states were all involved in a single review process where participants from multiple states would review items. The items were edited and then returned to the owning state for final approval.

A summary of the committee meetings appears in Exhibit D, with additional details about the participants in Appendix D, Fairness Committee Participant Details.

Exhibit D. Summary of Fairness Committee Meetings

State/Item Bank	Meeting	Number of Committee Members	Number of Items Reviewed	Number of Items Rejected
	February 2017	6	45	1
	December 2017	9	75	N/A
	December 2017	10	41	N/A
Connecticut	February 2018	3	42	N/A
Connecticut	November 2018	11	319	38
	December 2018	10	56	N/A
	January 2019	9	65	N/A
	September 2019	9	48	N/A <sup>a</sup>
	July 2017	22	25	2
11	September 2017	20	65	13
Hawaii	October 2018	29	85	6
	February 2019	21	44	0
IOOD	March 2018	13	152	N/A
ICCR	July 2021	С	124	5
ldaho	December 2018	15	111	1
Montana	July 2021	С	48	0
	January 2018	21	73	14
14004	March 2018	11	100	24
MSSA	January 2019	14	116	18
	July 2021	N/A	31	0
0	August 2017	5	110	5
Oregon	August 2018	9	256	56

State/Item Bank	Meeting	Number of Committee Members	Number of Items Reviewed	Number of Items Rejected
	December 2018	11	62	13
US Virgin Islands	October 2021	d	d	d
	August 2017	6	44	2
Utah	December 2017	6	48	1
	July/August 2021	С	56	2
	January 2017	28 <sup>b</sup>	34	N/A
West Virginia	January 2019	10	191	N/A
	July 2021	С	12	1
	December 2017	5	51	3
Wyoming	October 2018	5	37	N/A
	July 2021	С	41	0

Note. a Number of rejected items has not been finalized through client resolution at the time of writing this report.

# 2.5.4 Markup for Translation and Accessibility Features

After all approved state- and committee-recommended edits have been applied, the items are considered *locked* and ready for a portion of the accessibility tagging. TTS tagging is applied prior to field testing while Spanish translations and braille are applied post-field test. Accessibility markup is embedded into each item as part of the item development process rather than as a *post hoc* process applied to completed tests.

Accessibility markup, whether translations or for TTS, follow similar processes. One trained expert enters the markup, then a second expert reviews the work and recommends changes if necessary. If there is disagreement, a third expert is engaged to resolve the conflict.

Currently, science items are tagged with TTS. Spanish translations, including Spanish TTS and braille, are available for a subset of items.

#### 2.6 FIELD TESTING

A large pool of science field-test items was administered in the following nine states in spring 2018: Connecticut, Hawaii, New Hampshire, Oregon, Rhode Island, Utah, Vermont, West Virginia, and Wyoming. For Hawaii, Oregon, and Wyoming, items were embedded as field-test items in the legacy science test. Connecticut and Rhode Island conducted an independent field test in which all students participated, but no scores were reported. In New Hampshire, Utah, Vermont, and West Virginia, an operational field test was administered.

<sup>&</sup>lt;sup>b</sup>Number of Committee Members includes total committee members for ELA, mathematics, and science. The number for science only committee members is not available.

<sup>&</sup>lt;sup>c</sup>Multi-State review occurred over two weeks, with participants from multiple states involved. Items were reviewed by at least four participants.

<sup>&</sup>lt;sup>d</sup>U.S. Virgin Islands reviews were a review of previously accepted ICCR items by department staff.

In 2019, a second pool of field-test items was administered in the following nine states: Connecticut, Hawaii, Idaho, New Hampshire, Oregon, Rhode Island, Vermont, West Virginia, and Wyoming. For Hawaii, Idaho (elementary school), 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 Connecticut, New Hampshire, Oregon, Rhode Island, Vermont, and West Virginia, field-test items were administered as unscored items embedded within the operational items.

In 2021, a third wave of field-test items was administered in 12 states. An independent field test, in which students were administered a full set of items, was conducted for Idaho and Montana. Unscored field-test items were added as a separate segment to the operational (scored) legacy science test for Wyoming. 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 within the operational items.

CAI's field-test process is described in detail in Volume 1, Section 3.2.1, of this technical report.

#### 2.7 POST-FIELD-TEST REVIEW

Following the field test, items were subject to a substantial validation process. This included rubric validation and data review. These processes are described in Section 2.7.1, Rubric Validation, and Section 2.7.2, Data Review.

# 2.7.1 Rubric Validation

The validation process for the field-test items begins with rubric validation to verify and make any necessary revisions to the scoring rubrics. The rubric validation process occurs in two phases. During the first phase, CAI content experts work with the analysis team to prepare for the rubric validation meetings. The CAI content experts use the Rubric Evaluation and Verification for Items Scored Electronically (REVISE) system to generate student responses that are scientifically sampled to overrepresent responses most likely to have been mis-scored. Specifically, the sample overrepresents: (1) low-scored responses from otherwise high-scoring students, and (2) high-scored responses from otherwise low-scoring students. This process allows CAI to identify any potential scoring concerns before the rubric validation meeting, such as unanticipated (but accurate) responses, equivalent responses that were not originally considered, and responses receiving credit but should not (based on the content and the item rubric). At this point, the rubrics may be adjusted, and responses rescored.

The second phase of rubric validation involves committees of educators in each state. The committees review the response samples generated by CAI to make recommendations to change or to confirm the rubrics of each item. The committee recommendations are then discussed with the state of ownership to resolve any inconsistencies. The rubric is then edited or confirmed based on this resolution.

Exhibit E on the following page shows the features of REVISE.

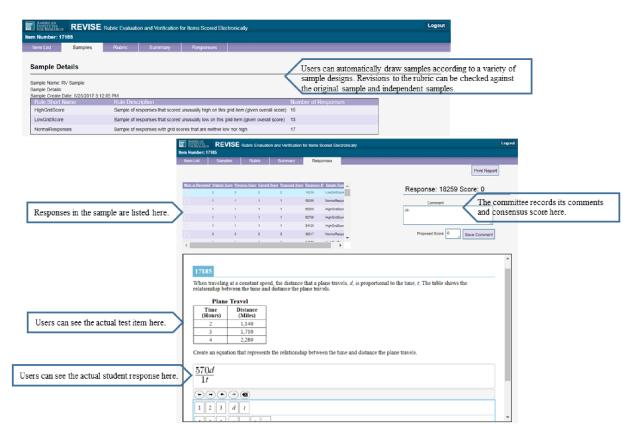


Exhibit E. Features of the REVISE Software

After the rubric validation meetings, CAI staff apply the approved revisions to the rubrics, and any items rejected as part of the process are rejected in ITS. ITS archives critical information regarding the scoring certification completed during the rubric validation process. This includes any rubric changes made during the scoring decision meetings and the sign-off completed by the senior content expert once the rubric has been changed, rescoring the entire sample, and the verification that the final rubric functioned as intended.

Following rubric validation, all items are subject to statistical checks, and flagged items are presented in data review committees.

#### 2.7.2 Data Review

Following rubric validation, all items are rescored and classical item statistics are computed for the scoring assertions, including item difficulty and item discrimination statistics, testing time, and differential item functioning (DIF) statistics. The states established standards for the statistics, and any items violating these standards are flagged for a second educator review. Even though the scoring assertions were the basic units of analysis to compute classical item statistics, the business rules to flag items for additional educator review were established at the item level, because assertions cannot be reviewed in isolation. A common set of business rules was defined for all the states participating in the field test. The classical item statistics were computed on the data of the students testing in the state that owned the item. For Rhode Island and Vermont, which share their item development, statistics were computed on the combined data of students testing in both states.

For ICCR items, the data from students testing in Connecticut, Hawaii, Idaho, Montana, New Hampshire, North Dakota, South Dakota, Rhode Island, Utah, Vermont, and West Virginia were combined (states that administered ICCR items and utilized either an independent field test or operational test).

Volume 1, Section 4, Annual Technical Report, describes in detail the statistical flags that send items to data review. The flags are designed to highlight potential content weaknesses, miskeys, or possible bias issues. Committee members are taught to interpret these flags and are given guidelines for examining the items for content or fairness issues.

For each of the states participating in the MOU, flagged items owned by the state were reviewed by a data review committee. The composition of the data review committees generally consisted of content experts from the state's department of education or state educators (in this case, the state educators were science teachers) and were supported by CAI content experts. ICCR field-test items were taken to committee members from several states participating in the MOU. The outcomes were decided by CAI science content leadership, taking the committees' recommendations into consideration.

At the start of each state-owned item data review meeting, CAI staff leads participants in a training session to familiarize them with the item development process, the purpose of the data review committee and the data review process, and the meaning of the various flags. Committee members are taught to interpret the various flags and are given guidelines for examining the items for content or fairness issues. The training includes a group review of item cards, which detail specific item attributes (including grade level and alignment to the science PEs, the content and rubric of the item, and various item statistics). A sample of the training materials used for these data review meetings is presented in Appendix E, Sample Data Review Training Materials. Participants use an online environment via laptop computers to review the items and interact with them in a manner similar to that of students, and to view the statistics associated with each item.

The items are then reviewed by the participants who are most familiar with the particular grade (band) level and the items' content domain. CAI content specialists, who are also well versed in item statistics, facilitate the discussion in each room with CAI psychometricians available to answer questions as they arise. At the end of each meeting day, CAI content specialists meet with the state content specialists to review the committee recommendations and decide whether to accept or reject the item for inclusion in the operational pool. Items that were rejected become eligible for potential changes and additional field test items.

Exhibit F summarizes the data review committee meetings. Details, including the composition of each committee, are presented in Appendix F, Data Review Committee Participant Details.

Number of Number of Number of Owner Meeting Committee **Item Type** Items Items **Members** Reviewed Rejected Connecticut August 2018 29 **Total** 18 11

Exhibit F. Summary of Data Review Committee Meetings

Owner	Meeting	Number of Committee Members	Item Type	Number of Items Reviewed	Number of Items Rejected
			Cluster	7	5
			Stand-Alone	11	6
	A 1 0040	00	Total	53	17
	August 2019	29	Cluster	14	6
			Stand-Alone	39	11
		0.5	Total	51	12
	August 2021 <sup>c</sup>	25	Cluster	8	2
			Stand-Alone	43	10
		40	Total	32	3
	August 2018	18	Cluster	7	1
			Stand-Alone	25	2
		18 25	Total	37	13
Hawaii	August 2019		Cluster	17	5
			Stand-Alone	20	8
	August 2021 °		Total	26	8
			Cluster	6	0
			Stand-Alone	20	8
			Total	84	8
	July 2018	18 N/Aª 25	Cluster	33	2
			Stand-Alone	51	6
	August 2019  August 2021d		Total	43	3
ICCR			Cluster	0	1
			Stand-Alone	43	2
			Total	75	6
			Cluster	11	2
			Stand-Alone	64	4
			Total	12	6
	August 2019	10	Cluster	4	3
Idaho			Stand-Alone	8	3
			Total	60	5
	August 2021 °	25	Cluster	26	1
			Stand-Alone	34	4
			Total	17	4
Montana	September 2021	4	Cluster	3	2
			Stand-Alone	14	2
		C.	Total	9	6
	August 2018	2 <sup>b</sup>	Cluster	2	0
MSSA			Stand-Alone	7	6
	August 2019	2 <sup>b</sup>	Total	14	4
		_	Cluster	2	1

Owner	Meeting	Number of Committee Members	Item Type	Number of Items Reviewed	Number of Items Rejected
			Stand-Alone	12	3
			Total	18	9
	August 2021 °	25	Cluster	4	4
			Stand-Alone	14	5
			Total	44	6
	September 2018	11	Cluster	28	5
Oregon			Stand-Alone	16	1
Oregon			Total	8	7
	August 2019	4	Cluster	1	1
			Stand-Alone	7	6
			Total	15	0
South Dakota <sup>d</sup>	September 2021	N/A <sup>e</sup>	Cluster	0	0
			Stand-Alone	15	0
	August 2018	16	Total	40	6
			Cluster	40	6
Utah			Stand-Alone	0	0
Otan	September 2021	6	Total	11	3
			Cluster	11	3
			Stand-Alone	0	0
			Total	3	1
	July 2018	4	Cluster	3	1
			Stand-Alone	0	0
		4	Total	7	6
West Virginia	September 2019		Cluster	1	1
			Stand-Alone	6	5
			Total	7	3
	August 2021 °	25	Cluster	1	1
			Stand-Alone	6	2
			Total	16	6
	October 2018	19	Cluster	6	1
			Stand-Alone	10	5
			Total	16	5
Wyoming	August 2019	10	Cluster	4	3
, 59	7 tagast 2010	70	Stand-Alone	12	2
		0-	Total	16	4
	August 2021 °	25	Cluster	3	1
			Stand-Alone	13	3

*Note.* <sup>a</sup>In summer 2019, ICCR field-test items were taken to Connecticut, Hawaii, and Idaho for committee review. <sup>b</sup>Conducted by Rhode Island Department of Education and Vermont Agency of Education science content experts. <sup>c</sup>Cross-state committee item data review.

### 3. SCIENCE ITEM BANK SUMMARY

Tests based on A Framework for K-12 Science Education (National Research Council, 2012) adopt a three-dimensional conceptualization of science understanding, including Science and Engineering Practices (SEPs), Crosscutting Concepts (CCCs), and Disciplinary Core Ideas (DCIs). Accordingly, the new science assessments are composed mostly of item clusters representing a series of interrelated student interactions directed towards describing, explaining, and predicting scientific phenomena. Some stand-alone items are added to increase the coverage of the test without increasing the testing time or testing burden.

CAI has built the Shared Science Assessment Item Bank in partnership with multiple states. The science item bank is robust and has been constructed to support multiple statewide science assessments. As described earlier, science items were written to the three-dimensional science standards. The Shared Science Assessment Item Bank is comprised of ICCR items and items developed for specific states, which are all shared with MOU partner states. These items follow the same specifications, test development processes, and review processes. 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. In 2019, 347 item clusters and stand-alone items were field tested, of which 268 were accepted and made available as operational items in 2020. In 2021, CAI field tested 545 item clusters and stand-alone items, of which 458 have passed rubric validation and item data review.

Each state using the Shared Science Assessment Item Bank selects items that are appropriately aligned and have passed required reviews (as described in Section Error! Reference source not found., Error! Reference source not found.) for use on its statewide assessment. The Shared Science Assessment Item Bank continues to grow as participating states continue to field test new items. Participating states collectively share the items and agree to field test new items each year.

# 3.1 CURRENT COMPOSITION OF THE SHARED SCIENCE ASSESSMENT ITEM BANK

The Shared Science Assessment Item Bank contains item clusters and stand-alone items. Item clusters represent a series of interrelated student interactions directed toward describing, explaining, and predicting scientific phenomena. Item clusters can consist of several item parts requiring the student to interact with the item in various ways. In addition, shorter items (stand-alone items) are included to increase the coverage of the assessments without also increasing testing time or testing burden.

Within each item (item cluster and stand-alone item), a series of explicit assertions is made about the knowledge and skills that a student has demonstrated based on specific features of the student's responses across multiple interactions. For example, a student may correctly graph data points indicating that they can construct a graph showing the relationship between two variables, but they may make an incorrect inference about the relationship between the two variables, therefore not supporting the assertion that the student can interpret relationships expressed graphically. Table 1

<sup>&</sup>lt;sup>d</sup>Legacy field-test items only.

<sup>&</sup>lt;sup>e</sup>State Department of Education review only.

lists the science interaction types. Examples of various interaction types can be found in Appendix G, Example Item Interactions.

Table 1. Science Interaction Types and Descriptions

Interaction Type	Associated Subtypes	Description
	Multiple-Choice	Traditional multiple-choice interaction allows the student to select a single option from several possible answer options.
Choice	Multi-Select	Traditional multi-select interaction (checkboxes) allows students to select one or more options from several possible answer choices.
	Simple Text Entry	Students type a response in a text box.
Text Entry	Embedded Text Entry	Students type their response in one or more text boxes that are embedded in a section of read-only text.
rext citily	Natural Language	Students are directed to provide a short, written response.
	Extended Response	Students are directed to provide a longer, written response in the form of an essay.
Table	Table Match	Interaction allows students to check a box to indicate if the information from a column header matches information from a row header.
	Table Input	Interaction solicits a student to complete tabular data.
	Edit Task	A student clicks a word and replaces it with another word that they type to revise a sentence.
Edit Task	Edit Task with Choice	A student clicks a word or phrase and chooses the replacement from several options.
	Edit Task Inline Choice	Drop-down menus are placed through the text, and a student chooses the correct option to complete the text.
	Selectable	Selectable hot-text interactions require students to select one or more text elements in the response area.
	Re-orderable	Re-orderable hot-text interactions require students to click and drag hot-text elements into a different order.
Hot-Text	Drag-from-Palette	Drag-from-palette hot-text interactions require students to drag elements from a palette into the available blank table cells or gaps (text boxes) in the response area.
	Custom	Custom hot-text interactions combine the functionality of the other hot-text interaction subtypes. Students responding to a custom hot-text interaction may need to select text elements, rearrange text elements, and/or drag text elements from a palette to blank table cells or drop targets in the response area.
Equation	N/A	Equation interactions require students to enter a response into input boxes. These boxes may stand alone, or they may be in

Interaction Type	Associated Subtypes	Description
		line with text or embedded in a table. The equation interaction may have an on-screen keypad that may consist of special mathematics characters. Students may also enter their response via a physical keyboard.
	Grid	Grid interactions require students to enter a response by interacting with a grid area in the answer space. The student may be required to draw a line or shape, plot a point, or create a graph. The student may also drag and drop or click selectable hot-spots.
Grid	Hot-Spot	Hot-spot interaction subtypes allow the student to create grid interactions with specific hot-spot functionality. These interactions require students to select hot-spot regions in the grid area.
	Graphic Gap Match	Graphic gap match interactions allow the student to create grid interactions with specific drag-and-drop functionality. These interactions require students to drag image objects from a palette to specified regions (gaps) in the grid area.
Simulation	N/A	Simulation interactions allow the student to investigate a phenomenon by selecting variables to get output data. Some simulations are accompanied by animations.

**Error!** Reference source not found.—

on the following pages provide the number of items in the Shared Science Assessment Item Bank available for use in the spring 2021 statewide assessments. Appendix H, Shared Science Assessment Item Bank provides the items available within the bank by grade band, performance expectation (PE), and origin.

Table 2. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank

Grade Band and Item Type	ICCR Items		MOU Items <sup>a</sup>	Total Bank Items	
Elementary School	130	24	285	439	
Cluster	41	13	165	219	
Stand-Alone	89	11	120	220	
Middle School	115	23	307	445	
Cluster	32	11	179	222	
Stand-Alone	83	12	128	223	
High School	122	16	232	370	
Cluster	43	6	96	145	
Stand-Alone	79	10	136	225	
Total	367	63	824	1254	

Note. <sup>a</sup>Other MOU states include Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming.

Table 3. Spring 2021 Shared Science Assessment Operational Item Bank

Grade Band and Item Type	ICCR Operational Items	MSSA Operational Items	MOU Operational Items <sup>a</sup>	Total Bank Operational Items
Elementary School	79	17	129	225
Cluster	32	9	72	113
Stand-Alone	47	8	57	112
Middle School	68	11	207	286
Cluster	24	5	133	162
Stand-Alone	44	6	74	124
High School	79	9	110	198
Cluster	28	4	56	88
Stand-Alone	51	5	54	110
Total	226	37	446	709

Note. <sup>a</sup>Other MOU operational item states include Connecticut, Hawaii, Idaho, Oregon, Utah, West Virginia, and Wyoming.

Table 4. Spring 2021 Shared Science Assessment Field-Test Item Bank

Grade Band and	ICCR Field-Test	MSSA Field-	MOU Field-Test	Total Bank
Item Type	Items	Test Items	Items <sup>a</sup>	Field-Test Items
Elementary School	51	7	156	214

Grade Band and Item Type	ICCR Field-Test Items	MSSA Field- Test Items	MOU Field-Test Items <sup>a</sup>	Total Bank Field-Test Items
Cluster	9	4	93	106
Stand-Alone	42	3	63	108
Middle School	47	12	100	159
Cluster	8	6	46	60
Stand-Alone	39	6	54	99
High School	43	7	122	172
Cluster	15	2	40	57
Stand-Alone	28	5	82	115
Total	141	26	378	545

Note. <sup>a</sup>Other MOU field-test item states include Connecticut, Hawaii, Idaho, Montana, Utah, West Virginia, and Wyoming.

Table 5. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank by Science Discipline

Grade Band	Science Discipline	Item Type	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items
	Earth and Space	Cluster	14	4	49	67
	Sciences	Stand-Alone	28	6	42	76
Elementary	Life Sciences	Cluster	14	4	51	69
School	Life Sciences	Stand-Alone	30	3	33	66
	Physical	Cluster	13	5	65	83
	Sciences	Stand-Alone	31	2	45	78
	Earth and Space Sciences	Cluster	11	3	47	61
		Stand-Alone	23	3	36	62
Middle	Life Sciences	Cluster	10	4	68	82
School		Stand-Alone	38	5	45	88
	Physical Sciences	Cluster	11	4	58	73
		Stand-Alone	22	4	46	72
	Earth and Space Sciences	Cluster	9	4	17	30
		Stand-Alone	12	4	29	45
High	Life Caianasa	Cluster	20	1	46	67
School	Life Sciences	Stand-Alone	49	3	55	107
	Physical	Cluster	14	1	32	47
	Sciences	Stand-Alone	18	3	52	73
Total	•	•	367	63	816 <sup>b</sup>	1246 <sup>b</sup>

*Note.* <sup>a</sup>Other MOU states include Hawaii, Idaho, Montana, MSSA (Rhode Island and Vermont), Oregon, Utah, West Virginia, and Wyoming . <sup>b</sup>Count excludes eight MOU items that do not align to the Next Generation Science Standards (NGSS).

Table 6. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank by Disciplinary Core Idea

Grade Band	Science Discipline	Disciplinary Core Idea	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items
	Earth and Space Sciences	ESS1	12	2	27	41
		ESS2	13	3	42	58
	Sciences	ESS3	17	5	22	44
		LS1	17	3	38	58
El	Life Sciences	LS2	5	1	15	21
Elementary School	Life Sciences	LS3	4	3	9	16
School		LS4	18	0	22	40
		PS1	12	4	31	47
	Dhysical Caianas	PS2	11	2	23	36
	Physical Sciences	PS3	17	1	37	55
		PS4	4	0	19	23
	Earth and Space Sciences	ESS1	15	1	23	39
		ESS2	9	2	31	42
		ESS3	10	3	29	42
	Life Sciences	LS1	10	5	40	55
		LS2	20	2	33	55
Middle School		LS3	4	0	14	18
		LS4	14	2	26	42
		PS1	9	3	32	44
	D	PS2	3	1	29	33
	Physical Sciences	PS3	14	3	24	41
		PS4	7	1	19	27
	Forth on 10	ESS1	7	3	15	25
	Earth and Space Sciences	ESS2	7	3	16	26
	Sciences	ESS3	7	2	15	24
Himb Californi		LS1	18	1	32	51
High School	Life Calerra	LS2	20	2	32	54
	Life Sciences	LS3	10	1	13	24
		LS4	21	0	24	45
	Physical Sciences	PS1	14	2	33	49

Grade Band	Science Discipline	Disciplinary Core Idea	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items
		PS2	8	1	19	28
		PS3	6	1	20	27
		PS4	4	0	12	16
Total			367	63	816 <sup>b</sup>	1246 <sup>b</sup>

Note. <sup>a</sup>Other MOU states include Hawaii, Idaho, Montana, MSSA (Rhode Island and Vermont), Oregon, Utah, West Virginia, and Wyoming. <sup>b</sup>Count excludes eight MOU items that do not align to the NGSS.

# 3.2 STRATEGY FOR BANK EVALUATION AND REPLENISHMENT

Both CAI and the participating MOU states continue to develop items to replenish and grow the Shared Science Assessment Item Bank. The general strategy for targeting item development gathers information from three sources:

- 1. Characteristics of released items to be replaced.
- 2. Characteristics of items that are overused.
- 3. Tabulations of content coverage and ranges of difficulty to identify gaps in the bank.

Before a test goes live, simulations are used to fine-tune the parameters of the algorithm that govern the item selection in a linear-on-the-fly test (LOFT) design. Among the many reports from the simulator are items that are seen by more than 20% of students. The characteristics of these items are the primary targets for development. Overused items become candidates for release in two years, once replacements have been introduced into the operational bank.

# 4. MULTI-STATE SCIENCE ASSESSMENT TEST CONSTRUCTION

# 4.1 TEST DESIGN

The Multi-State Science Assessment (MSSA) was administered online to students in grades 5, 8, and 11 using a linear-on-the-fly (LOFT) test design. Contrary to a fixed form, every student potentially sees a different set of items. Items are selected by an item selection algorithm so that the blueprint is met whenever possible. The algorithm that was used is the same algorithm that Cambium Assessment, Inc. (CAI) uses for the administration of adaptive tests. 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. Conversely, the content value decreases for items with content features that met the minimum. 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. By assigning a weight of zero to the information value of an item with respect to the underlying proficiency, the items are selected solely based on their contributions to meeting the blueprint. Details for CAI's adaptive testing algorithm are described in Appendix J, Adaptive Algorithm Design.

For the 2018 independent field test, a segmented design was used; items were administered grouped in four segments. The segments correspond to each of the three science disciplines and a (additional) field-test segment that could contain items from all three science disciplines.

In 2018, the order of the segments corresponding to the science disciplines was randomized over students. The additional field-test segment consisted of one item cluster and was always presented at the end of the test (segment four). The primary purpose was to collect additional student responses for the item clusters that had low exposure in the first three segments.

Starting from 2019, the scored operational part of the test consisted of the three segments corresponding to science disciplines. The embedded field-test segment consisted of two item clusters and four stand-alone items. In order to ensure that every student received exactly two item clusters and four stand-alone items as field-test items, the embedded field-test segment was split into two segments: one for field-test item clusters, and one for field-test stand-alone items.

The test was taken over two days. On the first day, half of the students received two operational segments, chosen at random from the three operational segments. The other half received one randomly chosen operational segment and the embedded field-test segments. The remaining segments were administered on the second day. Within one day, the order of the segments was randomized, with the restriction that the field-test segments for item clusters and stand-alone items were always administered right after each other.

# 4.2 TEST BLUEPRINTS

Test blueprints provide the following guidelines:

- Length of the test
- Science disciplines to be covered and the acceptable number of items across performance expectations (PEs) within each science discipline and Disciplinary Core Idea (DCI)

The blueprint for science is provided in

-Table 9.

Table 7. Science Test Blueprint, Grade 5

Grade 5	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
Discipline-Physical Sciences, PE Total = 17	2	2	4	4	6	6
DCI-Motion and Stability: Forces and Interactions	0	1	0	2	0	3
3-PS2-1: Forces-balanced and unbalanced forces	0	1	0	1	0	1
3-PS2-2: Forces-pattern predicts future motion	0	1	0	1	0	1
3-PS2-3: Forces-between objects not in contact	0	1	0	1	0	1
3-PS2-4: Forces-magnets*	0	1	0	1	0	1
5-PS2-1: Space Systems	0	1	0	1	0	1
DCI-Energy	0	1	0	2	0	3
4-PS3-1: Energy-relationship between speed and energy of object	0	1	0	1	0	1
4-PS3-2: Energy-transfer of energy	0	1	0	1	0	1
4-PS3-3: Energy-changes in energy when objects collide	0	1	0	1	0	1
4-PS3-4: Energy-converting energy from one form to another*	0	1	0	1	0	1
5-PS3-1: Matter and Energy	0	1	0	1	0	1
DCI–Waves and Their Applications in Technologies for Information Transfer	0	1	0	2	0	3
4-PS4-1: Waves-waves can cause objects to move	0	1	0	1	0	1
4-PS4-2: Structure, Function, Information Processing	0	1	0	1	0	1
4-PS4-3: Waves-using patterns to transfer information*	0	1	0	1	0	1
DCI-Matter and Its Interactions	0	1	0	2	0	3
5-PS1-1: Structure and Properties of Matter	0	1	0	1	0	1
5-PS1-2: Structure and Properties of Matter	0	1	0	1	0	1
5-PS1-3: Structure and Properties of Matter	0	1	0	1	0	1
5-PS1-4: Structure and Properties of Matter	0	1	0	1	0	1

Grade 5	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
Discipline-Life Sciences, PE Total = 12	2	2	4	4	6	6
DCI-From Molecules to Organisms: Structure and Function	0	1	0	2	0	3
3-LS1-1: Inheritance	0	1	0	1	0	1
4-LS1-1: Structure, Function, Information Processing	0	1	0	1	0	1
4-LS1-2: Structure, Function, Information Processing	0	1	0	1	0	1
5-LS1-1: Matter and Energy	0	1	0	1	0	1
DCI-Ecosystems: Interactions, Energy, and Dynamics	0	1	0	2	0	3
3-LS2-1: Ecosystems	0	1	0	1	0	1
5-LS2-1: Matter and Energy	0	1	0	1	0	1
DCI-Inheritance and Variation of Traits	0	1	0	2	0	3
3-LS3-1: Inheritance	0	1	0	1	0	1
3-LS3-2: Inheritance	0	1	0	1	0	1
DCI-Biological Evolution: Unity and Diversity	0	1	0	2	0	3
3-LS4-1: Ecosystems	0	1	0	1	0	1
3-LS4-2: Inheritance	0	1	0	1	0	1
3-LS4-3: Ecosystems	0	1	0	1	0	1
3-LS4-4: Ecosystems*	0	1	0	1	0	1
Discipline–Earth and Space Sciences, PE Total = 13	2	2	4	4	6	6
DCI-Earth's Systems	0	1	0	2	0	3
3-ESS2-1: Weather and Climate	0	1	0	1	0	1
3-ESS2-2: Weather and Climate	0	1	0	1	0	1
4-ESS2-1: Earth's Systems and Processes	0	1	0	1	0	1
4-ESS2-2: Earth's Systems and Processes	0	1	0	1	0	1
5-ESS2-1: Earth's Systems	0	1	0	1	0	1

Grade 5	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
5-ESS2-2: Earth's Systems	0	1	0	1	0	1
DCI-Earth and Human Activity	0	1	0	2	0	3
3-ESS3-1: Weather and Climate*	0	1	0	1	0	1
4-ESS3-2: Earth's Systems and Processes*	0	1	0	1	0	1
4-ESS3-1: Energy	0	1	0	1	0	1
5-ESS3-1: Earth's Systems	0	1	0	1	0	1
DCI-Earth's Place in the Universe	0	1	0	2	0	3
4-ESS1-1: Earth's Systems and Processes	0	1	0	1	0	1
5-ESS1-1: Space Systems	0	1	0	1	0	1
5-ESS1-2: Space Systems	0	1	0	1	0	1
PE Total = 42	6	6	12	12	18	18

*Note.* \* These PEs have an engineering component.

Table 8. Science Test Blueprint, Grade 8

Grade 8	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
Discipline-Physical Sciences, PE Total = 19	2	2	4	4	6	6
DCI-Matter and Its Interactions	0	1	0	2	0	3
MS-PS1-1: Structure and Properties of Matter	0	1	0	1	0	1
MS-PS1-2: Chemical Reactions	0	1	0	1	0	1
MS-PS1-3: Structure and Properties of Matter	0	1	0	1	0	1
MS-PS1-4: Structure and Properties of Matter	0	1	0	1	0	1
MS-PS1-5: Chemical Reactions	0	1	0	1	0	1
MS-PS1-6: Chemical Reactions*	0	1	0	1	0	1
DCI-Motion and Stability: Forces and Interactions	0	1	0	2	0	3
MS-PS2-1: Forces and Interactions*	0	1	0	1	0	1
MS-PS2-2: Forces and Interactions	0	1	0	1	0	1
MS-PS2-3: Forces and Interactions	0	1	0	1	0	1
MS-PS2-4: Forces and Interactions	0	1	0	1	0	1
MS-PS2-5: Forces and Interactions	0	1	0	1	0	1
DCI-Energy	0	1	0	2	0	3
MS-PS3-1: Energy	0	1	0	1	0	1
MS-PS3-2: Energy	0	1	0	1	0	1
MS-PS3-3: Energy*	0	1	0	1	0	1
MS-PS3-4: Energy	0	1	0	1	0	1
MS-PS3-5: Energy	0	1	0	1	0	1
DCI–Waves and Their Applications in Technologies for Information Transfer	0	1	0	2	0	3
MS-PS4-1: Waves and Electromagnetic Radiation	0	1	0	1	0	1
MS-PS4-2: Waves and Electromagnetic Radiation	0	1	0	1	0	1

Grade 8	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
MS-PS4-3: Waves and Electromagnetic Radiation	0	1	0	1	0	1
Discipline-Life Sciences, PE Total = 21	2	2	4	4	6	6
DCI–From Molecules to Organisms: Structures and Processes	0	1	0	2	0	3
MS-LS1-1: Structure, Function, Information Processing	0	1	0	1	0	1
MS-LS1-2: Structure, Function, Information Processing	0	1	0	1	0	1
MS-LS1-3: Structure, Function, Information Processing	0	1	0	1	0	1
MS-LS1-4: Growth, Development, Reproduction	0	1	0	1	0	1
MS-LS1-5: Growth, Development, Reproduction	0	1	0	1	0	1
MS-LS1-6: Matter and Energy	0	1	0	1	0	1
MS-LS1-7: Matter and Energy	0	1	0	1	0	1
MS-LS1-8: Structure, Function, Information Processing	0	1	0	1	0	1
DCI–Ecosystems: Interactions, Energy, and Dynamics	0	1	0	2	0	3
MS-LS2-1: Matter and Energy	0	1	0	1	0	1
MS-LS2-2: Interdependent Relationships in Ecosystems	0	1	0	1	0	1
MS-LS2-3: Matter and Energy	0	1	0	1	0	1
MS-LS2-4: Matter and Energy	0	1	0	1	0	1
MS-LS2-5: Interdependent Relationships in Ecosystems*	0	1	0	1	0	1
DCI-Hereditary: Inheritance and Variation of Traits	0	1	0	2	0	3
MS-LS3-1: Growth, Development, Reproduction	0	1	0	1	0	1
MS-LS3-2: Growth, Development, Reproduction	0	1	0	1	0	1
DCI-Biological Evolution: Unity and Diversity	0	1	0	2	0	3
MS-LS4-1: Natural Selection and Adaptation	0	1	0	1	0	1
MS-LS4-2: Natural Selection and Adaptation	0	1	0	1	0	1
MS-LS4-3: Natural Selection and Adaptation	0	1	0	1	0	1
MS-LS4-4: Natural Selection and Adaptation	0	1	0	1	0	1

Grade 8	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
MS-LS4-5: Growth, Development, Reproduction	0	1	0	1	0	1
MS-LS4-6: Natural Selection and Adaptation	0	1	0	1	0	1
Discipline–Earth and Space Sciences, PE Total = 15	2	2	4	4	6	6
DCI-Earth's Place in the Universe	0	1	0	2	0	3
MS-ESS1-1: Space Systems	0	1	0	1	0	1
MS-ESS1-2: Space Systems	0	1	0	1	0	1
MS-ESS1-3: Space Systems	0	1	0	1	0	1
MS-ESS1-4: History of Earth	0	1	0	1	0	1
DCI-Earth's Systems	0	1	0	2	0	3
MS-ESS2-1: Earth's Systems	0	1	0	1	0	1
MS-ESS2-2: History of Earth	0	1	0	1	0	1
MS-ESS2-3: History of Earth	0	1	0	1	0	1
MS-ESS2-4: Earth's Systems	0	1	0	1	0	1
MS-ESS2-5: Weather and Climate	0	1	0	1	0	1
MS-ESS2-6: Weather and Climate	0	1	0	1	0	1
DCI-Earth and Human Activity	0	1	0	2	0	3
MS-ESS3-1: Earth's Systems	0	1	0	1	0	1
MS-ESS3-2: Human Impacts	0	1	0	1	0	1
MS-ESS3-3: Human Impacts*	0	1	0	1	0	1
MS-ESS3-4: Human Impacts	0	1	0	1	0	1
MS-ESS3-5: Weather and Climate	0	1	0	1	0	1
PE Total = 55	6	6	12	12	18	18

Note. \* These PEs have an engineering component.

Table 9. Science Test Blueprint, Grade 11

Grade 11	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
Discipline-Physical Sciences, PE Total = 24	2	2	4	4	6	6
DCI-Matter and Its Interactions	0	1	0	2	0	3
HS-PS1-1: Structure and Properties of Matter	0	1	0	1	0	1
HS-PS1-2: Structure and Properties of Matter	0	1	0	1	0	1
HS-PS1-3: Structure and Properties of Matter	0	1	0	1	0	1
HS-PS1-4: Chemical Reactions	0	1	0	1	0	1
HS-PS1-5: Chemical Reactions	0	1	0	1	0	1
HS-PS1-6: Chemical Reactions*	0	1	0	1	0	1
HS-PS1-7: Chemical Reactions	0	1	0	1	0	1
HS-PS1-8: Nuclear Processes	0	1	0	1	0	1
DCI-Motion and Stability: Forces and Interactions	0	1	0	2	0	3
HS-PS2-1: Forces and Motion	0	1	0	1	0	1
HS-PS2-2: Forces and Motion	0	1	0	1	0	1
HS-PS2-3: Forces and Motion*	0	1	0	1	0	1
HS-PS2-4: Types of Interactions	0	1	0	1	0	1
HS-PS2-5: Types of Interactions	0	1	0	1	0	1
HS-PS2-6: Chemical Reactions*	0	1	0	1	0	1
DCI–Energy	0	1	0	2	0	3
HS-PS3-1: Energy	0	1	0	1	0	1
HS-PS3-2: Energy	0	1	0	1	0	1
HS-PS3-3: Energy*	0	1	0	1	0	1
HS-PS3-4: Energy	0	1	0	1	0	1
HS-PS3-5: Energy	0	1	0	1	0	1
DCI-Waves and Their Applications in Technologies for	0	1	0	2	0	3

Grade 11	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
Information Transfer						
HS-PS4-1: Wave Properties	0	1	0	1	0	1
HS-PS4-2: Wave Properties	0	1	0	1	0	1
HS-PS4-3: Wave Properties/Electromagnetic Radiation	0	1	0	1	0	1
HS-PS4-4: Electromagnetic Radiation	0	1	0	1	0	1
HS-PS4-5: Electromagnetic Radiation*	0	1	0	1	0	1
Discipline-Life Sciences, PE Total = 24	2	2	4	4	6	6
DCI–From Molecules to Organisms: Structures and Processes	0	1	0	2	0	3
HS-LS1-1: Structure and Function	0	1	0	1	0	1
HS-LS1-2: Structure and Function	0	1	0	1	0	1
HS-LS1-3: Structure and Function	0	1	0	1	0	1
HS-LS1-4: Growth and Development of Organisms	0	1	0	1	0	1
HS-LS1-5: Organization for Matter and Energy Flow in Organisms	0	1	0	1	0	1
HS-LS1-6: Organization for Matter and Energy Flow in Organisms	0	1	0	1	0	1
HS-LS1-7: Organization for Matter and Energy Flow in Organisms	0	1	0	1	0	1
DCI-Ecosystems: Interactions, Energy, and Dynamics	0	1	0	2	0	3
HS-LS2-1: Interdependent Relationships in Ecosystems	0	1	0	1	0	1
HS-LS2-2: Interdependent Relationships in Ecosystems	0	1	0	1	0	1
HS-LS2-3: Cycles of Matter and Energy Transfer in Ecosystems	0	1	0	1	0	1
HS-LS2-4: Cycles of Matter and Energy Transfer in Ecosystems	0	1	0	1	0	1
HS-LS2-5: Cycles of Matter and Energy Transfer in Ecosystems	0	1	0	1	0	1
HS-LS2-6: Ecosystem Dynamics, Functioning, and Resilience	0	1	0	1	0	1
HS-LS2-7: Ecosystem Dynamics, Functioning, and Resilience*	0	1	0	1	0	1
HS-LS2-8: Social Interactions and Group Behavior	0	1	0	1	0	1
DCI-Heredity: Inheritance and Variation of Traits	0	1	0	2	0	3

Grade 11	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
HS-LS3-1: Structure and Function	0	1	0	1	0	1
HS-LS3-2: Variation of Traits	0	1	0	1	0	1
HS-LS3-3: Variation of Traits	0	1	0	1	0	1
DCI-Biological Evolution: Unity and Diversity	0	1	0	2	0	3
HS-LS4-1: Evidence of Common Ancestry and Diversity	0	1	0	1	0	1
HS-LS4-2: Natural Selection	0	1	0	1	0	1
HS-LS4-3: Natural Selection	0	1	0	1	0	1
HS-LS4-4: Adaptation	0	1	0	1	0	1
HS-LS4-5: Adaptation	0	1	0	1	0	1
HS-LS4-6: Adaptation*	0	1	0	1	0	1
Discipline–Earth and Space Sciences, PE Total = 19	2	2	4	4	6	6
DCI-Earth's Place in the Universe	0	1	0	2	0	3
HS-ESS1-1: The Universe and Its Stars	0	1	0	1	0	1
HS-ESS1-2: The Universe and Its Stars	0	1	0	1	0	1
HS-ESS1-3: The Universe and Its Stars	0	1	0	1	0	1
HS-ESS1-4: Earth and the Solar System	0	1	0	1	0	1
HS-ESS1-5: The History of Planet Earth	0	1	0	1	0	1
HS-ESS1-6: The History of Planet Earth	0	1	0	1	0	1
DCI-Earth's Systems	0	1	0	2	0	3
HS-ESS2-1: Earth Materials and Systems	0	1	0	1	0	1
HS-ESS2-2: Earth Materials and Systems	0	1	0	1	0	1
HS-ESS2-3: Earth Materials and Systems	0	1	0	1	0	1
HS-ESS2-4: Weather and Climate	0	1	0	1	0	1
HS-ESS2-5: The Roles of Water in Earth's Surface Processes	0	1	0	1	0	1
HS-ESS2-6: Weather and Climate	0	1	0	1	0	1

Grade 11	Min Item Clusters	Max Item Clusters	Min Stand-Alone Items	Max Stand-Alone Items	Min Item Clusters + Min Stand-Alone Items	Max Item Clusters + Max Stand-Alone Items
HS-ESS2-7: Weather and Climate	0	1	0	1	0	1
DCI-Earth and Human Activity	0	1	0	2	0	3
HS-ESS3-1: Natural Resources	0	1	0	1	0	1
HS-ESS3-2: Natural Resources*	0	1	0	1	0	1
HS-ESS3-3: Human Impacts on Earth Systems	0	1	0	1	0	1
HS-ESS3-4: Human Impacts on Earth Systems*	0	1	0	1	0	1
HS-ESS3-5: Global Climate Change	0	1	0	1	0	1
HS-ESS3-6: Global Climate Change*	0	1	0	1	0	1
PE Total = 67	6	6	12	12	18	18

*Note.* \*These PEs have an engineering component.

Main characteristics of the blueprint were that any PE could be tested only once (indicated by the values of 0 and 1 for the Min and Max values of the individual PEs in					

-Table 9); in general, no more than one item cluster or two stand-alone items could be sampled from the same DCI, and no more than three total items could be sampled from the same DCI (as indicated by the Min and Max values in the rows representing DCIs).

While tests are not timed, the Rhode Island Department of Education (RIDE) and the Vermont Agency of Education (VT AOE) published estimated testing times for the MSSA. Combined percentile 85 of testing times are presented in **Error! Reference source not found.**, Rhode Island percentile 85 of testing times are presented in Table 11, and Vermont percentile 85 of testing times are presented in Table 12.

Subject	Grade	85th Percentile Testing
	5	119.18
Science	8	111.98
	11	108 12

Table 10. Combined Percentile 85 Testing Times by Grade

Table 11. Rhode Island Percentile 85 Testing Times by Grade

Subject	Grade	85th Percentile Testing
5		123.40
Science	8	112.25
	11	109.45

Table 12. Vermont Percentile 85 Testing Times by Grade

Subject	Grade	85th Percentile Testing
	5	110.63
Science	8	111.45
	11	105.31

#### 4.3 Online Test Construction

During fall 2020, CAI psychometricians and content experts worked with RIDE and VT AOE content specialists and leadership to build item pools for the spring 2021 administration. The MSSA test construction uses a structured test construction plan, explicit blueprints, and active collaborative participation from all parties.

The 2021 MSSA item pools were built by CAI test developers to match items exactly to the detailed test blueprints. Operational items were selected from nine item banks (ICCR, Connecticut, Hawaii, Idaho, MSSA, Oregon, Utah, West Virginia, and Wyoming) to fulfill the blueprint for that grade. Table 13—Table 17 on the following pages summarize the 2021 MSSA item pool. Appendix I, Multi-State Assessment Item Pool provides the 2021 MSSA item pool by grade, PE, and origin.

Table 13. MSSA Spring 2021 Operational and Field-Test Item Pool

Grade Band and Item Type	ICCR Items <sup>a</sup>	MSSA Items	MOU Items <sup>b</sup>	Total Pool Items
Elementary School	73	24	94	191
Cluster	26	13	49	88
Stand-Alone	47	11	45	103
Middle School	61	23	124	208
Cluster	19	11	68	98
Stand-Alone	42	12	56	110
High School	79	16	77	172
Cluster	33	6	39	78
Stand-Alone	46	10	38	94
Total	213	63	295	571

*Note.* <sup>a</sup>Includes 14 ICCR operational items only administered in Rhode Island. <sup>b</sup>Other MOU state items administered includes Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming.

Table 14. MSSA Spring 2021 Operational Item Pool

Grade Band and Item Type	ICCR Operational Items <sup>a</sup>	MSSA Operational Items	MOU Operational Items <sup>b</sup>	Total Operational Pool Items
Elementary School	60	17	59	136
Cluster	25	9	36	70
Stand-Alone	35	8	23	66
Middle School	49	11	94	154
Cluster	16	5	59	80
Stand-Alone	33	6	35	74
High School	58	9	56	123
Cluster	25	4	33	62
Stand-Alone	33	5	23	61
Total	167	37	209	413

*Note.* <sup>a</sup>Includes 14 ICCR operational items only administered in Rhode Island. <sup>b</sup>Other MOU state operational items administered includes Connecticut, Hawaii, Oregon, Utah, West Virginia, and Wyoming.

Table 15. MSSA Spring 2021 Field-Test Item Pool

Grade Band and Item Type	ICCR Field-Test Items	MSSA Field- Test Items	MOU Field-Test Items <sup>a</sup>	Total Field-Test Pool Items
Elementary School	13	7	35	55
Cluster	1	4	13	18
Stand-Alone	12	3	22	37
Middle School	12	12	30	54
Cluster	3	6	9	18
Stand-Alone	9	6	21	36
High School	21	7	21	49
Cluster	8	2	6	16
Stand-Alone	13	5	15	33
Total	46	26	86	158

*Note.* <sup>a</sup>Other MOU state field-test items administered includes Hawaii, Idaho, Montana, Utah, West Virginia, and Wyoming.

Table 16. MSSA Spring 2021 Operational and Field-Test Item Pool by Science Discipline

Grade	Science Discipline	Item Type	ICCR Items <sup>a</sup>	MSSA Items	MOU Items <sup>b</sup>	Total Pool Items
'	Earth and Space	Cluster	9	4	13	26
	Sciences	Stand-Alone	13	6	11	30
Grade 5	Life Sciences	Cluster	7	4	17	28
Grade 5	Life Sciences	Stand-Alone	15	3	12	30
	Physical	Cluster	10	5	19	34
	Sciences	Stand-Alone	19	2	22	43
	Earth and Space	Cluster	7	3	20	30
	Sciences	Stand-Alone	11	3	17	31
Grade 8	Life Sciences	Cluster	5	4	28	37
Graue o	Life Sciences	Stand-Alone	19	5	17	41
	Physical	Cluster	7	4	20	31
	Sciences	Stand-Alone	12	4	22	38
	Earth and Space	Cluster	8	4	8	20
	Sciences	Stand-Alone	11	4	11	26
Grade 11	Life Sciences	Cluster	15	1	16	32
Graue 11	Life Sciences	Stand-Alone	21	3	11	35
	Physical	Cluster	10	1	15	26
	Sciences	Stand-Alone	14	3	16	33
Total		-	213	63	295	571

<i>Note.</i> <sup>a</sup> Includes 14 ICCR operational items only administered in Rhode Island. <sup>b</sup> Other MOU state items administerincludes Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming.									

Table 17. MSSA Spring 2021 Operational and Field-Test Item Pool by Disciplinary Core Idea

Grade	Science Discipline	Disciplinary Core Idea	ICCR Items <sup>a</sup>	MSSA Items	MOU Items <sup>b</sup>	Total Pool Items
	Farth and Coasa	ESS1	6	2	9	17
	Earth and Space Sciences	ESS2	9	3	12	24
	Sciences	ESS3	7	5	3	15
		LS1	8	3	12	23
	Life Sciences	LS2	4	1	4	9
Grade 5	Life Sciences	LS3	2	3	6	11
		LS4	8	0	7	15
		PS1	7	4	9	20
	Physical	PS2	7	2	8	17
	Sciences	PS3	13	1	15	29
		PS4	2	0	9	11
	Earth and Space Sciences	ESS1	5	1	12	18
		ESS2	5	2	18	25
		ESS3	8	3	7	18
		LS1	5	5	18	28
	Life October	LS2	8	2	12	22
Grade 8	Life Sciences	LS3	2	0	6	8
		LS4	9	2	9	20
		PS1	5	3	15	23
	Physical	PS2	2	1	10	13
	Sciences	PS3	8	3	10	21
		PS4	4	1	7	12
		ESS1	7	3	6	16
	Earth and Space	ESS2	5	3	7	15
	Sciences	ESS3	7	2	6	15
Grade 11		LS1	9	1	7	17
	1.6.0	LS2	12	2	8	22
	Life Sciences	LS3	5	1	3	9
		LS4	10	0	9	19

Grade	Science Discipline	Disciplinary Core Idea	ICCR Items <sup>a</sup>	MSSA Items	MOU Items <sup>b</sup>	Total Pool Items
		PS1	11	2	11	24
	Physical	PS2	7	1	9	17
	Sciences	PS3	4	1	7	12
		PS4	2	0	4	6
Total			213	63	295	571

Note. <sup>a</sup>Includes 14 ICCR operational items only administered in Rhode Island. <sup>b</sup>Other MOU state items administered includes Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming.

More information about *p*-values, biserial correlations, and item response theory (IRT) parameters can be found in Volume 1, Annual Technical Report. The details on calibration, equating, and scoring of the MSSA can also be found in Volume 1.

#### 4.4 PAPER-PENCIL ACCOMMODATION FORM CONSTRUCTION

Student scores should not depend upon the mode of administration or type of test form. Because the MSSA was primarily administered in an online test system in spring 2021, only one student took the paper-pencil form in grade 5 and one in grade 8. Scores obtained via alternate modes of administration must be established as comparable to scores obtained through online testing. This section outlines the overall test development plans that ensured the comparability of online and paper-pencil tests.

To build paper-pencil forms, content specialists began with the online pool and removed any items that could not be rendered on paper. Next, content specialists constructed fixed forms adhering to the test blueprint. In spring 2021, the paper-pencil forms met all blueprint requirements.

#### 5. SIMULATION SUMMARY REPORT

This section describes the results of simulated test administrations used to configure and evaluate the adequacy of the item selection algorithm used to administer the 2020–2021 Multi-State Science Assessments (MSSA) for grades 5, 8, and 11. Simulations were carried out to configure the settings of the algorithm and to evaluate whether individual tests adhered to the test blueprint.

Some important settings included cset1 and cset2, which represent subsets of the item pool that were eligible for item selection. See Appendix J, Adaptive Algorithm Design, for more details of the current item selection algorithm. In spring 2021, cset1 and cset2 values were set to 5 and 1. Psychometricians reviewed the simulation results and configured settings based on some key diagnostics, including:

- Match-to-Test Blueprint. Determines that the tests have the correct number of test items overall and the appropriate proportion by content categories at each level of the content hierarchy, as specified in the test blueprints for every science grade.
- **Item Exposure Rate.** Evaluates the utility of item pools and identifies overexposed and underexposed items.

These diagnostics are interrelated. For example, if the test pool for a particular content category is limited (i.e., there are only a few test items available), achieving a 100% match to the blueprint for this content level will lead to a high item exposure rate, which means that a large number of students are sharing items. The software system that performs the simulation allows the adjustment of setting parameters to attain the best possible balance among these diagnostics. The simulation involves an iterative process that reviews initial results, adjusts these system parameters, runs new simulations, reviews the new results, and repeats the exercise until an optimal balance is achieved. The final setting would then be applied for the operational tests.

#### 5.1 FACTORS AFFECTING SIMULATION RESULTS

There are several factors that may influence simulation results for a linear-on-the-fly (LOFT) test administration. These include the following:

- The proportional relationship between the pool and the constraints to be met. Proportionally distributed pools tend to make better use of the pool (i.e., more uniform item exposure) and make it easier to meet blueprint and other constraints. For example, if the specifications call for at least one item cluster per Disciplinary Core Idea (DCI), but the pool has no item cluster for some DCIs, it may be impossible to meet this constraint.
- The correlational structure between constraints. It is easier to satisfy a constraint if there are instances of the constraint at all levels of another constraint. For example, if stand-alone items within a discipline are associated only with a specific DCI, it may be difficult to meet both the desired distribution of content and the desired distribution of item type.
- Whether or not there is a strict maximum on a given constraint. This means that the requirement must be met exactly in each test administration.

#### 5.2 RESULTS OF SIMULATED TEST ADMINISTRATIONS: ENGLISH

This section presents the simulation results for the English online tests, which is the test taken by the majority of all students (94.14%). Simulations were evaluated for all content areas using 1,000 simulated cases per grade.

#### 5.2.1 Summary of Blueprint Match

The simulation results showed no blueprint violations at all content levels for all three grades.

#### 5.2.2 Item Exposure

The simulator output also reports the degree to which the constraints set forth in the blueprints may yield greater exposure of items to students. This is reported by examining the percentage of test administrations in which an item appears. For instance, in a fixed paper-pencil form, 100% of the items appear on 100% of the test administrations because every test taker takes the same form. In an adaptive test or a LOFT test with a sufficiently large item pool, we would expect that most of the items would appear on a relatively small percentage of the test administrations only.

When this condition holds, it suggests that test administrations between students are more or less unique. Therefore, we calculated the item exposure rate for each item across by dividing the total number of test administrations in which an item appears by the total number of tests administered. Then we report the distribution of the item exposure rate (r) in six bins. The bins are r=0% (unused), 0%<r<=1%, 1%<r<=5%, 5%<r<=20%, 20%<r<=40%, 40%<r<=60%, 60%<r<=80%, and 80%<r<=100%. If global item exposure is minimal, we would expect the largest proportion of items to appear in the bins of 0%<r<=20%, an indication that most of the items appear on a very small percentage of the test forms.

Table 18 presents the percentage of items that falls into each exposure bin for all grades. Most test items (98% or more) are administered in 1%–40% of the test administrations. No item has an

exposure rate less than 1% and the minimum exposure rate is 3% in grade 5. A few items had an exposure rate higher than 60% because of the limitation of the current pool for some content categories.

Table 18. Item Exposure Rates by Grade: Percentage of Items by Exposure Rate,
Across All English Online Simulation Sessions

Grade	Total Items	[0,0]%	[0,1]%	[1,5]%	[5,20]%	[20,40] %	[40,60] %	[60,80] %	[80,100] %
5	130	-	-	6.15	80	12.31	0	0.77	0.77
8	146	-	-	6.16	84.25	8.22	1.37	0	0
11	118	-	-	10.17	62.71	25.42	0	0	1.69

#### 5.3 RESULTS OF SIMULATED TEST ADMINISTRATIONS: SPANISH

This section presents the simulation results for the Spanish tests. The Spanish item pool consists of a subset of Independent College and Career Readiness (ICCR) items and some Memorandum of Understanding (MOU) items that has a Spanish translation available. Table 1919 presents the numbers of items available for the Spanish tests.

Table 19. Spring 2019 Spanish Operational Item Pool

Grade	Item Type	Total Number of Items
	Cluster	11
5	Stand-Alone	23
8	Cluster	7
•	Stand-Alone	19
11	Cluster	8
	Stand-Alone	20
Total		88

Simulations were evaluated for all content areas using 1,000 simulated cases per grade.

#### 5.3.1 Summary of Blueprint Match

There was no blueprint violation at the discipline level for all three grades.

#### **5.3.2** Item Exposure

Table 20 presents the percentage of items that falls into each exposure bin for all grades. More than 90% of all test items were administered in more than 20% of the test administrations across the three grades. Some items had an exposure rate of 100% because of the limited Spanish item pool. Only those items were available to satisfy the blueprint constraints.

Table 20. Item Exposure Rates by Grade: Percentage of Items by Exposure Rate, Across All Spanish Simulation Sessions

Grade	Total Items	[0,0]%	[0,1]%	[1,5]%	[5,20]%	[20,40]%	[40,60]%	[60,80]%	[80,100]%
5	34	0	0	0	8.82	29.41	26.47	17.65	17.65
8	26	0	0	0	0	23.08	23.08	15.38	38.46
11	28	0	0	0	0	25	28.57	14.29	32.14

#### 6. OPERATIONAL TEST ADMINISTRATION SUMMARY REPORT

This section presents the blueprint match reports and item exposure rates for the spring 2021 operational test administrations.

#### **6.1 BLUEPRINT MATCH**

All tests in all grades met the blueprint specifications with a 100% match at all content levels.

#### **6.2** ITEM EXPOSURE

Table 21 and Table 22 present the item exposure rates of the spring 2021 test administration for Rhode Island and Vermont, respectively. The exposure rates were relatively similar to the simulation results described in Section 5.2.2, Item Exposure, for the English test administrations. The item exposure rate for field-test items ranged from 10% to 13% for all three grades. For the Spanish tests in Rhode Island, more items had high exposure rates compared to the English tests because of a smaller item pool. Also, the operational exposure rates were slightly different from the simulation results in some cases because of small population sizes in all three grades. In spring 2021, less than 200 students took the Spanish test in each grade in Rhode Island. The exposure rates are 100% for the Spanish test in Vermont because only one student took the Spanish test in Grade 8.

Table 21. Item Exposure Rates by Grade: Percentage of Items by Exposure Rate, Across All Spring 2019 Test Administrations in Rhode Island

Grade	Total Items	[0,0]%	[0,1]%	[1,5]%	[5,20]%	[20,40]%	[40,60]%	[60,80]%	[80,100]%			
					English							
5	135	0	0	5.19	85.93	6.67	0.74	0.74	0.74			
8	152	0	0	1.97	91.45	5.26	1.32	0	0			
11	121	0	0	2.48	74.38	20.66	1.65	0	0.83			
	Spanish											
5	34	0	0	0	8.82	26.47	35.29	8.82	20.59			

Grade	Total Items	[0,0]%	[0,1]%	[1,5]%	[5,20]%	[20,40]%	[40,60]%	[60,80]%	[80,100]%
8	26	0	0	0	0	19.23	26.92	11.54	42.31
11	28	0	0	0	0	28.57	28.57	10.71	32.14

Table 22. Item Exposure Rates by Grade: Percentage of Items by Exposure Rate, Across All Spring 2019 Test Administrations in Vermont

Grade	Total Items	[0,0]%	[0,1]%	[1,5]%	[5,20]%	[20,40]%	[40,60]%	[60,80]%	[80,100]%
					English				
5	130	0	0	4.62	81.54	11.54	0.77	0.77	0.77
8	146	0	0	1.37	90.41	6.85	1.37	0	0
11	119	0	0	8.47	63.56	26.27	0	0	1.69
	Spanish								
5	-	-	-	-	-	-	-	-	-
8	18	0	0	0	0	0	0	0	100
11	-	-	-	-	-	-	-	-	-

#### 7. REFERENCES

Council of Chief State School Officers (2015). Science Assessment Item Collaborative Assessment Framework for the Next Generation Science Standards. Washington, DC: Author. Retrieved from <a href="https://ccsso.org/sites/default/files/2017-12/SAICAssessmentFramework">https://ccsso.org/sites/default/files/2017-12/SAICAssessmentFramework</a> FINAL.pdf.

National Research Council. (2012). A Framework for K–12 Science Education: Practices, Crosscutting concepts, and Core Ideas. Washington, DC: The National Academies Press.

# Appendix A Item Writer Training Materials

Exhibit A-1. LABS Guidelines



## AMERICAN INSTITUTES FOR RESEARCH®

#### **LABS** Guidelines

#### 1. STEREOTYPING

Testing materials should not present persons stereotyped according to the following characteristics:

- Age
- Disability
- Gender
- Race/Ethnicity
- Sexual orientation

#### 2. SENSITIVE OR CONTROVERSIAL SUBJECTS

Controversial or potentially distressing subjects should be avoided or treated sensitively. For example, a passage discussing the historical importance of a battle is acceptable whereas a graphic description of a battle would not be. Controversial subjects include:

- Death and Disease
- Gambling\*
- Politics (Current)
- Race relations
- Religion
- Sexuality
- Superstition
- War

<sup>\*</sup>References to gambling should be avoided in mathematics items related to probability.

#### 3. ADVICE

Testing materials should not advocate specific lifestyles or behaviors except in the most general or universally agreed-upon ways. For example, a recipe for a healthful fruit snack is acceptable but a passage recommending a specific diet is not. The following categories of advice should be avoided:

- Religion
- Sexual preference
- Exercise
- Diet

#### 4. Dangerous Activity

Tests should not contain content that portrays people engaged in or explains how to engage in dangerous activities. Examples of dangerous activities include:

- Deep-sea diving
- Stunts
- Parachuting
- Smoking
- Drinking

#### 5. POPULATION DIVERSITY AND ETHNOCENTRISM

Testing materials should:

- Reflect the diversity of the testing population
- Use stimulus materials (such as works of literature) produced by members of minority communities
- Use personal names from different ethnic origin communities
- Use pictures of people from different ethnic origin communities
- Avoid *ethnocentrism*, or the attitude that all people should share a particular group's language, beliefs, culture, or religion

#### 6. DIFFERENTIAL FAMILIARITY AND ELITISM

Specialized concepts and terminology extraneous to the core content of test questions should be avoided. This caveat applies to terminology from the fields of:

- Construction
- Finance
- Sports
- Law
- Machinery
- Military topics
- Politics
- Science
- Technology
- Agriculture

#### 7. LANGUAGE USE

Language should be as inclusive as possible.

- Avoid masculine-coded words like mankind, manmade, and the generic "he"
- Use equal pairs such as husband and wife rather than man and wife

#### 8. LANGUAGE ACCESSIBILITY

The grammar and vocabulary should be clear, concise, and appropriate for the intended grade level. The following should be avoided or used with care:

- Passive constructions
- Idioms
- Multiple subordinate clauses
- Pronouns with unclear antecedents
- Multiple-meaning words
- Non-standard grammar
- Dialect
- Jargon

#### 9. ILLUSTRATIONS AND GRAPHICS

Illustrations and graphics should embody all of the previously referenced LABS Guidelines.

Exhibit A-2. LABS Checklist



## AMERICAN INSTITUTES FOR RESEARCH®

### **LABS-Checklist**

#### **STEREOTYPING CONSIDERATIONS**

	Does the material negatively represent, or stereotype people based on gender or sexual preference?
	Does the material portray one or more people with disabilities in a negative or stereotypical manner?
	Does the material portray one or more religious groups as aggressive or violent?
	Does the material romanticize or demean people based on socioeconomic status?
	Does the material portray one or more ethnic groups or cultures participating in certain stereotypical activities or occupations?
	Does the material portray one or more age groups in a negative or stereotypical manner?
	Does the material require a student to take a position that challenges authority?
SENSIT	IVE/CONTROVERSIAL MATERIAL CONSIDERATIONS
	Does the material present war or violence in an overly graphic manner?
	Does the material present sensitive or highly controversial subjects, such as death, war, abortion, euthanasia, or natural disasters, except where they are needed to meet State Content Standards?
	Does the material require test takers to disclose values that they would rather hold
	confidential?
	Does the material present sexual innuendoes?
_	Does the material present sexual innuendoes?

#### **ADVICE CONSIDERATIONS**

Does the material contain advice pertaining to health and well-being about which there
is not a universal agreement?

#### **POPULATION DIVERSITY**

Is the material written by members of diverse groups?
Does the material reflect the experiences of diverse groups?
Does the material portray people in positive nontraditional roles?
Does test material represent the racial and ethnic composition of the testing population?
Does the material reflect ethnocentrism?
Does the material refer to population subgroups accurately?
Does test material reflect diversity through the use of names, cultural references, pictures, and roles?

#### **DIFFERENTIAL FAMILIARITY/ELITISM**

Does the material contain phrases, concepts, and beliefs that are irrelevant to testing domain and are likely to be more familiar to specific groups that others?
Does the material require knowledge of individuals, events, or groups that is not familiar to all groups of students?
Does the material suggest that affluence is related to merit or intelligence?
Does the material suggest that poverty is related to increased negative behaviors in society?
Does the material use language, content, or context that is offensive to people of a particular economic status?
Does success with the material assume that the test taker has experience with a certain type of family structure?
Does the material favor one socioeconomic group over another?
Does the material assume values not shared by all test takers?

#### LINGUISTIC FEATURES/LANGUAGE ACCESSIBILITY/GRAPHICS

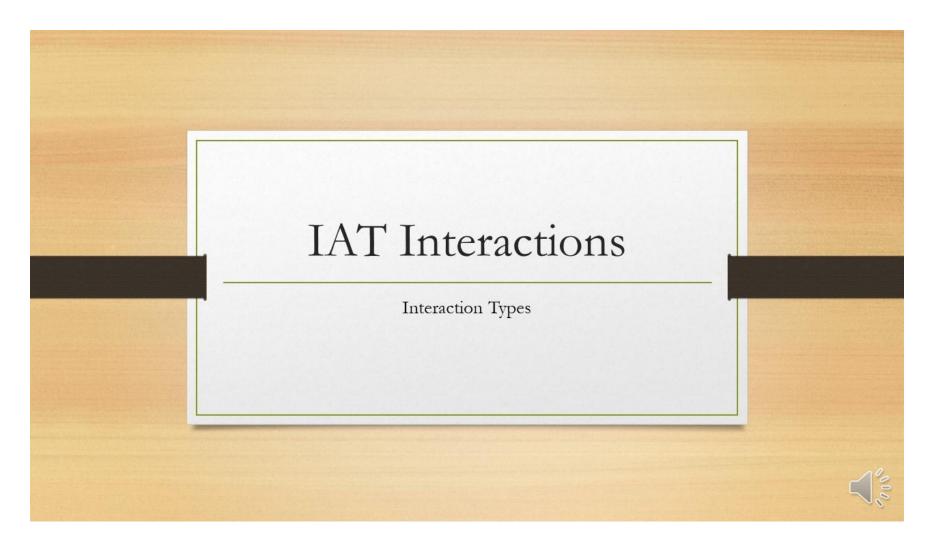
s grammar and vocabulary used in the items clear, concise, and appropriate for the
ntended grade level?

OTHER	QUESTIONS TO CONSIDER
	Does the material favor one age group over others except in a context where experience or maturation is relevant?
	Does the material use language, content, or context that is not accessible to one or more of the age groups tested?
	Does the material contain language or content that contradicts values held by a certain culture?
	Does the material favor one racial or ethnic group over others?
	Does the material degrade people based on physical appearance or any physical, cognitive, or emotional challenge?
	Does the material focus only on a person's disability rather than portraying the whole person?
	Does the material favor one religion and/or demean others?

□ Do the illustrations and graphics embody all of the previously referenced LABS

Guidelines?



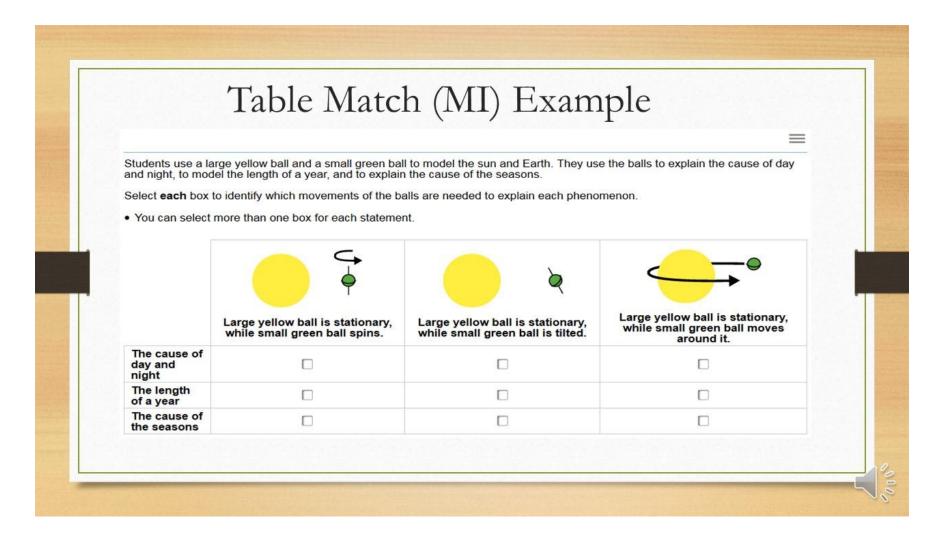


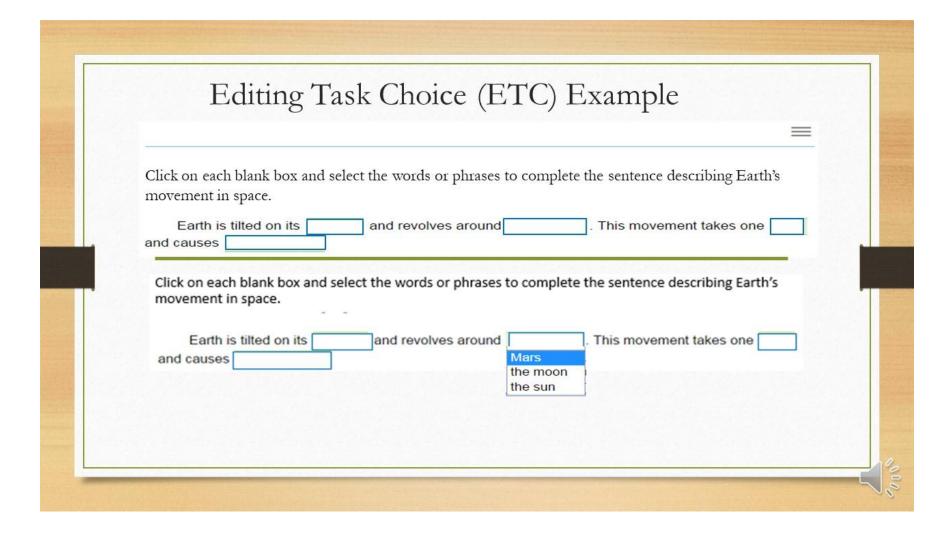
# Selected Response Interactions

- Selected Response interactions provide response options and the student selects the response(s). SR interaction types include:
  - Multiple Choice (MC)
  - Multi-Select (MS)
  - Table Match (MI)
  - Editing Task Choice (ETC)
  - Hot Text (HT)

These interactions are more accessible to all students!

# Multiple Select Example The hawksbill sea turtle builds nests on Hawaiian beaches. Female turtles lay their eggs in the nests. About two months later, the baby turtles hatch and crawl across the beaches to the ocean. Over the years, scientists have noticed a drop in the number of baby turtles making it to the ocean. Select the three observations that could explain the drop in the turtle population. Adult turtles get caught in nets. Baby turtles crawl quickly from the nests. Food left on the beach attracts predators of the turtles. The turtles mistake bright lights for the moon. Turtles eat plastic floating in the ocean.





# Hot Text (HT draggable) Example

A list of natural events is shown.

Click and drag the natural events to classify each natural event as either a fast or slow process that could shape and reshape Earth's surface.

#### **Fast and Slow Processes**

Slow Process
֡

- 1. A glacier melts, depositing sediment.
- 2. A mountain side collapses, causing a landslide.
- 3. A tsunami pushes sediment inland.
- An earthquake causes a crack along a road.
- 5. Waves carve an arch in a sea cliff.
- 6. Wind weathers a rock.

# Hot Text (HT selectable) Example

A list of natural events that could shape and reshape Earth's surface is shown.

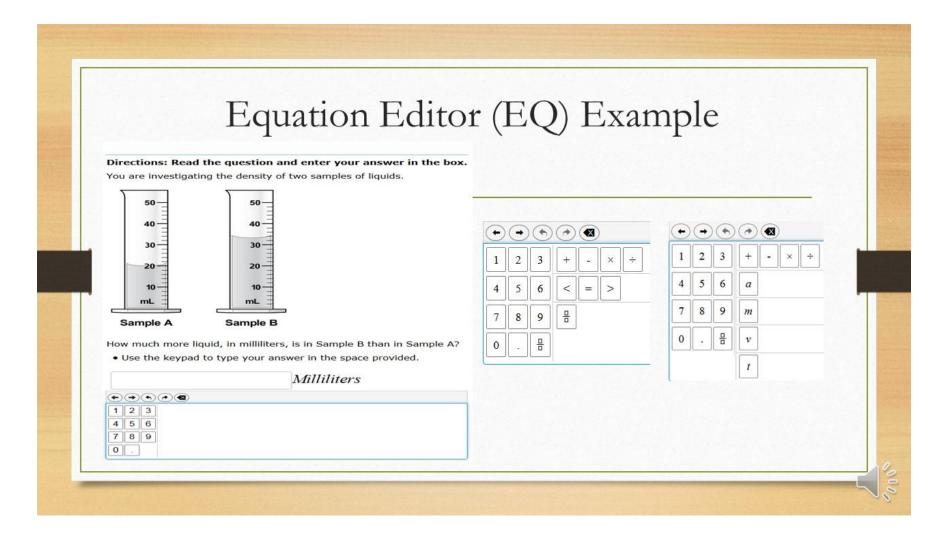
Click on each process below that happens slowly.

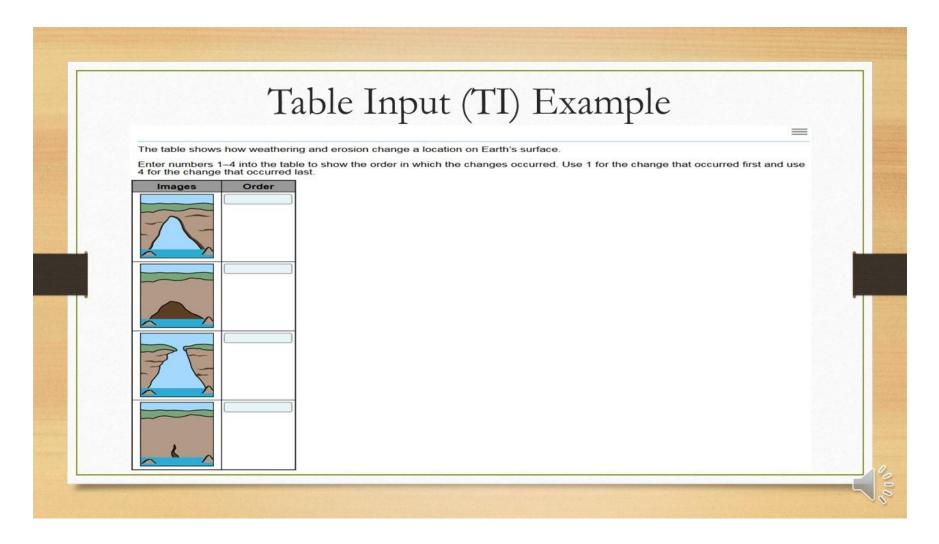
- · A glacier melts, depositing sediment.
- · A mountain side collapses, causing a landslide.
- A tsunami pushes sediment inland.
- · An earthquake causes a crack along a road.
- · Waves carve an arch in a sea cliff.
- · Wind weathers a rock.

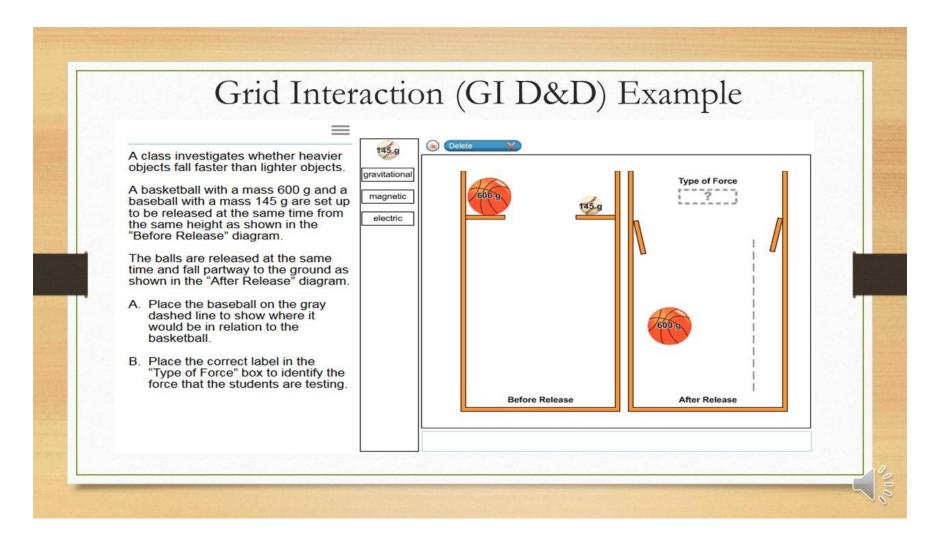
# Machine Scored Constructed Response Interactions

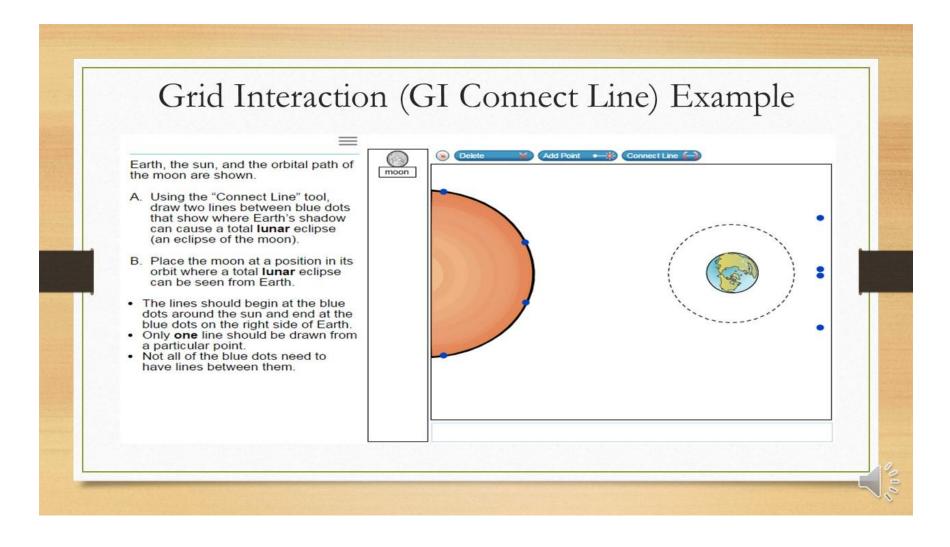
- Machine Scored Constructed Response interactions require scoring logic or a machine rubric within the interaction. MSCR interaction types include:
  - Equation Editor (EQ)
  - Table Interaction (TI)
  - Grid Interaction (GI)
  - Simulation (Sim)
  - Natural Language (NL)
  - Editing Task (ET)
  - Word Builder (WB)

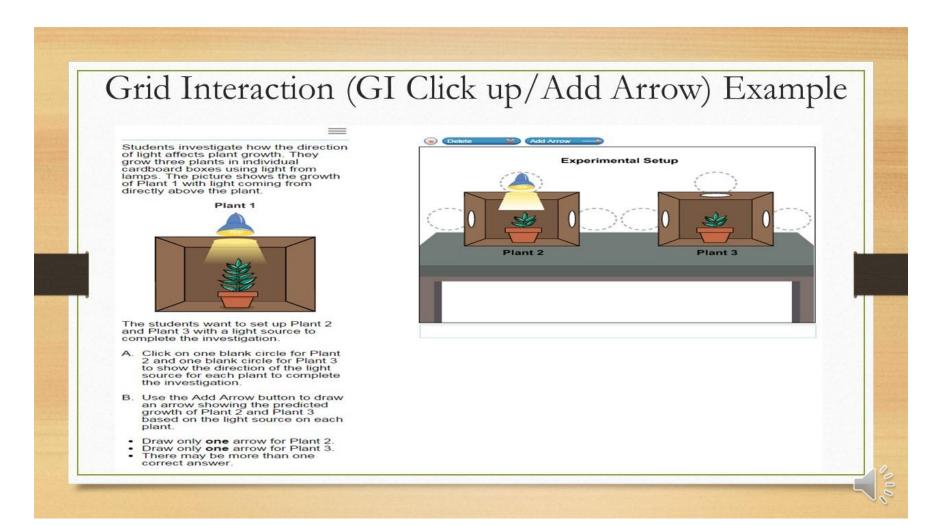
These interactions are less accessible to all students!

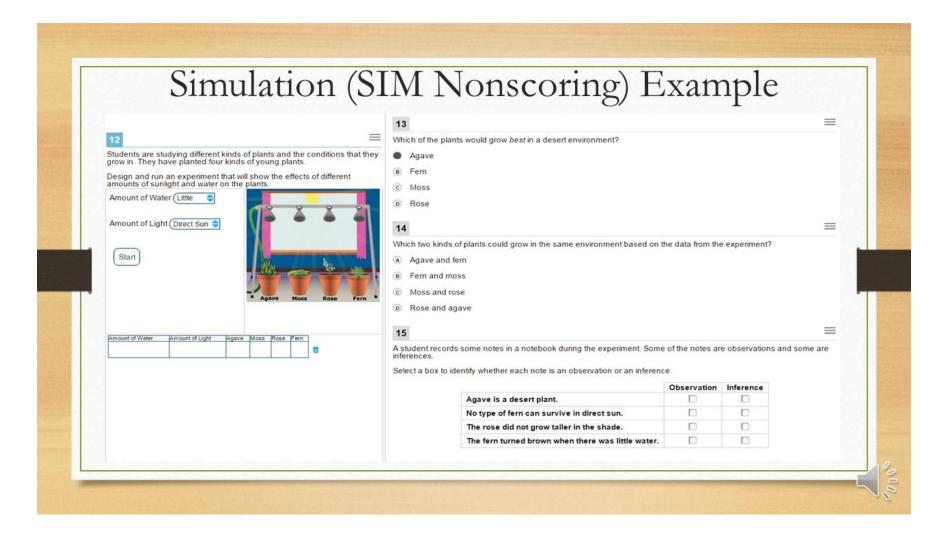


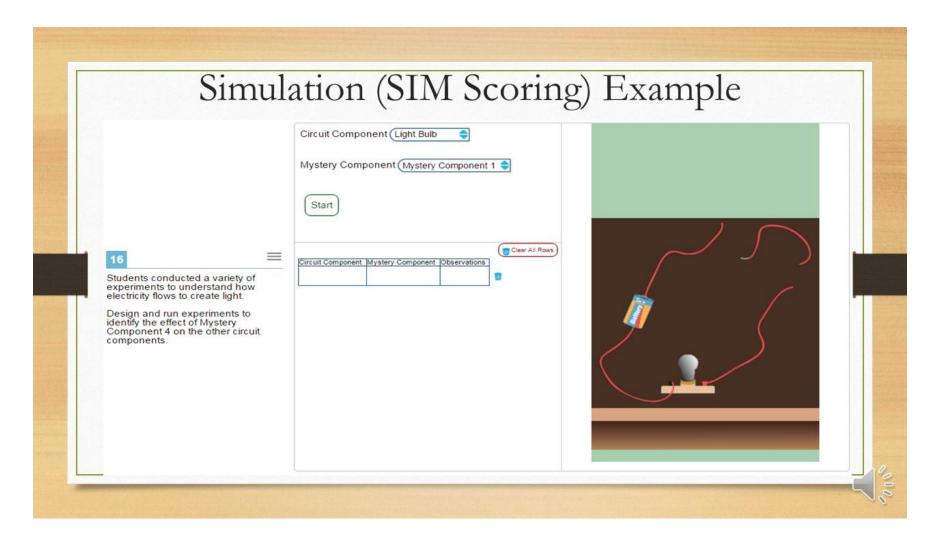


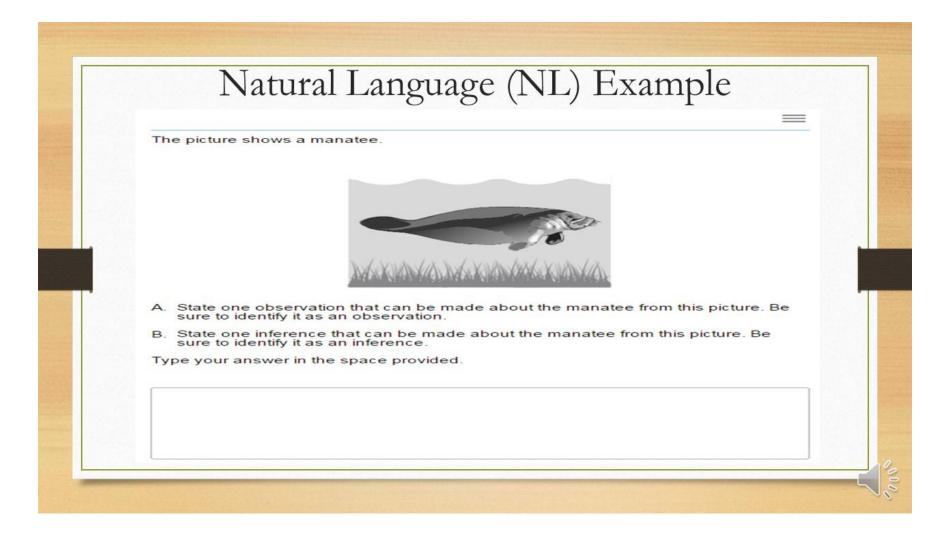












## Selected Response (SR) Interactions

Selected Response interactions provide response options and the student selects the response(s).

SR Interaction Type	Task Demands that can be Assessed				
Multiple Choice (MC)	Identify, Choose, Select, Label				
Multi Select (MS)	Identify, Choose, Select, Label				
Table Match (MI)	Classify, Categorize, Organize, Rank, Sort, Sequence				
Editing Task Choice (ETC)	Classify, Categorize, Organize, Sort, Sequence, Compare, Label, Construct an explanation/argument, Describe, Summarize, Complete				
Hot Text Selectable (HT)	Highlight, Identify, Select, Choose				

## Machine Scored Constructed Response (MSCR) Interactions

Machine Scored Constructed Response interactions require scoring logic or a machine rubric within the interaction. MSCR interaction types include:

Machine Scored Constructed Response Interaction Type	Task Demands that can be Assessed		
Equation Editor (EQ)	Calculate, Mathematically describe/represent/model, Identify		
Table Input (TI)	Calculate, Sequence, Identify, Organize, Chart		
Grid Interaction (GI)	Graph, Model, Represent, Show, Create		
Simulation Interaction (Sim)	Investigate, Experiment, Observe, Gather/collect data, Model		
Natural Language (NL)	Describe, Compare, Summarize, Explain		
Editing Task (ET)	Correct		
Word Builder (WB)	Identify		

Appendix B

Item Review Checklist

## **Item Review Checklist**

# Tier 1 – Sufficiency/Appropriateness of the Phenomenon to Assess the Performance Expectation

The elements	in	this	tier	are	critical
THE CICITOTICS	111	UIII		ui c	CITCICAL

Is the phenomenon based on a specific real-world scenario and focused enough to get the student to investigate what the Performance Expectation (PE) intends for them to investigate (i.e., the students' application of the Practice in the context of the Disciplinary Core Idea [DCI] and Crosscutting Concepts [CCC] as intended by the PE is sufficient to make sense of the phenomena)?
Is there an appropriate science-related activity that is puzzling and/or intriguing for students to engage in? Is the scenario focused on real-world observations that students can connect with or have direct experience with?
Is the context and complexity of the phenomenon grade-appropriate?
Cluster Task Statement: Does the "call to action" reflect the end goal of the interactions to be answered? Does the statement make sense? Is this an engaging and reasonable outcome to work towards?
Is the phenomenon presented in way(s) that all students can access and comprehend it based on information provided (including text, graphics, data, images, animations, etc.)? Is the phenomenon free of cultural bias, insensitivity or depreciation of unsafe situations?

#### Tier 2 – Review of Specific Elements by Component

#### **Stimulus**

Reading Load/Readability/Style

Is the reading load appropriate for the grade (i.e., the amount of text minimized to reduce cognitive load)?
Is the language and vocabulary appropriate for the grade?
Non-specific vocabulary should be one grade level lower than the tested grade.
Science vocabulary should be part of the "Science Vocabulary Students Are Expected to Know" in the item specifications.
Is all of the information in the stimulus necessary for the student to complete the item interactions?
Is language consistent throughout the cluster (i.e., does not switch between steam and vapor)?

[		Is everything in the active voice (i.e., avoids unnecessary and unclear passive construction)?
Measu	ren	nent/Units
[		Are the data in SI units? Check style guide for exceptions.
[		Are units of measurement introduced or defined before they are used in graphs/tables?
[		Are the dependent/independent variables on the correct axes or in the correct columns?
[		Are the graphs/tables/pictures free of extraneous information and appropriate for the grade level?
[		Is there information included in graphs/pictures/tables that is not necessary and can be removed?
]		Do the graphs/tables/pictures depend on color? Is there another way to represent the difference in the data other than by color (e.g., using patterns)?
Data S	oui	rce and Scientific Reference
[		Is content both accurate and appropriate in its context?
[		Are the data sources appropriate for the subject/grade and taken from reliable academic sources?
[		Does the item use the most up-to-date explanation?
Forma	ttir	g
[		Is everything presented within the browser dimensions (1024x768) without horizontal scrolling?
[		Are the tables/graphs/etc. laid out in a way that is easy to read?
[		Are details and text in animations easy to see? Are labels in diagrams easy to read?
[		Is the average file size appropriate for test delivery (approximately $100KB$ , $250KB$ maximum)?
<u>Item</u>		
Interac	ctio	n and Alignment to Specifications
[		Does the item make sense if you are responding to the interactions as if you are the student in the intended grade-level?
[		Does the interaction require the student to demonstrate the science practice and/or content that the PE is assessing them on?
[		Are the interactions grade level/developmentally appropriate and do they follow a logical progression? Do the interactions use appropriate scaffolding to guide students in making sense of the phenomena?
[		Do the interactions align with the task demands?

Do the interactions avoid redundancy? Do the student interactions follow a coherent progression?
Do the student interactions follow a coherent progression? Does the order of the interactions allow students to make sense of the phenomenon or problem?
Is the item stem worded in a way that makes the intent of the interaction clear to the student?
Is it clear to the student what they will be scored on in the interaction?
Is the language (e.g., words, phrases) consistent throughout the stimulus and items?

#### Grade Appropriate

- Is the content within the item accurate and grade appropriate?
- Are the correct units used? Are the units grade appropriate? Where necessary, are the abbreviations of the units introduced?
- Is the number of item parts/scoring assertions appropriate for the grade level?
- Is the mathematics level appropriate for the grade being tested?

#### Formatting

- Is everything presented within the browser frame without horizontal scrolling?
- Are the tables/graphs/etc. easy to read? Are the images created in an appropriate color palette per the Style Guide?
- Are details and text in animations easy to see?

#### Tier 3 – Review of the Scoring and Assertion(s)

#### Scoring Accuracy

Do the interactions/task provide clear guidance on how student responses will be scored/interpreted?
Are scores assigned appropriately as correct or incorrect?
Are the dependencies logical?
Are any of the scoring assertions exclusive (i.e., the student can get only one assertion correct and not another at any given time)?
Is the correct answer clear and distinct from the distractors?
Does the scoring result in an appropriate distribution of points?

#### **Scoring Assertions**

Is the appropriate wording used for each scoring assertion (e.g., <feature of="" response=""> providing some evidence of <what about="" infer="" student="" the="" to="" want="" we="">)?</what></feature>
Does the inference follow from the data?
Are the assertions specific to the individual interactions (i.e., does not just repeat the PE)?
Are the scoring assertions in the same order as the interactions?
Does the wording of the scoring assertion make it very clear which interaction and action it refers to?

#### Strategies for Editing Text to Produce Plain Language

- Reduce excessive length
- Use common words
- Avoid ambiguous words
- Limit irregularly spelled words
- Avoid inconsistent naming and graphic conventions
- Avoid multiple terms for the same concept
- Limit the use of embedded clauses and phrases
- Avoid the passive voice

# Appendix C Content Advisory Committee Participant Details

## **Content Advisory Committee Participant Details**

Table C-1. Content Advisory Committee Participants, Science

State	Date	Location	Grade Band	Number of Teachers in Each Group	Teacher Demographic Summary	Number of Science Items Reviewed	Number of Science Items Approved by Teacher Committees
ICCR	March 2018	Virtual	Elementary School		Gender: Male 27%, Female 73% Ethnicity: Not collected		ı
			Middle School	26ª	26ª	State: Connecticut 46%, Hawaii 8%, Maryland 4%, Oregon 12%, West Virginia 27%, Utah 4% Teaching Experience: General Education 31%, General Ed and Other 12%, Science Curriculum Specialist 15%, Science	152
			High School		Department Head 8%, STEM Consultant 8%, No response 27%		
	February 2017	Cromwell, Connecticut	Elementary School	11	Gender: Male 22%, Female 78% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	45	31
			Middle School	14			
			High School	16			
		New Britain, Connecticut	Elementary School	12	Gender: Male 26%, Female 74% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	40	N/A <sup>b</sup>
	May 2017		Middle School	15			
			High School	15			
Connecticut	October 2017	New Britain.	Elementary School	11	11 Gender: Male 20%, Female 80% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	75	64
		Connecticut	Middle School	12			
			High School	18			
	November	er New Britain, Connecticut	Elementary School	7	Gender: Male 17%, Female 83%  Ethnicity: Not collected  Region: Not collected  Teaching Experience: Not collected	41	32
	2017		Middle School	11			
			High School	17			
	January 2018	New Britain, Connecticut	Elementary School	11	Gender: Male 18%, Female 82% Ethnicity: Not collected Region: Not collected	42	25
			Middle School	14			

State	Date	Location	Grade Band	Number of Teachers in Each Group	Teacher Demographic Summary	Number of Science Items Reviewed	Number of Science Items Approved by Teacher Committees
			High School	8	Teaching Experience: Not collected		
	October	New Britain,	Elementary School	13	Gender: Male 16%, Female 84% Ethnicity: Not collected		54
	2018	Connecticut	Middle School	16	Region: Not collected	84	
			High School	16	Teaching Experience: Not collected		
	November	New Britain.	Elementary School	10	Gender: Male 14%, Female 86% Ethnicity: Not collected		
	2018	Connecticut	Middle School	18	Region: Not collected	235	200
			High School	21	Teaching Experience: Not collected		
	December	New Britain.	Elementary School	10	Gender: Male 19%, Female 81% Ethnicity: Not collected		
	2018	Connecticut	Middle School	7	Region: Not collected	56	55
			High School	15	Teaching Experience: Not collected		
	January	New Britain.	Elementary School	13	Gender: Male 18%, Female 82% Ethnicity: Not collected	65	59
	2019	Connecticut	Middle School	13	Region: Not collected		
			High School	18	Teaching Experience: Not collected		
	September	Rocky Hill,	Elementary School	14	Gender: Male 18%, Female 82% Ethnicity: Not collected	60	57
	2019	Connecticut	Middle School	16	Region: Not collected		
			High School	20	Teaching Experience: Not collected		
		July 2017 Honolulu, Hawaii	Elementary School	7	Gender: Male 36%, Female 64% Ethnicity: Black 5%, Chinese and White 5%, Filipino 9%, Hawaiian 14%, Hispanic 9%, Japanese 14%, White 41%, No response 5% Region: Not collected	25	N/A <sup>b</sup>
	July 2017		Middle School	8			
Hawaii			High School	7	Teaching Experience: General Education 64%, General Education w/SPED Certification 5%, SPED Teacher 5%, Other 23%, No response 5%		
	September 2017	Honolulu, Hawaii	Elementary School	6	Gender: Male 25%, Female 75%	65	N/A <sup>b</sup>

State	Date	Location	Grade Band	Number of Teachers in Each Group	Teacher Demographic Summary	Number of Science Items Reviewed	Number of Science Items Approved by Teacher Committees
			Middle School	8	Ethnicity: Black 5%, Filipino 10%, Hispanic 10%, Japanese 15%, White 50%, No response 10%		
			High School	6	Region: Not collected Teaching Experience: General Education 65%, General Education w/SPED Certification 15%, Other 20%		
			Elementary School	10	Gender: Male 17.24%, Female 82.76% Ethnicity: White 27.59%, N/A 10.34%, Hispanic 10.34%, Asian 31.03%, Hawaiian		
	October 2018	Hawaii	Middle School	6	3.45%, Asian Pacific Islander 6.9%, Two or More: 10.34%  Region: Not collected  Teaching Experience: General Education	85	79
			High School	12	82.76%, SPED Teacher 0%, ELL Teacher 0%, General Education w/ SPED Certification 0%, Other 24.14%		
		February Honolulu, 2019 Hawaii	Elementary School	8	Gender: Male 20%, Female 80% Ethnicity: White 35%, Asian 50%, Two or More: 15%		
				Middle School	6	Region: Not collected Teaching Experience: General Education	44
			High School	7	65%, SPED Teacher 5%, General Education w/ SPED Certification 5%, Other 25%		
Idaho	December 2018	KI/Au	Elementary School	21ª	Gender: Not collected Ethnicity: Not collected Region: Not collected	111	N/A <sup>b</sup>
	2010		Middle School	Teaching Experience: Not collected			
MSSA <sup>c</sup>			Elementary School	15	Gender: Not collected Ethnicity: Not collected		
	January 2018	N/A <sup>d</sup>	Middle School	14	State: 90% Rhode Island, 10% Vermont Teaching Experience: General Education	73	N/A <sup>b</sup>
			High School	13	69%, Bilingual Education 2%, Science Coordinator 14%, Other 14%		
	March 2018	N/A <sup>d</sup>	Elementary School	12	Gender: Not collected Ethnicity: Not collected	100	N/A <sup>b</sup>
			Middle School	13	State: Rhode Island 25%, Vermont 75%	73	

State	Date	Location	Grade Band	Number of Teachers in Each Group	Teacher Demographic Summary	Number of Science Items Reviewed	Number of Science Items Approved by Teacher Committees
			High School	9	Teaching Experience: Not collected		
	January		Elementary School Middle School	-	Gender: Male 25.71%, Female 74.29% Ethnicity: Not collected Region: Not collected		
	2019	N/A <sup>d</sup>	High School  21a  Teaching Experions 68.57%, Special Education 0%,	Teaching Experience: General Education 68.57%, Special Education 2.86%, Bilingual Education 0%, Administration 0%, Other 28.57%, N/A 5.71%	116	N/A <sup>b</sup>	
			Elementary School	4	Gender: Male 10%, Female 90% Ethnicity: Not collected		
	August 2017	Salem, Oregon	Middle School High School	3	Region: Urban 50%, Suburban 0%, Rural 50% Teaching Experience: Regular Education 100%, Bilingual Education 10%, Special Education 10%, Administration 20%, Other 0%	235	142
	August	Salem,	Elementary School	4	Gender: Male 20%, Female 80% Ethnicity: Other 5%, White 95%	257	200
Oregon	2018	Oregon	Middle School High School	8	Region: Urban 56%, Suburban 0%, Rural 44% Teaching Experience: Bilingual Education		
			Elementary School	6	65%, Special Education 65%, Other 55%  Gender: Male 38%, Female 63%  Ethnicity: Asian 6%, White 94%		48
	December	Virtual	Middle School	5	Region: Urban 50%, Suburban 50%, Rural 0%	62	
	2018	2018	High School	5	Teaching Experience: General Education 38%, Bilingual Education 63%, Special Education 25%		
			Grade 6	6	Gender: Male 26.09%, Female 73.91% Ethnicity: White 91.3%, Native American 4.35%, Other 4.35%		
Utah	July 2017	Park City, Utah	Grade 7	6	Region: Not collected Teaching Experience: General Education	55	51
			Grade 8	6	100%, Special Education 4.35%, Bilingual Education 0%, Administration 0%, Other 4.35%		
	December	Salt Lake	Grade 6	12	Gender: Male 16%, Female 83.87% Ethnicity: American Indian or Alaska Native	64	60
	2017	City, Utah	Grade 7	12	and White 3.23%, Other 3.23%, White 93.55% Region: Not collected	64	62

State	Date	Location	Grade Band	Number of Teachers in Each Group	Teacher Demographic Summary	Number of Science Items Reviewed	Number of Science Items Approved by Teacher Committees
			Grade 8	12	Teaching Experience: General Education 87.09%, General Education and Other 9.68%, General Education and ESOL 3.23%		
	January 2017	N/A <sup>d</sup>	Elementary School	- 28 <sup>a, e</sup>	Gender: Not collected Ethnicity: Not collected Region: Not collected	39	N/A <sup>b</sup>
			Middle School  Elementary School		Teaching Experience: Not collected  Gender: Male 11.11%, Female 88.89% Ethnicity: White 88.89%, Black 11.11%		
West Virginia	October 2018	N/A <sup>d</sup>	Middle School	10ª	Region: Rural 100%, Urban 0%, Suburban 0% Teaching Experience: General Education 100%, Special Education 0%, Bilingual Education 0%, Administration 0%, Other 0%	191	N/A <sup>b</sup>
			Elementary School	6	Gender: Male 13.04%, Female 86.96% Ethnicity: White 86.96%, Asian 4.35%, Black 4.35%, N/A 4.35%  Region: Rural 69.57%, Urban 30.43%, Suburban 0%, N/A 4.35%  Teaching Experience: General Education 71.74%, Special Education 4.35%, Bilingual Education 0%, Administration 0%, Other 13.04%, N/A 13.04%	50	N/A <sup>b</sup>
	July 2019	N/A <sup>d</sup>	Middle School	6		50	
	<b>December</b> Ch	Cheyenne,	Elementary School	6	Gender: Not collected Ethnicity: Not collected		
Wyoming	2017	Wyoming	Middle School	8	Region: Not collected  Teaching Experience: Not collected	51	N/A <sup>b</sup>
	October		High School	4	reaching Experience. Not collected		N/A <sup>b</sup>
		Cheyenne,	Elementary School		Gender: Not collected Ethnicity: Not collected		
	2018	Wyoming	Middle School	14 <sup>a</sup>	Region: Not collected	37	
			High School		Teaching Experience: Not collected		

*Note.* <sup>a</sup>Number of Committee Members by grade band is not available.

<sup>&</sup>lt;sup>b</sup>Number of science items approved by teacher committees is unavailable at the time of writing this report.

<sup>°</sup>MSSA = Rhode Island and Vermont's Multi-State Science Assessment.

<sup>&</sup>lt;sup>d</sup>Location of Content Advisory Committee Meeting is unavailable at the time of writing this report.

eNumber of Committee Members includes total committee members for ELA, math, and science. The number for science only committee members is not available.

# Appendix D Fairness Committee Participant Details

## **Fairness Committee Participant Details**

Table D-1. Fairness Committee Participants, Science

State	Date	Location	Total Number of Committee Members	Teacher Demographic Summary	Number of Items Reviewed
ICCR	March 2018	Virtual	13	Gender: Male 15%, Female 85% Ethnicity: Not collected State: Connecticut 46%, Indiana 8%, Utah 15%, West Virginia 23%, Wyoming 8% Teaching Experience: General Education 8%, General Education and Other 15%, EL Instructional Coach 8%, No response 69%	152
	February 2017	Cromwell, Connecticut	6	Gender: Male 17%, Female 83% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	45
	December 2017	New Britain, Connecticut	9	Gender: Male 22%, Female 78% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	75
	December 2017	Cromwell, Connecticut	10	Gender: Male 30%, Female 70% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	41
Connecticut	February 2018	New Britain, Connecticut	3	Gender: Male 33%, Female 67% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	42
	November 2018	New Britain, Connecticut	11	Gender: Male 9%, Female 91% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	319
	December 2018	New Britain, Connecticut	10	Gender: Male 20%, Female 80% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	56
	January 2019	New Britain, Connecticut	9	Gender: Male 22%, Female 78% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	65

State	Date	Location	Total Number of Committee Members	Teacher Demographic Summary	Number of Items Reviewed
	September 2019	Cromwell, Connecticut	9	Gender: Male 11%, Female 89% Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	48
Hawaii	July 2017	Honolulu, Hawaii	22	Gender: Male 36%, Female 64% Ethnicity: Black 5%, Chinese and White 5%, Filipino 9%, Hawaiian 14%, Hispanic 9%, Japanese 14%, White 41%, No response 5% Region: Not collected Teaching Experience: General Education 64%, General Education w/SPED Certification 5%, SPED Teacher 5%, Other 23%, No response 5%	25
	September 2017	Honolulu, Hawaii	20	Gender: Male 25%, Female 75%  Ethnicity: Black 5%, Filipino 10%, Hispanic 10%, Japanese 15%, White 50%, No response 10%  Region: Not collected  Teaching Experience: General Education 65%, General Education w/SPED Certification 15%, Other 20%	65
	October 2018	Honolulu, Hawaii	29	Gender: Male 20.69%, Female 79.31% Ethnicity: White 27.59%, Japanese 10.34%, N/A 10.34%, Hispanic 10.34%, Chinese 6.9%, Asian 6.9%, Hawaiian 3.45%, Asian Pacific Islander 6.9%, Filipino 3.45%, Multi-Racial/Ethnic 13.8% Region: Not collected	85
	February 2019	Honolulu, Hawaii	21	Gender: Male 20%, Female 80% Ethnicity: White 35%, Asian 50%, Two or More: 15% Region: Not collected Teaching Experience: General Education 65%, SPED Teacher 5%, General Education w/ SPED Certification 5%, Other 25%	44
ldaho	December 2018	N/Aª	15	Gender: Not collected Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	111
MSSA♭	January 2018	N/Aª	21	Gender: Not collected Ethnicity: Not collected State: Rhode Island 100%, Vermont 0% Teaching Experience: General Education 67%, Bilingual Education 14%, Special Equation 5%, Science Coordinator 5%, Other 10%	73
	March 2018	N/Aª	11	Gender: Not collected Ethnicity: Not collected State: Rhode Island 55%, Vermont 45% Teaching Experience: Not collected	100

State	Date	Location	Total Number of Committee Members	Teacher Demographic Summary	Number of Items Reviewed
	January 2019	N/Aª	14	Gender: Male 22.86%, Female 62.86% Ethnicity: Not collected Region: Not collected Teaching Experience: General Education 68.57%, Special Education 2.86%, Bilingual Education 0%, Administration 0%, Other 17.14%, Coach 11.43%	116
Oregon	August 2017	Salem, Oregon	5	Gender: Male 0%, Female 100% Ethnicity: Not collected Region: Urban 80%, Suburban 20%, Rural 0% Teaching Experience: Regular Education 40%, Bilingual Education 20%, Special Education 20%, Administration 60%, Other 20%	110
	August 2018	Salem, Oregon	39	Gender: Male 26%, Female 74% Ethnicity: Asian 3%, Hispanic 8%, Native American 3%, White 82%, Other 10% Region: Urban 56%, Suburban 0%, Rural 44% Teaching Experience: General Education 15%, Bilingual Education 72%, Special Education 33%, Other 33%	257
	December 2018	Virtual	11	Gender: Male 9%, Female 91% Ethnicity: Hispanic 9%, White 91% Region: Urban 55%, Suburban 0%, Rural 45% Teaching Experience: General Education 27%, Bilingual Education 64%, Special Education 18%, Administration 9%, Other 64%	62
Utah	August 2017	Park City, Utah	6	Gender: Male 0%, Female 100%  Ethnicity: American Indian or Alaska Native 33%, Hispanic 33%, White 33%  Region: Urban 0%, Suburban 0%, Rural 17%, Unknown/No response/Not applicable 83%  Teaching Experience: General Education 17%, Special Education 17%, Administrator 33%, Other 33%	44
	December 2017	Salt Lake City, Utah	6	Gender: Male 16.67%, Female 83.33% Ethnicity: Black 33.33%, Native American 33.33%, Hispanic 16.67, White 0%, N/A 16.67% Region: Not collected Teaching Experience: General Education 0%, Special Education 0%, Bilingual Education 0%, Administration 33.33%, Other 83.33%	48
West Virginia	January 2017	N/Aª	28°	Gender: Not collected Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	34

State	Date	Location	Total Number of Committee Members	Teacher Demographic Summary	Number of Items Reviewed
	January 2019	N/Aª	10	Gender: Male 11.11%, Female 88.89% Ethnicity: Black 11.11%, White 88.89% Region: Rural 100%, Urban 0%, Suburban 0% Teaching Experience: General Education 100%, Special Education 0%, Bilingual Education 0%, Administration 0%, Other 0%	191
W	December 2017	Cheyenne, Wyoming	5	Gender: Not collected Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	51
Wyoming	October 2018	Cheyenne, Wyoming	5	Gender: Not collected Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	37

*Note.* <sup>a</sup>Location of Fairness Committee Meeting is unavailable at the time of writing this report. <sup>b</sup>MSSA = Rhode Island and Vermont's Multi-State Science Assessment.

<sup>&</sup>lt;sup>c</sup>Number of Committee Members includes total committee members for ELA, math, and science. The number for science only committee members is not available.

# Appendix E Sample Data Review Training Materials

## **Sample Data Review Training Materials**

# Data Review for NGSS, 2019

AMERICAN INSTITUTES FOR RESEARCH

# Read and Sign Non-Disclosure

- Read and Sign Non-Disclosure
- Turn in to AIR Facilitator

American Institutes for Research

# Overview of Training

- Steps in the Development Process
- Describe the structure of Three-Dimensional clusters
- Describe scoring assertions
- Role of the Data Review Committee
- Data Review Process
- Participant Guidelines

AMERICAN INSTITUTES FOR RESEARCH

3

# Steps in the Development Process

- AIR Writes Clusters & Standalones
- AIR Internal Review (Content & Fairness)
- Client & Educator Review (Content & Fairness)
- Field Test with Students
- Rubric Validation Process
- Update Scores & Generate Field Test Data
- Review of Field Test Data

AMERICAN INSTITUTES FOR RESEARCH

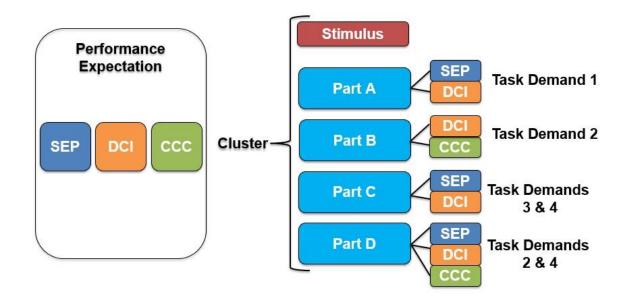
4

## NGSS in Hawaii

- A new Science Assessment has been developed to assess how well students master the NGSS
- The items of the new assessment look very different
  - Focus on item clusters
    - » Aligned to a single performance expectation
    - » Consisting of multiple interactions

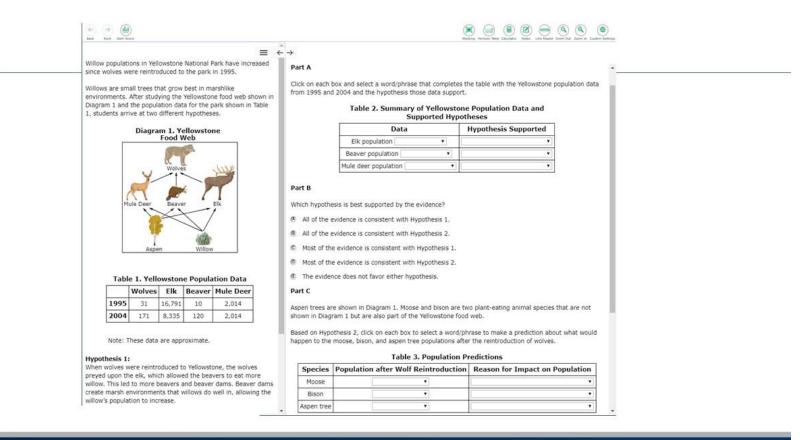
Astronous Incriminate and Resease.

## Structure of AIR Clusters



American Institutes for Research

Ь



American Institutes for Research

## Scoring Assertions

#### Scoring

- Within each item cluster, a series of explicit assertions can be made about the knowledge and skills that a student has demonstrated based on specific features of the student's responses
- Scoring assertions can be supported based on students' responses in one or more interactions within an item cluster.
- For example:
  - » A student correctly graphs data points indicating that (s)he can construct a graph showing the relationship between two variables,
  - » Makes an incorrect inference about the relationship between the two variables, thereby not supporting the assertion that the student can interpret relationships expressed graphically

American Institutes for Research

### 2019 Test Administrations and Rubric Validation

- Items were embedded as field test items in Spring 2019
- This past June, Items went through rubric validation
  - To check whether assertions were scored correctly
    - » Looking at actual student responses
  - It was determined for two items that student facing changes were necessary, so they will be updated and re-field tested this next year
  - Some assertions were modified (deleted/added)

AMERICAN INSTITUTES FOR RESEARCH

9

#### Data Review

- After rubric validation, statistics were computed at the assertionlevel
  - Assertions can only be evaluated in the context of the entire item
  - Inclusion in data review will be decided at the item level, not at the assertion level
  - Inclusion is based on statistical flags that rely on assertion level statistics but are evaluated for the entire item

#### Data Review

- Flagging is based on business rules related to:
  - Difficulty of the cluster
  - Relation between the score on cluster and the overall student's score
  - Response time of the cluster
  - Statistical flags for differential item functioning
- These items may be perfectly fine, but we want your input
  - Is this a good item and set of assertions?
  - Do you see any reason for why the item was flagged from a content perspective?

## p-value

- The p-value is the proportion of students for which the assertion is TRUE
- -Corresponds to the difficulty of an item in a traditional assessment
- Across an item bank, we want to see assertions with p-values across the full range to be able to precisely measure proficiency across all proficiency levels
  - » A low p-value is not bad per se
- However, we want to make sure the low p-value is not a result of an item being misleading

### p-value

- Criteria for clusters:
  - » average *p*-value < .30 (across the assertions within a cluster)
  - » average *p*-value > .85 (across the assertions within a cluster)
- Criteria for stand-alone items (typically has 1-3 assertions):
  - » average *p*-value < . 15 (across the assertions within a stand-alone)
  - » average *p*-value > .95 (across the assertions within a stand-alone)

- Item-total correlation
  - We expect students who do well on the test overall to have a higher probability of doing well on individual assertions
  - —The item-total correlation describes that relation
  - -Criterion
    - » Average item-total (biserial) correlation < .25
    - » One or more assertions with an item-total correlation < 0

## Differential item functioning

- Fair items behave similar across groups
- Probability of answering correctly is the same for all students of similar ability regardless of group membership

## Groups are defined by

- Gender
- Ethnicity
- Economically disadvantaged vs. not
- LEP vs. not
- Special Education vs. not

- Severity of possible bias
  - "A" No statistical evidence of DIF
  - "B" Evidence for potential mild DIF
  - "C" Evidence for potential severe DIF
- Direction of possible bias
  - "-" assertion favors reference groups (whites/females/non LEPs)
  - "+" assertion favors focal group

- Criteria
  - For clusters: 2 or more assertions show 'C' DIF in the same direction
  - For stand-alone items: 1 or more assertions show 'C' DIF in the same direction

## Timing

 We want a good balance between the amount of information an item provides, and the time students spend on the item

#### Criteria

- For clusters: percentile 80 > 15 minutes
  - » A percentile 80 of x minutes: 80% of the students spent x minutes or less on the cluster
- For stand-alone items: percentile 80 > 3 minutes
- Assertions per minute < .5 for clusters and stand-alone items

## **Data Review Process**

American Institutes for Research

#### **Process**

- Item is presented with information on
  - » Grade
  - » Discipline
  - » Topic
  - » Performance Expectation
- Facilitator will present the cluster or stand-alone item
- Statistics on the assertions of the item are presented
  - Including the reason for flagging
- Evaluation of item (Stimulus, interactions, assertions)
- For every item, one of the following decisions is made
  - Reject
  - Accept as is

## Participant Guidelines

- Keep phones turned off & stowed while in the meeting room.
  - If needed, please take the call outside of meeting room
- Keep your name tent visible.
- Do not keep personal items on the table with secure materials.
  - No personal laptop or tablet use is allowed in the meeting rooms.
- Do not speak to other panelists about specific items outside of the meeting rooms.
- To limit disruptions, try to take breaks at designated break times.
- If you have any questions about the review or procedures, feel free to talk to AIR or DOE staff during breaks or at lunch.

American Institutes for Research 21

## Questions?

American Institutes for Research

## Appendix F Data Review Committee Participant Details

### **Data Review Committee Participant Details**

Table F-1. Data Review Committee Participants, Science

State	Date	Location	Grade Band	Number of Teachers in Each Group	Teacher Demographic Summary	Number of Items Reviewed
			Elementary School		Gender: Not collected	
	July 2018	Virtual	Middle School	18ª	Ethnicity: Not collected Region: Not collected	84
IOOD			High School		Teaching Experience: Not collected	
ICCR			Elementary School			
	August 2019	N/A	Middle School	N/A <sup>b</sup>	N/A <sup>b</sup>	43
			High School			
	August 2018	New Britain, Connecticut	Elementary School	10	Gender: Male 12%, Female 88%	18
			Middle School	8	Ethnicity: Not collected Region: Not collected	
0			High School	8	Teaching Experience: Not collected	
Connecticut	August 2019	Cromwell, Connecticut	Elementary School	7	Gender: Male 17%, Female 83%	53
			Middle School	10	Ethnicity: Not collected Region: Not collected	
			High School	6	Teaching Experience: Not collected	
	August 2018	Honolulu, Hawaii	Elementary School	18ª	Gender: Not collected	
Hawaii			Middle School		18ª	Ethnicity: Not collected  Region: Not collected
			High School		Teaching Experience: Not collected	
	August 2019	Honolulu, Hawaii	Elementary School	6	Gender: Male 29%, Female 71%	37

State	Date	Location	Grade Band	Number of Teachers in Each Group	Teacher Demographic Summary	Number of Items Reviewed
			Middle School	7	Ethnicity: American Indian and White 12%, Asian 41%, Asian and White 6%, Hispanic and White 12%, Native Hawaiian or Pacific Islander 18%, White 12%	
			High School	5	Region: Not collected Teaching Experience: General Education 53%, General Education with SPED Certification 6%, Bilingual Education 0%, Administration 0%, Other 29%, Special Education 12%	
Idaho	August 2019	N/A°	Elementary School	- 10ª	Gender: Male 20%, Female 70%, Did not specify 1% Ethnicity: White 100% Region: Rural 60%, Urban 0%, Suburban 40%	12
			Middle School		Teaching Experience: General Education 60%, Administration 2%, Coach 20%	
MSSAd	August 2018	N/A <sup>e</sup>	N/A <sup>e</sup>	N/A <sup>e</sup>	N/A <sup>e</sup>	9
MISSA	August 2019	N/A <sup>e</sup>	N/A <sup>e</sup>	N/A <sup>e</sup>	N/A <sup>e</sup>	14
	September 2018	Salem, Oregon	Elementary School	3	Gender: Male 18%, Female 82% Ethnicity: White 100%	44
			Middle School	4	Region: Urban 27%, Suburban 0%, Rural 73% Teaching Experience: Regular Education 64%, Bilingual	
0			High School	4	Education 55%, Special Education 36%, Administration 18%, Other 18%	
Oregon	August 2019	019 Remote	Elementary School	1	Gender: Male 50%, Female 50% Ethnicity: White 100%	
			Middle School	2	Region: Urban 50%, Suburban 0%, Rural 50% Teaching Experience: Regular Education 50%, Bilingual	8
			High School	1	Education 25%, Special Education 25%, Administration 25%, Other 75%	
		August 2018 Salt Lake City, Utah	Grade 6	6	Gender: Male 7%, Female 93% Ethnicity: White 87%, Unknown 13% Region: Urban 0%, Suburban 13%, Rural 27%, Unknown/no	40
Utah	August 2018		Grade 7	5		
			Grade 8	5	response 60%  Teaching Experience: General Education 100%	
West	July 2049	y 2018 N/A°	Elementary School		Gender: Not collected Ethnicity: Not collected	3
Virginia	July 2018		Middle School	Region: Not collected Teaching Experience: Not collected	ა 	

State	Date	Location	Grade Band	Number of Teachers in Each Group	Teacher Demographic Summary	Number of Items Reviewed
	September 2019	N/A°	Elementary School  Middle School	4 <sup>a</sup>	Gender: Not collected Ethnicity: Not collected Region: Not collected	7
			Wildule Scribbi		Teaching Experience: Not collected	
	October 2018	Cheyenne, Wyoming	Elementary School	11ª	Gender: Not collected Ethnicity: Not collected Region: Not collected Teaching Experience: Not collected	
			Middle School			16
Wyoming			High School			
	August 2019	Cheyenne, Wyoming	Elementary School	3	Gender: Male 10%, Female 90% Ethnicity: N/A Region: Urban 0% Suburban 40%, Rural 60% Teaching Experience: 90% Regular Education, 10% Administration	
			Middle School	4		16
			High School	3		

Note. <sup>a</sup>Number of Committee Members by grade band is not available.

<sup>&</sup>lt;sup>b</sup>In summer 2019, ICCR field-test items were taken to Connecticut, Hawaii, and Idaho for committee review.

<sup>&</sup>lt;sup>c</sup>Location of Data Review Committee Meeting is unavailable at the time of writing this report.

<sup>&</sup>lt;sup>d</sup>MSSA = Rhode Island and Vermont's Multi-State Science Assessment.

<sup>&</sup>lt;sup>e</sup>Conducted by the Rhode Island Department of Education and the Vermont Agency of Education science content experts.

# Appendix G Example Item Interactions

## Interaction Types Available in the Multi-State Science Assessment (MSSA)

#### **Review of Different Interaction Types**

Interaction Type	Associated Sub-Types	Legacy Item Types Supported
<u>Choice</u>	Multiple Choice	MC
	Multiple Select	MS
	Scaffolding	ASI2, ASI3
Text Entry	Simple Text Entry	EA, ECR, LA, OE, SA, SR, WCR, RW, SCR
	Embedded Text Entry	CL, FI
	Natural Language	NL
	Extended Response	ER
<u>Table</u>	Table Match	MI
	Table Input	ТІ
	Column Match	MI
Edit Task	Edit Task	ET
	Edit Task with Choice	ETC
	Edit Task Inline Choice	ETC
Hot Text	<u>Selectable</u>	HTQ
	Re-orderable	НТ
	Drag-from-Palette	DnD
	Custom	HTQ, HT, DnD
<u>Equation</u>	N/A	EQN
<u>Grid</u>	Grid	GI
	Hot Spot	GI
	Graphic Gap Match	GI
Simulation*	N/A	SIM

Note. the abbreviations correlate to the attributes used in CAI's Item Tracking System

#### **Multiple-Choice Interactions**

Multiple-Choice (MC) interactions require students to select a single option from a list of possible answer options. The number and orientation of answer options in a multiple-choice interaction are configurable. Answer options may appear vertically, horizontally, vertically-stacked (in a specified number of columns), or horizontally-stacked (in a specified number of rows).

What is the product of 68 and 90?				
A	612			
®	1,260			
	6,120			
<b>©</b>	6,300			

#### **Multiple-Select Interactions**

Multiple-Select interactions require students to select one or more options from a list of possible answer options. The number and orientation of answer options in a multiple-select interaction are configurable. Answer options may appear vertically, horizontally, horizontally-stacked (in a specified number of rows), or vertically-stacked (in a specified number of columns).

Select the values that are greater than or equal to ½.					
□ 0.6	□ .45				
□ 2/6	☐ One Fifth				
□ 5/8	□ 2/10				

#### **Text Entry Interactions**

The Text Entry Interaction Editor allows you to create content for the following interaction types:

- Error! Reference source not found.

#### **Simple Text Entry Interactions**

Simple Text Entry interactions require students to type a response in a text box. For Simple Text Entry interactions, we can allow you to specify the maximum response length for the text box and the type of text editor available to students.

Select a sentence in the passage that does not fit with the overall structure and explain why it is disruptive to the organization of the passage.
Type your answer in the space provided.

#### **Embedded Text Entry Interactions**

Embedded Text Entry interactions require students to type their response in one or more text boxes that are embedded in a section of read-only text.

Fill in the blanks in th	e sentence below.
The quick	fox jumps over the lazy

#### **Extended Response Interactions**

Extended Response interactions require students to type a response in a text box. Extended Response interactions are scored by an uploaded essay scoring model that analyzes the student's response to identify variations of acceptable key words and phrases. For Extended Text Entry interactions, we can allow you to specify the maximum response length for the text box and the type of text editor available to students.

Select a sentence in the passage that does not fit with the overall structure and explain why it is disruptive to the organization of the passage.		
Type your answer in the space provided.		



Alert: Extended Response interactions cannot be combined with any other interactions in the item.

#### **Table Entry Interaction**

The Table Entry Interaction Editor allows you to create content for the following interaction types:

- Error! Reference source not found.
- Error! Reference source not found.
- Error! Reference source not found.

#### **Table Match Interactions**

Table Match interactions arrange two sets of match options in a table, with one set listed in columns and the other set listed in rows. Students match options in the columns to options in the rows by marking checkboxes in the cells where the columns and rows intersect.

For each number listed in the rows of the table, mark the checkboxes for each column that describes that number.								
	Perfect Prime Odd Even Square Number Number Number							
5								
12								
9	9 🗆 🗆 🗆							

Table Match interactions allow you to customize the number of match options in each set and enter the content for each match option. You can also set restrictions on the number of matches students can make. By default, the panel includes a basic table consisting of three rows and columns (including the row header and column header).

#### **Table Input Interactions**

Table Input interactions provide students with a table that includes one or more blank cells. Each blank cell displays a text box in which students can type their response.

in the table below.	Enter a stage direction that you might give to each theater technician listed in the table below.  The first one has been done for you.					
Theater technicians						
Set designer	A circular bench around a small obelisk					
Props manager						
Sound technician						
Lighting technician						

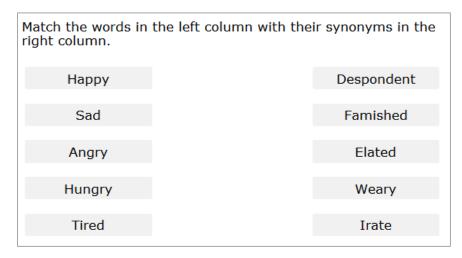
Table Input interactions allows you to customize the number of rows and columns in the table, specify which cells display text boxes, and enter content for the read-only cells. By default, the panel includes a basic table consisting of three rows and columns (including the row header and column header).



**Alert:** If a table does not include row headers, then it must include column headers. If a table does not include column headers, then it must include row headers.

#### **Column Match Interactions**

Column Match interactions provide students with two columns that each contain a set of match options. Students respond to the interaction by selecting a match option in the left column and then selecting the corresponding match option in the right column. A match option in one set may have one, multiple, or no matches in the other set.



Column Match interactions allows you to customize the number of match options in each set and enter the content for each match option. By default, the panel includes two single-column tables, each of which includes two match options. You can also set restrictions on the number of matches students can make.

#### **Edit Task Interactions**

The Edit Task Interaction Editor allows you to create content for the following interaction types:

- Error! Reference source not found.
- Error! Reference source not found.
- Error! Reference source not found.

#### **Edit Task Interactions**

Edit Task interactions provide students with a sentence or paragraph containing one or more tagged text elements. Tagged elements usually contain an error, such as improper spelling or grammar.

To respond to these interactions, students click a tagged element and enter corrected text in an editing window. The entered text replaces the original tagged text.

The sentence below contains several grammatical mistakes. Click the highlighted words to correct the grammar.

The quick foxes jumps over the lazy, dogs.

Edit Task interactions allow you to enter the text that appears in the response area and tag elements within the text that students can edit.



Warning: You cannot include hand-scored and machine-scored interactions in the same item.

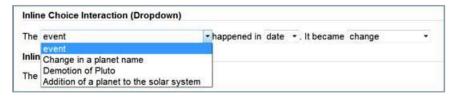
#### **Edit Task with Choice Interactions**

Edit Task with Choice interactions are similar to Edit Task interactions. The only difference is that when responding to Edit Task with Choice interactions, students replace the tagged text elements with options selected from a drop-down list.

Edit Task with Choice interactions allow you to enter the text that appears in the response area and tag elements within the text that students can edit.

#### **Edit Task Inline Choice Interactions**

Edit Task Inline Choice interactions are similar to Edit Task with Choice interactions. The only difference is that students select replacement options from a drop-down list embedded within the read-only text, rather than accessing the drop-down list via a pop-up window.



#### **Hot Text Interactions**

The Hot Text Interaction Editor allows you to create content for the following interaction types:

- Error! Reference source not found.

#### **Selectable Hot Text Interactions**

Selectable Hot Text interactions require students to select one or more text elements in the response area.

Select the sentences that support the inference that the area is in danger of losing its moose population. Select **all** that apply.

A similar boom-and-bust cycle occurs between predator and prey. Ten times the size of a wolf, a moose has long, strong legs and a dangerous kick. So wolves prey mainly on old and weak animals. Good hunting means food for the whole pack. Wolves then raise lots of pups, and their numbers increase. More wolves mean more mouths to feed and more moose get eaten. However, when the moose population decreases, wolves starve.

Selectable Hot Text interactions allows you to set the minimum and maximum number of elements students can select, enter the text that appears in the response area, and tag the text elements that will be selectable.

#### Re-orderable Hot Text Interactions

Re-orderable Hot Text interactions require students to click and drag hot text elements into a different order.

Place the following sentences in the correct order.

Hey Jude. And make it better. Don't be afraid. Take a sad song.

Re-orderable Hot Text interactions allow you to enter the re-orderable text elements in the response area. You can specify the elements' orientation and set them to appear in random order to students.

#### Drag-from-Palette Hot Text Interactions a.k.a. Hot Text Gap Match

Drag-from-Palette Hot Text interactions require students to drag elements from a palette into the available blank table cells or "gaps" (text boxes) in the response area. Palette elements may consist of text and/or images. Students may be able to drag the same palette element into multiple gaps, depending on the interaction's configuration.

Drag and drop the characteristics into the appropriate table cells be					
Fortunato's character Montressor's character					
Sinister and calculating					
Cowardly and irreverent					
Egotistical and rude					
Lazy and inconsiderate					

Drag-from-Palette Hot Text interactions allow you to enter the elements that appear in the palette, enter static text for the response area, and create the gap targets where students can drag the text elements. You can enter all of the elements in a single text box or enter each segment in its own text box.

- Can set a minimum/maximum number of times a student is required/allowed to use a specific palette object
- Only supports drag-and-drop of palette items (images or plain text) onto pre-defined drop targets ("gaps" or "blanks") in the body text
  - These palette items are always confined to a special palette region (no "preplacing" them)
  - There is some control over palette placement
  - o The items can only be placed in predefined "target" regions

#### **Custom Hot Text Interactions**

Custom Hot Text interactions combine the functionality of the other Hot Text interaction subtypes. Students responding to a Custom Hot Text interaction may need to select text elements, rearrange text elements, and/or drag text elements from a palette to blank table cells or drop targets in the response area. In many ways, this is the grid of the text-interaction world. In practice, it is typically used to do drag-and-drop with text, but it can technically do more:

- Supports dragging and dropping text elements onto drop target areas
  - Text elements can originally be placed anywhere in the interaction (there's no dedicated palette)
  - Multiple elements can be dropped onto a target
    - this constitutes a "group"
    - much like grid hotspots, you can set constraints on the group

- Supports selectable text elements
- Like grid hotspots, these too can be grouped

Use the word bank to fill in the blank in the sentence below. Then, select all the words in the sentence that are nouns.
Word bank:
young dull good rich
Sentence:
All work and no play makes Jack a boy.

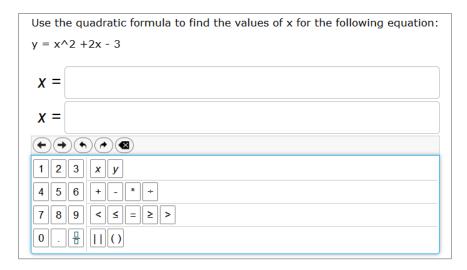
Custom Hot Text interactions allow you to create groups of text elements, as well as the drop targets and static text that appear in the response area. When you create a group of text elements, you must assign a Hot Text functionality to that group. The following functionalities are available:

- **Selectable:** When you assign this functionality to a group, the text elements in the group behave like elements in a Selectable Hot Text interaction. You cannot add drop target elements to this kind of group.
- **Draggable:** When you assign this functionality to a group, the text elements in the group behave like elements in a Re-Orderable Hot Text interaction. If you assign this functionality to a group and also add drop targets to the group, the text elements in the group behave like elements in a Drag-from-Palette Hot Text interaction.

You can create as many groups as you wish, but you can only assign one Hot Text functionality to each group.

#### **Equation Interaction Editor**

The Equation Interaction Editor allows you to create content for Equation interactions only. Equation interactions require students to enter a response into input boxes using an on-screen keypad, which may consist of special mathematics characters. Students can also enter their response via a physical keyboard, but they cannot enter any characters that are not included in the on-screen keyboard.



Equation interactions allow you to select the buttons to include in the on-screen keypad, enter static text in the response area, and specify the number of input boxes to include in the response area. When selecting buttons to include in the keypad, you can add individual buttons or an entire row or tab of buttons.

#### **Grid Interactions**

The Grid Interaction Editor allows you to create content for the following interaction types:

- Grid Interactions
- Hot Spot Interactions
- Graphic Gap Match Interactions

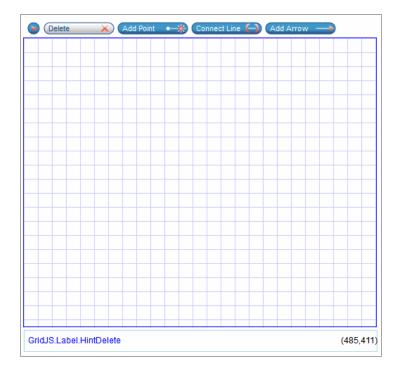


**Note:** Although there are three options available in the **Interaction Type** drop-down list, the generic **Grid** option allows you to create interactions with functionality similar to Hot Spot and Graphic Gap Match sub-types.

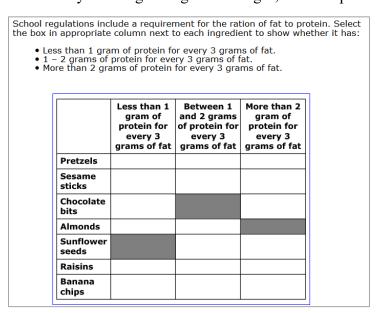
#### **Grid Interactions Types**

Grid interactions require students to enter a response by interacting with a grid area in the answer space. There are three general ways in which students can interact with the grid area.

• **Graphing Functionality:** Students can use various tool buttons to add points, lines, and other geometric shapes to the grid area. Only the Grid interaction sub-type allows you to create interactions with this functionality.



- **Hot Spot Functionality:** Students can click or hover over interactive regions in the grid area (hot spots) in order to activate them. Activated hot spots become highlighted, become outlined, or display an image. The Grid and Hot Spot interaction sub-types allow you to create interactions with this functionality.
  - o Hotspots can be defined in groups, each of which can have its own selection constraints
  - These regions support events so clicking a hotspot might change the appearance of the interaction by showing/hiding other images, for example



- **Drag-and-Drop Functionality:** Students can click image or text objects and drag them into various locations in the grid area. The objects for these interactions are either provided in a palette beside the grid area or pre-placed within the grid area itself. The Grid and Graphic Gap Match interaction sub-types allow you to create interactions with this functionality; however, only Graphic Gap Match interactions allow text objects.
  - o These palette items can be "preplaced" on the canvas or listed in a separate palette
  - The items can be placed anywhere on the canvas or guided to specific regions with snap points





**Note:** The functionalities of these interaction types are not mutually exclusive. A single Grid interaction may require students to select hot spots and place objects, or graph lines and select hot spots, and so on. However, a Grid interaction cannot include preplaced objects if it also includes the **Delete** tool button above the grid area.

#### **Grid Hot Spot Interactions**

Hot Spot interaction sub-types allow you to create Grid interactions with hot spot functionality. These interactions require students to select hot spot regions in the grid area.

- Only supports click-to-select "hotspots"
  - o No visual side-effect events are supported
  - No hotspot groups are supported

#### **Grid Graphic Gap Match Interactions**

Graphic Gap Match interactions allow you to create Grid interactions with both hot spot and dragand-drop functionality. These interactions require students to drag image objects from a palette to hot spot regions (gaps) in the grid area.

• Only supports drag-and-drop of palette items (images or plain text) onto the canvas/background

- These palette items are always confined to a special palette region (no "preplacing" them on the canvas)
- o The items can only be placed in predefined "target" regions



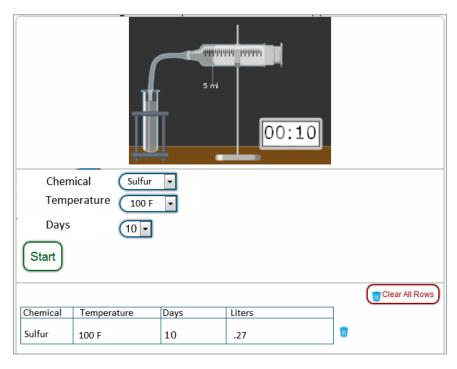
**Alert:** Graphic Gap Match interactions do not allow you to enable Snap-to-Point or Snap-to-Grid mode. You cannot pre-place image or text objects in the grid area with Graphic Gap Match Interactions.

Basically, graphic gap match and hotspot are dedicated interactions that don't support all the features of a grid. The trade-off here is:

- Graphic gap match and hotspot interactions are rendered differently (more simplistically)
- In some ways, graphic gap match and hotspot are easier to author and maintain
- Grid interactions need to use the "grid rubric tool," which is quite complicated

#### **Simulation Interaction Editor**

The Simulation Interaction Editor allows you to create content for Simulation interactions only. Simulation interactions consist of an animation tool, a set of input tools, and an output table. Students select parameters from the input tools to influence the animation. After the animation runs, the simulation results appear in the output table. Students can run multiple trials with different parameters to insert additional rows into this table.



## Appendix H Shared Science Assessment Item Bank

# **Shared Science Assessment Item Bank**

Table H-1. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank by Performance Expectation, Elementary School

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items
		4-ESS1-1	3	0	9	12
	ESS1	5-ESS1-1	2	2	10	14
		5-ESS1-2	7	0	8	15
		3-ESS2-1	4	1	5	10
		3-ESS2-2	2	0	7	9
	ESS2	4-ESS2-1	2	0	8	10
Earth and Space Sciences	E552	4-ESS2-2	2	0	9	11
		5-ESS2-1	0	1	8	9
		5-ESS2-2	3	1	5	9
		3-ESS3-1	3	1	6	10
	F662	4-ESS3-1	5	1	3	9
	ESS3	4-ESS3-2	6	2	5	13
		5-ESS3-1	3	1	8	12
		3-LS1-1	3	2	5	10
		4-LS1-1	10	0	9	19
	LS1	4-LS1-2	2	1	11	14
		5-LS1-1	2	0	13	15
	1.60	3-LS2-1	4	0	8	12
Life Ocionese	LS2	5-LS2-1	1	1	7	9
Life Sciences	1.00	3-LS3-1	3	2	5	10
	LS3	3-LS3-2	1	1	4	6
		3-LS4-1	2	0	8	10
	104	3-LS4-2	8	0	4	12
	LS4	3-LS4-3	5	0	5	10
		3-LS4-4	3	0	5	8
Dhysical Caionas	DC4	5-PS1-1	4	0	9	13
Physical Sciences	PS1	5-PS1-2	3	1	8	12

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items
		5-PS1-3	4	1	7	12
		5-PS1-4	1	2	7	10
		3-PS2-1	4	1	6	11
		3-PS2-2	3	0	3	6
	PS2	3-PS2-3	1	0	5	6
		3-PS2-4	1	1	4	6
		5-PS2-1	2	0	5	7
		4-PS3-1	4	0	9	13
		4-PS3-2	5	0	4	9
	PS3	4-PS3-3	3	0	8	11
		4-PS3-4	3	0	9	12
		5-PS3-1	2	1	7	10
		4-PS4-1	1	0	8	9
	PS4	4-PS4-2	1	0	8	9
		4-PS4-3	2	0	3	5
otal	•		130	24	285	439

Note. aMOU state item sources include Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming.

Table H-2. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank by Performance Expectation, Middle School

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items <sup>b</sup>
		MS-ESS1-1	7	0	4	11
	ESS1	MS-ESS1-2	3	0	5	8
	E331	MS-ESS1-3	4	0	7	11
		MS-ESS1-4	1	1	7	9
		MS-ESS2-1	1	0	7	8
		MS-ESS2-2	3	1	5	9
	ESS2	MS-ESS2-3	2	0	7	9
Earth and Space Sciences	E332	MS-ESS2-4	1	0	5	6
		MS-ESS2-5	1	0	5	6
		MS-ESS2-6	1	1	2	4
		MS-ESS3-1	2	0	4	6
		MS-ESS3-2	2	0	8	10
	ESS3	MS-ESS3-3	0	1	5	6
		MS-ESS3-4	3	1	8	12
		MS-ESS3-5	3	1	4	8
		MS-LS1-1	0	0	5	5
		MS-LS1-2	2	1	6	9
		MS-LS1-3	0	0	6	6
	1.04	MS-LS1-4	2	0	2	4
	LS1	MS-LS1-5	0	2	4	6
		MS-LS1-6	3	1	5	9
Life Sciences		MS-LS1-7	1	1	6	8
Life Sciences		MS-LS1-8	2	0	6	8
		HS-LS2-4	0	0	1	1
		MS-LS2-1	3	0	9	12
	1.00	MS-LS2-2	3	0	4	7
	LS2	MS-LS2-3	2	1	6	9
		MS-LS2-4	8	0	7	15
		MS-LS2-5	4	1	6	11

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Banl Items <sup>b</sup>
	LS3	MS-LS3-1	1	0	6	7
	LSS	MS-LS3-2	3	0	8	11
		MS-LS4-1	5	0	5	10
		MS-LS4-2	1	0	7	8
	LS4	MS-LS4-3	2	0	4	6
		MS-LS4-4	2	0	4	6
		MS-LS4-5	3	1	1	5
		MS-LS4-6	1	1	5	7
		MS-PS1-1	1	0	5	6
		MS-PS1-2	3	1	5	9
	PS1	MS-PS1-3	1	1	4	6
	P31	MS-PS1-4	1	0	8	9
		MS-PS1-5	1	0	8	9
		MS-PS1-6	2	1	2	5
		MS-PS2-1	1	0	4	5
		MS-PS2-2	1	0	5	6
	PS2	MS-PS2-3	1	1	4	6
Physical Sciences		MS-PS2-4	0	0	8	8
-		MS-PS2-5	0	0	8	8
		MS-PS3-1	2	1	4	7
		MS-PS3-2	1	0	7	8
	PS3	MS-PS3-3	3	1	5	9
		MS-PS3-4	3	1	3	7
		MS-PS3-5	5	0	5	10
		MS-PS4-1	1	0	8	9
	PS4	MS-PS4-2	5	0	6	11
		MS-PS4-3	1	1	5	7
			115	23	300	438

*Note.* <sup>a</sup>MOU state item sources include Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming. <sup>b</sup>Count excludes seven middle school MOU items that do not align to the NGSS.

Table H-3. Spring 2021 Shared Science Assessment Operational and Field-Test Item Bank by Performance Expectation, High School

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items <sup>b</sup>
		HS-ESS1-1	1	1	3	5
		HS-ESS1-2	2	0	1	3
	ESS1	HS-ESS1-3	1	0	2	3
	E331	HS-ESS1-4	1	2	3	6
		HS-ESS1-5	1	0	3	4
		HS-ESS1-6	1	0	3	4
		HS-ESS2-1	0	1	1	2
		HS-ESS2-2	2	1	3	6
		HS-ESS2-3	1	0	1	2
Earth and Space Sciences	ESS2	HS-ESS2-4	1	0	5	6
•		HS-ESS2-5	0	0	2	2
		HS-ESS2-6	2	1	2	5
		HS-ESS2-7	1	0	2	3
	5000	HS-ESS3-1	2	0	3	5
		HS-ESS3-2	1	1	2	4
		HS-ESS3-3	0	1	2	3
	ESS3	HS-ESS3-4	1	0	2	3
		HS-ESS3-5	2	0	4	6
		HS-ESS3-6	1	0	2	3
		HS-LS1-1	3	0	7	10
		HS-LS1-2	3	0	7	10
		HS-LS1-3	0	0	3	3
	LS1	HS-LS1-4	5	0	3	8
		HS-LS1-5	1	1	5	7
Life Sciences		HS-LS1-6	4	0	2	6
		HS-LS1-7	2	0	5	7
		HS-LS2-1	2	0	3	5
	LS2	HS-LS2-2	2	1	6	9
		HS-LS2-3	1	0	5	6

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items <sup>b</sup>
		HS-LS2-4	5	1	4	10
		HS-LS2-5	2	0	4	6
		HS-LS2-6	3	0	4	7
		HS-LS2-7	4	0	5	9
		HS-LS2-8	1	0	1	2
		HS-LS3-1	3	0	6	9
	LS3	HS-LS3-2	4	1	3	8
		HS-LS3-3	3	0	4	7
		HS-LS4-1	8	0	5	13
		HS-LS4-2	4	0	4	8
		HS-LS4-3	2	0	5	7
	LS4	HS-LS4-4	2	0	4	6
		HS-LS4-5	5	0	4	9
		HS-LS4-6	0	0	2	2
		HS-PS1-1	2	0	4	6
		HS-PS1-2	3	0	5	8
		HS-PS1-3	2	1	5	8
		HS-PS1-4	1	0	2	3
	PS1	HS-PS1-5	1	0	7	8
		HS-PS1-6	2	0	3	5
		HS-PS1-7	3	0	4	7
		HS-PS1-8	0	1	3	4
		HS-PS2-1	2	0	4	6
Physical Sciences		HS-PS2-2	1	1	4	6
		HS-PS2-3	0	0	4	4
	PS2	HS-PS2-4	3	0	3	6
		HS-PS2-5	1	0	1	2
		HS-PS2-6	1	0	3	4
		HS-PS3-1	1	0	3	4
		HS-PS3-2	1	0	4	5
	PS3	HS-PS3-3	1	0	7	8
		HS-PS3-4	1	0	3	4

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Bank Items <sup>b</sup>
		HS-PS3-5	2	1	3	6
		HS-PS4-1	2	0	3	5
		HS-PS4-2	0	0	1	1
	PS4	HS-PS4-3	0	0	4	4
		HS-PS4-4	0	0	3	3
		HS-PS4-5	2	0	1	3
Total	_		122	16	231	369

*Note*. <sup>a</sup>MOU state item sources include Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming. <sup>b</sup>Count excludes one high school MOU item that does not align to the NGSS.

# Appendix I Multi-State Science Assessment (MSSA) Item Pool

# **Multi-State Science Assessment Item Pool**

Table I-1. Spring 2021 MSSA Operational and Field-Test Item Pool by Performance Expectation, Grade 5

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Item Poo
		4-ESS1-1	2	0	2	4
	ESS1	5-ESS1-1	1	2	4	7
		5-ESS1-2	3	0	3	6
		3-ESS2-1	4	1	1	6
		3-ESS2-2	2	0	2	4
Fauth and Occasi	ESS2	4-ESS2-1	1	0	2	3
Earth and Space Sciences	E552	4-ESS2-2	1	0	2	3
Sciences		5-ESS2-1	0	1	3	4
		5-ESS2-2	1	1	2	4
	ESS3	3-ESS3-1	1	1	0	2
		4-ESS3-1	2	1	1	4
		4-ESS3-2	3	2	0	5
		5-ESS3-1	1	1	2	4
	LS1	3-LS1-1	1	2	3	6
		4-LS1-1	5	0	2	7
		4-LS1-2	1	1	3	5
		5-LS1-1	1	0	4	5
	LS2	3-LS2-1	3	0	2	5
Life Sciences	LSZ	5-LS2-1	1	1	2	4
Life Sciences	LS3	3-LS3-1	1	2	3	6
	LS3	3-LS3-2	1	1	3	5
		3-LS4-1	2	0	2	4
	LS4	3-LS4-2	3	0	2	5
	L54	3-LS4-3	2	0	2	4
		3-LS4-4	1	0	1	2
		5-PS1-1	3	0	4	7
<b>Physical Sciences</b>	PS1	5-PS1-2	2	1	3	6
		5-PS1-3	2	1	1	4

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Item Pool
		5-PS1-4	0	2	1	3
		3-PS2-1	2	1	3	6
		3-PS2-2	2	0	1	3
	PS2	3-PS2-3	1	0	1	2
		3-PS2-4	1	1	1	3
		5-PS2-1	1	0	2	3
		4-PS3-1	4	0	4	8
		4-PS3-2	3	0	2	5
	PS3	4-PS3-3	3	0	3	6
		4-PS3-4	1	0	4	5
		5-PS3-1	2	1	2	5
		4-PS4-1	0	0	3	3
	PS4	4-PS4-2	1	0	5	6
		4-PS4-3	1	0	1	2
otal	•		73	24	94	191

Note. aMOU state items administered includes Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming.

Table I-2. Spring 2021 MSSA Operational and Field-Test Item Pool by Performance Expectation, Grade 8

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Item Poo
		MS-ESS1-1	1	0	2	3
	E004	MS-ESS1-2	1	0	3	4
	ESS1	MS-ESS1-3	2	0	3	5
		MS-ESS1-4	1	1	4	6
		MS-ESS2-1	0	0	3	3
		MS-ESS2-2	2	1	3	6
Fauth and Once	ESS2	MS-ESS2-3	1	0	4	5
Earth and Space Sciences	E332	MS-ESS2-4	1	0	3	4
Sciences		MS-ESS2-5	1	0	4	5
		MS-ESS2-6	0	1	1	2
		MS-ESS3-1	2	0	0	2
		MS-ESS3-2	1	0	5	6
	ESS3	MS-ESS3-3	0	1	0	1
		MS-ESS3-4	2	1	2	5
		MS-ESS3-5	3	1	0	4
		MS-LS1-1	0	0	3	3
		MS-LS1-2	0	1	3	4
		MS-LS1-3	0	0	3	3
	LS1	MS-LS1-4	2	0	2	4
	LSI	MS-LS1-5	0	2	0	2
		MS-LS1-6	1	1	2	4
		MS-LS1-7	1	1	3	5
Life Sciences		MS-LS1-8	1	0	2	3
		MS-LS2-1	3	0	4	7
	LS2	MS-LS2-3	0	1	3	4
	LOZ	MS-LS2-4	1	0	4	5
		MS-LS2-5	4	1	1	6
	LS3	MS-LS3-1	0	0	2	2
	LSS	MS-LS3-2	2	0	4	6
	LS4	MS-LS4-1	3	0	3	6

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Item Poo
		MS-LS4-2	0	0	3	3
		MS-LS4-3	2	0	1	3
		MS-LS4-4	2	0	1	3
		MS-LS4-5	2	1	0	3
		MS-LS4-6	0	1	1	2
		MS-PS1-1	1	0	2	3
		MS-PS1-2	2	1	3	6
	PS1	MS-PS1-3	1	1	1	3
	P31	MS-PS1-4	0	0	2	2
		MS-PS1-5	0	0	6	6
		MS-PS1-6	1	1	1	3
		MS-PS2-1	1	0	0	1
		MS-PS2-2	0	0	3	3
	PS2	MS-PS2-3	1	1	3	5
Physical Sciences		MS-PS2-4	0	0	2	2
		MS-PS2-5	0	0	2	2
		MS-PS3-1	0	1	1	2
		MS-PS3-2	1	0	3	4
	PS3	MS-PS3-3	1	1	3	5
		MS-PS3-4	2	1	0	3
		MS-PS3-5	4	0	3	7
		MS-PS4-1	1	0	4	5
	PS4	MS-PS4-2	3	0	2	5
		MS-PS4-3	0	1	1	2
otal			61	23	124	208

Note. aMOU state items administered includes Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming.

Table I-3. Spring 2021 MSSA Operational and Field-Test Item Pool by Performance Expectation, Grade 11

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Item Poo
		HS-ESS1-1	1	1	2	4
		HS-ESS1-2	2	0	1	3
	ESS1	HS-ESS1-3	1	0	1	2
	E331	HS-ESS1-4	1	2	0	3
		HS-ESS1-5	1	0	2	3
		HS-ESS1-6	1	0	0	1
		HS-ESS2-1	0	1	0	1
		HS-ESS2-2	2	1	2	5
Earth and Space Sciences		HS-ESS2-3	1	0	0	1
	ESS2	HS-ESS2-4	0	0	3	3
		HS-ESS2-5	0	0	1	1
		HS-ESS2-6	1	1	0	2
		HS-ESS2-7	1	0	1	2
		HS-ESS3-1	2	0	0	2
		HS-ESS3-2	1	1	1	3
	ESS3	HS-ESS3-3	0	1	1	2
		HS-ESS3-4	1	0	0	1
		HS-ESS3-5	2	0	3	5
		HS-ESS3-6	1	0	1	2
		HS-LS1-1	0	0	1	1
		HS-LS1-2	2	0	2	4
		HS-LS1-3	0	0	2	2
	LS1	HS-LS1-4	2	0	1	3
		HS-LS1-5	0	1	1	2
Life Sciences		HS-LS1-6	4	0	0	4
		HS-LS1-7	1	0	0	1
		HS-LS2-1	2	0	2	4
	1.00	HS-LS2-2	2	1	1	4
	LS2	HS-LS2-3	0	0	1	1
		HS-LS2-4	2	1	1	4

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Item Poo
		HS-LS2-5	0	0	1	1
		HS-LS2-6	3	0	1	4
		HS-LS2-7	2	0	1	3
		HS-LS2-8	1	0	0	1
		HS-LS3-1	1	0	1	2
	LS3	HS-LS3-2	3	1	1	5
		HS-LS3-3	1	0	1	2
		HS-LS4-1	4	0	1	5
		HS-LS4-2	2	0	2	4
	LS4	HS-LS4-3	2	0	2	4
		HS-LS4-4	2	0	2	4
		HS-LS4-5	0	0	2	2
	PS1	HS-PS1-1	1	0	2	3
		HS-PS1-2	2	0	2	4
		HS-PS1-3	2	1	3	6
		HS-PS1-4	1	0	0	1
		HS-PS1-5	0	0	1	1
		HS-PS1-6	2	0	1	3
		HS-PS1-7	3	0	1	4
		HS-PS1-8	0	1	1	2
	PS2	HS-PS2-1	1	0	4	5
		HS-PS2-2	1	1	2	4
Physical Sciences		HS-PS2-3	0	0	1	1
		HS-PS2-4	3	0	2	5
		HS-PS2-5	1	0	0	1
		HS-PS2-6	<u>.</u> 1	0	0	1
	PS3	HS-PS3-1	<u>.</u> 1	0	0	1
		HS-PS3-2	 1	0	2	3
		HS-PS3-3	0	0	2	2
		HS-PS3-4	0	0	1	1
		HS-PS3-5	2	1	2	5
	PS4	HS-PS4-1	0	0	1	1

Science Discipline	Disciplinary Core Idea	Performance Expectation	ICCR Items	MSSA Items	MOU Items <sup>a</sup>	Total Item Pool
		HS-PS4-3	0	0	2	2
		HS-PS4-4	0	0	1	1
		HS-PS4-5	2	0	0	2
Total			79	16	77	172

Note. aMOU state items administered includes Connecticut, Hawaii, Idaho, Montana, Oregon, Utah, West Virginia, and Wyoming.

# Appendix J Adaptive Algorithm Design

# TABLE OF CONTENTS

1.	Introduction, Background, and Definition	ONS2
1.1	Blueprint	
1.2	Content Value	
	1,2.1 Content Value for Single Items	
	1,2,2 Content Value for Sets of Items	
1.3	Information Value	7
	1.3.1 Individual Information Value	7
	1,3,2 Binary Items	
	1.3.3 Polytomous Items	
	1.3.4 Item Group Information Value	10
2.	Entry and Initialization	11
2.1	Item Pool	11
2.2	Adjust Segment Length	
2.3	Initialization of Starting Theta Estimates	
2.4	Insertion Of Embedded Field-Test Items	12
3.	ITEM SELECTION	13
3.1	Trimming The Custom Item Pool	13
3.2	Recycling Algorithm	
3.3	ADAPTIVE ITEM SELECTION	14
3.4	Selection of The Initial Item	15
3.5	Exposure Control	15
4.	TERMINATION	16
APPI	PENDIX 1. DEFINITIONS OF USER-SETTABLE PARAM	METERS17
APPI	PENDIX 2. SUPPORTING DATA STRUCTURES	19
ADD	DENDUM. ADJUSTMENTS TO THE USE OF ITEM CLU	STERS20

# **Adaptive Item Selection Algorithm**

## 1. Introduction, Background, and Definitions

This document describes the adaptive item selection algorithm. The item selection algorithm is designed to cover a standards-based blueprint, which may include content, cognitive complexity, and item type constraints. The item selection algorithm will also include:

- the ability to customize an item pool based on access constraints and screen items that have been previously viewed or may not be accessible for a given individual;
- a mechanism for inserting embedded field-test items; and
- a mechanism for delivering "segmented" tests in which separate parts of the test are administered in a fixed order.

This document describes the algorithm and the design for its implementation for the test delivery system (TDS). The implementation builds extensively on the algorithm implemented in the Cambium Assessment, Inc (CAI)'s TDS and incorporates substantial CAI intellectual property. CAI will release the algorithm and the implementation described here under the same open-source license under which the rest of the open-source system is released.

The general approach described here is based on a highly parameterized multiple-objective utility function. The objective function includes:

- a measure of content match to the blueprint;
- a measure of overall test information; and
- measures of test information for each reporting category on the test.

We define an objective function that measures an item's contribution to each of these objectives, weighting them to achieve the desired balance among them. Equation (1) sketches this objective function for a single item.

$$f_{ijt} = w_2 \frac{1}{\sum_{r=1}^{R} d_{rj}} \sum_{r=1}^{R} s_{rit} p_r d_{rj} + w_1 \sum_{k=1}^{K} q_k h_{1k}(v_{kijt}, V_{kit}, t_k) + w_0 h_0(u_{ijt}, U_{it}, t_0)$$
(1)

where the term w represents user-supplied weights that assign relative importance to meeting each of the objectives  $d_{rj}$  indicates whether item j has the blueprint-specified feature r, and  $p_r$  is the user-supplied priority weight for feature r. The term  $s_{rit}$  is an adaptive control parameter that is described. In general,  $s_{rit}$  increases for features that have not met their designated minimum as the end of the test approaches.

The remainder of the terms represents an item's contribution to measurement precision:

- $v_{kijt}$  is the value of item j toward reducing the measurement error for reporting category k for examinee i at selection t; and
- $u_{ijt}$  is the value of item j in terms of reducing the overall measurement error for examinee i at selection t.

The terms  $U_{it}$  and  $V_{kit}$  represent the total information overall and on reporting category k, respectively.

The term  $q_k$  is a user-supplied priority weight associated with the precision of the score estimate for reporting category k. The terms t represent precision targets for the overall score  $(t_0)$  and each score reporting category score. The functions h(.) are given by:

$$h_0(u_{ijt}, U_{it}, t_0) = \begin{cases} au_{ijt} & \text{if } U_{it} < t_0 \\ bu_{ijt} & \text{otherwise} \end{cases}$$

$$h_{1k}(v_{kijt}, V_{kit}, t_k) = \begin{cases} c_k v_{kijt} & \text{if } V_{kit} < t_k \\ d_k v_{kijt} & \text{otherwise} \end{cases}$$

Items can be selected to maximize the value of this function. This objective function can be manipulated to produce a pure, standards-free adaptive algorithm by setting  $w_2$  to zero or a completely blueprint-driven test by setting  $w_1 = w_0 = 0$ . Adjusting the weights to optimize performance for a given item pool will enable users to maximize information subject to the constraint that the blueprint is virtually always met.

We note that the computations of the content values and information values generate values on very different scales, and that the scale of the content value varies as the test progresses. Therefore, we normalize both the information and content values before computing the value of Equation (1).

This normalization is given by  $x = \begin{cases} 1 & \text{if } min = max \\ \frac{v - min}{max - min} & \text{otherwise} \end{cases}$ , where min and max represent the

minimum and maximum, respectively, of the metric computed over the current set of items or item groups.

The remainder of this section describes the overall program flow, the form of the blueprint, and the various value calculations employed in the objective function. Subsequent sections describe the details of the selection algorithm.

#### 1.1 BLUEPRINT

Each test will be described by a single blueprint for each segment of the test and will identify the order in which the segments appear. The blueprint will include:

- an indicator of whether the test is adaptive or fixed form;
- termination conditions for the segment, which are described in a subsequent section;
- a set of nested content constraints, each of which is expressed as:

- o the minimum number of items to be administered within the content category;
- o the maximum number of items to be administered within the content category;
- o an indication of whether the maximum should be deterministically enforced (a "strict" maximum);
- o a priority weight for the content category  $p_r$ ;
- o an explicit indicator as to whether this content category is a reporting category; and
- o an explicit precision-priority weight  $(q_k)$  for each group identified as a reporting category.
- a set of non-nested content constraints, which are represented as:
  - o a name for the collection of items meeting the constraint;
  - o the minimum number of items to be administered from this group of items;
  - o the maximum number of items to be administered from this group of items;
  - o an indication of whether the maximum should be deterministically enforced (a "strict" maximum);
  - o a priority weight for the group of items  $p_r$ ;
  - o an explicit indicator as to whether this named group will make up a reporting category; and
  - o an explicit precision-priority weight  $(q_k)$  for each group identified as a reporting category.
  - O The priority weights,  $p_r$  on the blueprint, can be used to express values in the blueprint match. Large weights on reporting categories paired with low (or zero) weights on the content categories below them may allow more flexibility to maximize information in a content category covering fewer fine-grained targets, while the reverse would mitigate toward more reliable coverage of finer-grained categories, with less content flexibility within reporting categories.

An example of a blueprint specification appears in Appendix J-1.

#### 1.2 CONTENT VALUE

Each item or item group will be characterized by its contribution to meeting the blueprint, given the items that have already been administered at any point. The contribution is based on the presence or absence of features specified in the blueprint and denoted by the term d in Equation (1). This section describes the computation of the content value.

## 1.2.1 Content Value for Single Items

For each constraint appearing in the blueprint (r), an item i either does or does not have the characteristic described by the constraint. For example, a constraint might require a minimum of four and a maximum of six algebra items. An item measuring algebra has the described characteristic, and an item measuring geometry, but algebra does not. To capture this constraint, we define the following:

- $d_j$  is a feature vector in which the elements are  $d_{rj}$ , summarizing item j's contribution to meeting the blueprint. This feature vector includes content categories such as claims and targets as well as other features of the blueprint, such as Depth of Knowledge (DOK) and item type.
- $S_{it}$  is a diagonal matrix, the diagonal elements of which are the adaptive control parameters  $S_{rit}$ .
- p is the vector containing the user-supplied priority weights  $p_r$ .

The scalar content value for an item is given by  $C_{ijt} = d_i S_{it} p$ .

Letting  $z_{rit}$  represent the number of items with feature r administered to student i by iteration t, the value of the adaptive control parameters is:

$$s_{rit} = \begin{cases} m_{it} \left( 2 - \frac{z_{rit}}{Min_r} \right) & \text{if } z_r < Min_r \\ 1 - \frac{z_{rit} - Min_r}{Max_r - Min_r} & \text{if } Min_r < z_{rit} < Max_r \\ \left( Max_r - z_{rit} \right) - 1 & \text{if } Max_r \le z_{rit} \end{cases}$$

The blueprint defines the minimum  $(Min_r)$  and maximum  $(Max_r)$  number of items to be administered with each characteristic (r).

The term  $m_{it} = \frac{T}{T-t}$  where T is the total test length. This has the effect of increasing the algorithm's preference for items that have not yet met their minimums as the end of the test nears and the opportunities to meet the minimum diminish.

This increases the likelihood of selecting items for content that has not met its minimum as the opportunities to do so are used up. The value s is highest for items with content that has not met its minimum, declines for items representing content for which the minimum number of items has been reached but the maximum has not, and turns negative for items representing content that has met the maximum.

#### 1.2.2 Content Value for Sets of Items

Calculation of the content value of sets of items is complicated by two factors:

- 1. The desire to allow more items to be developed for each set and to have the most advantageous set of items administered.
- 2. The design objective of characterizing the information contribution of a set of items as the expected information over the working theta distribution for the examinee.

The former objective is believed to enhance the ability to satisfy highly constrained blueprints while still adapting to obtain good measurement for a broad range of students. The latter arises from the recognition that English Language Arts (ELA) tests will select one set of items at a time, without an opportunity to adapt once the passage has been selected.

The general approach involves successive selection of the highest content value item in the set until the indicated number of items in the set have been selected. Because the content value of an item changes with each selection, a temporary copy of the already-administered content vector for the examinee is updated with each selection such that subsequent selections reflect the items selected in previous iterations.

Exhibit A on the following page presents a flowchart for this calculation. Readers will note the check to determine whether  $w_0 > 0$  or  $w_1 > 0$ . These weights, defined with Equation (1), identify the user-supplied importance of information optimization relative to blueprint optimization. In cases such as independent field tests, this weight may be set to zero, as it may not be desirable to make item administration dependent on the match to student performance. In more typical adaptive cases where item statistics will not be recalculated, favoring more informative items is generally better. The final measure of content value for the set of selected set of items is divided by the number of items selected to avoid a bias toward selection of sets with more items.

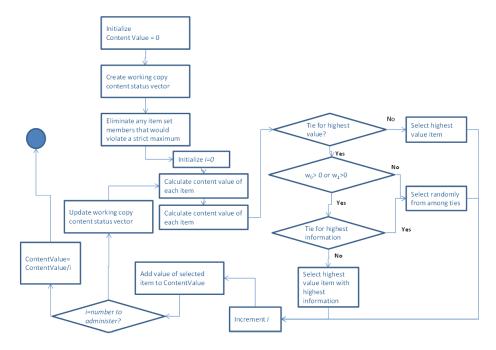


Exhibit A. Content Value Calculation for Item Sets

## 1.3 Information Value

Each item or item group also has value in terms of maximizing information, both overall and on reporting categories.

#### 1.3.1 Individual Information Value

The information value associated with an item will be an approximation of information. The system will be designed to use generalized Item Response Theory (IRT) models; however, it will treat all items as though they offer equal measurement precision. This is the assumption made by the Rasch model, but in more general models, items known to offer better measurement are given preference by many algorithms. Subsequent algorithms are then required to control the exposure of the items that measure best. Ignoring the differences in slopes serves to eliminate this bias and help equalize exposure.

# 1.3.2 Binary Items

The approximate information value of a binary item will be characterized as  $I_j(\theta) = p_j(\theta)(1 - p_j(\theta))$ , where the slope parameters are artificially replaced with a constant.

# 1.3.3 Polytomous Items

In terms of information, the best polytomous item in the pool is the one that maximizes the expected information,  $I_i(\theta)$ . Formally,  $I_i(\theta) > I_k(\theta)$  for all items  $k \neq j$ . The true value  $\theta$ ,

however, remains unknown and is accessed only through an estimate,  $\hat{\theta} \sim N(\bar{\theta}, \sigma_{\theta})$ . By definition of an expectation, the expected information  $I_i(\theta) = \int I_i(t) f(t|\bar{\theta}, \sigma_{\theta}) dt$ .

The intuition behind this result is illustrated in Exhibit B. In Exhibit B, each panel graphs the distribution of the estimate of  $\theta$  for an examinee. The top panel assumes a polytomous item in which one step threshold (A1) matches the mean of the  $\theta$  estimate distribution. In the bottom panel, neither step threshold matches the mean of the  $\theta$  estimate distribution. The shaded area in each panel indicates the region in which the hypothetical item depicted in the panel provides more information. We see that approximately 2/3 of the probability density function is shaded in the lower panel, while the item depicted in the upper panel dominates in only about 1/3 of the cases. In this example, the item depicted in the lower panel has a much greater probability of maximizing the information from the item, despite the fact that the item in the upper panel has a threshold exactly matching the mean of the estimate distribution and the item in the lower panel does not.

Exhibit B. Two Example Items, with the Shaded Region Showing the Probability that the Item Maximizes Information for the Examinee Depicted

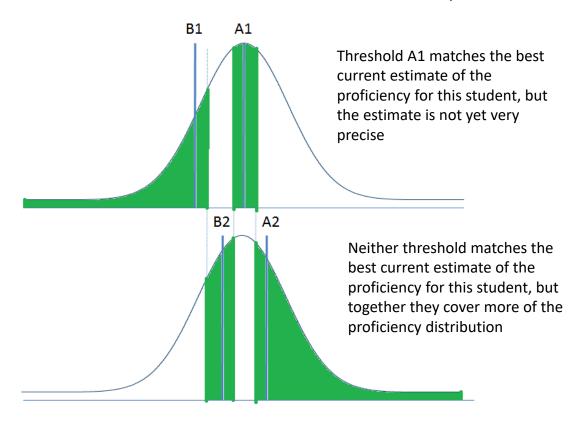
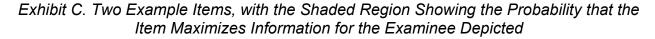
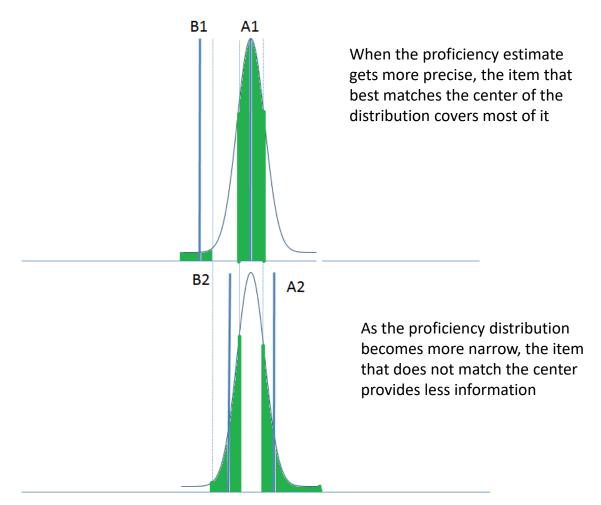


Exhibit C on the following page shows what happens to information as the estimate of this student's proficiency becomes more precise (later in the test). In this case, the item depicted in the top panel maximizes information about 65 to 70 percent of the time, compared to about 30 to 35 percent for the item depicted in the lower panel. These are the same items depicted in the Exhibit B, but in this case, we are considering information for a student with a more precise current proficiency estimate.





The approximate information value of polytomous items will be characterized as the expected information, specifically  $E[I_j(\theta)|m_i,s_i]=\int \sum_{k=1}^K I_{jk}(t)\,p_j(k|t)\phi(t;m_i,s_i)dt$ , where  $I_{jk}(t)$  represents the information at t of response k to item j,  $p_j(k|t)$  is the probability of response k to item j (artificially holding slope constant), given proficiency t,  $\phi(.)$  represents the normal probability density function, and  $m_i$  and  $s_i$  represent the mean and standard deviation of examinee i's current estimated proficiency distribution.

We propose to use Gauss-Hermite quadrature with a small number of quadrature points (approximately five). Experiments show that we can complete this calculation for 1,000 items in fewer than 5 milliseconds, making it computationally reasonable.

As with the binary items, we propose to ignore the slope parameters to even exposure and avoid a bias toward the items with better measurement.

# 1.3.4 Item Group Information Value

Item groups differ from individual items in that a set of items will be selected for administration. Therefore, the goal is to maximize information across the working theta distribution. As with the polytomous items, we propose to use Gauss-Hermite quadrature to estimate the expected information of the item group.

In the case of multiple-item groups

$$E[I_g(\theta)|m_i, s_i] = \frac{1}{J_g} \int \sum_{i=1}^{J_g} I_{g(i)}(t) \, \phi(t; m_i, s_i) dt$$

Where  $I_g(.)$  is the information from item group g,  $I_{g(j)}$  is the information associated with item  $j \in g$ , for the  $J_g$  items in set g. In the case of polytomous items, we use the expected information, as described above.

### 2. Entry and Initialization

At startup, the system will

- create a custom item pool;
- initialize theta estimates for the overall score and each score point; and
- insert embedded field-test items.

#### 2.1 ITEM POOL

At test startup, the system will generate a *custom item pool*, a string of item IDs for which the student is eligible. This item pool will include all items that

- are active in the system at test startup; and
- are not flagged as "access limited" for attributes associated with this student.

The list will be stored in ascending order of ID.

#### 2.2 ADJUST SEGMENT LENGTH

Custom item pools run the risk of being unable to meet segment blueprint minimums. To address this special case, the algorithm will adjust the blueprint to be consistent with the custom item pool. This capability becomes necessary when an accommodated item pool systematically excludes some content.

Let

S be the set of top-level content constraints in the hierarchical set of constraints, each consisting of the tuple (name, min, max, n);

**C** be the custom item pool, each element consisting of a set of content constraints **B**;

f, p integers represent item shortfall and pool count, respectively; and

**t** be the minimum required items on the segment.

For each s in S, compute n as the sum of active operational items in C classified on the constraint.

```
f = summation over S (min - n)

p = summation over S (n)

if t - f < p, then t = t - f
```

#### 2.3 Initialization of Starting Theta Estimates

The user will supply five pieces of information in the test configuration:

1. A default starting value if no other information is available

- 2. An indication whether prior scores on the same test should be used, if available
- 3. Optionally, the test ID of another test that can supply a starting value, along with
- 4. Slope and intercept parameters to adjust the scale of the value to transform it to the scale of the target test
- 5. A constant prior variance for use in calculation of working EAP scores

#### 2.4 INSERTION OF EMBEDDED FIELD-TEST ITEMS

Each blueprint will specify

- the number of field-test items to be administered on each test;
- the first item position into which a field-test item may be inserted; and
- the last item position into which a field-test item may be inserted.

Upon startup, select randomly from among the field-test items or item sets until the system has selected the specified number of field-test items. If the items are in sets, the sets will be administered as a complete set, and this may lead to more than the specified number of items administered.

The probability of selection will be given by  $p_j = \frac{\sum_{j=1}^K K_j}{\sum_{j=1}^K a_j K_j} a_j K_j \frac{m}{N_j}$ , where

 $p_j$  represents the probability of selecting the item;

*m* is the targeted number of field-test items;

 $N_i$  is the total number of active items in the field-test pool;

 $K_i$  is the number of items in item set j; and

 $a_j$  is a user-supplied weight associated with each item (or item set) to adjust the relative probability of selection.

The  $a_j$  variables are included to allow for operational cases in which some items must complete field testing sooner or enter field testing later. While using this parameter presents some statistical risk, not doing so poses operational risks.

For each item set, generate a uniform random number  $r_j$  on the interval  $\{0,1\}$ . Sort the items in ascending order by  $\frac{r_j}{p_j}$ . Sequentially select items, summing the number of items in the set. Stop the selection of field-test items once  $FTNMin \leq m \leq FTNMax = \sum_{j=0}^{\infty} K_j$ .

Next, each item is assigned to a position on the test. To do so, select a starting position within f - FTMax - FTMin positions from FTMin, where FTMax is the maximum allowable position for field-test items and FTMin is the minimum allowable position for field-test items. FTNMin and FTNMax refer to the minimum and maximum number of field-test items, respectively. Distribute the items evenly within these positions.

#### 3. ITEM SELECTION

Exhibit D summarizes the item selection process. If the item position has been designated for a field-test item, administer that item. Otherwise, the adaptive algorithm kicks in.

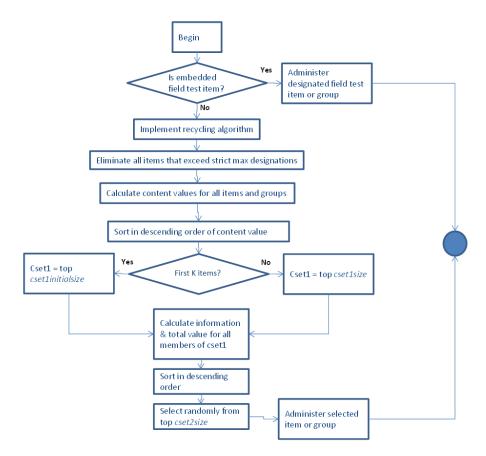


Exhibit D. Summary of Item Selection Process

This approach is a "content first" approach designed to optimize match to blueprint. An alternative, "information first" approach, is possible. Under an information first approach, all items within a specified information range would be selected as the first set of candidates, and subsequent selection within that set would be based, in part, on content considerations. The engine is being designed so that future development could build such an algorithm using many of the calculations already available.

#### 3.1 TRIMMING THE CUSTOM ITEM POOL

At each item selection, the active item pool is modified in four steps:

- 1. The custom item pool is intersected with the active item pool, resulting in a custom active item pool.
- 2. Items already administered on this test are removed from the custom active item pool.

- 3. Items that have been administered on prior tests are tentatively removed (see Section 3.2, Recycling Algorithm).
- 4. Items that measure content that has already exceeded a strict maximum are tentatively removed from the pool, removing entire sets containing items that meet this criterion.

## 3.2 RECYCLING ALGORITHM

When students are offered multiple opportunities to test, or when prior tests have been started and invalidated, students will have seen some of the items in the pool. The trimming of the item pool eliminates these items from the pool. It is possible that in such situations, the pool may no longer contain enough items to meet the blueprint.

Hence, items that have been seen on previous administrations may be returned to the pool. If there are not enough items remaining in the pool, the algorithm will recycle items (or item groups) with the required characteristic that is found in insufficient numbers. Working from the least recently administered group, items (or item groups) are reintroduced into the pool until the number of items with the required characteristics meets the minimum requirement. When item groups are recycled, the entire group is recycled rather than an individual item. Items administered on the current test are never recycled.

## 3.3 ADAPTIVE ITEM SELECTION

Selection of items will follow a common logic, whether the selection is for a single item or an item group. Item selection will proceed in the following three steps:

- 1. Select Candidate Set 1 (cset1).
  - a. Calculate the content value of each item or item group.
  - b. Sort the item groups in descending order of content value.
  - c. Select the top *cset1size*, a user-supplied value that may vary by test.
- 2. Select Candidate Set 2 (*cset2*).
  - a. Calculate the information values for each item group in *cset1*.
  - b. Calculate the overall value of each item group in *cset1* as defined in Equation (1).
  - c. Sort *cset2* in descending order of value.
  - d. Select the top *cset2size* item groups, where *cset2size* is a user-supplied value that may vary by test.
- 3. Select the item or item group to be administered.
  - a. Select randomly from *cset2* with uniform probability.

Note that a "pure adaptive" test, without regard to content constraints, can be achieved by setting cset I size to the size of the item pool and  $w_2$ , the weight associated with meeting content constraints

in Equation (1), to zero. Similarly, linear on-the-fly tests can be constructed by setting  $w_0$  and  $w_1$  to zero.

#### 3.4 SELECTION OF THE INITIAL ITEM

Selection of the initial item can affect item exposure. At the start of the test, all tests have no content already administered, so the items and item groups have the same content value for all examinees. In general, it is a good idea to spread the initial item selection over a wider range of content values. Therefore, we define an additional user-settable value, *cset1initialsize*, which is the size of Candidate Set 1 on the first *K* items only, where *K* is the number of reporting categories. Similarly, we define *cset2initialisize*.

#### 3.5 EXPOSURE CONTROL

This algorithm uses randomization to control exposure and offers several parameters that can be adjusted to control the tradeoff between optimal item allocation and exposure control. The primary mechanism for controlling exposure is the random selection from CSET2, the set of items or item groups that best meet the content and information criteria. These represent the "top k" items, where k can be set. Larger values of k provide more exposure control at the expense of optional selection.

In addition to this mechanism, we avoid a bias toward items with higher measurement precision by treating all items as though they measured with equal precision by ignoring variation in the slope parameter. This has the effect of randomizing over items with differing slope parameters. Without this step, it would be necessary to have other *post hoc* explicit controls to avoid the overexposure of items with higher slope parameters, an approach that could lead to different test characteristics over the course of the testing window.

#### 4. TERMINATION

The algorithm will have configurable termination conditions. These may include

- administering a minimum number of items in each reporting category and overall;
- achieving a target level of precision on the overall test score;
- achieving a target level of precision on all reporting categories; and
- achieving a score insufficiently distant from a specified score with sufficient precision (e.g., less than two standard errors below proficient). Cambium Assessment, Inc (CAI) envisions this being used in conjunction with other termination conditions to allow very high or very low achieving students to continue on to a segment that contains items from adjacent grades but barring other students from those segments.

We will define four user-defined flags indicating whether each of these is to be considered in the termination conditions (*TermCount*, *TermOverall*, *TermReporting*, *TermTooClose*). A fifth user-supplied value will indicate whether these are taken in conjunction or if satisfaction of any one of them will suffice (*TermAnd*). Reaching the minimum number of items is always a necessary condition for termination.

In addition, two conditions will each individually and independently cause termination of the test:

- 1. Administering the maximum number of items specified in the blueprint
- 2. Having no items in the pool left to administer

# **APPENDIX 1. DEFINITIONS OF USER-SETTABLE PARAMETERS**

This appendix summarizes the user-settable parameters in the adaptive algorithm.

Parameter Name	Description	Entity Referred to by Subscript Index
$w_0$	Priority weight associated with match to blueprint	N/A
$w_1$	Priority weight associated with reporting category information	
$w_2$	Priority weight associated with overall information	N/A
$q_k$	Priority weight associated with a specific reporting category	reporting categories
$p_r$	Priority weight associated with a feature specified in the blueprint (These inputs appear as a component of the blueprint.)	features specified in the blueprint
а	Parameter of the function $h(.)$ that controls the overall information weight when the information target has not yet been hit	N/A
b	Parameter of the function $h(.)$ that controls the overall information weight after the information target has been hit	N/A
$c_k$	Parameter of the function $h(.)$ that controls the information weight when the information target has not yet been hit for reporting category $k$	reporting categories
$d_k$	Parameter of the function $h(.)$ that controls the information weight after the information target has been hit for reporting category $k$	
cset1size	Size of candidate pool based on contribution to blueprint match	N/A
cset1initialsize	set1initialsize Size of candidate pool based on contribution to blueprint match for the first <i>K</i> items or item sets selected	
cset2size	Size of final candidate pool from which to select randomly	N/A
cset2initialsize	Size of candidate pool based on contribution to blueprint match and information for the first item or item set selected	
$t_0$	Target information for the overall test	N/A
$t_k$	Target information for reporting categories	reporting categories
startTheta	A default starting value if no other information is available	N/A
startPrevious	startPrevious An indication of whether previous scores on the same test should be used, if available	
startOther	The test ID of another test that can supply a starting value, along with startOtherSlope	N/A
startOtherSlope	startOtherSlope Slope parameter to adjust the scale of the value to transform it to the scale of the target test	

Parameter Name	Description	Entity Referred to by Subscript Index
startOtherInt	Intercept parameter to adjust the scale of the value to transform it to the scale of the target test	N/A
FTMin	Minimum position in which field-test items are allowed	N/A
FTMax	Maximum position in which field-test items are allowed	N/A
FTNMin	Target minimum number of field-test items	N/A
FTNMax	Target maximum number of field-test items	N/A
$a_j$	Weight adjustment for individual embedded field-test items used to increase or decrease their probability of selection	field-test items
AdaptiveCut	The overall score cutscore, usually proficiency, used in consideration of <i>TermTooClose</i>	
TooCloseSEs	The number of standard errors below which the difference is considered "too close" to the adaptive cut to proceed. In general, this will signal proceeding to a final segment that contains off-grade items.	
TermOverall	Flag indicating whether to use the overall information target as a termination criterion	N/A
TermReporting	Flag to indicate whether to use reporting category information target as a termination criterion	N/A
TermCount	Flag to indicate whether to use minimum test size as a termination condition	N/A
TermTooClose	Terminate if you are not sufficiently distant from the specified adaptive cut	
TermAnd	Flag to indicate whether the other termination conditions are to be taken separately or conjunctively	N/A

## **APPENDIX 2. SUPPORTING DATA STRUCTURES**

## Cambium Assessment, Inc (CAI) Cautions and Caveats

- Use of standard error termination conditions will likely cause inconsistencies between the blueprint content specifications, and the information criteria will cause unpredictable results, likely leading to failures to meet blueprint requirements.
- The field-test positioning algorithm outlined here is very simple and will lead to deterministic placement of field-test items.

## ADDENDUM. ADJUSTMENTS TO THE USE OF ITEM CLUSTERS

Cambium Assessment, Inc (CAI) adjusted the adaptive algorithm to the use of item clusters as follows:

- Using marginal maximum likelihood estimator (MMLE) to update proficiency estimates, marginalizing out cluster effects.
- Normalizing the information by the number of assertions within an item, to avoid overselection of item clusters and stand-alone items with more assertions.