



Bureau of Indian Education 2024–25 Technical Report

Grades 5, 8, and 11 Science

**Prepared by Cognia and the
Bureau of Indian Education**

Table Of Contents

Chapter 1. Introduction to the Assessment Program	5
1.1 Purposes and Uses of the Assessment Program.....	5
1.2 Statements of Intended Score Interpretations and Uses (SIUs)	5
1.2.1 Primary Intended BIE Science Assessment Score Interpretations and Uses	6
1.3 Introduction to Validity Arguments for the Program: Rationales for the Approach.....	7
Chapter 2. Overview of the Program.....	8
2.1 History of the Program.....	8
2.2 Stakeholder Involvement	8
2.2.1 Educator Committees	8
2.3 Student Participation	8
Chapter 3. Test Content.....	10
3.1 Content Standards.....	10
3.1.1 Eligible Standards	10
3.2 Assessment Design.....	10
3.2.1 Content Coverage Blueprint.....	12
3.2.2 Operational Section	12
3.2.3 Field-Test Sections	13
3.2.4 Item Types	13
3.2.5 Stimulus Types.....	14
Chapter 4. Test Development	15
4.1 Overview of General Approach	15
4.2 Item Specifications	15
4.3 Item Review Committees and Processes.....	16
4.3.1 Content and Item Reviews.....	16
4.3.2 Bias and Sensitivity Review	16
4.4 Test Forms Construction	17
4.4.1 Item and Stimulus Selection	17
4.4.2 Selection Specifications to Meet Blueprint Requirements	18
4.4.3 Cognitive Processes	19
Chapter 5. Test Administration.....	20
5.1 Roles and Responsibilities	20
5.2 Test Administrator's Manual.....	20
5.3 TA and Proctor Training Requirements and 2025 Test Administrations	20
5.4 Testing Irregularity Reports	20
5.5 Test Security.....	21
Chapter 6. Scoring: Scope of Work, Processes, and Procedures	22

6.1 Scope of Work	22
6.2 Operational Scoring: Processes and Procedures	22
6.2.1 Score Verification of Multiple-Choice Items	22
6.2.2 Benchmarking Field-Test Items	22
6.2.3 Scoring of Open-Ended Response Items	22
Chapter 7. Classical Item and Test Analysis.....	28
7.1 Classical Item Statistics	28
7.2 Total Test and Subscore Intercorrelations	29
Chapter 8. Psychometrics: Item Response Theory (IRT) Scaling and Equating.....	31
8.1 IRT Models	31
8.2 Dimensionality Analysis.....	32
8.3 Item Response Theory Results	34
8.4 Equating.....	35
8.5 Reported Total Test and Subtest Scale Scores	35
8.6 Performance Levels.....	36
8.6.1 Percentages of Students in Each Performance Level	36
8.7 Differential Item Functioning.....	37
Chapter 9. Score Reliability.....	39
9.1 IRT Marginal Reliabilities	39
9.2 Decision Accuracy and Consistency	41
Chapter 10. Score Reporting.....	44
10.1 Individual Student Reports	44
10.2 Scale Score	44
10.3 Achievement Level	44
10.4 Standards Performance Indicators	45
10.5 Comparison of Student Performance at the School, District, and Overall Level.....	45
10.6 Student Roster.....	45
Chapter 11. Validity Arguments to Support Intended Score Interpretation and Uses.....	47
11.1 Primary Intended Score Interpretation	50
11.2 Primary Intended Score Uses.....	54
11.2.1 Intended Score Use for Individual Students	54
11.2.2 Intended Score Use for Groups of Students.....	54
11.3 Conclusions and Next Steps	56
11.3.1 Research Agenda	57
References	58
Appendices	60

Appendix A	List of Acronyms
Appendix B	Performance Level Descriptors
Appendix C	Accommodation Frequencies
Appendix D	IRC and Bias Review Meeting Participants
Appendix E	Scorer Qualification Rates
Appendix F	Classical Item Statistics
Appendix G	Item Response Theory Parameters
Appendix H	Test Characteristic Curves & Test Information Functions
Appendix I	Raw to Scaled Score Lookup Tables
Appendix J	Decision Accuracy and Consistency Results
Appendix K	BIE Science Reporting Business Requirements
Appendix L	Score Report Interpretation Quick Guide
Appendix M	Cumulative Scaled Score Distributions
Appendix N	Scaled Score Descriptive Statistics

List of Tables and Figures

Table 2-1. Participation Number by Grade, as a Function of Demographic Variable.....	9
Table 2-2. Number of Participating Students, as a Function of Grade	9
Table 3-1. Student Testing Experience	11
Table 3-2. Grades 5, 8, 11 BIE Science Assessment Operational Test Blueprint.....	12
Table 3-3. Distribution of Emphasis Across Content Standards in Terms of Percentage of Total Test Points by Grade—Science Grades 5, 8, 11—Spring 2025	12
Table 3-4. Distribution of Raw-Score Points Across Reporting Categories by Grade—Science Grades 5, 8, 11	13
Table 3-5. Science Item Types	14
Table 4-1. Criteria for Flagged Items Based on Classical Test Theory (CTT) Statistics	19
Table 4-2. Criteria for Flagged Items Based on Item Response Theory (IRT) Statistics	19
Table 6-1. Overview of BIE Science Assessment Scope-of-Work	22
Figure 6-1. Cognia Scoring Staff.....	23
Table 6-2. Summary of Interrater Reliability Statistics for BIE Science—OP Items by Grade	27
Table 6-3. Score of Record Rules.....	27
Table 7-1. Item Summary by Grade.....	29
Table 7-2. Pearson Correlations of Total Test and Subtest Raw Scores on BIE Science Grade 5, as a Function of Operational Form	29
Table 7-3. Pearson Correlations of Total Test and Subtest Raw Scores on BIE Science Grade 8, as a Function of Operational Form	29
Table 7-4. Pearson Correlations of Total Test and Subtest Raw Scores on BIE Science Grade 11, as a Function of Operational Form	30
Table 8-1. DIMTEST Hypothesis Testing Results, as a Function of Grade*	34
Table 8-2. Spring 2025 Scaled Score Slopes and Intercepts by Subject and Grade	36
Table 8-3. Spring 2025 Cutpoints on the Theta Metric and Reporting Scale by Grade	36
Table 8-4. Performance Level Distribution as a Function of Grade*	36
Table 8-5. Number of Items with Low DIF	38
Table 9-1. IRT Marginal Reliability by Grade and Form	40
Table 9-2. IRT Marginal Reliability by Grade and Subgroup	40
Figure 9-1a. Overall Decision Accuracy for Science by Grade for Form 1 (Set A).....	42
Figure 9-1b. Overall Decision Accuracy for Science by Grade for Form 2 (Set B).....	43
Figure 11-1. BIE Validity Argument Model.....	48
Table 11-1. Relationships Among Score Interpretations and Uses, Claims, and Sub-Claims, and Supporting Evidence.....	49
Table 11-2. Relevance and Completeness or Completeness of Evidence in Support of SIUs and Claims Underlying Validity Arguments for BIE Score Interpretations and Uses	50
Table 11-3. Status of Evidence for All SIUs, Claims, and Subclaims.....	56

Chapter 1. Introduction to the Assessment Program

1.1 Purposes and Uses of the Assessment Program

The Bureau of Indian Education (BIE—see Appendix A for a list of acronyms) Science Assessment is a summative assessment for science, administered to students in grades 5, 8, and 11. It is designed to provide evidence to determine a student's grade-level proficiency and progress toward college and/or career readiness, as defined by the BIE, by showing BIE students have mastered the Next Generation Science Standards. The BIE Science Assessment is a key component of BIE's Every Student Succeeds Act (ESSA) plan to meet ESSA's general assessment requirements.

As the BIE Science Assessment is an end-of-grade/end-of-course single measure, interpretations and uses of its scores should be supplemented with additional measures, including information from classroom summative, interim, and formative assessments in science. In keeping with the practices outlined in *Standards for Educational and Psychological Testing*, each student's score should be used as part of a body of evidence regarding mastery and should not be used in isolation to make high-stakes decisions (AERA, APA, & NCME, 2014). Hence, the aggregation of student scores on the BIE Science Assessment at the school or district levels, or for the entire BIE, is generally a more reliable indicator of program success, particularly when monitored over several years.

The BIE Science Assessment is designed to provide point-in-time information about the academic achievement and progress of students. Student results are reported according to academic achievement descriptors utilizing scale scores for each of the four performance levels: Novice, Nearing Proficiency, Proficient, and Advanced. The results from these assessments provide educators and the public with information to guide the creation of future educational practices to meet the needs of students, while monitoring the continuous improvement efforts of schools, districts, and the BIE in achieving a world-class education system for all students.

1.2 Statements of Intended Score Interpretations and Uses (SIUs)

The phrase “intended score interpretations for uses” appears several times in *Standards for Educational and Psychological Testing* and is the core of the field's views on validity and validation. For the BIE Science Assessment and other assessment programs, the phrase refers broadly to test scores (e.g., total test scale scores, aggregations of test scores, the percentages of students at or above standard), and other test performance information elements, such as the definition of “at or above standard” in the performance level descriptors (PLDs—see Appendix B).

1.2.1 Primary Intended BIE Science Assessment Score Interpretations and Uses

Primary Score Interpretations:

- For Elementary and Middle Schools, performance on the BIE Science Assessment indicates student mastery of grade levels 3–5 and 6–8 expectations for integration of Science and Engineering Practices (SEPs), Disciplinary Core Ideas (DCIs), and Crosscutting Concepts (CCCs) as presented in the standards, which is the progression for the next level of science curriculum and is a predictor of being on track for college and career readiness.
- For High School, the BIE Science Assessment is designed to measure whether students are on track to be ready for college or career, as defined by the BIE, by showing they have mastered the Next Generation Science Standards, which require the integration of SEPs, DCIs, and CCCs to explain phenomena and solve problems.

Primary Score Uses:

Educators, administrators, and other stakeholders at the school level can use the Science Assessment and its results to (a) monitor trends in student performance, (b) design professional development for teachers, and (c) drive accountability results.

- Teachers can use the BIE Science Assessment and its results to better integrate assessment with their instructional planning.
- Parents can use the BIE Science Assessment and its results to get information about (a) what their child knows and can do, and (b) their child's progress over time.

The intended score interpretation and uses stated here align with the original statements of intended score interpretations and uses in the National Center and State Collaborative 2015 Operational Assessment Technical Manual.

The BIE Science Assessment is designed, developed, and implemented to support three intended SIUs, according to the broad interpretation of the phrase above. These interpretations and uses are applicable to assessments in general and to specific applications with individual students and groups of students, as described below.

SIU 1: Intended Score Interpretation

The BIE Science Assessment provides reliable and valid information about important knowledge and skills in grade-level science attained by general education students.

- Claim 1.1: The content of the tests represents the content of the standards.
- Claim 1.2: The test items are construct-relevant.
- Claim 1.3: Test scores on the BIE Science Assessments provide reliable information about student performance and accurate classifications into performance levels.
- Claim 1.4: Item and test scoring are implemented accurately; approved scoring rules are implemented accurately.

SIU 2: Intended Score Use for Individual Students

Scale scores can be used to compare an individual student's performance to the performance of other students in BIE.

- Claim 2.1: Educators and other stakeholders at the school level can use results from the BIE Science Assessment to describe and monitor student achievement status with respect to mastery of the content standards.

SIU 3: Intended Score Use for Groups of Students

SIU statements for groups of students are applicable to aggregate reporting of student subgroups (e.g., English learners, students with disabilities, racial/ethnic subgroups) within those levels of aggregation.

- Claim 3.1: Educators can use results from the BIE Science Assessment to support instructional planning for groups of students.
- Claim 3.2: School and BIE stakeholders can use results from the BIE Science Assessment to make comparisons between organizations.

Claims, subclaims, and evidence that support the intended interpretations and uses of BIE Science Assessment scores are provided in Chapter 11.

1.3 Introduction to Validity Arguments for the Program: Rationales for the Approach

This report documents test development procedures and psychometric outcomes for the 2025 BIE Science Assessment. These technical aspects of the program contribute to the accumulation of validity evidence to support BIE Science Assessment score interpretations and uses. Because the interpretations of test scores, not the test itself, are evaluated for validity, this report presents documentation to substantiate intended interpretations (AERA et al., 2014). Subsequent chapters of this report discuss test development, test alignment, test administration, scoring, equating, item analyses, reliability, scaled scores, performance levels, and reporting. Each of these topics contributes important information toward establishing the validity of the assessment program. Note, however, that this report does not include certain aspects of a comprehensive validity argument that could be important to consider when drawing conclusions about validity. For example, additional sources of validity evidence might speak to the extent to which BIE Science Assessment scores converge with other measures of the same or similar constructs and diverge from measures of different constructs and consequences that arise from scores at the student, school, and district levels, as well as from scores for BIE as a whole.

The *Standards for Educational and Psychological Testing* (AERA et al., 2014) also gives a framework for describing sources of evidence that should be considered when constructing a validity argument. These sources include evidence based on the following five areas: test content, response processes, internal structure, relationship to other variables, and consequences of testing. These sources address different aspects of supporting evidence for validity arguments; they are not distinct types of validity. Instead, each contributes to a body of evidence about the comprehensive validity of score interpretations and uses and the intended interpretations and uses. Moreover, these sources represent only a partial list of sources of evidence from the BIE Science Assessment design, development, test administration, analysis, and reporting processes that are relevant to the overall validity arguments for intended interpretations and uses of BIE Science Assessment scores and other information.

Chapter 2. Overview of the Program

2.1 History of the Program

The BIE’s mission is to provide quality education opportunities from early childhood through life in accordance with a tribe’s needs for cultural and economic well-being, in keeping with the wide diversity of Indian tribes as distinct cultural and governmental entities. To better support this mission, the BIE began the process in 2020 of administering a unique science assessment. BIE Science was scheduled to have its first administration in Spring 2021, and due to unforeseen constraints, used the New Mexico Assessment of Science Readiness assessments. Beginning in 2022, the BIE administered a unique science assessment separate from any other state, the BIE Science Assessment. After its first administration in Spring 2022, BIE Science had its second administration in Spring 2023. The third administration took place in Spring 2024, and the fourth in Spring 2025. BIE Science continues to work with Cognia to develop test materials that are culturally relevant. BIE schools have the option to test both online and on paper, depending on their specific technological access and needs.

The window for the 2025 BIE Science administration with Cognia was March 10–April 18, 2025.

2.2 Stakeholder Involvement

Cognia and the BIE work together on all aspects of the implementation of the BIE Science Assessment program. The BIE also works with several stakeholder groups for input into the implementation of the program.

2.2.1 Educator Committees

Educator committees are periodically convened for the purpose of content development. The committees include those listed below, with the details of each committee found in Chapter 4.

- **National Item Review Committee:** Cognia convenes a national item review committee to review the content of the items that are created. BIE educators comprise two seats per grade/content span for those committees.
- **National Bias Review Committee:** Cognia convenes a national bias committee to look for bias and sensitivity concerns in the content that is created. BIE educators comprise two seats on that committee.
- **BIE-specific Review Committee:** Cognia facilitates a committee comprised only of BIE educators to review content and to look for bias and sensitivity concerns in custom BIE items that are created. These items focus on culturally relevant phenomena for BIE students.

2.3 Student Participation

BIE policy defines student participation on a BIE Science Assessment as attempting five or more items on the given assessment. Table 2-1 provides participation rates, as a function of assessment subject (science), and background/demographic variable. The number of students participating in the BIE Science Assessment in Spring 2025 per grade ranged from approximately 900 to 1,900.

Table 2-1. Participation Number by Grade, as a Function of Demographic Variable

Type	Grade		
	5	8	11
All	1,928	1,677	1,095
Female	961	825	542
Male	967	852	553
Gender Undefined	0	0	0
American Indian/Alaska Native	1,906	1,662	1,065
Non-American Indian/Alaska Native	21	13	15
Asian	6	6	5
Black/African American	0	0	0
Multi-racial	10	4	4
Caucasian/White	5	3	2
Pacific Islander/Native Hawaiian	0	0	4
Currently receiving LEP services	321	202	108
Not receiving LEP services	394	405	226
Special Ed	351	299	152
Non-special Ed	1,573	1,374	925
Economically Disadvantaged Students	1,840	1,625	1,052
Non-economically Disadvantaged Students	21	28	8
Gifted Students	81	116	63
Non-gifted Students	103	109	44
Title1 Students	1,870	1,623	1,042
Non-title1 Students	15	13	6

The BIE Science Assessments were administered in either computer-based or paper-based formats. Most students utilized computer-based administration. Table 2-2 contains the number of students utilizing computer-based or paper-based administrations. Additionally, the usage statistics on accommodation(s) and accessibility feature(s), as a function of grade are available in Appendix C. Only students who met the attemptedness rule (i.e., attempted 5 or more items) contributed to the frequencies in the aforementioned tables. Of the students that participated in the Spring 2025 administration, the table below indicates the numbers of students who were assessed in each mode.

Table 2-2. Number of Participating Students, as a Function of Grade

Grade	Computer-Based	Paper-Based
5	1,905	23
8	1,655	22
11	1,062	33

The small numbers of students per grade in Spring 2025 (the total number of BIE students who participated in the 2025 science assessments was less than 2,000 per grade) impacted the analyses reported in this document in a variety of ways, such as increased levels of sampling, measurement, and/or estimation error in test reliabilities (overall and subgroup), decision accuracies and consistencies, classical item statistics, dimensionality effect sizes, scaled score distributions, and performance level distributions.

Chapter 3. Test Content

3.1 Content Standards

Test content, including items and stimuli, for the BIE Science Assessment was developed according to the *Next Generation Science Standards*. These standards are the basis for the test designs developed for each grade and are used to inform the development of items. Each item is designed to measure a specific standard, or performance expectation, and align to multiple dimensions of the standard (Disciplinary Core Idea, Science and Engineering Practice, Crosscutting Concept).

The specific content standards were subsequently grouped into categories for the purpose of communicating with students, families, and educators. The content standards that are eligible to be included in the BIE Science Assessment are described in the following sections.

3.1.1 Eligible Standards

The BIE Science Assessment assesses the Next Generation Science Standards at grades 5, 8, and 11. All standards are eligible for assessment. However, because of the number of standards per grade, not all standards can be tested every year. The design of the BIE Science Assessment allows for all assessable standards to be included in the assessment at least once within a three-to-five-year period.

3.2 Assessment Design

Table 3-1 provides a summary of the number of items and points by item type, usage (i.e., operational items or field-test items), and estimated testing time for each grade level and content area of the BIE Science Assessment. The test is administered in three sessions. Test forms contain core operational items, matrix operational items, and matrix field-test items. There are two sets of operational items, set A and set B, differing in the standalone MS-2 items that are in the set (but still following the same content blueprint), to support sufficient assessment of all content standards over time. All operational items count toward student score, with the core operational items being common across both forms and the matrix operational items being administered across the two different operational forms¹. Matrix field-test items are items administered to subsets of students to “try out” performance (with different students receiving different field-test items) and therefore do not count toward student score. The types of items on the BIE Science Assessment are 1-point machine-scored items (MS-1), 2-point machine-scored items (MS-2), and 4-point constructed-response/open-ended items (OE-4). Additional item-type descriptions can be found in section 3.2.4.

¹ In 2024, the BIE Grade 5 Science operational Forms A and B were identical, constructed with identical clusters and OE4 items, and the same MS2 items so all items were core operational. All students received the correct lookup tables for scoring and reporting. Standards coverage was consistent with the intended blueprint design, as shown in the tables in this section.

Table 3-1. Student Testing Experience

Grades 5, 8	Cluster/Passage Items			Standalone Items		Total Items	Total Points
	Stimulus/Passage	MS-1	MS-2	MS-2	OE		
Core Operational Items	6	12	12	0	3	27	48
Matrix Operational Items	0	0	0	8	0	8	16
Matrix Field-Test Items	2	4	4	4	1	13	24
Total Student Experience	8	16	16	12	4	48	88
Estimated Testing Time (min)							150 (50/50/50)
Grade 11	Cluster/Passage Items			Standalone Items		Total Items	Total Points
	Stimulus/Passage	MS-1	MS-2	MS-2	OE		
Core Operational Items	6	12	12	0	3	27	48
Matrix Operational Items	0	0	0	10	0	10	20
Matrix Field-Test Items	2	4	4	5	1	14	26
Total Student Experience	8	16	16	15	4	51	94
Estimated Testing Time (min)							165 (55/55/55)

BIE Science Assessment Specifications

The reporting categories for the BIE Science Assessment are based on the science domains in the Next Generation Science Standards. Target percentages for the distribution of operational test points for each of the reporting categories reflect the distribution in the standards, so as not to over- or under-represent content. Because of their small percentage of standards overall, items aligned to standards in Engineering, Technology, and Applications of Science are reported under the reporting category domain that matches the context of the design problem presented. Percentages for the three reported domains are shown in the tables in the next section.

Specifications for the full test blueprints for the construction of the operational forms reflect the reporting category specifications. These constructs represent key aspects of the standards to which items are aligned; as such, the percentage of operational test points for each should be maintained from year to year. Note that some of the points for each reporting category come from clusters (a grouping of four items: 2 MS-1 and 2 MS-2—all associated with a common stimulus), and some points come from standalone/discrete items.

Form 1 is used to create print and accommodated forms. On the paper-based test (PBT) form, any technology-enhanced item is replaced with an equivalent multiple-choice or multi-select item for administration. All other items on the PBT form remain the same as on the online form.

Table 3-2. Grades 5, 8, 11 BIE Science Assessment Operational Test Blueprint

Reporting Category	Grade 5				
	Ideal Number of Clusters	Ideal Number of Standalone MS-2	Ideal Number of Standalone OE	Ideal Number of Core Points	Ideal Percent of Core Points (+/- 4%)
Practices and Crosscutting Concepts in Physical Sciences	2	4–6	1	24–28	40%
Practices and Crosscutting Concepts in Life Sciences	2	1–3	1	18–22	30%
Practices and Crosscutting Concepts in Earth and Space Sciences	2	1–3	1	18–22	30%
	Grade 8				
	Ideal Number of Clusters	Ideal Number of Standalone MS-2	Ideal Number of Standalone OE	Ideal Number of Core Points	Ideal Percent of Core Points (+/- 4%)
Practices and Crosscutting Concepts in Physical Sciences	2	2–4	1	20–24	35%
Practices and Crosscutting Concepts in Life Sciences	2	2–4	1	20–24	35%
Practices and Crosscutting Concepts in Earth and Space Sciences	2	1–3	1	18–22	30%
	Grade 11				
	Ideal Number of Clusters	Ideal Number of Standalone MS-2	Ideal Number of Standalone OE	Ideal Number of Core Points	Ideal Percent of Core Points (+/- 4%)
Practices and Crosscutting Concepts in Physical Sciences	2	3–5	1	22–26	35%
Practices and Crosscutting Concepts in Life Sciences	2	3–5	1	22–26	35%
Practices and Crosscutting Concepts in Earth and Space Sciences	2	1–3	1	18–22	30%

3.2.1 Content Coverage Blueprint

The distribution of emphasis for BIE content standards in science for the Spring 2025 assessment is shown in Table 3-3. Assessable standards cover Physical Sciences, Life Sciences, Earth and Space Sciences, and Engineering, Technology, and Applications of Science (ETS).

Table 3-3. Distribution of Emphasis Across Content Standards in Terms of Percentage of Total Test Points by Grade—Science Grades 5, 8, 11—Spring 2025

Standards Category	Grade 5		Grade 8		Grade 11	
	Total Points	% of Total Core Points	Total Points	% of Total Core Points	Total Points	% of Total Core Points
Physical Sciences	21	32.81%	23	35.94%	16	23.53%
Life Sciences	20	31.25%	17	26.56%	24	35.29%
Earth and Space Sciences	17	26.56%	21	32.81%	22	32.35%
ETS	6	9.38%	3	4.69%	6	8.83%
Grand Total	64	100.00%	64	100.00%	68	100.00%

3.2.2 Operational Section

As mentioned at the start of section 3.2, there are two sets of operational items, set A and set B, differing in the standalone MS-2 items that are in the set (but still following the same content blueprint)². Table 3-4 shows the reporting categories in the BIE Science Assessment test design, and the maximum possible number of raw-score points students could earn in each reporting category on the Spring 2025 assessment for both operational forms. Note: Because only operational items are counted toward

² In 2024, the BIE Grade 5 Science operational Forms A and B were identical, constructed with identical clusters and OE4 items, and the same MS2 items. All students received the correct lookup tables for scoring and reporting. Standards coverage was consistent with the intended blueprint design, as shown in the tables in this section.

students' scaled scores, only operational items are reflected in this table. The number of items and item types that are used to achieve these distributions are provided in the tables at the beginning of section 3.2. Any items aligned to standards in ETS are reported under the reporting category domain that matches the context of the design problem presented.

Table 3-4. Distribution of Raw-Score Points Across Reporting Categories by Grade—Science Grades 5, 8, 11

Reporting Category	Grade 5		Grade 8		Grade 11	
	Total Points	% of Total Core Points	Total Points	% of Total Core Points	Total Points	% of Total Core Points
Practices and Crosscutting Concepts in Physical Sciences	24	37.50%	23	35.94%	22	32.35%
Practices and Crosscutting Concepts in Life Sciences	20	31.25%	20	31.25%	24	35.30%
Practices and Crosscutting Concepts in Earth and Space Sciences	20	31.25%	21	32.81%	22	32.35%
Grand Total	64	100.00%	64	100.00%	68	100.00%

3.2.3 Field-Test Sections

All items are appropriately field tested prior to operational use. Items for the BIE Science Assessment came primarily from a national item bank, and these items were field tested with BIE students and/or other students using the same item bank content. For the field testing done within the BIE Science Assessment forms, the test utilizes a matrix design that embeds the field-test items within each form. Matrix field-test items are items administered to subsets of students to “try out” performance (with different students receiving different field-test items) and therefore do not count toward student scores. In grades 5 and 8, the BIE Science Assessment contains a total of 13 field-test items per form: two clusters (with four items each), four MS-2 standalones, and one OE-4. The grade 11 test contains 14 field-test items per form: two clusters (with four items each), five MS-2 standalones, and one OE-4.

Note that starting in 2024, a BIE-specific cluster was developed to be culturally responsive to BIE students. A grade 8 cluster was field tested in Spring 2025. A grade 11 cluster will be field tested in Spring 2026. This process will be repeated with the development of a grade 5 cluster for field testing in Spring 2027. These items are not part of the national item bank.

3.2.4 Item Types

Item types are chosen to best balance the desires for making efficient use of limited testing time and providing coverage of a broad range of knowledge and skills. The item types used on the BIE Science Assessment and the functions of each are listed below.

The item types on the BIE Science Assessment include machine-scored 1-point items (MS-1), machine-scored 2-point items (MS-2), and open-ended items (OE-4). Some of the MS-1 and MS-2 items are grouped together in clusters.

MS-1 items may be multiple choice, multiple select, or technology enhanced (e.g., drag-and-drop, hot spot, drop-down selections). MS-1 items are only found in clusters. They are all machine-scored as correct only; partial credit is not awarded.

MS-2 items have two parts (Part a and Part b) for students to answer. These items may combine multiple choice, multiple select, and/or technology-enhanced interactions across the two parts. MS-2 items are included in clusters and as standalone items. They are all machine-scored, and students may earn 2, 1, or 0 points across Part a and Part b.

An item cluster is a set of items all associated with a common stimulus. Clusters contain four items, with two of the items being worth 1 point (MS-1) and two of the items being worth 2 points (MS-2). The clusters typically align to two PEs, and all clusters measure all three dimensions of the PEs being assessed.

Open-ended items (OE-4) are worth 4 points. These items require students to write an extended response to a prompt. The prompt may be a single prompt, or more typically, the items are written with multiple, scaffolded parts for students to respond to. These items are hand-scored, with scorers using a rubric and scoring notes to evaluate responses on a scale from 0–4.

Each type of item on the assessment is worth a specific number of points in the student’s total Science score, as shown in Table 3-5.

Table 3-5. Science Item Types

Item Type	Maximum Number of Points Available
MS-1	1
MS-2	2
OE-4	4

3.2.5 Stimulus Types

On the BIE Science Assessment, all clusters are written with an extended, rich stimulus that may be one or more paragraphs in length. The stimulus must present a single, rich science phenomenon or engineering design problem aligned to the standards/performance expectations being assessed. The phenomenon or problem must launch and support a single storyline, or sequence of sense-making, which is carried out in the items.

The stimulus may present any variety of elements to provide the necessary information related to the phenomenon or problem and the storyline: text paragraphs, graphs, data tables, models, drawings, etc. All information in the stimulus should be necessary, but not conceptually sufficient, for students to respond (i.e., students must also use their own knowledge of the constructs in the standards to answer the items, rather than simply identify given information), and the stimulus must provide enough information to allow students to engage in the SEPs, DCIs, and CCCs of the targeted standards as they respond to items.

Chapter 4. Test Development

4.1 Overview of General Approach

This chapter provides an overview of the development of the BIE Science Assessment, including test and item specifications, item reviews, and test construction.

According to the Standards for Educational and Psychological Testing (AERA, APA, & NCME, 2014), “important validity evidence can be obtained from an analysis of the relationship between a test’s content and the construct it is intended to measure” (p. 14). Accordingly, the descriptions of the test development procedures included in this chapter provide evidence that supports both the content and construct validity of the assessments.

4.2 Item Specifications

The test design for the BIE Science Assessment is based on the three content domains of Physical Sciences, Life Sciences, and Earth and Space Sciences. Items are expected to align to the multiple dimensions of the standards (Disciplinary Core Ideas, Science and Engineering Practices, Crosscutting Concepts) in each domain, such that every item is at least two-dimensional, if not three-dimensional. To emphasize this multi-dimensional nature of the items, the names of the reporting categories incorporate the three dimensions (Practices and Crosscutting Concepts in Physical Sciences, Practices and Crosscutting Concepts in Life Sciences, Practices and Crosscutting Concepts in Earth and Space Sciences). Students are expected to demonstrate sense-making by using core ideas, practices, and crosscutting concepts together to respond to items.

Items assessing Engineering, Technology, and Applications of Science are reported within the Physical, Life, or Earth and Space Sciences category, depending on the content match of the design problem presented in the item.

As content support, students taking the grade 11 test are provided with a Periodic Table reference sheet. No items on the assessment require a calculator or other mathematical tools to answer.

Cognitive Complexity

In addition to being created according to subject-area content standards, each item on the BIE Science Assessment is assigned information about its cognitive complexity.

Because the items on the BIE Science Assessment are NGSS-aligned, the cognitive complexity of the items is evaluated with a different framework than Depth of Knowledge (DOK). This framework, Cognitive Complexity Framework for SSIB, is based on *Achieve’s A Framework to Evaluate Cognitive Complexity in Science* (September 2019).

Under the Cognitive Complexity Framework for SSIB, four indicators are used to classify the cognitive complexity of each item: stimulus, science and engineering practice, disciplinary core idea, and crosscutting concept. For each indicator, the classification in terms of high, medium, or low complexity is based on how the students are using the indicator to respond to the item—specifically, to what degree does students’ engagement with the indicator contribute to the level of sense-making required by the item.

The evaluation of cognitive complexity is done at the individual item level. For an operational BIE Science Assessment form, after summing the operational points that reflect cognitive complexity at each

complexity level across all four indicators, the target distribution is that at least 10% of the points should be high cognitive complexity and no more than 35% of the points should be low cognitive complexity.

4.3 Item Review Committees and Processes

Items used on the BIE Science Assessment are developed to measure achievement on BIE's science standards, the Next Generation Science Standards. Cognia content specialists, in collaboration with BIE, ensure this alignment, and ongoing independent evaluations are held to verify alignment. In addition, independent reviews are scheduled to ensure that items and passages conform to bias and sensitivity guidelines.

4.3.1 Content and Item Reviews

The test developers at Cognia review newly developed items for

- alignment to the intended content standard;
- item integrity, including content and structure, format, clarity, and possible ambiguity;
- desired correct responses;
- appropriateness and quality of graphics;
- appropriateness of scoring-guide descriptions and distinctions;
- completeness of associated item documentation (e.g., scoring guide, content codes, key, grade level, cognitive complexity); and
- appropriateness for the designated grade level.

Newly developed stimuli and items for the BIE Science Assessment also undergo review by nationally representative panels of content and assessment experts, including educators from across many states. The purpose of these reviews is to evaluate items and determine their suitability for assessment by answering the following four questions:

- Does the item align with the assigned content standard(s)?
- Is the content accurate?
- Are the content and context grade-level appropriate?
- Does the item provide maximum accessibility for all students?

Additionally, Cognia has started to develop custom, culturally responsive items for use on the BIE Science Assessment. In 2024, the first custom cluster, at grade 5, was field tested. This cluster was developed in 2023 and reviewed by BIE as well as a specific committee of BIE teachers for both content and bias/sensitivity feedback. The same process was repeated with development of a cluster in grade 8 that was field tested in 2025, and development of a cluster in grade 11 that will be field tested in 2026. This year a grade 5 cluster will be developed to be field tested in 2027.

4.3.2 Bias and Sensitivity Review

Bias and sensitivity review is an essential component of the passage- and item-review process. All Cognia content specialists receive training in bias and sensitivity issues. Controversial and biased topics are avoided in the test development process. Internal reviews include review of not only content but context, with a particular awareness of bias and sensitivity issues that are specific to BIE students.

Since no one person is well-versed in the full spectrum of possible concerns, the bias and sensitivity review committee helps to ensure that all potential issues are identified. All stimuli and items undergo bias and sensitivity review prior to field testing.

The bias and sensitivity review committee comprises a diverse group of people who represent a variety of national student subgroups. The people currently serving on this committee span a variety of U.S. racial and ethnic groups, with an emphasis placed on representation for minoritized populations. These representatives have also been selected to represent specific subgroups of the student population, such as students who are English learners, students with disabilities, students within the LGBTQ+ community, and students impacted by special familial or home environment considerations, as well as having varied experiences with urban/suburban/rural environments and economically disadvantaged students.

Educator committees will periodically convene for the purpose of content development. Committees for both item content review (one committee per grade) and bias/sensitivity review (one committee for all grades) will be held, with BIE educators able to participate with two seats per committee. See Appendix D for additional information concerning the IRC and Bias Review Committees.

4.4 Test Forms Construction

The Cognia content specialists and psychometricians work collaboratively to produce operational test forms using sequential and iterative procedures that support both the content and construct validity of the assessments.

4.4.1 Item and Stimulus Selection

After field testing and item data review, Cognia test developers carefully select the items that will appear in the operational tests. In consultation with Cognia psychometricians, test developers consider the following criteria in selecting sets of items for the operational tests:

- **Content coverage/match to test design and blueprints.** The test designs and blueprints stipulate a specific number of items by item type and content distribution.
- **Item difficulty and complexity.** Item statistics are evaluated to ensure quality psychometric characteristics, as well as similar levels of difficulty and complexity from year to year.
- **“Clueing” items.** Items are reviewed for any information that might “clue” or provide information that would help to answer another item.

Test developers sort and lay out stimuli and items into test forms. During assembly of the test forms, the following criteria are considered:

- **Key patterns.** The sequence of keys (correct answers) is reviewed to ensure that their order appears random.
- **Option balance.** Multiple-choice (MC) items are balanced across forms so that key options are not markedly disproportionate.
- **Page fit.** Items always appear one per screen for online testing. Common science stimuli always appear to the left of the associated item.
- **Visual appeal.** Every effort is made to make each item as accessible as possible. Each item’s presentation may differ slightly depending on the delivery method and size of the screen.

A reviewer designated by the BIE reviews the test form and, prior to approval, specifically considers the following criteria:

- **Construct validity.** The test content is evaluated to determine the degree to which the test measures what it claims, or purports, to be measuring and items/tasks are aligned to the appropriate indicator/standard/measurable outcome.
- **Key accuracy.** Item keys (and the number of designated keys) are reviewed to ensure accuracy.

- **Positive phrasing in item stems.** Items are checked to ensure that negative words such as “not” and/or “except” are rarely, if ever, used.
- **Specific determiners.** Words such as “always,” “never,” “totally,” and “absolutely” are avoided whenever possible to prevent inadvertent cueing of correct or incorrect answers.
- **Clueing/clanging item associations.** The items on the test are reviewed to ensure that the answer to an item is not given away within another item on the same form (clueing) or that an item context is not too similar to another item on the same form (clanging).
- **Bias/sensitivity concerns.** The test is reviewed by all appropriate stakeholders within the BIE and assessment bureaus to ensure that the content is appropriate for students.
- **Errors or typos.** The test is reviewed to verify that the content and metadata are accurate and there appear to be no obvious human errors.

4.4.2 Selection Specifications to Meet Blueprint Requirements

All BIE Science Assessment items are appropriately field tested prior to operational use. Once stimuli have been field tested with a set of items, content specialists evaluate the statistics from the items associated with each stimulus. Often, items associated with a stimulus demonstrate a range of student performance, which is largely dependent upon factors inherent to each item. However, if a circumstance is encountered where many items associated with a stimulus are not performing as expected, this is evaluated carefully. While this scenario does not automatically mean the stimulus contains content that is not comprehensible or accessible, it does signal the need to thoroughly review the stimulus in relation to the item content and reevaluate the acceptability of the stimulus. Cognia assessment content specialists also review all the aspects of item content, and this is especially important when data indicate that further scrutiny is warranted.

The process for item data includes the following information for all field-tested items:

- classical item difficulty for all items (i.e., p -value)
- score distributions for polytomous items
- item option selection distribution for multiple-choice items
- 10 most frequent student responses for non-MC machine-scored items
- item-total and option-total correlations
- Item Response Theory (IRT) statistics

Differential item functioning (DIF) is generated using the standardization DIF procedure (Dorans & Kulick, 1986)³ to produce classifications among student subgroups, such as gender, economically disadvantaged status, and Limited English Proficiency (LEP) status. The flags listed in Tables 4-1 and 4-2 are used to identify those items that require an additional level of scrutiny.

³ DIF occurs when an item has difficulty measures that vary across contexts for similarly able subgroups of examinees. DIF procedures are designed to identify items on which the performances of certain subgroups of interest differ from each other after controlling for construct-relevant achievement. In order to ensure meaningful results, DIF statistics are not computed for populations containing less than 200 students in both subgroups. Analysis was conducted using field-test data to detect potential DIF at the item level. The standardization DIF procedure (Dorans & Kulick, 1986) was employed to evaluate subgroup differences. The computed DIF indices have a theoretical range of -1.0 to 1.0 for multiple-choice items. Critical values are defined as 0.05 and 0.10 and the values are flagged as statistically significant using $\alpha = 0.05$. If the absolute value of standardized DIF is equal to or greater than 0.10, the item is classified “C” DIF; items with absolute values greater than or equal to 0.05 are classified as “B” DIF; otherwise, items are classified as “A” DIF.

Table 4-1. Criteria for Flagged Items Based on Classical Test Theory (CTT) Statistics

Item-Flagging Criteria	Concern
If p -value of keyed response < 0.10	Item too difficult
If p -value of keyed response > 0.90	Item too easy
If p -value of distractor* $> p$ -value of keyed response	Possible mis-key
If p -value of distractor* > 0.35	Possible second correct option
If item-total correlation < 0.15	Poorly discriminating item
If item-total correlation < 0.00	Non-discriminating or negatively discriminating item
If DIF analysis is B or C	Possible bias in item (B, B-, C, C-)

*Note: These analyses examine item score and item option selection distribution for polytomous and selected-response items, respectively.

Table 4-2. Criteria for Flagged Items Based on Item Response Theory (IRT) Statistics

Item-Flagging Criteria	Concern
If IRT a -parameter < 0.30	Poorly discriminating
If IRT b -parameter < -3.00	Easy item
If IRT b -parameter > 3.00	Hard item
If IRT c -parameter > 0.35	Low ability students answer correctly (i.e., guessing)
If IRT b -parameter standard error of estimation > 0.3	Uncertainty around the item difficulty

The item content of each flagged item is reviewed and discussed by Cognia content specialists to make a decision regarding acceptability of the item. At the end of the process, all field-tested items are designated with a status of “Accept,” “Rework,” or “Reject.” Accepted items become eligible for operational testing. Rework items are eligible to be edited and field tested again so new item data can be generated. Rejected items are removed from the pool of items eligible for operational testing.

Cognia understands that item-level data review must be conducted thoroughly and carefully because of the impacts on test construction, which need to be consistent from administration to administration. Being experts in their respective content areas, Cognia’s content specialists also understand that some assessed standards are typically more challenging for students than others, and the specialists can simultaneously make good decisions about both content and data in accepting or rejecting items for operational use based on the item statistics. Finally, Cognia understands that items with DIF statistic flags need to be scrutinized for potential sources of bias. While a flag does not automatically mean the item contains biased content, it does signal the need to thoroughly review the item content and evaluate the ways in which the different focal groups would have access and ability to answer the item to ensure it is fair for all students.

4.4.3 Cognitive Processes

Cognitive process assesses whether students’ cognitive skills and processes align with those specified in the construct domains defined by test developers for all students and subgroups.

Subject matter experts, including educators, pay close attention to cognitive processes through test review. A team of science content specialists developed and internally reviewed items for BIE Science (items within Cognia’s Science Secure Item Bank). Subsequently, state/entity department of education content and/or assessment specialists, along with science educators, reviewed the items. Their input ensured that the tasks effectively measured the intended cognitive ability. Science educators participated in the review of all items developed and field tested, in addition to Cognia content expert reviews for alignment, content accuracy, and cognitive complexity levels. Committee review meetings included discussions and confirmation of target students’ characteristics, such as various learning targets, use of different science dimensions, and the degree of cognitive processing elicited by the stimulus and application of each dimension. Educators and specialists with expertise in science content provided feedback in these committee meetings, indicating that the test content did not require extraneous cognitive processes for engagement.

Chapter 5. Test Administration

Orderly and secure test administrations are necessary to protect secure test content and ensure that test data are validity-interpretable to meet score reporting and accountability reporting requirements.

5.1 Roles and Responsibilities

The Test Coordinator's Manual emphasizes that School Test Coordinators (STCs) play a crucial role as the primary source of assessment information for school stakeholders. STCs are responsible for keeping their schools informed about assessment policies, changes, and providing resources to teachers.

Manuals, including the Test Coordinator's Manual and the Test Administrator's Manual, ensure consistent administration procedures across schools, emphasizing test security and ethical administration. These documents are accessible on Cognia's [BIE Science Help and Support Website](#).

Additional staff, such as Test Administrators (TAs), are vital for a successful BIE Science Assessment. TAs must follow procedures outlined in the Test Administrator's Manual and attend training provided by the STC. If additional staff are needed, those who have received training and signed the Confidentiality Agreement may assist with one-to-one accommodations. STCs are appointed at the local level.

5.2 Test Administrator's Manual

For Spring 2025, the Test Administrator's Manual (TAM) outlined the steps to follow before, during, and after administration of the Spring 2025 BIE Science Assessment. Understanding of and compliance with each of these steps is vital for successful administration.

The TAM covers administration policies such as security guidelines and administration information, accessibility features and accommodations including requirements for computer-based tests (CBT) and paper-based tests (PBT), preparing for CBTs and PBTs, administering CBTs and PBTs, directions and scripts for use during CBT and PBT administrations, and what to do at the completion of CBTs and PBTs.

5.3 TA and Proctor Training Requirements and 2025 Test Administrations

All TAs and proctors involved in test administration, preparation, and security are required to attend training provided by the STC. Training should include information on test security policies and procedures, test administration procedures, documentation and provision of testing accommodations, and the importance of strictly following all directions in the manuals.

5.4 Testing Irregularity Reports

During the Spring 2025 BIE Science testing window, test administrators and coordinators were trained to report test administration irregularities. The BIE defines a testing irregularity as any incident in the handling or administration of a test that results in questioning the accuracy of the data or security of the test that may or may not result in an invalidation.

Seven test irregularities were submitted to BIE for the following reasons:

- One grade 5 student claimed test interruption and completed the test. The school invalidated the scores for the student.
- A grade 5 student got sick and went home. The student was rescheduled for a make-up test.

- A grade 5 student checked out for an appointment a few minutes after starting and paused/exited test. The student completed the session as a makeup on another day.
- A grade 11 student had his left wireless earbud in his ear before testing. The school confiscated the earbud before starting the test.
- A grade 11 student was reassessed on all sessions instead of a missed session only. BIE requested Cognia to count the first attempt and merge the test sessions.
- TA did not follow the scheduled time for a grade 5 student and started session 3 before all students completed session 2.
- Testing time was shortened for a grade 8 student due to inclement weather. The student was given time (50 mins) to complete and start off at the last question answered.

5.5 Test Security

The BIE Assessment Program requires that the BIE Science Assessment be treated with the highest level of test security and accountability. The security of BIE materials must be maintained before, during, and after the test administration. TAs, proctors, and STCs are required to follow the guidelines in the TAM for distributing, collecting, and returning testing materials. All testing personnel are required to have access to a central, locked storage space for safekeeping of test materials until print materials are returned to Cognia.

To maintain the validity of the tests administered in the assessment program, keeping all test questions secure is absolutely necessary. If security is breached or compromised, the assessment results may not be valid. If one student or school has advantages not awarded to another, the test administration is no longer standardized and loses the important distinction of being appropriate for program accountability.

TAs must follow these security guidelines before, during, and after testing:

- Receive training on test security and administration by the STC.
- Complete the BIE Confidentiality Agreement and return it to the STC.
- Follow the testing schedule established by the school.
- Ensure the TA is not assigned to a classroom in which a relative is being tested.
- Carry out standard examination procedures.
- Ensure secure test materials are secured in a central and locked area when not in use.
- Use the security checklist or a similar tracking tool daily, as provided by the STC, during test administration to check in and check out all test materials.
- Report any possible breaches of security to the STC immediately. Examples of security breaches include but are not limited to
 - improper handling of test materials, such as
 - someone reproducing any student responses,
 - allowing any unauthorized access to test materials before, during, or after testing, or
 - leaving test materials (including computers being used for CBTs) unsecured when the TA or a proctor is not in the classroom, and
 - improper test administration procedures, such as
 - coaching students during testing,
 - altering student responses in any way, or
 - stray-mark cleanup, including but not limited to erasing double-marks or lightly marked answers.

School staff are prohibited from studying or discussing online test questions in any manner, either among themselves or with students before, during, or after testing.

Chapter 6. Scoring: Scope of Work, Processes, and Procedures

6.1 Scope of Work

The operational and field-test items on the 2024–25 BIE Science Assessment included 4-point open-ended response items. Table 6-1 outlines the number and type of each item per grade.

Table 6-1. Overview of BIE Science Assessment Scope-of-Work

Grade 5	Grade 8	Grade 11
OP – 3 OE4	OP – 3 OE4	OP – 3 OE4
FT – 2 OE4	FT – 2 OE4	FT – 2 OE4

OP=Operational, OE4= 4-point open-ended response item

FT=Field Test, OE4= 4-point open-ended response item

6.2 Operational Scoring: Processes and Procedures

6.2.1 Score Verification of Multiple-Choice Items

For both computer-based tests (CBTs) and paper-based tests (PBTs), responses to multiple-choice items were compared to scoring keys using item analysis software. This robust software compared each student response to multiple-choice items to the respective answer key and assigned a maximum score of 1 point for correct responses and 0 points for incorrect answers. In PBTs, if students filled in multiple bubbles in response to one item, the response was assigned 0 points. At the end of an administration, a second independent validation of all the student responses was conducted to compare and validate results to ensure accurate machine scoring.

6.2.2 Benchmarking Field-Test Items

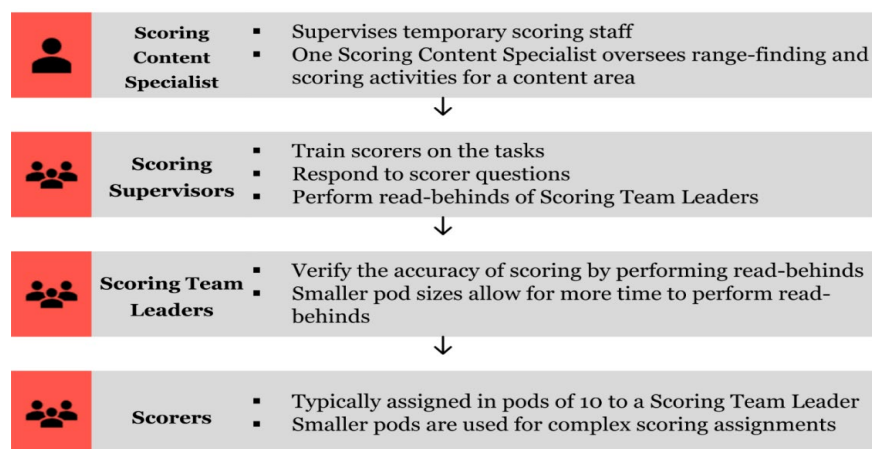
Field-test items were benchmarked by scoring leadership, who selected a range of responses to be reviewed. Benchmarking meetings were held between the Scoring Content Specialist and the Content Development Specialist for each grade to discuss the way the students engaged with each item and to review the suggested scores assigned to the benchmarked materials. Each of the field-test items were reviewed to determine their scorability and to set the scoring standards using exemplar student responses. Anchor and practice sets were created to be used in the training of scorers. One qualification set of ten responses was also identified before field-test scoring began.

6.2.3 Scoring of Open-Ended Response Items

6.2.3.1 Personnel Structure

Cognia’s personnel structure for scoring responses consisted of four hierarchical levels, as shown in Figure 6-1.

Figure 6-1. Cognition Scoring Staff



All responses were scored by fully vetted scorers who were supervised by Scoring Team Leaders (STLs). The Scoring Supervisors monitored the work of the STLs assigned to them. The Scoring Content Specialist monitored the work of the Scoring Supervisors, STL, and scorers. Scoring Content Specialists report to the Scoring Content Group Manager. This hierarchical structure, whereby each level monitors the one below, ensures reliable quality and consistency in scoring.

Scoring Content Specialist

The Scoring Content Specialist functioned as the primary lead for his or her designated content area and as a liaison between scoring activities and the Scoring Project Manager to ensure that established quality standards and production schedules were met.

During scoring, the Scoring Content Specialist was responsible for supervising all scoring staff working on the project, including Scoring Supervisors and STLs. The Scoring Content Specialist was also responsible for assuring the consistency and accuracy of scoring work performed by individual scorers and across groups of scorers.

Scoring Supervisors

Scoring Supervisors managed the scorer training and supervised the STLs and scorers working on a designated item and/or content. Scoring Supervisors worked closely with the STLs to ensure consistency and provide counsel and retraining to scorers, as necessary. In addition, Scoring Supervisors engaged in supervisory oversight and performed quality-control checks to ensure the consistency and accuracy of the STLs. Scoring Supervisors, who were responsible for monitoring training and conducting the retraining of scorers, were selected for their ability to instruct and for their level of expertise in their respective disciplines.

Scoring Team Leaders

The STLs were responsible for supervising and monitoring the group of scorers assigned to them. STLs worked closely with their scorers to maintain consistently accurate scoring. They provided quality checks, and they counseled scorers, as necessary. STLs were responsible for monitoring and maintaining accurate scoring of their assigned scorers. This included performing read-behinds on scorers and monitoring other quality-control measures. STLs were responsible for arbitrating responses scored by multiple scorers when the assigned scores varied by more than one score point. The arbitration process ensured that such responses received the necessary attention by providing an additional review before

assigning a third and final resolution score. In addition to the essential quality control, the arbitration process provided continued opportunities for scorer training.

Because the read-behinds that the STLs performed moderated the scoring process and thus maintained the integrity of the scores, individuals chosen to fill STL positions were selected for their accuracy and content knowledge.

Scorers

Scorers are individuals who evaluate student responses and assign scores.

6.2.3.2 Scorer Recruitment

Cognia actively sought a diverse pool of scorers with a broad range of backgrounds, including teachers, scientists, business professionals, graduate school students, and retired educators.

The minimum requirement to assume a position as a scorer or Scoring Team Leader is 48 college credits, which include classes related to the content area being scored. Scoring Supervisors must hold a bachelor's degree with classes related to the content area being scored. All potential scorers and leadership staff submitted documentation (e.g., résumés and/or transcripts) as evidence of meeting the education and experience requirements. Each scorer and leadership staff signed a binding non-disclosure/confidentiality agreement as well.

6.2.3.3 Scoring Platform

For the scoring of the 2024–25 BIE Science test administration, images were loaded into iScore, the proprietary image-based scoring system used by Cognia to view and record scores submitted by scorers for each open-ended item. The iScore system ensures the security of student responses and test items. During scoring, no student names or schools/districts associated with viewed student work are visible to scorers, and all Scoring Services temporary associates are subject to the same nondisclosure requirements as full-time Cognia staff. Cognia maintained security during scoring by using a highly secure, server-to-server interface, ensuring that access to all student response images was limited only to scorers and appropriate Cognia staff.

Scorers evaluated most student responses from images rendered by the online testing platform and a small number of responses from scanned images of paper-based tests. Whether administered in an online or a paper/pencil environment, all responses were scored applying the same scoring criteria.

Prior to the beginning of scoring, databases were created for each grade to receive submitted student responses for each item to be scored. To provide maximum security for all testing and scoring materials, each scorer was required to log on to the scoring systems using a unique combination of an assigned username and password.

6.2.3.4 Leadership Training

Scoring Supervisors reviewed training materials and consulted with the Scoring Content Specialist in advance of scorer training to ensure full understanding of the scoring parameters and decisions for each item. Scoring Supervisors then trained STLs during a separate session prior to scorer training. The STLs, responsible for scoring student responses, were required to achieve the same standard as scorers on item qualification sets: a minimum accuracy scoring rate of 70 percent exact, and 90 percent exact plus adjacent agreement (70/90).

6.2.3.5 Scorer Training

For the scoring of BIE science operational and field-test items, all scorer training was conducted during live training sessions via secure Zoom videoconferencing. Scorer training at the item level followed this process:

- Brief overview of the assessment program and the training process
- Review of the student prompt, passage, and the scoring rubric
- Guided review of the Anchor Set responses
- Analysis and discussion of the score for each anchor response and the scoring rationale
- Independent review and scoring of responses in the Practice Set to replicate the actual scoring process
- Discussion of each practice response, revealing the actual score assigned to the student response and the scoring rationale
- Methodical review of all scoring criteria while paying particular attention to the fine lines that determine the cut-points between adjacent score points
- Question and answer segment addressing any remaining scorer questions

After training, scorers were given two opportunities to qualify for operational items. If scorers were unable to attain a score match of at least 70 percent exact and 90 percent exact plus adjacent agreement on the first qualifying set, they were retrained by discussing the responses contained in the first qualification set with respect to the score-point descriptions of the rubric and by comparing them to the responses of the Anchor Set. Following this retraining, scoring leadership would administer a second qualification set. If scorers achieved a scoring accuracy rate of at least 70 percent exact and 90 percent exact plus adjacent agreement on the second qualification set, they could score student responses. Scorers who failed to pass the minimum threshold were not allowed to score that item. They were either trained on another item or dismissed from the project. Appendix E captures the qualification rates for all operational grades.

For field-test scoring, one qualification set was administered for each item. If scorers did not meet the standard of at least 70 percent exact and 90 percent exact plus adjacent, they were not permitted to score that field-test item.

6.2.3.6 Monitoring Scoring Quality

Scorers were required to demonstrate and maintain their ability to score student responses accurately and consistently throughout the scoring process. The iScore system enabled scoring leadership to measure and monitor individual and group performance on each scored item in terms of accuracy and consistency, and in terms of read rate (scoring speed) and overall production rate on a constant, real-time basis. Scoring tools employed to measure scoring quality were as follows:

- Read-behind scoring
- Double-blind scoring
- Embedded validity responses
- Recalibration sets

Each scorer's performance on the above procedures was monitored and recorded by the scoring system, and scoring leadership could review data related to the accuracy, consistency, and overall quality of scoring. Scoring leadership was always available to answer scorer questions. They also counseled and retrained scorers as needed to determine whether a scorer should continue scoring. If a scorer's performance did not meet the prescribed quality standards, scoring leadership initiated a process through

which that scorer's work was invalidated and returned to the scoring queue of unscored responses to be re-scored by those scorers who demonstrated scoring accuracy at or above standard.

Read-Behind Scoring

Read-behind scoring allowed scoring leadership to monitor each scorer's scoring performance by way of an immediate real-time snapshot of the scorer's accuracy. The data generated by read-behind scoring presented leadership with opportunities to answer questions and to provide counsel to scorers who may have had trouble maintaining the scoring standards. If the scores assigned by the scorer and the STL were discrepant (more than one score point apart) or if there were a significant number of adjacent scores (one score point apart) between the scorer and the STL, scoring leadership then counseled and retrained the scorer. Scoring leadership determined when or whether these scorers were given access to resume operational scoring. Retrained scorers were subject to additional monitoring and read-behinds.

The number of read-behinds for each scorer varied depending on the accuracy of the scorer. BIE scoring specifications require a read-behind target of 2% of responses scored. Consistently accurate scorers would only receive the minimum number of read-behinds, whereas scorers who exhibited difficulties in maintaining accuracy or consistency received additional read-behinds.

In addition to scorers, scoring leadership was also subject to quality assurance reviews, which were administered by the Scoring Content Specialist. They monitored scoring leadership's accuracy and consistency by reviewing the read-behind results.

Double-Blind Scoring and Arbitration Resolution

Double-blind scoring refers to the process of two scorers independently scoring the same response. During this process, neither scorer has any knowledge of the other scorer's score. The double-blind process helps inform scoring leadership about the consistency of scoring among peer scorers who actively score an item. At least 2 percent of responses in all grades require a double-blind score.

During double-blind scoring, the scoring system distributes randomly selected responses assigned for double-blind scoring to different scorers without alerting either scorer. The scoring system then records each scorer's score and routes any scoring discrepancies of more than one point between the two scores to an arbitration response queue for resolution by the STL. The percentage of double-blind responses sent to arbitration by a scorer because of a difference in actual scores (i.e., not including blank or unreadable responses) should not have exceeded 10 percent. If a scorer's arbitration percentage exceeded this threshold, scoring leadership counseled, retrained, and/or dismissed the scorer.

Embedded Validity Responses

Validity responses are prescored responses that serve calibration purposes at the onset of scoring an item. Ten validity responses were embedded in the first 100 live student responses and distributed to each scorer in randomized order. Scorers were not aware when they were scoring an embedded validity response as compared to a live student response. Scorers who demonstrated an accuracy rate of less than 70% exact on each composite score were counseled and the STL increased the number of read-behinds to ensure accuracy.

Recalibration Sets

A set of five calibration papers was administered starting with the second day of scoring an item. This set of five responses, selected by scoring leadership, was inserted into the scoring queue, and served as a refresher. It was used to gauge the scorers' ability to maintain accurate scoring of the item on days following their initial item training.

Interrater Reliability

Section 6.2.3.6 of this report describes in detail the processes that were implemented to monitor the quality of the hand-scoring of student responses for open-ended items. One of these processes was double-blind scoring: A minimum of 2 percent of student responses in all grades was randomly selected and routed to two different scorers to be scored independently. Results of the double-blind scoring were used to identify scorers who required retraining or other interventions and are presented here as evidence of the reliability of the BIE Science Assessment.

Table 6-2 shows the number of included scores and the percentage of exact and adjacent agreement, as well as percent third reads for adjudication.

Table 6-2. Summary of Interrater Reliability Statistics for BIE Science—OP Items by Grade

Grade	Total # of Responses Scored	Total # of Double-Blind Responses Scored	Total % Double-Blind Responses Scored	Score Categories	Score Point Ranges	% Exact	% Adjacent	% Third Reads
5	5,714	266	5	1	0-4	75	22	2.6
8	4,933	231	5	1	0-4	83	16	1.8
11	3,197	231	7	1	0-4	84	11	4

6.2.3.7 Score-of-Record Rules

Per scoring specifications, the final score-of-record (SOR) was determined as shown in Table 6-3.

Table 6-3. Score of Record Rules

Score-of-Record Rules	
Condition	Resolution
Only one (first) read is provided	First read is SOR
Only one (first) read is provided	If the first read and the second read are equal, the first read is SOR.
First and second reads are provided	If the first and second read differ by 1 point (i.e., they are adjacent to each other), the first read is SOR.
First and second reads are provided	If the first and second read differ by 2 or more points (i.e., they are discrepant), the resolution score provided by leadership is SOR.
One or two reads and one read-behind score	Read-behind score (provided by leadership) is SOR.
Two reads and two read-behind scores	The later read-behind is SOR.

Chapter 7. Classical Item and Test Analysis

As noted in the *Principles of Educational and Psychological Testing* (Brown, 1983), “A test is only as good as the items it contains.” A complete evaluation of a test’s quality must include an evaluation of each item. Both *Standards for Educational and Psychological Testing* (AERA et al., 2014) and *Code of Fair Testing Practices in Education* (Joint Committee on Testing Practices, 2004) include standards for identifying quality items. Items should assess only knowledge or skills that are identified as part of the domain being tested and should avoid assessing irrelevant factors. Items should also be unambiguous and free of grammatical errors, potentially insensitive content or language, and other confounding characteristics. In addition, items must not unfairly disadvantage students in particular racial, ethnic, or gender groups.

7.1 Classical Item Statistics

All operational items were evaluated in terms of classical item difficulty, which under classical test theory practices is defined as the average scored response on an item, divided by the maximum possible score for the item. Although this index is traditionally described as an estimate of item difficulty, it is properly interpreted as an easiness index. The greater in value a classical item difficulty is, the easier the item.

Items that are answered correctly by almost all students provide little information about differences in student abilities, but they do indicate knowledge or skills that have been mastered by most students. Similarly, items that are correctly answered by very few students provide little information about differences in student abilities, but they may indicate knowledge or skills that have not yet been mastered by most students. In general, to provide adequate measurement, classical difficulty indices should range from near-chance performance (e.g., 0.25 for four-option multiple-choice items) to 0.90, with a majority of items generally falling between around 0.4 to 0.7. However, on a standards-referenced assessment such as the BIE science, it is appropriate to include items with very low or very high item difficulty values to ensure sufficient content coverage.

A desirable characteristic of an item is for higher-ability students to perform better on the item than lower-ability students do. The correlation between student performance on a single item and total test score is a commonly used measure of this characteristic of the item. Within classical test theory, the item-total correlation is referred to as the item’s classical discrimination because it indicates the extent to which successful performance on an item discriminates between high and low scores on the test. Each of the item-total correlations reported here is the Pearson correlation between scored responses on a given item and total raw scores. This Pearson correlation is commonly referred to as the point-biserial correlation (for a dichotomously scored item) and a point-polyserial correlation (for a polytomously scored item). The theoretical range of these correlations is -1.0 to $+1.0$, with a typical observed range from 0.2 to 0.6. Discrimination indices can be thought of as measures of how closely an item assesses the same knowledge and skills assessed by other items contributing to the criterion total score. That is, the discrimination index can be thought of as a measure of construct consistency.

A comparison of indices across grade levels is complicated because these indices are population-dependent. Direct comparisons would require that either the items or the students were common across groups. Since that is not the case, it cannot be determined whether differences in performance across grade levels are due to differences in student abilities, differences in item difficulties, or both.

Classical item difficulties and item-total correlations are provided in Appendix F, and a summary is provided in Table 7-1. These statistics were calculated based on a national student sample from all the states and entities that used Cognia Secure Science Item Bank (SSIB).

Table 7-1. Item Summary by Grade

Grade	Item Type	Item Mean
5	MC	0.33
5	OR	0.74
8	MC	0.40
8	OR	0.61
11	MC	0.36
11	OR	0.57

7.2 Total Test and Subscore Intercorrelations

When subscores are strongly related to each other, it implies a high internal consistency between subscores. The Pearson correlation matrices among the individual reporting categories (i.e., subscores) are shown in Tables 7-2 to 7-4. These correlations are based on BIE student performance, which includes their total test scale score and their three subscores.

Table 7-2. Pearson Correlations of Total Test and Subtest Raw Scores on BIE Science Grade 5, as a Function of Operational Form

Subtest	Number of Items	Total Test	Earth and Space Sciences	Physical Sciences	Life Sciences
Operational Form 1					
Earth and Space Sciences	12	0.87	1	--	--
Physical Sciences	12	0.84	0.62	1	--
Life Sciences	11	0.91	0.70	0.64	1
Operational Form 2					
Earth and Space Sciences	12	0.86	1	--	--
Physical Sciences	12	0.84	0.58	1	--
Life Sciences	11	0.91	0.67	0.65	1

Table 7-3. Pearson Correlations of Total Test and Subtest Raw Scores on BIE Science Grade 8, as a Function of Operational Form

Subtest	Number of Items	Total Test	Earth and Space Sciences	Physical Sciences	Life Sciences
Operational Form 1					
Earth and Space Sciences	13	0.79	1	--	--
Physical Sciences	12	0.88	0.55	1	--
Life Sciences	10	0.84	0.51	0.59	1
Operational Form 2					
Earth and Space Sciences	13	0.80	1	--	--
Physical Sciences	12	0.88	0.55	1	--
Life Sciences	10	0.82	0.49	0.58	1

Table 7-4. Pearson Correlations of Total Test and Subtest Raw Scores on BIE Science Grade 11, as a Function of Operational Form

Subtest	Number of Items	Total Test	Earth and Space Sciences	Physical Sciences	Life Sciences
Operational Form 1					
Earth and Space Sciences	12	0.83	1	--	--
Physical Sciences	13	0.85	0.55	1	--
Life Sciences	12	0.87	0.61	0.59	1
Operational Form 2					
Earth and Space Sciences	12	0.75	1	--	--
Physical Sciences	13	0.83	0.41	1	--
Life Sciences	12	0.84	0.49	0.54	1

Chapter 8. Psychometrics: Item Response Theory (IRT) Scaling and Equating

This chapter describes the procedures used to scale the BIE Science Assessment, where all test forms were pre-equated for the Spring 2025 administration.

8.1 IRT Models

All BIE items were previously calibrated using item response theory (IRT). IRT uses mathematical models to define a relationship between an unobserved measure of student proficiency, usually referred to as theta (θ), and the probability (p) of getting a dichotomous item correct or of getting a particular score on a polytomous item. In IRT, all items are assumed to be independent measures of the same construct (i.e., of the same θ). Another way to think of θ is as a mathematical representation of the latent trait of interest. Several common IRT models are used to specify the relationship between θ and p (Hambleton & van der Linden, 1997; Hambleton & Swaminathan, 1985). The process of determining the specific mathematical relationship between θ and p is called item calibration. After items are calibrated, they are defined by a set of parameters that specify a nonlinear, monotonically increasing relationship between θ and p . Once the item parameters are known, an estimate of θ for each student can be calculated. This estimate, $\hat{\theta}$, is considered to be an estimate of the student's performance. It has characteristics that may be preferable to those of raw scores for equating and scaling purposes.

For the BIE Science Assessments, the three-parameter logistic (3PL) model was used for dichotomous (selected-response) items and the Graded-Response Model (GRM) was used for polytomous (constructed-response) items. The 3PL model for dichotomous items can be defined as:

$$P_i(\theta_j) = P(U_i = 1|\theta_j) = c_i + (1 - c_i) \frac{\exp[D a_i(\theta_j - b_i)]}{1 + \exp[D a_i(\theta_j - b_i)]}$$

Where U indexes the scored response on an item,

i indexes items,

j indexes students,

α represents item discrimination,

b represents item difficulty,

c is the lower asymptote parameter, and

D is a normalizing constant equal to 1.701.

In the GRM for polytomous items, an item is scored in a $k + 1$ graded category that can be viewed as a set of k dichotomies. At each point of dichotomization (i.e., at each threshold), a two-parameter model can be used. This implies that a polytomous item with $k + 1$ categories can be characterized by k Item Category Threshold Curves (ICTCs) of the two-parameter logistic form:

$$P_{ik}^*(\theta_j) = P(U_i \geq k|\theta_j) = \frac{\exp[Da_i(\theta_j - b_i + d_{ik})]}{1 + \exp[Da_i(\theta_j - b_i + d_{ik})]},$$

where U indexes the scored response on an item,

i indexes the items,

j indexes students,

k indexes threshold ($k = 0, 1, \dots, m$),

a represents item discrimination,

b represents item difficulty,

d represents item category threshold, and

D is a normalizing constant equal to 1.701.

After computing k ICTCs in the GRM, $k + 1$ Item Category Characteristic Curves (ICCCs), which indicate the probability of responding to a particular category given θ , are derived by subtracting adjacent ICTCs:

$$P_{ik}(\theta_j) = P(U_i = k|\theta_j) = P_{ik}^*(\theta_j) - P_{i(k+1)}^*(\theta_j),$$

where P_{ik} represents the probability that the score on item i falls in category k , and

P_{ik}^* represents the probability that the score on item i falls above the threshold k .

Note that $P_{i0}^* = 1$ and $P_{i(m+1)}^* = 0$.

The GRM is also commonly expressed as:

$$P_{ik}(\theta_j) = \frac{\exp[Da_i(\theta_j - b_i + d_{ik})]}{1 + \exp[Da_i(\theta_j - b_i + d_{ik})]} - \frac{\exp[Da_i(\theta_j - b_i + d_{i(k+1)})]}{1 + \exp[Da_i(\theta_j - b_i + d_{i(k+1)})]}.$$

The Item Characteristic Curve (ICC) for polytomous items is computed as a weighted sum of ICCCs, where each ICCC is weighted by a score assigned to a corresponding category:

$$E(U_i|\theta_j) = \sum_{k=0}^{m+1} w_{ik} P_{ik}(\theta_j),$$

where w_{ik} is the weighting constant and is equal to the number of score points for score category k on item i .

See Lord and Novick (1968), Hambleton and Swaminathan (1985), and Baker and Kim (2004) for more information about item calibration and parameter estimation.

8.2 Dimensionality Analysis

Tests are constructed with multiple content-area subcategories and their associated knowledge and skills. Hence, the potential exists for dimensions being invoked beyond the common primary dimension. Generally, the content-area subcategories are highly correlated with each other, and the primary dimension they share typically explains an overwhelming majority of the variance in test scores. The

presence of just such a dominant primary dimension is the psychometric assumption that provides the foundation for the unidimensional item response theory (IRT) models that are used for scaling and equating of the BIE tests.

The purpose of dimensionality analysis is to investigate whether violation of the assumption of test unidimensionality is statistically detectable and, if so, (a) the degree to which unidimensionality is violated and (b) the nature of the multidimensionality. Findings from dimensionality analyses performed on the BIE operational items for science are reported below. (Note: Only operational items were analyzed since they are used for score reporting.)

The dimensionality analyses were conducted using the nonparametric IRT-based methods DIMTEST (Stout, 1987; Stout, Froelich, & Gao, 2001) and DETECT (Zhang & Stout, 1999). Both methods use the estimated average conditional covariances for item pairs as their basic statistical building block. A conditional covariance is the covariance between two items conditioned on expected total score for the rest of the test, and the average conditional covariance is obtained by averaging across every possible conditioning score. When a test is strictly unidimensional, all conditional covariances are expected to take on values within random noise of zero, indicating statistically independent item responses for examinees with equal expected total test scores. Nonzero conditional covariances are essentially violations of the principle of local independence, and local dependence implies multidimensionality. Thus, nonrandom patterns of positive and negative conditional covariances indicate multidimensionality.

DIMTEST is a hypothesis-testing procedure for detecting violations of local independence. The data are first divided into a training sample and a cross-validation sample. Then an exploratory analysis of the conditional covariances is conducted on the training sample data to find the cluster of items that display the greatest evidence of local dependence. The cross-validation sample is then used to test whether the conditional covariances of the selected cluster of items display local dependence, conditioned on total score on the non-clustered items. The DIMTEST statistic follows a standard normal distribution under the null hypothesis of unidimensionality.

The DETECT statistic is an effect-size measure of multidimensionality. As with DIMTEST, the data are first divided into a training sample and a cross-validation sample. The training sample is used to find a set of mutually exclusive and collectively exhaustive clusters of items that best fit a systematic pattern of positive conditional covariances for pairs of items from the same cluster and negative conditional covariances from different clusters. Next, the clusters from the training sample are used with the cross-validation sample data to average the conditional covariances. Within-cluster conditional covariances are summed; from this sum, the between-cluster conditional covariances are subtracted. This difference is divided by the total number of item pairs, and this average is multiplied by 100 to yield an index of the average violation of local independence for an item pair. DETECT values less than 0.2 indicate very weak multidimensionality (or near unidimensionality); values of 0.2 to 0.4, weak multidimensionality; values of 0.4 to 1.0, moderate multidimensionality; and values greater than 1.0, strong multidimensionality (Roussos & Ozbek, 2006).

DIMTEST and DETECT were applied to the BIE science tests, Set A and Set B for each of grades 5, 8, and 11. First, each dataset was split into a training sample and a cross-validation sample. DIMTEST was then applied to each sample. As shown in Table 8-1, the DIMTEST null hypothesis of unidimensionality was not rejected for any of the datasets. Because the DIMTEST statistic did not find any evidence for rejecting the null hypothesis of unidimensionality, no DETECT analyses were needed or conducted for estimating the sizes of the violations of local independence.

In summary, the dimensionality analyses conducted for the BIE Science Assessments found no evidence of any violations of the assumption of unidimensionality and, thus, strongly support the use of the unidimensional IRT models whose implementation is detailed in this chapter.

Table 8-1. DIMTEST Hypothesis Testing Results, as a Function of Grade*

Content Area	Grade	P-value	Interpretation
Science (operational Set A)	5	0.26	Does not detect violation of unidimensionality
	8	0.29	Does not detect violation of unidimensionality
	11	0.25	Does not detect violation of unidimensionality
Science (operational Set B)	5	0.11	Does not detect violation of unidimensionality
	8	0.15	Does not detect violation of unidimensionality
	11	0.33	Does not detect violation of unidimensionality

* Since the DIMTEST did not reject H_0 (unidimensionality), we don't need to run DETECT.

8.3 Item Response Theory Results

The tables in Appendix G give the IRT item parameters of all common items on the Spring 2025 BIE Science Assessment by grade.

Test characteristic curves (TCCs) are based on the IRT item parameters and display the expected (average) raw score associated with each θ_j value between -4.0 and 4.0 , or equivalently the expected (average) raw score associated with each observable scaled score (see Section 8.4 for details on scaled scores). Mathematically, the TCC is computed by summing the ICCs of all items that contribute to the raw score. Using the notation introduced in Section 7.1, the expected raw score at a given value of θ_j is

$$E(X|\theta_j) = \sum_{i=1}^n E(U_i|\theta_j),$$

Where i indexes the items (and n is the number of items contributing to the raw score),

j indexes students (here, θ_j runs from -4 to 4), and

$E(X|\theta_j)$ is the expected raw score for a student of ability θ_j ,

U indexes the scored response on an item.

The expected raw score monotonically increases with θ_j , consistent with the notion that students of high ability tend to earn higher raw scores than do students of low ability. Most TCCs are “S-shaped”—flatter at the ends of the distribution and steeper in the middle.

Test information functions (TIFs) display the amount of statistical information the test provides at each value of θ_j or equivalently display the amount of statistical information the test provides at each observable scaled score. TIFs depict test score precision across the entire latent trait continuum. There is an inverse relationship between the information of a test and its conditional standard error of measurement (CSEM). The CSEM at a given θ_j [$CSEM(\theta_j)$] is equal to the inverse of the square root of the statistical information at θ_j (e.g., Hambleton, Swaminathan, & Rogers, 1991). That is, the $CSEM(\theta_j)$ is equal to the inverse of the square root of the TIF at a given θ_j [$TIF(\theta_j)$] the expression for which can be written as follows:

$$CSEM(\theta_j) = \frac{1}{\sqrt{TIF(\theta_j)}}$$

Compared to the tails, TIFs are often higher near the middle of the θ distribution, where most students are located and where most items are sensitive by design.

Appendix H contains graphs of the TCC and TIF for the BIE science tests, Set A and Set B for each of grades 5, 8, and 11. Each TCC graph displays the expected raw score (on the vertical axis) for the entire theta scale (on the horizontal axis). Each TCC graph also has a set of vertical lines that indicate the values of the theta cut scores for the given grade. Each TIF graph displays test information value (on the vertical axis) at the entire theta scale (on the horizontal axis). Each TIF graph also has a set of vertical lines that indicate the values of the scaled score cut scores for the given subject and grade.

8.4 Equating

The purpose of equating is to ensure that scores obtained from different forms of a test are equivalent to each other. Equating may be used if multiple test forms are administered in the same year or to equate one year's forms to those given in the previous year.

The cut scores for 2025 BIE Science Assessments were established in 2022 by Cognia with a standard setting using previously calibrated items from the Cognia SSIB. As new SSIB items are field tested each year (for BIE and other SSIB contracts), the operational items on the forms serve as equating items to bring the field-test items onto the SSIB scale.

The pre-equating process uses item bank values of the IRT item parameters to place the pre-equated test form onto the established IRT scale. Equating ensures that students are not given an unfair advantage or disadvantage because the test form they took is easier or harder than those taken by other students.

8.5 Reported Total Test and Subtest Scale Scores

The θ scale used in IRT calibrations is not readily understood by most stakeholders. As such, reporting scales are used for BIE reporting. The reporting scales are linear transformations of the underlying θ scale. To obtain a student's scaled score on a given assessment, the student's raw score (i.e., total number of points earned) is translated into a value on the underlying θ scale using TCC mapping. The student's θ value is translated into a scaled score (SS) using the following linear equation:

$$SS = \beta_0 + \beta_1\theta$$

where β_0 is an intercept constant and

β_1 is a slope constant,

The CSEM can also be translated into a scaled CSEM. Whereas values of the CSEM are on the θ scale, values of the scaled CSEM are on the reporting scale. The scaled CSEM is obtained via the following equation:

$$\text{Scaled CSEM} = \beta_1 \times \text{CSEM}(\theta)$$

Table 8-2 shows the slope and intercept terms used for the Spring 2025 BIE Science Assessment to calculate the scaled scores. See Appendix I for Raw to Scale Score Lookup Tables.

Table 8-2. Spring 2025 Scaled Score Slopes and Intercepts by Subject and Grade

Grade	Slope	Intercept
5	12.5	553.57
8	10.0	855.10
11	7.5	1159.72

It is important to note that converting from raw scores to θ values to scaled scores does not change students' achievement-level classifications. Given the relative simplicity of raw scores, it is fair to question why scaled scores are reported instead of raw scores. Scaled scores make the reporting of grade-level results consistent across test forms and administrations. It is this uniformity across scaled scores that facilitates the understanding of student performance. The psychometric advantage of scaled scores over raw scores comes from there being linear transformations of θ . Since the θ scale is used for pre- or post-equating, scaled scores are comparable from one year to the next. Raw scores are not.

8.6 Performance Levels

The cut scores used for the Spring 2025 BIE Science Assessments are the cut scores that were established for the Cognia SSIB in 2022, on which the BIE Science summative assessments are based. BIE conducted a standard validation in August 2022; however, due to low participation, BIE determined the cut scores set for the SSIB would be used. A standards validation was held in 2023 to verify these cut scores against the BIE Science Assessment. The cut scores on the theta scale and the reporting scale, used for the Spring 2025 BIE Science Assessments, are presented in Table 8-3.

Table 8-3. Spring 2025 Cutpoints on the Theta Metric and Reporting Scale by Grade

Theta				Scale Score			
Cut1	Cut2	Cut3	Min	Cut1	Cut2	Cut3	Max
-0.75048	0.51466	1.70117	500	544	560	574	590
-0.96101	0.48988	2.73095	800	845	860	882	890
-0.76114	0.03716	2.91134	1100	1154	1160	1181	1190

8.6.1 Percentages of Students in Each Performance Level

The empirical performance level distributions for the Spring 2025 administration of BIE Science Assessments are shown in Table 8-4.

Table 8-4. Performance Level Distribution as a Function of Grade*

Grade	Number of Students	Novice	Nearing Proficiency	Proficient	Advanced	% Novice	% Nearing Proficiency	% Proficient	% Advanced
5	1,928	801	877	225	25	42	45	12	1
8	1,677	404	1,029	241	3	24	61	14	0
11	1,095	641	288	166	0	59	26	15	0

*Calculations based on those students attempting 5 or more items.

8.7 Differential Item Functioning

The *Code of Fair Testing Practices in Education* (Joint Committee on Testing Practices, 2004) explicitly states that subgroup differences in performance should be examined when sample sizes permit and that actions should be taken to ensure that differences in performance are due to construct-relevant, rather than irrelevant, factors. Chapter 3 of *Standards for Educational and Psychological Testing* (AERA et al., 2014) includes similar guidelines. As part of the effort to identify such problems, BIE Science items were evaluated in terms of DIF statistics.

For items in the SSIB, the standardization DIF procedure (Dorans & Kulick, 1986) was employed to evaluate subgroup differences. The standardization DIF procedure is designed to identify items for which subgroups of interest perform differently, beyond the impact of differences in overall achievement. The DIF procedure calculates the difference in item performance for two groups of students (at a time) matched for achievement on the total test. Specifically, average item performance is calculated for students conditional on scale score. Then an overall average is calculated, weighting by the pooled scale score distribution so that it is the same for the two groups.

When differential performance between two groups occurs on an item (i.e., a DIF index in the “low” or “high” categories, explained below), it may or may not indicate item bias, e.g., caused by construct-irrelevant factors. Other construct-relevant reasons could also lead to DIF, such as course-taking patterns or differences in school curricula. On the other hand, if subgroup differences in performance can be traced to differential experience (such as geographical living conditions or access to technology), the inclusion of such items should be reconsidered.

For items in the SSIB, four subgroup comparisons were evaluated for DIF:

- Male compared with Female
- White compared with Black
- White compared with Hispanic
- White compared with Asian
- White compared with Multiracial
- Not economically disadvantaged status compared with economically disadvantaged

The DIF statistics were calculated based only on the members of the subgroup in question in the computations; values were calculated only for subgroups with 100 or more students. Computed DIF indices have a theoretical range from -1.0 to 1.0 for selected-response items. Dorans and Holland (1993) suggested that index values between -0.05 and 0.05 should be considered negligible. The preponderance of the BIE Science items fell within this range. Dorans and Holland further stated that items with values between -0.10 and -0.05 and those with values between 0.05 and 0.10 (i.e., “low” DIF) should be inspected to ensure that no possible effect is overlooked, and that items with values outside the -0.10 to 0.10 range (i.e., “high” DIF) are more unusual and should be examined very carefully; thus, content experts conducted a review of items flagged for DIF.

No items used in BIE Science Assessments showed high DIF; a handful of items showed low DIF, as indicated in Table 8-5. These results indicate that the content bias reviews for science were conducted thoroughly.

Table 8-5. Number of Items with Low DIF

Subject	Grade	Demographic	DIF	Group	Favor focal low	Number of	Items	Favor reference high
			Reference	Focal		Favor reference low	Favor focal high	
Science	5	ELL	0	1	0	1	0	0
		Gender	M	F	1	4	0	0
		IEP	Y	N	0	4	0	1
	8	ELL	0	1	3	1	0	0
		Gender	M	F	0	3	0	0
		IEP	N	Y	2	0	0	0
	11	Gender	M	F	1	0	0	0

Chapter 9. Score Reliability

9.1 IRT Marginal Reliabilities

IRT marginal reliability estimation is based on applying the standard classical test theory (CTT) formula, relating variances of true score, observed score, and measurement error, in the IRT setting. In CTT, the relationship between these variances is given by:

$$\sigma_X^2 = \sigma_T^2 + \sigma_E^2$$

where σ_X^2 is the observed-score variance,

σ_T^2 is the true-score variance, and

σ_E^2 is the error variance.

Starting from this equation, it can be shown that the formula for CTT reliability can be expressed by:

$$CTT \text{ Reliability} = 1 - \frac{\sigma_E^2}{\sigma_X^2}.$$

IRT marginal reliability is based on extending the CTT model to an IRT framework (Samejima, 1994) and provides an IRT-based estimate of the overall test reliability. Error variance is estimated as the mean squared conditional standard error of measurement (CSEM) of the theta estimates across students within a grade. Observed-score variance is estimated as the variance of the theta estimates across students within a grade. Equivalently, the mean squared CSEM of the scale scores and the variance of the scale scores can be used in place of the CSEM of the theta estimates and the variance of the theta estimates, respectively. IRT marginal reliability is then given by the following formula:

$$IRT \text{ Marginal Reliability} = 1 - \frac{\overline{CSEM(\theta)^2}}{Var(\hat{\theta})} = 1 - \frac{\overline{CSEM(SS)^2}}{Var(SS)},$$

where $\overline{CSEM(\theta)^2}$ is the mean squared CSEM,

$\overline{CSEM(SS)^2}$ is the mean squared scaled CSEM,

$Var(\hat{\theta})$ is the variance of theta estimates, and

$Var(SS)$ is the variance of scaled scores.

Using this formula, IRT marginal reliability estimates were calculated for each grade. The reliability of a test can also be evaluated by simply examining the CSEMs themselves. CSEMs facilitate the interpretation of individual scale scores. With any given scale score estimate for a student, the reasonable limits of the true scale score for the student can be calculated by using the CSEM for the scale score.

Table 9-1 presents descriptive scale score statistics, IRT-based reliability, and mean scale score CSEMs by grade based on BIE student population. As shown in the table, most of the values reached levels associated with adequate reliability (0.70 or higher).

Table 9-1. IRT Marginal Reliability by Grade and Form

Grade	Scale Form	Min	Max	Mean	SD	IRT Marginal Reliability	Mean Scaled CSEM
5	1	519	589	548.68	11.11	0.87	4.05
	2	517	590	547.75	10.62	0.86	3.98
8	1	800	889	852.00	8.19	0.80	3.52
	2	801	883	850.64	8.92	0.82	3.71
11	1	1128	1178	1153.70	6.77	0.81	2.87
	2	1117	1178	1152.69	6.89	0.76	3.25

At the total test level and per grade, IRT marginal reliability ranged from 0.76 to 0.87. Note that IRT marginal reliability is partially dependent upon the variance in scaled scores. When present, range restriction in smaller samples can reduce the variance in scaled scores and therefore reduce the resulting value of IRT marginal reliability. Thus, while not necessarily ideal, the observed values of IRT marginal reliability may have been impacted due to the somewhat small sample sizes from the BIE student population taking the Spring 2025 administration.

While subgroup reliability results are included in Table 9-2 for subgroups with at least 100 students, many of the subgroups have fewer than 100 students per subject and grade. Because the subgroup reliabilities are based on very small samples, no interpretations ought to be made about the adequacy of these subgroup reliabilities.

Given that, the results in Table 9-2 should be interpreted with appropriate levels of caution. Reliabilities are dependent not only on the measurement properties of a test, but also on the statistical distribution of the studied subgroup. Additionally, reliability estimates can be artificially depressed for subgroups with little variability in test scores (Draper & Smith, 1998).

Table 9-2. IRT Marginal Reliability by Grade and Subgroup

Grade	Type	Number of Student	Min	Max	Mean	SD	IRT Marginal Reliability	Mean Scaled CSEM
5	All	1,928	517	590	547.61	10.90	0.86	4.04
	Female	961	517	590	547.62	10.34	0.85	4.01
	Male	967	519	590	547.60	11.43	0.87	4.08
	Currently receiving LEP services	321	521	583	547.84	9.54	0.82	3.98
	Not receiving LEP services	394	517	590	548.42	11.48	0.87	4.05
	Economically Disadvantaged Students	1,840	517	590	547.58	10.88	0.86	4.04
	Non-gifted Students	103	527	572	546.30	10.02	0.83	4.07
	Title1 Students	1,870	517	590	547.55	10.87	0.86	4.04
8	All	1,677	800	889	851.01	8.77	0.81	3.66
	Female	825	800	883	850.73	8.40	0.79	3.66
	Male	852	801	889	851.28	9.11	0.83	3.66
	Currently receiving LEP services	202	832	869	850.23	6.63	0.70	3.60
	Not receiving LEP services	405	827	883	852.58	8.78	0.83	3.54
	Economically Disadvantaged Students	1,625	800	889	851.08	8.80	0.81	3.66
	Gifted Students	116	836	883	858.46	9.69	0.88	3.33
	Non-gifted Students	109	827	877	850.05	7.99	0.78	3.66
11	Title1 Students	1,623	800	889	850.95	8.71	0.81	3.66
	All	1,095	1117	1178	1153.29	6.78	0.79	3.03
	Female	542	1127	1178	1152.91	6.33	0.76	3.03
	Male	553	1117	1178	1153.67	7.18	0.81	3.03
	Currently receiving LEP services	108	1136	1169	1151.71	6.18	0.73	3.12
	Not receiving LEP services	226	1117	1175	1153.56	7.20	0.80	3.04
	Economically Disadvantaged Students	1,052	1117	1178	1153.30	6.77	0.79	3.03
	Title1 Students	1,042	1117	1178	1153.27	6.74	0.79	3.03

9.2 Decision Accuracy and Consistency

While related to reliability, the accuracy and consistency of classifying students into achievement categories are even more important statistics in a standards-based reporting framework (Livingston & Lewis, 1995). After the achievement levels were specified and students were classified into those levels, empirical analyses were conducted to estimate the statistical accuracy and consistency of the classifications.

Accuracy refers to the extent to which decisions based on test scores match decisions that would have been made if the scores did not contain any measurement error. Evaluation of decision accuracy is essential, considering all test scores contain measurement error. Consistency measures the extent to which classification decisions based on test scores match the decisions based on scores from a second, parallel form of the same test. Consistency can be evaluated directly from actual responses to test items if two complete and parallel forms of the test are given to the same group of students. In operational test programs, however, such a design is usually impractical.

However, techniques have been developed to estimate both the accuracy and the consistency of classification decisions based on a single administration of a test. The Rudner (2001, 2005) technique was used for the BIE Science assessments because it can be easily applied to data that are scored in the IRT theta metric or any linear transformation of this metric, such as the scale scores. The applicability of the Rudner technique to IRT-based metrics distinguishes this method from methods based on observed scores, such as the Lewis and Livingston (1995) method.

For details of the Rudner method, refer to Rudner (2001, 2005); given here is a brief review of the basic idea behind the method. Using an examinee's estimated scale score and standard error, assuming a normal probability distribution, the method first calculates for all examinees at a fixed value of true scale score, the expected proportion whose observed scale score is in an interval [a,b]. Then, by summing over all examinees whose true scale scores are in an interval [c,d], the method yields the expected proportion of all examinees whose true scale score is in [c,d] and whose observed scale score is in [a,b]. Setting [a,b] and [c,d] to correspond to the true score intervals defined by the cut scores yields the elements of a classification table that shows the expected proportion of all examinees with observed and true scale scores in each cell. These proportions can then be used to calculate both classification accuracy and classification consistency estimates.

For the classification accuracy tables, cell $[i, j]$ represents the estimated proportion of students whose true scale score fell into classification i (where $i = 1$ to 4, for the four achievement levels) and whose observed scale score fell into classification j (where $j = 1$ to 4). The sum of the diagonal entries (i.e., the proportion of students whose true and observed classifications matched) signified overall accuracy.

Another way to measure consistency is to use κ (kappa; Cohen, 1960), which indicates the proportion of consistent classifications after removing the proportion of consistent classifications that would be expected by chance. It is calculated using the following formula:

$$\kappa = \frac{(\text{Observed agreement}) - (\text{Chance agreement})}{1 - (\text{Chance agreement})} = \frac{\sum_i C_{ii} - \sum_i C_{i.} C_{.i}}{1 - \sum_i C_{i.} C_{.i}},$$

where $C_{i.}$ is the proportion of students whose observed achievement level would be Level i (where $i = 1-3$) on the first hypothetical parallel form of the test;

$C_{.i}$ is the proportion of students whose observed achievement level would be Level i (where $i = 1-3$) on the second hypothetical parallel form of the test; and

C_{ii} is the proportion of students whose observed achievement level would be Level i (where $i = 1-3$) on both hypothetical parallel forms of the test.

Because κ is corrected for chance, its values are lower than are other consistency estimates.

Figure 9-1 shows the overall decision accuracy for science by grade level. More details on decision accuracy and consistency (DAC) are provided in Appendix J. Table J-1 in Appendix J includes overall accuracy and consistency indices, along with kappa. Accuracy and consistency values conditional on performance level are also provided in Table J-1. For these calculations, the denominator is the proportion of students associated with a given performance level. Following is an example from Table J-1, looking at Level 1 for grade 5 for scale form 1.

- The *conditional accuracy* value was 0.85. This indicates that among the students whose true scale scores placed them in Level 1, 85% would be expected to be in this same level again when categorized according to their observed scale scores.
- The *consistency* value was 0.79. This indicates that among the students whose *observed scale scores* placed them in Level 1, 79% would be expected to be in this same level again if a second parallel test form were used.

For some testing situations, the greatest concern may be decisions regarding level thresholds. For example, in testing done for the Every Student Succeeds Act (ESSA) accountability purpose, the primary concern is distinguishing between students who are proficient and those who are not yet proficient. Table J-2 in Appendix J provides accuracy and consistency estimates at each cutpoint, as well as false positive and false negative decision rates. A false positive rate is the proportion of students whose observed scores were above the cut and whose true scores were below the cut. A false negative rate is the proportion of students whose observed scores were below the cut and whose true scores were above the cut.

Figure 9-1a. Overall Decision Accuracy for Science by Grade for Form 1 (Set A)

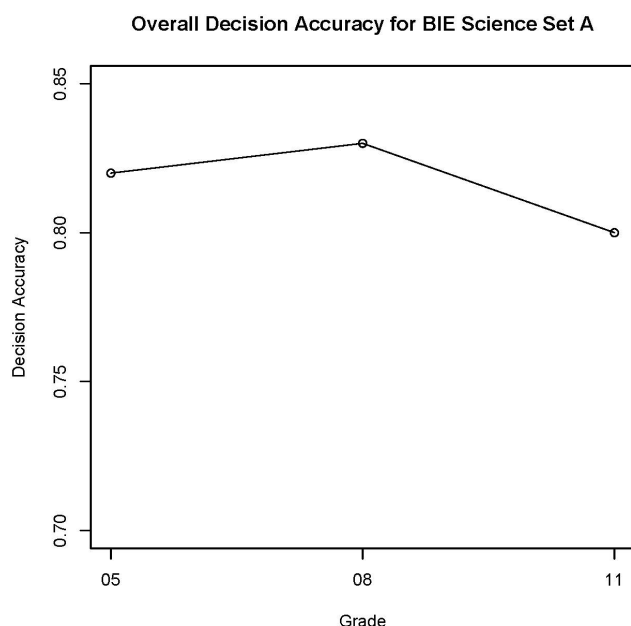
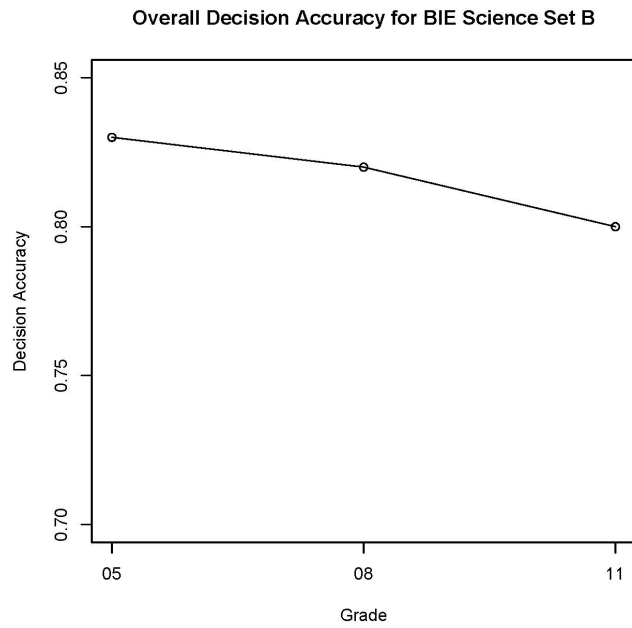


Figure 9-1b. Overall Decision Accuracy for Science by Grade for Form 2 (Set B)



Chapter 10. Score Reporting

10.1 Individual Student Reports

The individual student report (ISR) is distributed to schools and parents. The report is also available in a downloadable format to schools in the Download Center in LightHouse™. Schools can log into the secure portal to access their students' reports. The report is produced in color. One copy of the report is sent home to parents by the school and one copy of the report is retained by the school. Each report is printed duplex on one 8 ½ x 11-inch sheet of paper. The front of the report contains demographic identifying information about the student. It also contains a letter to parents and guardians. A list of resources is presented on the front page as well. A graphical display of the scale scores and achievement levels and achievement level descriptors is shown. No design changes were made from the previous administration.

The back page of the report contains all the students' performance information. Each student's performance on the BIE Science Assessment is described in the ISR. Based on the student's earned score, the student is assigned to one of the following four achievement levels: Novice, Nearing Proficiency, Proficient, or Advanced. The report contains scale scores, achievement levels, standard error of measurement, and reporting category performance indicators. The report also contains the percentage of students in each performance indicator for each reporting category. The student's scale score is compared to the average scale score at the school, district, and overall BIE levels. The report provides a descriptor for the achievement level earned and the level above if the achievement level earned is not Advanced. Examples are provided along with the descriptors. For additional information concerning the individual student report, see Appendix K (BIE Science Reporting Business Requirements) and Appendix L (Score Report Interpretation Quick Guide) for the sample reports.

10.2 Scale Score

A scale score is a numerical value that summarizes student performance. Not all students respond to the same set of test items, so each student's scaled score accounts for the slight differences in difficulty among the various forms and administrations of the test. The resulting scale score allows for an appropriate comparison across test forms and administration years within a grade or course and content area. BIE reports provide overall scale scores for science, which determine a student's achievement level for the content area. The scale score range is XX00–XX90 where XX = the student's tested grade. For example, in grade 5 the range is 500–590.

For example, a student who earns an overall scale score of 800 on one form of the grade 8 assessment would be expected to earn an overall scale score of 800 on any other form of the grade 8 assessment. Furthermore, the student's overall scale score and level of mastery of concepts and skills would be comparable to a student who took the same assessment the previous year or the following year. For cumulative scaled-score distributions, see Appendix M; for scaled score descriptive statistics, see Appendix N.

10.3 Achievement Level

Each BIE achievement level is a broad category that is defined by a student's overall scale score and is used to report overall student performance by describing how well students met the expectations for their grade levels. There are four achievement levels for the Spring 2025 BIE Science Assessments:

Novice. Students demonstrate evidence of **emerging** understanding and use of college and career readiness knowledge, skills, and abilities.

Nearing Proficiency. Students demonstrate evidence of **partial** understanding and use of college and career readiness knowledge, skills, and abilities.

Proficient. Students demonstrate evidence of **satisfactory** understanding and use of college- and career-readiness knowledge, skills, and abilities.

Advanced. Students demonstrate evidence of **thorough** understanding and use of college and career readiness knowledge, skills, and abilities.

These PLDs are referred to as Policy Definitions for reporting BIE performance in science. Range PLDs describe the knowledge and skills that students throughout each proficiency level's range are expected to demonstrate in each grade. For example, in line with the nature of the science standards, the science range PLDs combine science and engineering practices, disciplinary core ideas, and crosscutting concepts that students in grades 5, 8, and 11 are expected to integrate and demonstrate. The range PLDs appear in Appendix B.

10.4 Standards Performance Indicators

Standards performance for BIE assessments is reported using words and colors that indicate whether the student performed above standard, at/near standard, or below standard in each standard. Additional information about standard performance indicators is in the Score Report Interpretation Quick Guide, Appendix L in this document. Performance indicators for each standard are summarized at the overall level and presented as whole number percentages in a horizontal stacked bar that is color coded to represent the performance indicators.

10.5 Comparison of Student Performance at the School, District, and Overall Level

The last section of the back page is a table containing the Achievement levels, Achievement Level Descriptors, scale score range for each achievement level and the student's overall science scale score on a bar graph compared to the scale score average at the school, district, and overall (BIE) levels. The student's scale score bar is colored orange. The other bars are colored in the same shade of blue.

10.6 Student Roster

A student roster is produced for each participating school. The roster lists all students in the reporting dataset. The student is either reported as tested or with a not-tested reason. See the Reporting Business Requirements document for more information on participation statuses. The Roster is printed duplex in color on 8 ½- x 11-inch sheets of paper. The report is oriented landscape. The report lists the following information about each student:

Student Name

Student's NASIS ID

Student's Grade

Student's Gender

Student's Testing Status

Student's Scale Score

Student's Achievement Level

The achievement level is named and has a background color that matches the color coding on the ISR.

The report is marked as confidential. The roster is printed and shipped to the schools.

Chapter 11. Validity Arguments to Support Intended Score Interpretation and Uses

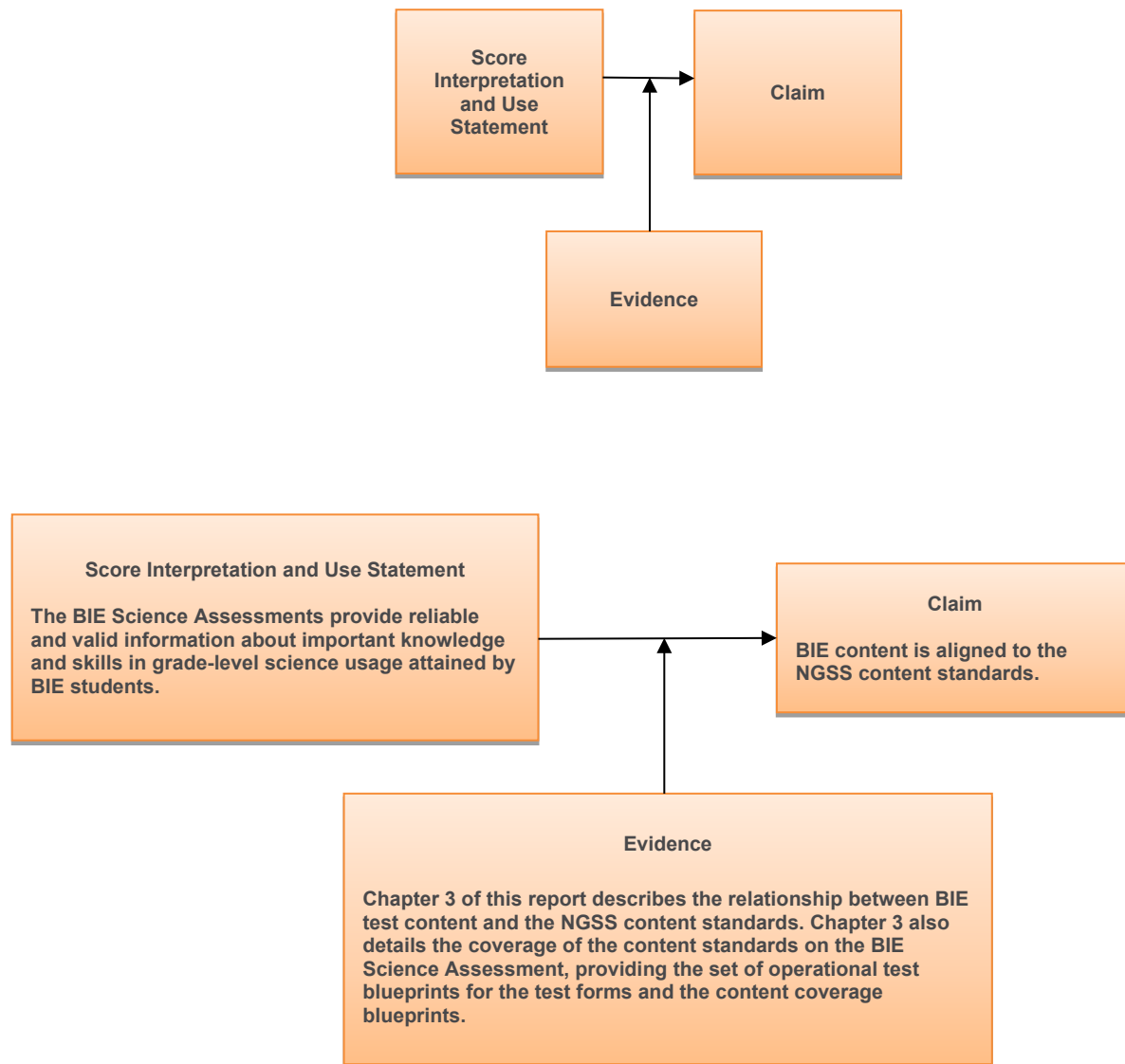
This chapter presents the primary intended score interpretation and three primary intended score uses. This chapter also presents the claims and subclaims that underlie these four score interpretations and uses (SIUs) and the evidence that supports the claims and subclaims. The BIE validity argument model is introduced and applied to develop validity arguments to support the four SIUs.

The *Standards for Educational and Psychological Testing* (2014) defines validity as “the degree to which evidence and theory support the interpretations of test scores for proposed uses of tests” (p. 11). Elaborating on that definition, the *Standards* asserts that “it is the interpretations of test scores for proposed uses that are evaluated, not the test itself” (p. 11) and that “validation logically begins with an explicit statement of the proposed interpretation of test scores, along with a rationale for the relevance of the interpretation to the proposed use” (p. 11). This definition applies specifically to intended interpretations and uses of test scores, rather than to the broader program of curriculum and instruction in which a testing program is embedded or to the surrounding education and school improvement policies and aspirations for student learning.

Further, the *Standards* states that “a sound validity argument integrates various strands of evidence into a coherent account of the degree to which existing evidence and theory support the intended interpretations of test scores for specific uses” (p. 21). We use these views in the *Standards*, that evidence must be used to support score interpretation and use claims, as the basis for the BIE validity argument model, which we describe next.

Emerging practice in state assessment programs is to construct validity arguments based on Toulmin’s model of argumentation (Toulmin, 1958), Chapelle’s proposed practice-oriented adaptation (2021), and Kane’s formulation of validity arguments (2013). A model for BIE validity arguments, derived from these three conceptualizations, is shown in Figure 11-1. The first panel shows the BIE model; the second panel is an illustration for a validity argument for a score interpretation and use statement.

Figure 11-1. BIE Validity Argument Model



Adapted from Chapelle (2021) Figures 2.1-2.3, Kane (2013) Figure 1, and Toulmin (1958).

The first panel in Figure 11-1, displays a generic representation of a BIE SIU statement and how it connects directly to a corresponding validity argument and claims and subclaims, shown at the right. The evidence that supports the BIE SIUs, claims, and validity arguments connects directly to the claims- validity argument pairings. The second panel displays a specific example of this relationship.

In the BIE validity argument model, the overall validity argument is that the existing design, procedural, and psychometric evidence supports the intended score interpretations and uses. Each of the interpretations and uses represents a set of claims, subclaims, and validity arguments that require supporting evidence to connect the evidence to the score interpretation and use. This line of reasoning and argumentation creates supported validity arguments. The remaining sections of this chapter describe the a) intended SIUs; b) claims, subclaims, and validity arguments, which connect the BIE design, procedural, psychometric, and other program information to the SIUs; and c) evidence that supports each SIU and validity argument (for which detail is provided in the previous chapters).

The relationships among the score interpretations and uses, claims, and validity arguments appear in Table 11-1. Each entry in the table is presented following the table, with descriptions and summaries of the supporting evidence.

Table 11-1. Relationships Among Score Interpretations and Uses, Claims, and Sub-Claims, and Supporting Evidence

Claims	Claims and Subclaims that Support Score Interpretations and Uses
SIU 1: Primary Intended Score Interpretation	
The BIE Science Assessment provides reliable and valid information about important knowledge and skills in grade-level science attained by BIE general education students.	
Claim 1.1: The content of the tests represents the content of the standards.	<ul style="list-style-type: none"> 1.1.1. BIE content is aligned to Next Generation Science Standards. 1.1.2. BIE items are aligned to Next Generation Science Standards.
Claim 1.2: The test items are construct-relevant.	<ul style="list-style-type: none"> 1.2.1. Items require application of the knowledge, skills, and abilities (KSAs) of the targeted construct. 1.2.2. Items are free of bias and sensitivity issues. 1.2.3. Students' cognitive skills and processes match those identified in the construct domains for all students and for each subgroup.
Claim 1.3: Test scores on the BIE Science Assessments provide reliable information about student performance and accurate classifications into performance levels.	<ul style="list-style-type: none"> 1.3.1. BIE scores and performance level categorizations are adequately reliable for their intended purpose. 1.3.2. Item characteristics support intended interpretations about all students who take the BIE Science Assessment. 1.3.3. Test characteristics support intended interpretations about all students who take the BIE Science Assessment.
Claim 1.4: Item and test scoring are implemented accurately.	<ul style="list-style-type: none"> 1.4.1. Machine-scored items were scored accurately. 1.4.2. Constructed-response item scoring training and monitoring procedures met industry standards.
SIU 2: Intended Score Use for Individual Students	
Claim 2.1: Educators, schools, and district administrators can use results from the BIE Science Assessment to describe and monitor student achievement status with respect to mastery of the content standards.	<ul style="list-style-type: none"> 2.1.1. BIE test scores and performance level categorizations of individual students are adequately reliable and valid measures of student achievement status with respect to mastery of the content standards.
SIU 3: Intended Score Use for Groups of Students	
Claim 3.1: Educators can use results from the BIE Science Assessment to support instructional planning for groups of students.	<ul style="list-style-type: none"> 3.1.1. Teachers find the performance level descriptors and their students' performance levels useful for planning instruction, especially for students whose test scores fall within performance levels 1 and 2. 3.1.2. Teachers find their students' scale score information useful for planning instruction, especially for students whose test scores fall within performance levels 1 and 2.
Claim 3.2: Schools, districts, and state-level stakeholders can use results from the BIE Science Assessment to make comparisons between organizations (e.g., schools, districts).	<ul style="list-style-type: none"> 3.2.1. BIE scores and performance levels for groups of students are adequately reliable and valid to enable school, district, and state leaders to monitor changes in means, standard deviations, and performance level percentages for classroom, school, district, and state groups. 3.2.2. BIE scores and proficiency level categorizations of groups of students are adequately reliable and valid to enable monitoring of grade-level performance and student-cohort performance.

Evidence that supports SIUs and claims in BIE validity arguments is summarized below, using a relevance-rating scale, with rating levels defined in Table 11-2.

Table 11-2. Relevance and Completeness or Completeness of Evidence in Support of SIUs and Claims Underlying Validity Arguments for BIE Score Interpretations and Uses

Complete Evidence	When all required pieces of relevant evidence are provided to support a validity argument
Moderate to Substantial Evidence	When several pieces of relevant evidence are provided, but not all required pieces of evidence are provided
Limited Evidence	When only one or two pieces of evidence are provided, where the evidence may be only marginally relevant, or where more than one or two pieces of evidence are required
No Evidence	When no relevant evidence exists

11.1 Primary Intended Score Interpretation

The primary intended score interpretation for the BIE Science Assessment (SIU 1) states that the BIE Assessments provide reliable and valid information about important knowledge and skills in grade-level science usage attained by general education students.

Claim 1.1. The content of the tests represents the content of the standards.

Items used on the BIE Science Assessment are developed to measure achievement on the Next Generation Science Standards in the assessed content areas. Cognia content specialists, in collaboration with the BIE, ensure this alignment, and ongoing independent evaluations are held to verify alignment. In addition, independent reviews are scheduled to ensure that items and passages conform to bias and sensitivity guidelines.

Subclaim 1.1.1. BIE Science Assessment content is aligned to Next Generation Science Standards.

Evidence: Chapter 3 of this report describes the relationship between BIE test content and the Next Generation Science Standards. Chapter 3 also details the coverage of the content standards on the BIE Science Assessment, providing the set of operational test blueprints for test forms and the content coverage blueprints.

The BIE Science Assessment uses Cognia Secure Science Item Bank (SSIB) and has the exact same test design as well as uses many of the same items on the test as the other state that had the alignment study. Theoretically, the alignment study results on the other state should still hold for the BIE Science Assessment and no additional alignment study is needed.

Summary of evidence: Complete evidence.

Subclaim 1.1.2. BIE Science Assessment items are aligned to the Next Generation Science Standards.

Evidence: Chapter 4 describes the item specifications and standardized item writer training in support of new item development. Chapter 4 also details the item review process performed by item review committees, to ensure item content alignment to the intended content standard. Content alignment for the BIE science item bank and test forms has been verified by an independent alignment study in other states. Therefore, an additional alignment study that is specific to BIE assessments is not needed.

Summary of evidence: Complete evidence.

Claim 1.2: The test items are construct-relevant.

Subclaim 1.2.1. Items require the application of the KSAs of the targeted construct.

Evidence: The 2025 operational items are aligned with the Next Generation Science Standards. The evidence for element 1.2.1 is directly linked to the subclaims 1.1.1 and 1.1.2 above.

Summary of evidence: Complete evidence.

Subclaim 1.2.2. Items are free of bias and sensitivity issues.

Evidence: During the item development process, the items follow a rigorous development cycle that includes reviews by subject matter experts and by Item Content and Bias and Sensitivity panelists. The item development process also includes data reviews, during which item-level statistics—including differential item functioning (DIF) statistics—are reviewed. See Chapter 8 for a detailed description of the item review process.

Summary of evidence: Complete evidence, based on current Cognia procedures for the Spring 2025 testing season.

Subclaim 1.2.3. Students' cognitive skills and processes match those identified in the construct domains for all students and for each subgroup.

Evidence: Cognitive-process evidence examines the extent to which students' cognitive skills and processes match those identified in the construct domains defined by test developers for all students and for each subgroup.

Evidence of cognitive processes can be obtained through subject matter experts (including educators) review or a study. The items for BIE Science (Cognia's Science Secure Item Bank) were developed and internally reviewed by a team of science content specialists, and then the items were reviewed by state/entity department of education content and/or assessment specialists and science educators from the field. Their input helped ensure the tasks would measure the intended cognitive ability. In addition to Cognia content expert reviews of the items for alignment, content accuracy, and cognitive complexity levels, science educators were involved in the review of all items developed and field tested. When reviewing items in committee review meetings, target students' characteristics (such as various learning targets, use of the different science dimensions, and the degree of cognitive processing elicited by the stimulus and application of each dimension) were all discussed and confirmed. The feedback from educators/specialists with expertise with science content in these committee meetings indicated that the test content did not require extraneous cognitive processes for engagement.

Summary of evidence: Complete evidence, based on current Cognia procedures for the Spring 2025 testing season.

Claim 1.3: Test scores on the BIE Science Assessment provide reliable information about student performance and accurate classifications into performance levels.

Subclaim 1.3.1. BIE scores and performance level categorizations are adequately reliable for their intended purpose.

Evidence:

Score Reliability: Score reliability is supported through multiple sources of evidence. Chapter 9 provides a description of IRT reliability theory, interpretation guidelines and the

relevant equations. Table 9-1 contains the reliability estimates by grade, while Table 9-2 also contains reliability results disaggregated by student subgroups. These reliability estimates are consistent with industry standards, which can be observed in technical reports posted online by other state assessment programs. In addition, Chapter 6 details the scoring processes implemented to ensure score quality across item types. For open-ended items requiring hand-scoring, double-blind scoring procedures were employed to minimize scorer bias and enhance score reliability. Interrater reliability results presented in Table 6-2, further confirm that scoring consistency was effectively maintained.

Scale Score Standard Errors: Chapter 8 provides a description of calculation and interpretation of the scaled scores and Chapter 9 provides a description of the calculation of the standard error for a scaled score. The average standard error for reported scaled scores is reported in Table 9-1. The scale score standard error can be compared to the scale score range and the scale score standard deviation to provide some context for interpretation. These standard error estimates are consistent with industry standards, which can be observed in technical reports posted online by other state assessment programs.

Decision Consistency and Accuracy Estimates: Decision accuracy is an estimate of the probability that the observed classification is the true classification. Decision consistency is an estimate of the probability that students would receive the same classification if they tested twice on parallel forms. Chapter 9 describes the theory and equations underlying the estimation of classification accuracy and consistency. Decision accuracy and consistency results are provided in Appendix J. These decision consistency and accuracy estimates are consistent with industry standards, which can be observed in technical reports posted online by other state assessment programs.

Summary of evidence: Complete evidence.

Subclaim 1.3.2. Item characteristics support intended interpretations about all students who take the BIE Science Assessment.

Evidence: The psychometric characteristics most pertinent to evaluating the adequacy of individual items are the estimated item parameters. The item parameter estimates are provided in Appendix G. For dichotomously scored items, the item parameters include the discrimination, difficulty, and lower asymptote parameters. For polytomously scored items, the item parameter estimates include the discrimination, location, and item-category parameters. All items undergo statistical analyses at the time of field testing, including classical, DIF, and IRT analyses. As stated in Chapter 4, the results of these analyses are reviewed in Data Review meetings with national subject matter experts. After field testing and prior to operational administration, items from the previous operational administration are reviewed for their item information function (IIF) contributions at the performance level cuts to evaluate and rate the quality of each item. After each operational administration, dimensionality analyses are also conducted to determine the adequacy of the unidimensional IRT model used for scaling, equating, and scoring the BIE students.

Summary of evidence: Complete evidence.

Subclaim 1.3.3. Test characteristics support intended interpretations about all students who take the BIE Science Assessment.

Evidence: High correlations (e.g., greater than or equal to 0.7) among content area subdomain indicators and the relatively low reliability of these indicators demonstrate that such indicators

must be interpreted and used cautiously, and in conjunction with other information about student achievement and learning needs in these areas.

Dimensionality: Dimensionality analysis was conducted on each grade-level test. Chapter 7, section 7.2, provides a detailed description of the dimensionality hypothesis testing and effect-size estimation methods and provides dimensionality results. No statistically significant violations of local independence were noted.

Conditional Standard Errors of Measurement: Chapter 9 provides a detailed description of the psychometric model that was fitted to the data, the test information function (TIF), and the inverse transformation of the TIF into the Conditional Standard Error of Measurement (CSEM). The TIF and CSEM are inverse transformations of each other. Whereas the TIF indicates test score precision, the CSEM indicates the converse, i.e., test score imprecision or measurement error. The TIF and its analogue, the CSEM, are the most pertinent products of the psychometric model in evaluating the adequacy of a test (form). By examining the value of CSEM at each of the performance level cut scores, the psychometric appropriateness and accuracy of each test can be evaluated.

Content Coverage: Subclaims 1.1.1, 1.1.2, and 1.2.1 above detail the evidence in support of the content coverage and the alignment of the content to the BIE standards.

Scoring: Subclaims 1.4.1 and 1.4.2 detail the evidence in support of accurate item and test scores.

Summary of evidence: Complete evidence.

Claim 1.4: Item and test scoring are implemented accurately.

Subclaim 1.4.1. Machine-scored items were scored accurately.

Evidence: As described in section 6.2.1 of Chapter 6 and in Chapter 7, a classical item analysis on the set of machine-scored items is performed prior to scaling and equating. This ensures that for each machine-scored item, the response designated as the correct response was indeed the correct response.

Summary of evidence: Complete evidence.

Subclaim 1.4.2. Constructed-response item scoring training and monitoring procedures met industry standards.

Evidence: As detailed in Chapter 6, scorer recruitment, training, and qualification and scoring-monitoring procedures follow industry standards. Section 6.2.3, Scoring of Open-Ended Response Items, describes all the procedures that are used to ensure the accuracy of the scoring for the open-ended (constructed) response items, including administrator training and monitoring, benchmarking and identification of scoring materials, scorer recruitment and qualifications, scoring leadership, qualification, specific scoring rules to ensure accuracy, monitoring of quality control, quality reports, and interrater reliability.

Summary of evidence: Complete evidence.

11.2 Primary Intended Score Uses

11.2.1 Intended Score Use for Individual Students

Claim 2.1: Educators, schools, and district administrators can use results from the BIE Science Assessment to describe and monitor student achievement status with respect to mastery of the content standards.

Subclaim 2.1.1. BIE test scores and performance level categorizations of individual students are adequately reliable and valid measures of student achievement status with respect to mastery of the content standards.

Evidence:

Scaled Score Standard Errors: Chapter 8 provides a description of calculation and interpretation of the scaled scores and Chapter 9 provides a description of the calculation of the standard error for a scaled score. The average standard error for reported scaled scores is reported in Appendix I. The scaled score standard error can be compared to the scaled score range and the scaled score standard deviation to provide some context for interpretation.

Decision Consistency and Accuracy Estimates: Decision accuracy is an estimate of the probability that the observed classification is the true classification. Decision consistency is an estimate of the probability that students would receive the same classification if they tested twice on parallel forms. Chapter 9 describes the theory and equations underlying the estimation of classification accuracy and consistency. Decision accuracy and consistency results are provided in Appendix J.

Content Coverage: Subclaims 1.1.1, 1.1.2, and 1.2.1 above detail the evidence in support of the content coverage and the alignment of the content to the BIE standards.

Scoring: Subclaims 1.4.1 and 1.4.2 detail the evidence in support of accurate item and test scores.

Summary of evidence: Complete evidence.

11.2.2 Intended Score Use for Groups of Students

Claim 3.1: Educators can use results from the BIE Science Assessment to support instructional planning for groups of students.

Subclaim 3.1.1. Teachers find the performance level descriptors and their students' performance levels useful for planning instruction, especially for students whose test scores fall within performance levels 1 and 2.

Evidence: None.

Summary of evidence: No evidence. An example of a source of evidence could be a survey of teachers to begin to understand the degree to which teachers find BIE performance level descriptors and their students' performance levels useful for planning instruction.

Subclaim 3.1.2. Teachers find their students' scale score information useful for planning instruction, especially for students whose test scores fall within performance levels 1 and 2.

Evidence: None.

Summary of evidence: No evidence. An example of a source of evidence could be a survey of teachers to begin to understand the degree to which teachers find BIE scores useful for planning instruction.

Claim 3.2: Schools, districts, and state-level stakeholders can use results from the BIE Science Assessment to make comparisons between organizations (e.g., schools, districts).

Subclaim 3.2.1. BIE scores and performance levels for groups of students are adequately reliable and valid to enable school, district, and state leaders to monitor changes in means, standard deviations, and performance level percentages for classroom, school, district, and state groups.

Evidence: Evidence for the reliability and validity of the scores and the corresponding scoring processes is presented above under Claim 1.3, which cites Chapter 6 on scoring, Chapter 8 on IRT scaling and equating, and Chapter 9 on IRT reliability and decision accuracy and consistency. The reliability of aggregated scores (e.g., means) is typically as high as or higher than individual score reliabilities (e.g., Brennan, 1995). Appendix J contains the decision accuracy and consistency results for the overall test as well as by performance level and by cut score. Subclaims 1.1.1, 1.1.2, and 1.2.1 above detail the evidence in support of the content coverage and the alignment of the content to the BIE standards. Subclaims 1.4.1 and 1.4.2 detail the evidence in support of accurate item and test scores.

Summary of evidence: Moderate to substantial evidence.

Subclaim 3.2.2. BIE scores and proficiency level categorizations of groups of students are adequately reliable and valid to enable monitoring of grade-level performance and student-cohort performance.

Evidence: Evidence for the reliability and validity of the scores and the corresponding scoring processes is presented above under Claim 1.3, which cites Chapter 6 on scoring, Chapter 8 on IRT scaling and equating, and Chapter 9 on IRT reliability and decision accuracy and consistency. The reliability of aggregated scores (e.g., means) is typically as high as or higher than individual score reliabilities (e.g., Brennan, 1995). Appendix J contains the decision accuracy and consistency results for the overall test as well as by performance level and by cut score. Subclaims 1.1.1, 1.1.2, and 1.2.1 above detail the evidence in support of the content coverage and the alignment of the content to the BIE standards. Subclaims 1.4.1 and 1.4.2 detail the evidence in support of accurate item and test scores.

Summary of evidence: Moderate to substantial evidence.

11.3 Conclusions and Next Steps

The majority of the claims and subclaims that support the four claims—that is, the primary intended score interpretations and three intended score uses—are supported by solid evidence. These claims and subclaims and their supporting evidence comprise the validity arguments for BIE scores. Table 11-3 summarizes the relevance ratings for each claim and subclaim.

Table 11-3 indicates the following:

Primary Score Intended Score Interpretation

Of the three claims and nine subclaims that support the intended score interpretation, eight sets of evidence are complete, and one set of evidence is moderate to substantial.

Intended Score Use for Individual Students

The one claim has one supporting subclaim that is moderate to substantial.

Intended Score Use for Groups of Students

Of the two claims and four supporting subclaim evidence, two sets of evidence are moderate to substantial and two subclaims do not currently have evidence.

Table 11-3. Status of Evidence for All SIUs, Claims, and Subclaims

SIUs, Claims, and Subclaims	Relevance of the Evidence to the Validity Argument			
	No Evidence Exists Currently	Limited	Moderate to Substantial	Complete
SIU 1: Primary Intended Score Interpretation				
The BIE Science Assessment provides reliable and valid information about important knowledge and skills in grade-level science usage attained by general education students.				
1.1.1. BIE content is aligned to Next Generation Science Standards.				X
1.1.2. BIE items are aligned to Next Generation Science Standards.				X
1.2.1. Items require application of the KSAs of the targeted construct.				X
1.2.2. Items are free of bias and sensitivity issues.				X
1.2.3. Students' cognitive skills and processes match those identified in the construct domains for all students and for each subgroup.				X
1.3.1. BIE scores and performance level categorizations are adequately reliable for their intended purpose.				X
1.3.2. Item characteristics support intended interpretations about all students who take the BIE Science Assessment.				X
1.3.3. Test characteristics support intended interpretations about all students who take the BIE Science Assessment.				X

continued

SIUs, Claims, and Subclaims	Relevance of the Evidence to the Validity Argument			
	No Evidence Exists Currently	Limited	Moderate to Substantial	Complete
1.4.1. Machine-scored items were scored accurately.				X
1.4.2. Constructed-response item scoring training and monitoring procedures met industry standards.				X
SIU 2: Intended Score Use for Individual Students				
Scale scores can be used to compare an individual student's performance to the performance of other students in BIE.				
2.1.1. BIE test scores and performance level categorizations of individual students are adequately reliable and valid measures of student achievement status with respect to mastery of the content standards.				X
SIU 3: Intended Score Use for Groups of Students				
SIU statements for groups of students are applicable to aggregate reporting of student subgroups (e.g., English learners, students with disabilities, racial/ethnic subgroups) within those levels of aggregation.				
3.1.1. Teachers find the performance level descriptors and their students' performance levels useful for planning instruction, especially for students whose test scores fall within performance levels 1 and 2.	X			
3.1.2. Teachers find their students' scale score information useful for planning instruction, especially for students whose test scores fall within performance levels 1 and 2.	X			
3.2.1. BIE scores and performance levels for groups of students are adequately reliable and valid to enable school, district, and state leaders to monitor changes in means, standard deviations, and performance level percentages for classroom, school, district, and state groups.			X	
3.2.2. BIE scores and proficiency level categorizations of groups of students are adequately reliable and valid to enable monitoring of grade-level performance and student-cohort performance.			X	

11.3.1 Research Agenda

The Score Card ratings provide a road map for a research agenda for the BIE science program. Specifically, the BIE and Cognia can work together to identify the highest priority claims and subclaims for which *No Evidence Exists Currently* and the evidence is *Limited* and plan studies to gather relevant evidence and strengthen validity arguments. This will be a topic of discussion and planning for more immediate and longer-term efforts during the 2025–2026 school year.

References

- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for educational and psychological testing*. Washington, DC: American Educational Research Association.
- Baker, F. B., & Kim, S. H. (2004). *Item response theory: Parameter estimation techniques* (2nd ed.). New York: Marcel Dekker, Inc.
- Brennan, R. L. (1995). The conventional wisdom about group mean scores. *Journal of Educational Measurement*, 32(4), 385–396.
- Brown, F. G. (1983). *Principles of educational and psychological testing* (3rd ed.). Fort Worth: Holt, Rinehart, and Winston.
- Chapelle, C. A. (2021). *Argument-based validation in testing and assessment*. Sage.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20, 37–46.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16, 297–334.
- Dorans, N. J., & Kulick, E. (1986). Demonstrating the utility of the standardization approach to assessing unexpected differential item performance on the Scholastic Aptitude Test. *Journal of Educational Measurement*, 23, 355–368.
- Dorans, N. J., & Holland, P. W. (1993). DIF detection and description: Mantel-Haenszel and standardization. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 35–66). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Draper, N. R., & Smith, H. (1998). *Applied regression analysis* (3rd ed.). New York, NY: John Wiley and Sons.
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Boston, MA: Kluwer Academic Publishers.
- Hambleton, R. K., Swaminathan, H., & Rogers, J. H. (1991). *Fundamentals of item response theory*. Newbury Park, CA: Sage.
- Hambleton, R. K., & van der Linden, W. J. (1997). *Handbook of modern item response theory*. New York, NY: Springer-Verlag.

- Joint Committee on Testing Practices. (2004). *Code of fair testing practices in education*. Washington, D.C.: National Council on Measurement in Education.
- Kane, M. T. (2013). Validating the interpretations and uses of test scores. *Journal of Educational Measurement*, 50(1), 1–73.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32, 179–197.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theories of meta test scores*. Reading, MA: Addison-Wesley.
- Persuasiveness. (2019). In *Vocabulary.com* Retrieved from <https://www.lexico.com/en/definition/persuasive>
- Plausibility. (2019). In *Lexico.com* Retrieved from <https://www.lexico.com/en/definition/plausibility>
- Relevance. (2019). In *Lexico.com* Retrieved from <https://www.lexico.com/en/definition/relevance>
- Roussos, L. A., & Ozbek, O. Y. (2006). Formulation of the DETECT population parameter and evaluation of DETECT estimator bias. *Journal of Educational Measurement*, 43, 215–243.
- Rudner, L. M. (2001). Computing the expected proportions of misclassified examinees. *Practical Research & Evaluation*, 7 (14). Retrieved from <http://PARE.online.net/getvn.asp?v=7&n=14>
- Rudner, L. M. (2005). Expected classification accuracy. *Practical Research & Evaluation*, 10 (13). Retrieved from <http://pareonline.net/pdf/v10n13.pdf>
- Samejima, F. (1994). Estimation of reliability coefficients using the test information function and its modifications. *Applied Psychological Measurement*, 18, 229–244.
- Stout, W. F. (1987). A nonparametric approach for assessing latent trait dimensionality. *Psychometrika*, 52, 589–617.
- Stout, W. F., Froelich, A. G., & Gao, F. (2001). Using resampling methods to produce an improved DIMTEST procedure. In A. Boomsma, M. A. J. van Duign, & T. A. B. Snijders (Eds.), *Essays on item response theory* (pp. 357–375). New York, NY: Springer-Verlag.
- Toulmin, S. (1958). *The uses of argument*. Cambridge, UK: Cambridge University Press.
- Zhang, J., & Stout, W. F. (1999). The theoretical DETECT index of dimensionality and its application to approximate simple structure. *Psychometrika*, 64, 213–249.

Appendices

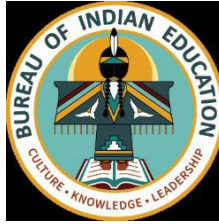
APPENDIX A

LIST OF ACRONYMS

Common Acronyms Used in Assessment Reports	
3PL	three-parameter logistic
AERA	American Educational Research Association
APA	American Psychological Association
BIE	Bureau of Indian Education
CBT	computer-based test
CR	constructed response items
CSEM	conditional standard error of measurement
CTT	classical test theory
DAC	decision accuracy and consistency
DETECT	Dimensionality Evaluation to Enumerate Contributing Traits
DIF	differential item functioning
DIMTEST	computer program used by Cognia
DOK	depth of knowledge
ESSA	Every Student Succeeds Act
ETS	Engineering, Technology, and the application of science
GRM	Graded-Response Model
ICC	item characteristic curve
ICCC	Item Category Characteristic Curve
ICTC	Item Category Threshold Curve
IIF	item information function
IRT	item response theory
ISR	individual student report
KSA	knowledge, skills, and abilities
LEP	limited English proficiency
MS	machine scored items
NCME	National Council on Measurement in Education
OE	open-ended items
PADDI	Principled Assessment Design, Development, and Implementation
PBT	paper-based test
PE	performance expectation
PLD	performance level descriptor
SEM	standard error of measurement
SIU	score interpretation and use
SR	selected response items
SSIB	Cognia's Summative Science Item Bank
STC	School Test Coordinator
STL	Scoring Team Leader
TA	Test Administrator
TAM	Test Administrator's Manual
TCM	Test Coordinator's Manual
TCC	test characteristic curve
TIF	test information function

APPENDIX B

PERFORMANCE LEVEL DESCRIPTORS



BIE Science Proficiency Level

Descriptors Grade 5

Policy PLDs

Policy PLDs define the knowledge and skill level expectations for all grades and content areas for the BIE Science Assessment.

Level 4. Advanced

Students demonstrate evidence of **thorough** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 3. Proficient

Students demonstrate evidence of **satisfactory** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 2. Nearing Proficiency

Students demonstrate evidence of **partial** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 1. Novice

Students demonstrate evidence of **emerging** understanding and use of college and career readiness knowledge, skills, and abilities.

Range PLDs

Range PLDs describe the knowledge and skills that students throughout the range of each proficiency level are expected to be able to demonstrate. In line with the nature of the science standards, the statements combine science and engineering practices, disciplinary core ideas, and crosscutting concepts that students are expected to integrate and demonstrate.

Advanced

Students at the **Advanced** level demonstrate evidence of thorough understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. In addition to demonstrating the skills and understandings at the Proficient level, students performing at the Advanced level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of thorough understanding and use of the NGSS Standards:

- Develop, use, and analyze a model to describe and explain phenomena using an understanding of matter as tiny particles, describe quantities that should be measured to explain phenomena using an understanding of conservation of matter during physical and chemical changes, describe observations and measurements that can be used to identify materials based on their properties, and plan and conduct an investigation to determine whether a new substance with different properties is formed when two substances are mixed. (PS1)
- Plan and conduct an investigation to provide multiple pieces of evidence about phenomena using an understanding of the effects of balanced and unbalanced forces on the motion of an object, predict the future motion of an object based on complex patterns in observations and measurements, ask detailed questions to describe phenomena using an understanding of cause and effect of electric and magnetic interactions between objects not in contact with each other, thoroughly define a simple design problem that can be solved using magnets, and support an argument with multiple pieces of evidence about phenomena using an understanding that the gravitational force of Earth on objects is directed down. (PS2)
- Construct an explanation supported by multiple pieces of evidence about phenomena using an understanding of the relationship between the speed and energy of an object; provide and analyze evidence that energy can be transferred from place to place; predict and explain outcomes for the changes in energy that occur when objects collide; thoroughly design, test, and refine a device that converts energy from one form to another; and use models to explain phenomena using an understanding that food energy was once energy from the Sun. (PS3)
- Develop models to explain phenomena using an understanding that waves can cause objects to move and that light allows objects to be seen; and compare

and explain multiple solutions that use patterns to transfer information. (PS4)

- Develop models to explain phenomena using an understanding of the diversity and commonalities of the life cycles of organisms, construct an argument supported by multiple pieces of evidence about phenomena using an understanding that plants and animals have internal and external structures that support life functions, explain phenomena using an understanding that animals receive, process, and respond to information from their senses, and support an argument with multiple pieces of evidence about phenomena using an understanding that plants get the materials they need for growth chiefly from air and water. (LS1)
- Construct an argument supported by multiple pieces of evidence about phenomena using an understanding that some animals form groups that help members survive; and develop a model to explain phenomena using an understanding of the movement of matter among plants, animals, decomposers, and the environment. (LS2)
- Analyze and interpret data to provide multiple pieces of evidence about phenomena using an understanding that plants and animals have inherited traits and that variation of these traits exists in groups of similar organisms; and support an explanation with multiple pieces of evidence about phenomena using an understanding that traits can be influenced by the environment. (LS3)
- Analyze and interpret fossil data to provide multiple pieces of evidence of organisms and the environments in which they lived; construct an explanation supported by multiple pieces of evidence of phenomena using an understanding that variation among individuals of the same species is a survival advantage and that in a particular habitat some organisms survive well, some survive less well, and some cannot survive at all; and make a claim supported by multiple pieces of evidence about the merit of a solution to a problem caused by changes to the environment and the types of plants and animals that live there. (LS4)
- Explain phenomena using an understanding that patterns in rock formations and fossils in rock layers provide evidence to support an explanation for changes in a landscape over time, support an argument about phenomena using multiple pieces of evidence and an understanding that differences in the apparent brightness of the Sun compared to other stars is due to their relative distances from Earth, and represent and explain data in graphical displays to reveal patterns of daily changes in shadows, day and night, and the seasonal appearance of stars. (ESS1)
- Represent data to explain typical seasonal weather conditions, combine and synthesize information to describe climates in different regions of the world, provide multiple pieces of evidence for phenomena using an understanding of the effects of weathering and the rate of erosion, analyze and interpret data from maps to explain patterns of Earth's features, develop a model to explain

phenomena using an understanding of how multiple systems on Earth interact, and describe and graph the percentages of water to provide multiple pieces of evidence about the distribution of water on Earth. (ESS2)

- Make a claim supported by multiple pieces of evidence about the merit of a design solution that reduces the impacts of a weather-related hazard, combine and synthesize information to explain that energy and fuels are derived from natural resources, how their uses affect the environment, and ways individual communities protect Earth's resources and the environment, and generate and compare multiple solutions to reduce the impacts of several natural Earth processes on humans. (ESS3)
- Define a simple design problem including detailed criteria for success and constraints; generate and compare multiple detailed solutions to a problem; and plan and carry out fair tests to identify more than one way to improve a model or prototype. (ETS1)

Proficient

Students at the **Proficient** level demonstrate evidence of satisfactory understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. In addition to demonstrating the skills and understandings at the Nearing Proficiency level, students performing at the Proficient level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of satisfactory understanding and use of the NGSS Standards:

- Develop and use a model to describe phenomena using an understanding of matter as tiny particles, measure and graph quantities to describe phenomena using an understanding of conservation of matter during physical or chemical changes, make observations and measurements to identify materials based on their properties, and conduct an investigation to determine whether a new substance with different properties is formed when two substances are mixed. (PS1)
- Plan and conduct an investigation to provide one piece of evidence about phenomena using an understanding of the effects of balanced and unbalanced forces on the motion of an object, predict the future motion of an object based on patterns in observations and measurements, ask questions to describe phenomena using an understanding of cause and effect of electric or magnetic interactions between objects not in contact with each other, define a simple design problem that can be solved using magnets, and support an argument with one piece of evidence about phenomena using an understanding that the gravitational force of Earth on objects is directed down. (PS2)
- Construct an explanation supported by one piece of evidence about phenomena using an understanding of the relationship between the speed and energy of an object; provide evidence that energy can be transferred from place to place;

predict outcomes for the changes in energy that occur when objects collide; design, test, and refine a device that converts energy from one form to another; and use models to describe phenomena using an understanding that food energy was once energy from the Sun. (PS3)

- Develop models to describe phenomena using an understanding that waves can cause objects to move and that light allows objects to be seen; and compare multiple solutions that use patterns to transfer information. (PS4)
- Develop models to describe phenomena using an understanding of the diversity and commonalities of the life cycles of organisms, construct an argument supported by one piece of evidence about phenomena using an understanding that plants and animals have internal and external structures that support life functions, describe phenomena using an understanding that animals receive, process, and respond to information from their senses, and support an argument with one piece of evidence about phenomena using an understanding that plants get the materials they need for growth chiefly from air and water. (LS1)
- Construct an argument supported by one piece of evidence about phenomena using an understanding that some animals form groups that help members survive; and develop a model to describe phenomena using an understanding of the movement of matter among plants, animals, decomposers, and the environment. (LS2)
- Analyze and interpret data to provide one piece of evidence about phenomena using an understanding that plants and animals have inherited traits and that variation of these traits exists in groups of similar organisms; and support an explanation with one piece of evidence about phenomena using an understanding that traits can be influenced by the environment. (LS3)
- Analyze and interpret fossil data to provide one piece of evidence of organisms and the environments in which they lived; construct an explanation supported by one piece of evidence of phenomena using an understanding that variation among individuals of the same species is a survival advantage and that in a particular habitat some organisms survive well, some survive less well, and some cannot survive at all; and make a claim supported by one piece of evidence about the merit of a solution to a problem caused by changes to the environment and the types of plants and animals that live there. (LS4)
- Describe phenomena using an understanding that patterns in rock formations and fossils in rock layers provide evidence to support an explanation for changes in a landscape over time, support an argument about phenomena using one piece of evidence and an understanding that differences in the apparent brightness of the Sun compared to other stars is due to their relative distances from Earth, and represent data in graphical displays to reveal patterns of daily changes in shadows, day and night, and the seasonal appearance of stars. (ESS1)

- Represent data to describe typical seasonal weather conditions, combine information to describe climates in different regions of the world, provide one piece of evidence for phenomena using an understanding of the effects of weathering or the rate of erosion, analyze and interpret data from maps to describe patterns of Earth's features, develop a model to describe phenomena using an understanding of how Earth's systems interact, and describe and graph the percentages of water to provide one piece of evidence about the distribution of water on Earth. (ESS2)
- Make a claim supported by one piece of evidence about the merit of a design solution that reduces the impacts of a weather-related hazard, combine information to describe how energy and fuels are derived from natural resources, how their uses affect the environment, and ways individual communities protect Earth's resources and the environment, and generate and compare multiple solutions to reduce the impacts of natural Earth processes on humans. (ESS3)
- Define a simple design problem including criteria for success and constraints; generate and compare multiple solutions to a problem; and plan and carry out fair tests to identify one way to improve a model or prototype. (ETS1)

Nearing Proficiency

Students at the **Nearing Proficiency** level demonstrate evidence of partial understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. Students performing at the Nearing Proficiency level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of partial understanding and use of the NGSS Standards:

- Use a model to describe phenomena using an understanding of matter as tiny particles, graph quantities to describe phenomena using an understanding of conservation of matter during physical or chemical changes, make observations to identify materials based on their properties, and use data to determine whether a new substance with different properties is formed when two substances are mixed. (PS1)
- Conduct an investigation to provide one piece of evidence about phenomena using an understanding of the effects of balanced or unbalanced forces on the motion of an object, predict the future motion of an object based on simple patterns in observations and measurements, ask questions using an understanding of cause and effect of electric or magnetic interactions between objects not in contact with each other, partially define a simple design problem that can be solved using magnets, and make a claim about phenomena using an understanding that the gravitational force of Earth on objects is directed down. (PS2)

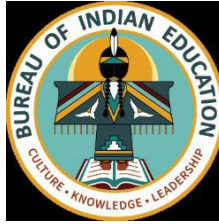
- Describe phenomena using an understanding of the relationship between the speed and energy of an object, explain that energy can be transferred from place to place, describe the changes in energy that occur when objects collide, describe elements of a device that converts energy from one form to another, and describe phenomena using an understanding that food energy was once energy from the Sun. (PS3)
- Use models to describe phenomena using an understanding that waves can cause objects to move or that light allows objects to be seen; and describe a solution that uses patterns to transfer information. (PS4)
- Use models to describe phenomena using an understanding of the diversity and commonalities of the life cycles of organisms, make a claim about phenomena using an understanding that plants and animals have internal and external structures that support life functions, describe phenomena using an understanding that animals receive or process or respond to information from their senses, and make a claim about phenomena using an understanding that plants get the materials they need for growth chiefly from air and water. (LS1)
- Make a claim about phenomena using an understanding that some animals form groups that help members survive; and use a model to describe phenomena using an understanding of the movement of matter among plants, animals, decomposers, and the environment. (LS2)
- Use data to describe phenomena using an understanding that plants and animals have inherited traits OR that variation of these traits exists in groups of similar organisms; and make a claim about phenomena using an understanding that traits can be influenced by the environment. (LS3)
- Use fossil data to describe organisms and the environments in which they lived; describe phenomena using an understanding that variation among individuals of the same species is a survival advantage or that in a particular habitat some organisms survive well, some survive less well, and some cannot survive at all; and describe a solution to a problem caused by changes to the environment and the types of plants and animals that live there. (LS4)
- Describe phenomena using an understanding that patterns in rock formations or fossils in rock layers provide evidence to support changes in a landscape over time, describe phenomena using an understanding that differences in the apparent brightness of the Sun compared to other stars is due to their relative distances from Earth, and use data in graphical displays to reveal patterns of daily changes in shadows or day and night or the seasonal appearance of stars. (ESS1)
- Use data to describe typical seasonal weather conditions, describe climates in different regions of the world, describe phenomena using an understanding of the effects of weathering or the rate of erosion, use maps to describe patterns of Earth's features, use a model to describe phenomena using an understanding of

how Earth's systems interact, and provide a description of the distribution of water on Earth. (ESS2)

- Describe a design solution or component of a design solution that reduces the impacts of a weather-related hazard, use information to determine that energy and fuels are derived from natural resources or how their uses affect the environment or ways individual communities protect Earth's resources and the environment, and describe a solution or component of a solution to reduce the impacts of natural Earth processes on humans. (ESS3)
- Define a simple design problem including at least one criterion for success or one constraint; generate a solution to a problem or compare two solutions to a problem; and use results of a fair test to identify one way to improve a model or prototype. (ETS1)

Novice

Students at the **Novice** level demonstrate evidence of emerging understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences.



BIE Science Proficiency Level

Descriptors Grade 8

Policy PLDs

Policy PLDs define the knowledge and skill level expectations for all grades and content areas for the BIE Science Assessment.

Level 4. Advanced

Students demonstrate evidence of **thorough** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 3. Proficient

Students demonstrate evidence of **satisfactory** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 2. Nearing Proficiency

Students demonstrate evidence of **partial** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 1. Novice

Students demonstrate evidence of **emerging** understanding and use of college and career readiness knowledge, skills, and abilities.

Range PLDs

Range PLDs describe the knowledge and skills that students throughout the range of each proficiency level are expected to be able to demonstrate. In line with the nature of the science standards, the statements combine science and engineering practices, disciplinary core ideas, and crosscutting concepts that students are expected to integrate and demonstrate.

Advanced

Students at the **Advanced** level demonstrate evidence of thorough understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. In addition to demonstrating the skills and understanding at the Proficient level, students performing at the Advanced level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of thorough understanding and use of the NGSS Standards:

- Develop, use, and analyze models to describe or explain phenomena using an understanding of the structure of matter; to predict, describe, and explain phenomena using an understanding of changes in particle motion and state; and to provide evidence for and describe phenomena using an understanding of conservation of mass during multiple physical and chemical changes. (PS1)
- Plan, carry out, and refine investigations to provide evidence to explain phenomena using an understanding of the effects of forces, interactions, and mass on the motion of objects, as well as analyze various data to evaluate claims about phenomena using an understanding of gravitational interactions in systems, and to design and/or compare multiple solutions to a problem using an understanding of systems of colliding objects. (PS2)
- Plan, carry out, and refine investigations, use and analyze data, and develop, use, and analyze models to explain phenomena using an understanding of various relationships involving kinetic and potential energy in systems, as well as apply such understanding to design, test, and evaluate devices to solve problems related to energy transfer and to support and/or evaluate claims about phenomena related to energy transfer. (PS3)
- Develop, use, and apply mathematical representations, patterns, and models to explain phenomena using an understanding of wave properties and relationships and wave interactions with various materials, and synthesize multiple sources of information using an understanding of signal types to evaluate claims about phenomena related to reliability of digital and analog signals. (PS4)
- Use multiple pieces of evidence from investigations of phenomena to explain that living things are made of cells; develop and use models of phenomena to describe the function of a cell and its parts and to describe how food is

rearranged in organisms through chemical reactions to support growth and/or release energy; support arguments about phenomena using multiple pieces of evidence and an understanding that the body is a system of interacting subsystems composed of cells; support an explanation for how several animal behaviors and several specialized plant structures affect the probability of reproductive success; use multiple pieces of evidence and an understanding of environmental and genetic factors to explain phenomena about how those factors influence the growth of organisms; construct an explanation using multiple pieces of evidence to explain the role of photosynthesis in the cycling of matter and flow of energy into and out of organisms; synthesize multiple sources of information about phenomena and an understanding of behaviors of organisms to determine that sensory receptors respond to stimuli by sending messages to the brain for immediate behavior or for storage as memories. (LS1)

- Analyze and interpret data about phenomena to provide evidence to explain multiple ways that resource availability affects populations, construct an argument supported by multiple pieces of evidence that populations are affected by changing physical and biological components of an ecosystem, develop and revise models to explain phenomena using an understanding of the cycling of matter and energy in an ecosystem, describe phenomena using an understanding of interactions among organisms and predict multiple patterns of interactions among organisms across multiple ecosystems, and evaluate multiple competing design solutions for phenomena that involve maintaining biodiversity in ecosystems. (LS2)
- Develop and use multiple models to explain phenomena using an understanding of how genetic mutations affect proteins resulting in harmful, beneficial, or neutral effects on an organism, and to explain phenomena using an understanding of how asexual reproduction results in offspring with identical genetic information and to explain how sexual reproduction results in offspring with genetic variation. (LS3)
- Analyze and interpret multiple pieces of data about phenomena using an understanding of the fossil record and modern organisms to show patterns in the change of life forms over time, apply multiple scientific ideas about phenomena to construct an explanation for the anatomical similarities and differences among modern and fossil organisms to infer evolutionary relationships, analyze pictorial data to compare similarities in embryological development across multiple familiar and unfamiliar species to identify evolutionary relationships, use multiple pieces of evidence and mathematical representations to explain phenomena using an understanding of how variation in genetic traits provides advantages to some individuals within a population and to support explanations of increases and decreases in specific traits over time, and explain phenomena by synthesizing multiple pieces of information about ways technologies have changed the way humans influence the inheritance of desired traits in organisms. (LS4)

- Develop, use, and revise a model of the Earth-Sun-Moon system to describe phenomena using an understanding of the cyclic pattern of the seasons and to describe phenomena using an understanding of the role of gravity in the motions within the solar system and galaxies, analyze and interpret data on multiple phenomena related to the scale properties of objects in the solar system, and use multiple pieces of evidence and an understanding of rock strata to explain phenomena about how the geologic time scale is used to organize Earth's history. (ESS1)
- Develop models to describe phenomena using an understanding of the flow of energy that drives the cycling of Earth's materials, use multiple pieces of evidence to explain phenomena using an understanding of how geoscience processes have changed Earth's surface at varying time and spatial scales, analyze and interpret multiple pieces of data to explain phenomena using an understanding of the evidence that supports past plate motions on Earth, develop models to describe phenomena using an understanding of the water cycle including energy and gravity, explain weather phenomena synthesizing and using evidence and an understanding of the interactions of air masses, and develop and use models to describe phenomena using an understanding of how unequal heating and Earth's rotation result in climate, atmospheric, and ocean circulation patterns. (ESS2)
- Use evidence to explain multiple phenomena using an understanding of how geoscience processes have resulted in uneven distribution of Earth's natural resources, analyze and interpret multiple pieces of data on natural hazard phenomena to forecast future catastrophic events and to inform the development of technologies to mitigate their effects, apply scientific principles to design a successful solution for monitoring and minimizing the human impacts on the environment, use multiple pieces of evidence to support an argument about phenomena using an understanding of how increases in human population impact Earth's systems, and ask multiple questions about phenomena to clarify evidence of multiple factors that have caused the rise in global temperatures. (ESS3)
- Define the criteria and constraints of a design problem with sufficient precision to ensure an optimal solution and using an understanding of scientific principles and potential impacts on people and the environment and an understanding of how those impacts may limit possible solutions, use a systematic process to evaluate how well multiple competing design solutions meet required criteria and constraints, analyze data from tests of multiple different design solutions to identify the best characteristics of each solution that can be combined into a new solution that will better meet criteria for success, and develop a realistic model of a proposed object, tool, or process that generates data while it is repeatedly tested and modified until an optimal design is achieved. (ETS1)

Proficient

Students at the **Proficient** level demonstrate evidence of satisfactory understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. In addition to demonstrating the skills and understanding at the Nearing Proficiency level, students performing at the Proficient level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of satisfactory understanding and use of the NGSS Standards:

- Develop and use models to describe phenomena using an understanding of the structure of matter; to predict and describe phenomena using an understanding of changes in particle motion and state; and to describe phenomena using an understanding of conservation of mass during one or two physical and chemical changes. (PS1)
- Plan and carry out investigations to produce data and/or provide evidence about phenomena using an understanding of the effects of forces, interactions, and mass on the motion of objects, as well as use direct data to support claims about phenomena using an understanding of gravitational interactions in systems and to design a solution to a problem using an understanding of systems of colliding objects. (PS2)
- Plan investigations, use data, and develop or use models to describe phenomena using an understanding of relationships involving kinetic and potential energy in systems, as well as apply such understanding to design and test devices to solve problems related to energy transfer and to support claims about phenomena related to energy transfer. (PS3)
- Use mathematical representations and patterns, and develop and use models, to describe phenomena using an understanding of wave properties and relationships and wave interactions with various materials and use multiple sources of information and an understanding of signal types to support claims about phenomena related to reliability of digital and analog signals. (PS4)
- Use one or two pieces of evidence from investigations of phenomena to explain that living things are made of cells, develop and use models of phenomena to describe the function of a cell and its parts and to describe how food is rearranged in organisms through chemical reactions to support growth and/or release energy, support arguments about phenomena using one or two pieces of evidence and an understanding that the body is a system of interacting subsystems composed of cells, support an explanation for how some animal behaviors and some specialized plant structures affect the probability of reproductive success, use one or two pieces of evidence and an understanding of environmental and genetic factors to explain phenomena about how those factors influence the growth of organisms, construct an explanation based on one or two pieces of evidence to explain the role of photosynthesis in the cycling

of matter and flow of energy into and out of organisms, use information about phenomena and an understanding of behaviors of organisms to determine that sensory receptors respond to stimuli by sending messages to the brain for immediate behavior or for storage as memories. (LS1)

- Analyze and interpret data about phenomena to provide evidence to explain one or two ways that resource availability affects populations, construct an argument supported by one or two pieces of evidence that populations are affected by changing physical or biological components of an ecosystem, develop models to describe phenomena using an understanding of the cycling of matter and energy in an ecosystem, describe phenomena using an understanding of interactions among organisms and predict one or two patterns of interactions among organisms across multiple ecosystems, and evaluate two competing design solutions for phenomena that involve maintaining biodiversity in ecosystems. (LS2)
- Develop and use one or two models to explain phenomena using an understanding of how genetic mutations affect proteins resulting in harmful, beneficial, or neutral effects on an organism and to explain phenomena using an understanding of how asexual reproduction results in offspring with identical genetic information and how sexual reproduction results in offspring with genetic variation. (LS3)
- Analyze and interpret one or two pieces of data about phenomena using an understanding of the fossil record and modern organisms to show patterns in the change of life forms over time, apply one or two scientific ideas about phenomena to construct an explanation for the anatomical similarities and differences among modern and fossil organisms to infer evolutionary relationships, analyze pictorial data to compare similarities in embryological development across multiple familiar species to identify evolutionary relationships, use one or two pieces of evidence or mathematical representations to explain phenomena using an understanding of how variation in genetic traits provides advantages to some individuals within a population and to support explanations of increases and decreases in specific traits over time, and explain phenomena by synthesizing one or two pieces of information about technologies that have changed the way humans influence the inheritance of desired traits in organisms. (LS4)
- Develop and use a model of the Earth-Sun-Moon system to describe phenomena using an understanding of the cyclic pattern of the seasons and to describe phenomena using an understanding of the role of gravity in the motions within the solar system and galaxies, analyze and interpret data on one or two phenomena related to the scale properties of objects in the solar system, and use one or two pieces of evidence and an understanding of rock strata to explain phenomena about how the geologic time scale is used to organize Earth's history. (ESS1)

- Develop a model to describe phenomena using an understanding of the flow of energy that drives cycling of Earth's materials, use one or two pieces of evidence to explain phenomena using an understanding of how geoscience processes have changed Earth's surface at varying time and spatial scales, analyze and interpret one or two pieces of data to explain phenomena using an understanding of the evidence that supports past plate motions on Earth, develop a model to describe phenomena using an understanding of the water cycle including energy and gravity, explain weather phenomena using evidence and an understanding of the interactions of air masses, and develop and use a model to describe phenomena using an understanding of how unequal heating and Earth's rotation result in climate, atmospheric, or ocean circulation patterns. (ESS2)
- Use evidence to explain one or two phenomena using an understanding of how geoscience processes have resulted in uneven distribution of Earth's natural resources, analyze and interpret one or two pieces of data on natural hazard phenomena to forecast future catastrophic events and to inform the development of technologies to mitigate their effects, apply scientific principles to design a solution for monitoring and minimizing human impacts on the environment, use one or two pieces of evidence to support an argument about phenomena using an understanding of how increases in human population impact Earth's systems, and ask one or two questions about phenomena to clarify evidence of one or two factors that have caused the rise in global temperatures. (ESS3)
- Define the criteria and constraints of a design problem with sufficient precision to ensure a successful solution and using an understanding of scientific principles and potential impacts on people and the environment and an understanding of how those impacts may limit possible solutions, use a systematic process to evaluate how well two competing design solutions meet required criteria and constraints, analyze data from tests of two different design solutions to identify the best characteristics of each solution that can be combined into a new solution that will better meet criteria for success, and develop a model of a proposed object, tool, or process that generates data while it is repeatedly tested and modified until an optimal design is achieved. (ETS1)

Nearing Proficiency

Students at the **Nearing Proficiency** level demonstrate evidence of partial understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. Students performing at the Nearing Proficiency level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of partial understanding and use of the NGSS Standards:

- Use models to identify the structure of matter as it relates to phenomena; to

describe basic phenomena using an understanding of changes in particle motion and state; and to identify that mass is conserved during physical or chemical changes that take place in various phenomena. (PS1)

- Identify or describe parts of investigations about phenomena using an understanding of the effects of forces, interactions, and mass on the motion of objects, describe aspects of phenomena using one or two pieces of data and an understanding of one gravitational interaction in a system, and identify a design or elements of a solution to a problem related to colliding objects. (PS2)
- Describe parts of investigations, use data, and use models to describe aspects of phenomena using an understanding of some relationships involving kinetic energy in systems; design and test a device for problems related to energy transfer; and identify principles that support claims about energy transfer in phenomena. (PS3)
- Use mathematical representations, patterns, and models to identify wave properties and wave interactions with various materials as they relate to phenomena and use one or two sources of information to identify that digital signals are more reliable than analog signals as demonstrated in various phenomena. (PS4)
- Use evidence from an investigation of a phenomenon to explain that living things are made of cells, use models of phenomena to describe the function of a cell and some of its parts and to describe that food is rearranged in organisms into new substances to support growth or to release energy, make claims about phenomena using evidence and a partial understanding that the body is a system of interacting subsystems composed of cells, describe some animal behaviors or specialized plant structures that may affect reproductive success, use evidence and a partial understanding of environmental or genetic factors to describe phenomena about how those factors influence the growth of organisms, construct an explanation based on one piece of evidence to describe the role of photosynthesis in the cycling of matter or flow of energy into and out of organisms, and use information about phenomena to describe that organisms use their senses to respond to stimuli immediately or to store information as memories. (LS1)
- Use data about phenomena to describe one way resource availability affects populations, make a claim supported by evidence that populations are affected by changing components of an ecosystem, use models to describe phenomena about the cycling of matter or energy in an ecosystem, use a partial understanding of interactions among organisms to describe one interaction between organisms, and describe a design solution or components of a solution for phenomena that involve maintaining biodiversity in ecosystems. (LS2)
- Use models to partially explain phenomena using an understanding of one way that genetic mutations affect organisms and to describe phenomena using an

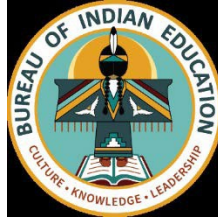
understanding of asexual reproduction resulting in offspring with identical genetic information, or about sexual reproduction resulting in offspring with genetic variation. (LS3)

- Use data about phenomena using a partial understanding of the fossil record and modern organisms to show patterns in the change of life forms over time, apply a scientific idea about phenomena to describe an anatomical similarity or difference between modern and fossil organisms, use pictorial data to describe similarities in embryological development across multiple species, use evidence or one mathematical representation to describe phenomena using a partial understanding of variation in genetic traits among individuals within a population or to describe increases and decreases in specific traits over time, and describe phenomena using a partial understanding about technologies that have changed the way humans influence the inheritance of desired traits in organisms. (LS4)
- Use a model of the Earth-Sun-Moon system to describe phenomena using a partial understanding of the cyclic pattern of the seasons and to describe phenomena using a partial understanding of the role of gravity in the motions within the solar system or galaxies, use data on phenomena related to the scale properties of objects in the solar system, and use evidence and a partial understanding of rock strata to describe some aspects about how the geologic time scale is related to Earth's history. (ESS1)
- Use a model to describe phenomena using a partial understanding of the flow of energy that drives the cycling of Earth's materials, use evidence to explain phenomena using a partial understanding of how geoscience processes have changed Earth's surface, use data to explain phenomena using a partial understanding of the evidence that supports past plate motions on Earth, use a model to describe phenomena using an understanding of the water cycle, describe weather phenomena using evidence and a partial understanding of the interactions of air masses, and use a model to describe phenomena using an understanding of how unequal heating and Earth's rotation result in some climate, atmospheric, or ocean circulation patterns. (ESS2)
- Use evidence to explain a phenomenon using a partial understanding of how geoscience processes have resulted in uneven distribution of some of Earth's natural resources, use data on natural hazard phenomena to forecast future catastrophic events or to inform the development of one technology that could be used to mitigate their effects, identify human impacts on the environment or design parts of a solution for monitoring or minimizing the human impacts, use evidence to make a claim about phenomena using a partial understanding of how increases in human population impact Earth's systems, and ask one question about a phenomenon to clarify evidence of one factor that has caused the rise in global temperatures. (ESS3)
- Define one criterion or constraint of a design problem using an understanding of scientific principles and/or potential impacts on people and the environment and

identify one way those impacts may limit possible solutions, use a systematic process to evaluate how well a design solution meets required criteria or constraints, analyze data from tests of a design solution to identify a characteristic of the solution that is necessary to meet the criteria for success, and develop a partial model of a proposed object, tool, or process that can be tested and modified. (ETS1)

Novice

Students at the **Novice** level demonstrate evidence of emerging understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences.



BIE Science Proficiency Level

Descriptors Grade 11

Policy PLDs

Policy PLDs define the knowledge and skill level expectations for all grades and content areas for the BIE Science Assessment.

Level 4. Advanced

Students demonstrate evidence of **thorough** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 3. Proficient

Students demonstrate evidence of **satisfactory** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 2. Nearing Proficiency

Students demonstrate evidence of **partial** understanding and use of college and career readiness knowledge, skills, and abilities.

Level 1. Novice

Students demonstrate evidence of **emerging** understanding and use of college and career readiness knowledge, skills, and abilities.

Range PLDs

Range PLDs describe the knowledge and skills that students throughout the range of each proficiency level are expected to be able to demonstrate. In line with the nature of the science standards, the statements combine science and engineering practices, disciplinary core ideas, and crosscutting concepts that students are expected to integrate and demonstrate.

Advanced

Students at the **Advanced** level demonstrate evidence of thorough understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. In addition to demonstrating the skills and understanding at the Proficient level, students performing at the Advanced level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of thorough understanding and use of the NGSS Standards:

- Use an understanding of the periodic table to predict multiple relative properties of elements, plan and conduct an investigation of phenomena related to multiple bulk scale properties of substances and explain how these relate to the strength of electrical forces between particles, explain phenomena about the release or absorption of energy by a chemical system by developing models to show changes in total bond energy, use multiple pieces of evidence to explain phenomena using an understanding of how changes in temperature and concentration affect reaction rate, explain phenomena about chemical systems at equilibrium using an understanding of how multiple conditions on the system could be changed to produce more or fewer products or more or fewer reactants, use multiple mathematical representations to support claims that mass is conserved during a chemical reaction, and develop multiple models to describe phenomena using an understanding of the changes in the nucleus of an atom and the energy released during fission, fusion, and radioactive decay. (PS1)
- Analyze and use multiple pieces of data from phenomena to show that $f = ma$; use multiple mathematical representations of phenomena and an understanding of momentum to support the claim that the total momentum of a system is conserved when there is no net force on the system; apply multiple scientific and engineering ideas to design, evaluate, and refine multiple devices that minimize force on an object during a collision; use mathematical representations of Newton's law of gravitation and Coulomb's law to describe and make predictions about familiar and unfamiliar phenomena using an understanding of gravitational and electrostatic forces between objects; plan and conduct an investigation of phenomena to produce multiple pieces of evidence that prove an electric current produces a magnetic field and that a changing magnetic field produces electric current; and communicate multiple pieces of information about

phenomena using an understanding of how the molecular structure of a material relates to its macroscopic properties and makes the material well suited for particular uses. (PS2)

- Create multiple computational models of phenomena to calculate changes in energy of a system when energy flows into and out of the system is known; develop and use multiple models to explain phenomena using an understanding of how energy at the macroscopic scale can be accounted for at the microscopic scale in energy associated with particle motion and relative position; design, build, and refine devices that convert one form of energy into another; plan and conduct an investigation of phenomena to provide multiple pieces of evidence using an understanding that when two components of different temperatures are combined within a closed system, both components eventually have the same temperature; and develop and use a model to explain phenomena related to the forces and the changes in energy between two objects interacting through electric fields and magnetic fields. (PS3)
- Explain phenomena using an understanding of multiple mathematical representations regarding relationships among frequency, wavelength, and speed of waves in various media; evaluate multiple questions about phenomena using an understanding of the advantages of using digital transmission and storage of information; use multiple phenomena to evaluate claims that electromagnetic radiation can be described using a wave or particle model; in the context of phenomena, evaluate multiple claims about the effects that different frequencies of electromagnetic radiation have on matter; and communicate technical information about phenomena using an understanding of how multiple specific technological devices use the principles of wave behavior and wave interactions to transmit and capture information and energy. (PS4)
- Use multiple pieces of evidence to explain phenomena using an understanding of how the structure of DNA determines the structure of proteins and how proteins carry out the functions of life through specialized cells; develop and use a complex model to describe phenomena using an understanding of the organization of interacting systems within multicellular organisms; plan and conduct an investigation to provide multiple pieces of evidence about phenomena that show that feedback mechanisms maintain homeostasis; use a complex model to describe phenomena using an understanding of how cell division and differentiation help produce and maintain complex organisms; use a complex model to describe phenomena using an understanding of how photosynthesis transforms light energy into stored chemical energy; use multiple pieces of evidence to construct and revise an explanation about phenomena using an understanding of how carbon, hydrogen, and oxygen from sugar molecules combine with other elements to form amino acids and other large carbon-based molecules; and use a complex model to describe phenomena using an understanding that cellular respiration is a chemical process that breaks the bonds in food and oxygen molecules and forms bonds in new

compounds, which results in a net transfer of energy. (LS1)

- Use mathematical and computational representations to support explanations of phenomena using an understanding of factors that affect carrying capacity of ecosystems at different scales; use multiple pieces of evidence and mathematical representations to support and revise explanations of phenomena using an understanding of factors affecting biodiversity and populations in ecosystems of different scales and to support claims for the cycling of matter and flow of energy among organisms in an ecosystem; use evidence to construct and revise an explanation of phenomena using an understanding of the cycling of matter and flow of energy in aerobic and anaerobic conditions; develop models to describe phenomena using an understanding of the role of photosynthesis and cellular respiration in the cycling of carbon among Earth's spheres; evaluate multiple claims, pieces of evidence, and reasoning about phenomena involving complex interactions in ecosystems using an understanding that these interactions maintain relatively consistent numbers and types of organisms under stable conditions, but changing conditions may result in a new ecosystem; design, evaluate, and refine solutions for reducing impacts of human activities on the environment or biodiversity; and evaluate multiple pieces of evidence about phenomena using an understanding of the role of group behavior on individual and species' chances to survive and reproduce. (LS2)
- Ask multiple questions about phenomena to clarify relationships surrounding the role of DNA in chromosomes in coding the instructions for traits passed from parents to offspring; use multiple pieces of evidence to make and defend a claim about phenomena using an understanding that inheritable genetic variations may result from new genetic combinations through meiosis, viable errors during replication, and/or mutations caused by environmental factors; and apply multiple concepts of statistics and probability to explain phenomena using an understanding of the variation and distribution of expressed traits in a population. (LS3)
- Communicate multiple pieces of scientific information about phenomena using an understanding that common ancestry and biological evolution are supported by multiple lines of empirical evidence; use multiple pieces of evidence to construct an explanation about phenomena using an understanding that the process of evolution primarily results from four factors: the potential for a species to increase in number, the heritable genetic variation of individuals in a species due to mutation and sexual reproduction, competition for limited resources, and the proliferation of those organisms that are better able to survive and reproduce in the environment; apply multiple concepts of statistics and probability to support explanations of phenomena using an understanding that organisms with an advantageous heritable trait tend to increase in proportion to organisms lacking the trait; use multiple pieces of evidence to construct an explanation of phenomena using an understanding of how natural

selection leads to adaptation of populations; evaluate multiple pieces of evidence supporting claims about phenomena using an understanding that changes in environmental conditions may result in: increases in numbers of individuals of some species, the emergence of new species over time, and the extinction of other species; and create and revise a simulation to test a solution to mitigate adverse impacts of human activity on biodiversity. (LS4)

- Use multiple pieces of evidence to develop a model to describe phenomena using an understanding of the life span of the Sun and the role of nuclear fusion in the Sun's core to release energy that reaches Earth in the form of radiation; use multiple pieces of evidence and an understanding of the phenomena of light spectra, motion of distant galaxies, and composition of matter in the universe to construct an explanation of the big bang theory; communicate multiple scientific ideas about phenomena using an understanding of the way stars produce different elements over their life cycles; use mathematical and computational representations of phenomena to predict the motion of orbiting objects in the solar system; evaluate multiple pieces of evidence of phenomena using an understanding of past and current movements of continental and oceanic crust and the theory of plate tectonics to explain the ages of crustal rocks; and apply scientific reasoning and an understanding of multiple pieces of evidence from ancient Earth materials, meteorites, and other planetary surfaces to describe phenomena about Earth's formation and early history. (ESS1)
- Develop models to describe phenomena using an understanding of how Earth's internal and surface processes operate at different spatial and temporal scales to form continental and ocean floor features; analyze multiple types of geoscience data about phenomena to make a claim that one change to Earth's surface can create feedback that causes changes to other Earth systems; use multiple pieces of evidence to develop a model of Earth's interior to describe phenomena using an understanding of the cycling of matter by thermal convection; use models to describe phenomena using an understanding of how variations in the flow of energy into and out of Earth's systems result in changes in climate; plan and conduct investigations of phenomena related to the properties of water using an understanding of water's effects on Earth materials and surface processes; develop quantitative models to describe phenomena using an understanding of the cycling of carbon among the hydrosphere, atmosphere, geosphere, and biosphere; and use multiple pieces of evidence to construct an argument about phenomena using an understanding of the simultaneous coevolution of Earth's systems and life on Earth. (ESS2)
- Use multiple pieces of evidence to construct an explanation about phenomena using an understanding of how the availability of natural resources, occurrence of natural hazards, and changes in climate have influenced human activity; using an understanding of cost-benefit ratios, evaluate multiple competing design solutions for developing, managing, and utilizing energy and mineral resources; create computational simulations of phenomena to show the

relationships among management of natural resources, the sustainability of human populations, and biodiversity; using an understanding of human impacts on natural systems, evaluate and refine a technological solution that reduces these impacts; analyze multiple pieces of geoscience data and multiple global climate models of phenomena to make a forecast of the current rate of climate change and associated future impacts to Earth systems; and use computational representations to describe phenomena using an understanding of the relationships among Earth systems and how those relationships are modified due to human activity. (ESS3)

- Analyze a major global challenge to specify multiple qualitative and quantitative criteria and constraints for solutions that account for multiple societal needs and wants; design an engineering solution to multiple complex real-world problems by breaking them down into smaller, more manageable problems; evaluate and refine a solution to a complex real-world problem based on prioritized criteria and trade-offs that account for a range of constraints, including cost, safety, reliability, and aesthetics, as well as social, cultural, and environmental impacts; and use a computer simulation to model the impact of multiple proposed solutions to a complex real-world problem with numerous criteria and constraints on interactions within and between systems relevant to the problem. (ETS1)

Proficient

Students at the **Proficient** level demonstrate evidence of satisfactory understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. In addition to demonstrating the skills and understanding at the Nearing Proficiency level, students performing at the Proficient level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of satisfactory understanding and use of the NGSS Standards:

- Use an understanding of the periodic table to predict one or two relative properties of elements, plan and conduct an investigation of phenomena related to one or two bulk scale properties of substances and explain how these relate to the strength of electrical forces between particles, explain phenomena about the release or absorption of energy by a chemical system by developing a model to show changes in total bond energy, use one or two pieces of evidence to explain phenomena using an understanding of how changes in temperature or concentration affect reaction rate, explain phenomena about chemical systems at equilibrium using an understanding of how one or two of the conditions on the system could be changed to produce more products, use one or two mathematical representations to support claims that mass is conserved during a chemical reaction, and develop one or two models to describe phenomena using an understanding of the changes in the nucleus of an atom and the energy released during fission, fusion, and radioactive decay. (PS1)

- Analyze and use one or two pieces of data from phenomena to show that $f = ma$; use one or two mathematical representations of phenomena and an understanding of momentum to support the claim that the total momentum of a system is conserved when there is no net force on the system; apply one or two scientific and engineering ideas to design, evaluate, and refine a device that minimizes force on an object during a collision; use mathematical representations of Newton's law of gravitation and Coulomb's law to describe and make predictions about familiar phenomena using an understanding of gravitational and electrostatic forces between objects; plan and conduct an investigation of phenomena to produce one or two pieces of evidence that prove an electric current produces a magnetic field and that a changing magnetic field produces electric current; and communicate one or two pieces of information about phenomena using an understanding of how the molecular structure of a material relates to its macroscopic properties and makes the material well suited for particular uses. (PS2)
- Create a computational model of phenomena to calculate changes in energy of a system when energy flows into and out of the system is known; develop and use one or two models to explain phenomena using an understanding of how energy at the macroscopic scale can be accounted for at the microscopic scale in energy associated with particle motion and relative position; design, build, and refine a device that converts one form of energy into another; plan and conduct an investigation of phenomena to provide one or two pieces of evidence using an understanding that when two components of different temperatures are combined within a closed system, both components eventually have the same temperature; and develop and use a model to explain phenomena related to the forces and the changes in energy between two objects interacting through electric or magnetic fields. (PS3)
- Explain phenomena using an understanding of one or two mathematical representations regarding relationships among frequency, wavelength, and speed of waves in various media; evaluate one or two questions about phenomena using an understanding of the advantages of using digital transmission and storage of information; use one or two phenomena to evaluate claims that electromagnetic radiation can be described using a wave or particle model; in the context of phenomena, evaluate one or two claims about the effects that different frequencies of electromagnetic radiation have on matter; and communicate technical information about phenomena using an understanding of how one or two specific technological devices use the principles of wave behavior and wave interactions to transmit and capture information and energy. (PS4)
- Use evidence to explain phenomena using an understanding of how the structure of DNA determines the structure of proteins and how proteins carry out the functions of life through specialized cells; develop and use a model to describe phenomena using an understanding of the organization of interacting

- systems within multicellular organisms; plan and conduct an investigation to provide evidence about phenomena that show that feedback mechanisms maintain homeostasis; use a model to describe phenomena using an understanding of how cell division and differentiation help produce and maintain complex organisms; use a model to describe phenomena using an understanding of how photosynthesis transforms light energy into stored chemical energy; use evidence to construct and revise an explanation about phenomena using an understanding of how carbon, hydrogen, and oxygen from sugar molecules combine with other elements to form amino acids and/or other large carbon-based molecules; and use a model to describe phenomena using an understanding that cellular respiration is a chemical process that breaks the bonds in food and oxygen molecules and forms bonds in new compounds, which results in a net transfer of energy. (LS1)
- Use mathematical or computational representations to support explanations of phenomena using an understanding of factors that affect carrying capacity of ecosystems at different scales; use one or two pieces of evidence and one or two mathematical representations to support and revise explanations of phenomena using an understanding of factors affecting biodiversity and populations in ecosystems of different scales and to support claims for the cycling of matter and flow of energy among organisms in an ecosystem; use evidence to construct or revise an explanation of phenomena using an understanding of the cycling of matter and flow of energy in aerobic and anaerobic conditions; develop a model to describe phenomena using an understanding of the role of photosynthesis and cellular respiration in the cycling of carbon among Earth's spheres; evaluate one or two claims, pieces of evidence, and reasoning about phenomena involving complex interactions in ecosystems using an understanding that these interactions maintain relatively consistent numbers and types of organisms under stable conditions, but changing conditions may result in a new ecosystem; design, evaluate, and refine a solution for reducing impacts of human activities on the environment or biodiversity; and evaluate one or two pieces of evidence about phenomena using an understanding of the role of group behavior on individual and species' chances to survive and reproduce. (LS2)
 - Ask one or two questions about phenomena to clarify relationships about the role of DNA in chromosomes in coding the instructions for traits passed from parents to offspring; use one or two pieces of evidence to make and defend a claim about phenomena using an understanding that inheritable genetic variations may result from new genetic combinations through meiosis, viable errors during replication, and/or mutations caused by environmental factors; and apply one or two concepts of statistics and probability to explain phenomena using an understanding of the variation and distribution of expressed traits in a population. (LS3)
 - Communicate one or two pieces of scientific information about phenomena

using an understanding that common ancestry and biological evolution are supported by multiple lines of empirical evidence; use one or two pieces of evidence to construct an explanation about phenomena using an understanding that the process of evolution primarily results from four factors: the potential for a species to increase in number, the heritable genetic variation of individuals in a species due to mutation and sexual reproduction, competition for limited resources, and the proliferation of those organisms that are better able to survive and reproduce in the environment; apply one or two concepts of statistics and probability to support explanations of phenomena using an understanding that organisms with an advantageous heritable trait tend to increase in proportion to organisms lacking the trait; use one or two pieces of evidence to construct an explanation of phenomena using an understanding of how natural selection leads to adaptation of populations; evaluate one or two pieces of evidence supporting claims about phenomena using an understanding that changes in environmental conditions may result in: increases in numbers of individuals of some species, the emergence of new species over time, and the extinction of other species; and create or revise a simulation to test a solution to mitigate adverse impacts of human activity on biodiversity. (LS4)

- Use one or two pieces of evidence to develop a model to describe phenomena using an understanding of the life span of the Sun and the role of nuclear fusion in the Sun's core to release energy that reaches Earth in the form of radiation; use one or two pieces of evidence and an understanding of the phenomena of light spectra, motion of distant galaxies, and composition of matter in the universe to construct an explanation of the big bang theory; communicate one or two scientific ideas about phenomena using an understanding of the way stars produce different elements over their life cycles; use mathematical or computational representations of phenomena to predict the motion of orbiting objects in the solar system; evaluate one or two pieces of evidence of phenomena using an understanding of past and current movements of continental and oceanic crust and the theory of plate tectonics to explain the ages of crustal rocks; and apply scientific reasoning and an understanding of one or two pieces of evidence from ancient Earth materials, meteorites, and other planetary surfaces to describe phenomena about Earth's formation and early history. (ESS1)
- Develop a model to describe phenomena using an understanding of how Earth's internal and surface processes operate at different spatial and temporal scales to form continental and ocean floor features; analyze one type of geoscience data about phenomena to make a claim that one change to Earth's surface can create feedback that causes changes to other Earth systems; use one or two pieces of evidence to develop a model of Earth's interior to describe phenomena using an understanding of the cycling of matter by thermal convection; use a model to describe phenomena using an understanding of how variations in the flow of energy into and out of Earth's systems result in changes in climate; plan

- and conduct an investigation of phenomena related to the properties of water using an understanding of water's effects on Earth materials and surface processes; develop a quantitative model to describe phenomena using an understanding of the cycling of carbon among the hydrosphere, atmosphere, geosphere, and biosphere; and use one or two pieces of evidence to construct an argument about phenomena using an understanding of the simultaneous coevolution of Earth's systems and life on Earth. (ESS2)
- Use one or two pieces of evidence to construct an explanation about phenomena using an understanding of how the availability of natural resources, occurrence of natural hazards, and changes in climate have influenced human activity; using an understanding of cost-benefit ratios, evaluate two competing design solutions for developing, managing, and utilizing energy and mineral resources; create a computational simulation of phenomena to show the relationships among management of natural resources, the sustainability of human populations, and biodiversity; using an understanding of human impacts on natural systems, evaluate or refine a technological solution that reduces these impacts; analyze one or two pieces of geoscience data and one or two global climate models of phenomena to make a forecast of the current rate of climate change and associated future impacts to Earth systems; and use a computational representation to describe phenomena using an understanding of the relationships among Earth systems and how those relationships are modified due to human activity. (ESS3)
 - Analyze a major global challenge to specify one or two qualitative and quantitative criteria and constraints for solutions that account for one or two societal needs and wants; design an engineering solution to a complex real-world problem by breaking it down into smaller, more manageable problems; evaluate a solution to a complex real-world problem based on prioritized criteria and trade-offs that account for a range of constraints, including cost, safety, reliability, and aesthetics, as well as social, cultural, and environmental impacts; and use a computer simulation to model the impact of two proposed solutions to a complex real-world problem with numerous criteria and constraints on interactions within and between systems relevant to the problem. (ETS1)

Nearing Proficiency

Students at the **Nearing Proficiency** level demonstrate evidence of partial understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences. Students performing at the Nearing Proficiency level can be expected to be able to demonstrate knowledge and skills like in the following examples, as evidence of partial understanding and use of the NGSS Standards:

- Use a partial understanding of the periodic table to predict one relative property of elements, conduct an investigation of phenomena related to bulk scale properties of substances or construct a partial explanation of how these relate to the strength of electrical forces between particles, describe phenomena about the release or absorption of energy by a chemical system, explain phenomena using a partial understanding of how changes in temperature or concentration affect reaction rate, explain phenomena about chemical systems at equilibrium using a partial understanding of how one condition on the system could be changed to produce more products, use a mathematical representation to support a claim that mass is conserved during a chemical reaction, and use models to describe phenomena using a partial understanding of the changes in the nucleus of an atom or the energy released during fission, fusion, or radioactive decay. (PS1)
- Use data from phenomena to show that $f = ma$, use a mathematical representation of a phenomenon and a partial understanding of momentum to support the claim that the total momentum of a system is conserved, apply a scientific or an engineering idea to design a device that minimizes force on an object during a collision, use mathematical representations of Newton's law of gravitation or Coulomb's law to describe phenomena using a partial understanding of gravitational and electrostatic forces between objects, conduct an investigation of phenomena to produce evidence that an electric current produces a magnetic field or that a changing magnetic field produces electric current, and communicate information about phenomena using a partial understanding of how the molecular structure of a material relates to its macroscopic properties or makes the material well suited for particular uses. (PS2)
- Use a computational model of phenomena to calculate changes in energy of a system when energy flows into and out of the system is known; use models to explain phenomena using a partial understanding of how energy at the macroscopic scale can be accounted for at the microscopic scale in energy associated with particle motion or relative position; design a device that converts one form of energy into another; conduct an investigation of phenomena that provides evidence using a partial understanding that when two components of different temperatures are combined within a closed system, both components eventually have the same temperature; and use a model to explain phenomena

related to the forces or the changes in energy between two objects interacting through electric or magnetic fields. (PS3)

- Explain phenomena using an understanding of a mathematical representation regarding one relationship among frequency, wavelength, and speed of waves in various media; ask questions about phenomena using a partial understanding of the advantages of using digital transmission and storage of information; use one phenomenon to support a claim that electromagnetic radiation may be described using a wave or particle model; make a claim about phenomena related to the effects that different frequencies of electromagnetic radiation have on matter; and communicate information about phenomena using a partial understanding of how specific technological devices use the principles of wave behavior or wave interactions to transmit or capture information and energy. (PS4)
- Use evidence to explain phenomena using a partial understanding of how the structure of DNA determines the structure of proteins; use a model to describe phenomena using a partial understanding of the organization of interacting systems within multicellular organisms; conduct an investigation to provide evidence about phenomena that show that some feedback mechanisms maintain homeostasis; use a model to describe phenomena using a partial understanding of how cell division or differentiation helps produce or maintain a complex organism; use a model to describe phenomena using a partial understanding of how photosynthesis transforms light energy into stored chemical energy; describe that carbon, hydrogen, and oxygen from sugar molecules combine with other elements to form amino acids or other large carbon-based molecules; and use a model to describe phenomena using a partial understanding that cellular respiration is a chemical process that breaks the bonds in food or oxygen molecules, forms bonds in new compounds, or results in a net transfer of energy. (LS1)
- Use mathematical or computational representations to support explanations of phenomena using a partial understanding of factors that affect carrying capacity of ecosystems; use evidence or mathematical representations to describe factors affecting biodiversity or populations in ecosystems and to support a claim for the cycling of matter or flow of energy among organisms in an ecosystem; use evidence to construct an explanation of phenomena using a partial understanding of the cycling of matter and flow of energy in aerobic or anaerobic conditions; use a model to describe phenomena using an understanding of the role of photosynthesis or cellular respiration in the cycling of carbon; evaluate a claim about a phenomenon involving interactions in ecosystems using a partial understanding that these interactions maintain relatively consistent numbers and types of organisms under stable conditions, but changing conditions may result in a new ecosystem; identify a solution for reducing impacts of human activities on the environment or biodiversity; and use evidence to describe the role of group behavior on individual and species' chances to survive and reproduce. (LS2)

- Ask a question about a phenomenon to clarify relationships about the role of DNA in chromosomes in coding the instructions for traits passed from parents to offspring; make a claim about phenomena using a partial understanding that inheritable genetic variations may result from new genetic combinations through meiosis, viable errors during replication, or mutations caused by environmental factors; and describe phenomena using an understanding of the variation and distribution of expressed traits in a population. (LS3)
- Communicate scientific information about phenomena using a partial understanding that common ancestry and biological evolution are supported by empirical evidence; use evidence to construct an explanation about phenomena using a partial understanding that the process of evolution primarily results from one or two of the following factors: the potential for a species to increase in number, the heritable genetic variation of individuals in a species due to mutation and sexual reproduction, competition for limited resources, and the proliferation of those organisms that are better able to survive and reproduce in the environment; apply concepts of statistics or probability to explain phenomena using a partial understanding that organisms with an advantageous heritable trait tend to increase in proportion to organisms lacking the trait; explain phenomena using a partial understanding of how natural selection leads to adaptation of populations; use evidence to support claims about phenomena using a partial understanding that changes in environmental conditions may result in one or two of the following: increases in numbers of individuals of some species, the emergence of new species over time, or the extinction of other species; and use a simulation to test a solution to mitigate adverse impacts of human activity on biodiversity. (LS4)
- Use evidence to develop a model to describe phenomena using a partial understanding of the life span of the Sun or the role of nuclear fusion in the Sun's core to release energy that reaches Earth in the form of radiation; use evidence and a partial understanding of the phenomena of light spectra, motion of distant galaxies, or composition of matter in the universe to construct an explanation of the big bang theory; communicate a scientific idea about a phenomenon using a partial understanding of the way stars produce different elements over their life cycles; use a mathematical or computational representation of phenomena to predict the motion of orbiting objects in the solar system; use evidence of phenomena using a partial understanding of past and current movements of continental and oceanic crust or the theory of plate tectonics to explain the ages of crustal rocks; and apply scientific reasoning or a partial understanding of evidence from ancient Earth materials, meteorites, and other planetary surfaces to describe phenomena about Earth's formation and early history. (ESS1)
- Use a model to describe phenomena using a partial understanding of how Earth's internal and surface processes operate to form continental and ocean floor features; use geoscience data about phenomena to make a claim that one

change to Earth's surface can create feedback that causes changes to other Earth systems; develop a model of Earth's interior to describe phenomena using an understanding of the cycling of matter by thermal convection; use a model to describe phenomena using a partial understanding of how variations in the flow of energy into and out of Earth's systems result in changes in climate; conduct an investigation of phenomena related to the properties of water using a partial understanding of water's effects on Earth materials or surface processes; use a quantitative model to describe phenomena using a partial understanding of the cycling of carbon among the hydrosphere, atmosphere, geosphere, and biosphere; and use evidence to make a claim about phenomena using a partial understanding of the simultaneous coevolution of Earth's systems and life on Earth. (ESS2)

- Use evidence to construct an explanation about phenomena using a partial understanding of how the availability of natural resources, occurrence of natural hazards, or changes in climate have influenced human activity; using a partial understanding of cost-benefit ratios, evaluate a design solution for developing, managing, or utilizing energy or mineral resources; use a computational simulation of phenomena to show some of the relationships among management of natural resources, the sustainability of human populations, and biodiversity; using a partial understanding of human impacts on natural systems, describe a technological solution that reduces these impacts; use geoscience data or global climate models of phenomena to make a forecast of the current rate of climate change or associated future impacts to Earth systems; and use a computational representation to describe phenomena using a partial understanding of the relationships among Earth systems or how those relationships are modified due to human activity. (ESS3)
- Analyze a major global challenge to specify a qualitative or quantitative criterion or constraint for a solution that accounts for a societal need or want; describe one or two ways a complex real-world problem could be broken down into smaller, more manageable problems that could be solved through engineering; explain how a solution to a complex real-world problem meets required criteria or explain one or two trade-offs of the solution; and use a computer simulation to model the impact of a proposed solution to a complex real-world problem with two or three criteria and constraints on interactions within or between systems relevant to the problem. (ETS1)

Novice

Students at the **Novice** level demonstrate evidence of emerging understanding and use of all three dimensions (science and engineering practices, crosscutting concepts, and disciplinary core ideas) to make sense of phenomena and/or to design solutions to problems in the physical, life, and Earth and space sciences.

APPENDIX C

ACCOMMODATION FREQUENCIES

Only students who met the attemptedness rule (i.e., attempted 5 or more items) contributed to the frequencies in these tables.

Table C-1. Number of Students Taking BIE Science Assessment and Utilizing Accommodation(s) /Accessibility Feature(s), as a Function of Grade*

Accommodation/Accessibility Feature	Accommodation		Grades	
	ID	5	8	11
AccAssistTech	3	1	0	0
AccReadAloudSelf	4	33	22	3
AccTactileGraphics	5	1	2	0
AccDirections	6	2	1	0
AccSpeechToText	7	51	64	5
AccComWordDictionary	13	3	4	0
AccSmallGroup	14	187	167	77
AccHumanReaderENG	19	36	13	3
AccAssistDevice	22	3	10	3
AccHumanSigner	29	0	2	0
AccHumanScribe	31	15	10	3
AccWordPredictionEmbed	32	26	1	1
AccColorContrast	33	12	17	0
AccModeTesting	34	8	7	3
AccAltSetting	35	78	55	9
AccSCI_TTSENG	38	172	120	45
AccWordPrediction	39	25	9	0
AccPictureDictionary	40	9	2	1
AccHeadphones	42	72	54	2
AccExtendedTime	45	228	188	87

**Only students who met the attemptedness rule (i.e., attempted 5 or more items) contributed to the frequencies in these tables.*

APPENDIX D
IRC AND BIAS REVIEW MEETING PARTICIPANTS

Grade 8 Science BIE-Specific Items Content & Bias Review Committee Final Attendance – June 2024 – Virtual

First Name	Last Name	School	Email Address	Meeting Status
Irish	Balunton	Na' Neelzhiin Ji Olta', Inc., NAV	irishb@naneelzhiin.org	Attended Meeting
Tyler	Bangert	Dishchii'bikoh Community School, TCS	tyler.bangert@dishchiibikoh.org	Attended Meeting
Conette	Crausos	Na' Neelzhiin Ji Olta', Inc., NAV	crausosconette91@gmail.com	Attended Meeting
Natalie	Hintz	Tate Topa Tribal School, TCS	natalie.hintz@k12.nd.us	Attended Meeting
Sherry	Johnson	Enemy Swim Day School/Tiospa Zina Tribal School, TCS	SherryJ@swo-nsn.gov	Attended Meeting
Leonila	Villaganas	Na Neelzhiin Ji Olta Inc., NAV	leonilav@naneelzhiin.org	Attended Meeting

Grade 11 Science BIE-Specific Items Content & Bias Review Committee Final Attendance – June 13, 2025 – Virtual

First Name	Last Name	School	Email Address	Meeting Status
Linda	Bimberg	Riverside Indian School, BOS	lindamap@hotmail.com	Applied and selected
Wendy	Fuller	Lac Courte Oreilles Ojibwe School, TCS	wendy.fuller@lcoosk12.org	Applied and selected
Rolanda	Morris	Wingate High School, BOS	Rolanda.morris@bie.edu	Applied and selected
Nate	Raynor	Mescalero Apache School	nate.raynor@mescalero.org	Applied and selected
Marlene	Murphy	Wingate High School	marlene.murphy@bie.edu	Applied and selected
Leny	Mendoza	Northwest High School-SASI, TCS	leny.mendoza@sasinm.com	Applied and selected
Aurelia	Shorty	Bureau of Indian Education	Assessment & Accountability	BIE-conducted review
Donald	Griffin	Bureau of Indian Education	Chief Academic Office	BIE-conducted review

APPENDIX E
SCORER QUALIFICATION RATES

Table E-1 summarizes the qualification rates for the 2025 operational assessment for BIE Science. Rates of success during qualification varied. Multiple factors determine the success of a scorer during qualification. These include familiarity with the assessment, grade level, and variation of item types. Please note that not all scorers who failed Qual 1 attempted Qual 2.

Table E-1. Qualification Summary for BIE Science

Grade 5	697063	697063	Scorers Qualified	697164	697164	Scorers Qualified	781454	781454	Scorers Qualified
	Qual 1	Qual 2		Qual 1	Qual 2		Qual 1	Qual 2	
Total Passed	22	7	29	11	4	15	15	0	15
Total Failed	7	0	0	4	0	0	0	0	0
Grade 8	664145	664145	Scorers Qualified	666113	666113	Scorers Qualified	787842	787842	Scorers Qualified
	Qual 1	Qual 2		Qual 1	Qual 2		Qual 1	Qual 2	
Total Passed	11	1	12	13	1	14	10	3	13
Total Failed	3	2	2	1	0	0	4	1	1
Grade 11	733046	733046	Scorers Qualified	663619	663619	Scorers Qualified	735361	735361	Scorers Qualified
	Qual 1	Qual 2		Qual 1	Qual 2		Qual 1	Qual 2	
Total Passed	7	6	13	12	3	15	13	2	15
Total Failed	8	2	2	3	0	0	2	0	0

APPENDIX F

CLASSICAL ITEM STATISTICS

Calculations are based on those students attempting 5 or more items on the form of the given BIE assessment. For 1-point items, the item-total correlation is the point-biserial. For 2 or more-point items, the item-total correlation is the point-polyserial.

All Classical Item Statistics listed in Appendix F were based on a national sample of students.

Table F-1. Classical Item Statistics for the Operational Items on Science Grade 5*

PsyltemNumber	Position	Item Type	Max Points	Item Mean	Item-Total Correlation
629699	1	MC	1	0.47	0.21
629707	2	MC	1	0.18	0.03
638349	3	MC	1	0.15	0.10
706738	4	MC	1	0.41	0.25
706792	5	MC	1	0.33	0.19
755474	6	MC	1	0.33	0.33
BIE100197	7	MC	1	0.45	0.34
638420	8	OR	1	0.17	0.40
706722	9	OR	1	0.17	0.42
706801	10	OR	1	0.32	0.53
755477	11	OR	1	0.32	0.46
BIE100191	12	OR	1	0.34	0.30
626442	13	OR	2	1.21	0.46
629703	14	OR	2	0.75	0.30
629711	15	OR	2	0.33	0.16
632444	16	OR	2	1.16	0.54
632570	17	OR	2	0.92	0.44
633861	18	OR	2	0.49	0.26
633993	19	OR	2	1.17	0.57
635886	20	OR	2	0.32	0.14
636880	21	OR	2	0.59	0.43
637796	22	OR	2	1.15	0.56
638324	23	OR	2	0.52	0.12
639510	24	OR	2	0.54	0.41
639571	25	OR	2	1.09	0.59
697027	26	OR	2	0.97	0.41
706726	27	OR	2	0.66	0.41
706747	28	OR	2	0.98	0.60
706765	29	OR	2	0.81	0.45
706847	30	OR	2	0.65	0.45
715507	31	OR	2	0.76	0.46
749674	32	OR	2	1.12	0.45
752774	33	OR	2	1.20	0.53
755479	34	OR	2	1.09	0.58
762478	35	OR	2	0.92	0.39
781469	36	OR	2	1.18	0.48
798694	37	OR	2	0.55	0.33
799604	38	OR	2	0.54	0.57
BIE100185	39	OR	2	0.42	0.20
BIE100193	40	OR	2	0.71	0.38
697063	41	OR	4	1.15	0.60
697164	42	OR	4	0.59	0.63
781454	43	OR	4	0.64	0.53

* Calculations are based on those students attempting 5 or more items on the form of the given BIE assessment. For 1-point items, the item-total correlation is the point-biserial. For 2 or more-point items, the item-total correlation is the point-polyserial.

Table F-2. Classical Item Statistics for the Operational Items on Science Grade 8*

PsyltemNumber	Position	Item Type	Max Points	Item Mean	Item-Total Correlation
636837	1	MC	1	0.33	0.36
636843	2	MC	1	0.36	0.28
711815	3	MC	1	0.45	0.23
711825	4	MC	1	0.39	0.40
741231	5	MC	1	0.47	0.22
779991	6	MC	1	0.50	0.30
780001	7	MC	1	0.32	0.28
710124	8	OR	1	0.28	0.17
713686	9	OR	1	0.37	0.37
713695	10	OR	1	0.16	0.38
741227	11	OR	1	0.06	0.18
750828	12	OR	1	0.12	0.41
631375	13	OR	2	0.78	0.37
631831	14	OR	2	1.04	0.37
635564	15	OR	2	0.75	0.39
636830	16	OR	2	1.02	0.46
636852	17	OR	2	0.70	0.38
637630	18	OR	2	0.59	0.51
639667	19	OR	2	0.77	0.46
641096	20	OR	2	0.95	0.28
697271	21	OR	2	0.33	0.37
697312	22	OR	2	0.45	0.13
697316	23	OR	2	0.80	0.47
707172	24	OR	2	0.72	0.39
711819	25	OR	2	0.61	0.35
711823	26	OR	2	0.52	0.23
713386	27	OR	2	0.65	0.42
713388	28	OR	2	0.82	0.28
741223	29	OR	2	0.58	0.29
741225	30	OR	2	1.06	0.54
741317	31	OR	2	0.68	0.26
758864	32	OR	2	0.44	0.36
763145	33	OR	2	0.66	0.28
763149	34	OR	2	0.61	0.32
779995	35	OR	2	0.65	0.46
780003	36	OR	2	0.61	0.32
780122	37	OR	2	0.40	0.26
799732	38	OR	2	0.83	0.48
800020	39	OR	2	0.49	0.24
850712	40	OR	2	0.92	0.33
664145	41	OR	4	0.55	0.55
666113	42	OR	4	0.49	0.53
787842	43	OR	4	0.48	0.43

* Calculations are based on those students attempting 5 or more items on the form of the given BIE assessment. For 1-point items, the item-total correlation is the point-biserial. For 2 or more-point items, the item-total correlation is the point-polyserial.

Table F-3. Classical Item Statistics for the Operational Items on Science Grade 11*

PsyltemNumber	Position	Item Type	Max Points	Item Mean	Item-Total Correlation
637964	1	MC	1	0.27	0.28
637994	2	MC	1	0.33	0.26
638025	3	MC	1	0.38	0.15
639319	4	MC	1	0.39	0.43
642634	5	MC	1	0.34	0.25
643598	6	MC	1	0.46	0.32
709172	7	MC	1	0.40	0.27
746281	8	MC	1	0.28	0.12
707594	9	OR	1	0.15	0.29
707606	10	OR	1	0.19	0.33
709184	11	OR	1	0.07	0.31
746259	12	OR	1	0.41	0.31
624850	13	OR	2	0.88	0.27
627081	14	OR	2	0.55	0.25
631404	15	OR	2	0.39	0.27
634611	16	OR	2	0.56	0.14
635856	17	OR	2	0.89	0.39
637979	18	OR	2	0.65	0.25
638054	19	OR	2	0.58	0.31
639344	20	OR	2	0.98	0.29
639346	21	OR	2	0.48	0.29
640122	22	OR	2	0.51	0.35
640447	23	OR	2	0.47	0.39
640580	24	OR	2	0.43	0.37
641189	25	OR	2	1.03	0.42
641194	26	OR	2	0.81	0.14
641199	27	OR	2	0.46	0.41
642454	28	OR	2	0.45	0.57
642552	29	OR	2	0.92	0.45
705054	30	OR	2	0.52	0.41
706121	31	OR	2	0.53	0.16
706123	32	OR	2	0.65	0.24
707610	33	OR	2	0.59	0.25
707617	34	OR	2	0.21	0.04
709258	35	OR	2	0.69	0.42
709262	36	OR	2	0.20	0.35
717815	37	OR	2	0.90	0.44
746267	38	OR	2	0.71	0.42
746295	39	OR	2	0.61	0.43
752118	40	OR	2	0.66	0.30
755818	41	OR	2	0.83	0.36
762012	42	OR	2	0.75	0.43
798143	43	OR	2	0.52	0.35
798455	44	OR	2	0.58	0.29
663619	45	OR	4	0.51	0.39
733046	46	OR	4	0.30	0.46
735361	47	OR	4	0.44	0.55

* Calculations are based on those students attempting 5 or more items on the form of the given BIE assessment. For 1-point items, the item-total correlation is the point-biserial. For 2 or more-point items, the item-total correlation is the point-polyserial.

APPENDIX G
ITEM RESPONSE THEORY PARAMETERS

Table G-1. IRT Parameters for Operational Items on the BIE Grade 5 Science Assessment

IREF	a	b	c	d0	d1	d2
629699	0.32902	0.01843	0			
629707	0.1823	4.63971	0			
638349	0.73025	1.79295	0.09797			
706738	0.70007	0.65128	0.18019			
706792	0.83975	1.37537	0.26479			
755474	0.62415	1.35728	0.05716			
BIE100197	0.65824	0.17106	0.1499			
638420	0.60293	0.95311	0			
706722	0.89085	0.9686	0			
706801	0.99637	0.32987	0			
755477	0.69609	0.20616	0			
BIE100191	0.57336	-0.22333	0			
626442	0.73417	-0.85181	0	0.68121	-0.68121	0
629703	0.38169	1.0554	0	1.76034	-1.76034	0
629711	0.31293	4.20951	0	3.13697	-3.13697	0
632444	0.69115	-0.74519	0	0.88335	-0.88335	0
632570	0.80927	-0.53795	0	0.77858	-0.77858	0
633861	0.44441	1.50359	0	1.34344	-1.34344	0
633993	0.81904	-0.91899	0	0.46635	-0.46635	0
635886	0.24564	4.63068	0	3.04072	-3.04072	0
636880	0.7288	0.50971	0	0.85312	-0.85312	0
637796	0.83618	-1.03236	0	0.33649	-0.33649	0
638324	0.17996	4.31528	0	2.13197	-2.13197	0
639510	0.74656	0.6709	0	0.97387	-0.97387	0
639571	0.97911	-0.60144	0	0.77809	-0.77809	0
697027	0.64913	-0.42395	0	0.96675	-0.96675	0
706726	0.66281	0.44392	0	0.9297	-0.9297	0
706747	0.67866	-0.25679	0	0.994	-0.994	0
706765	0.79008	-0.08108	0	0.71768	-0.71768	0
706847	0.66793	0.7281	0	0.98715	-0.98715	0
715507	0.65414	-0.22633	0	0.73017	-0.73017	0
749674	0.63499	-1.00514	0	1.17411	-1.17411	0
752774	0.98878	-0.72263	0	0.60063	-0.60063	0
755479	1.24091	-0.69001	0	0.4044	-0.4044	0
762478	0.64481	-0.03493	0	0.9875	-0.9875	0
781469	0.70224	-0.80778	0	1.05637	-1.05637	0
798694	0.52071	1.3203	0	1.55999	-1.55999	0
799604	1.05607	0.18789	0	0.56868	-0.56868	0
BIE100185	0.37478	3.1181	0	2.5059	-2.5059	0
BIE100193	0.50915	0.59051	0	1.35874	-1.35874	0
697063	0.63796	1.88408	0	3.15837	1.15146	-1.02845
697164	0.93395	1.37909	0	1.89904	0.64813	-0.51044
781454	0.78062	2.06197	0	2.27376	0.93157	-0.48343

Table G-2. IRT Parameters for Operational Items on the BIE Grade 8 Science Assessment

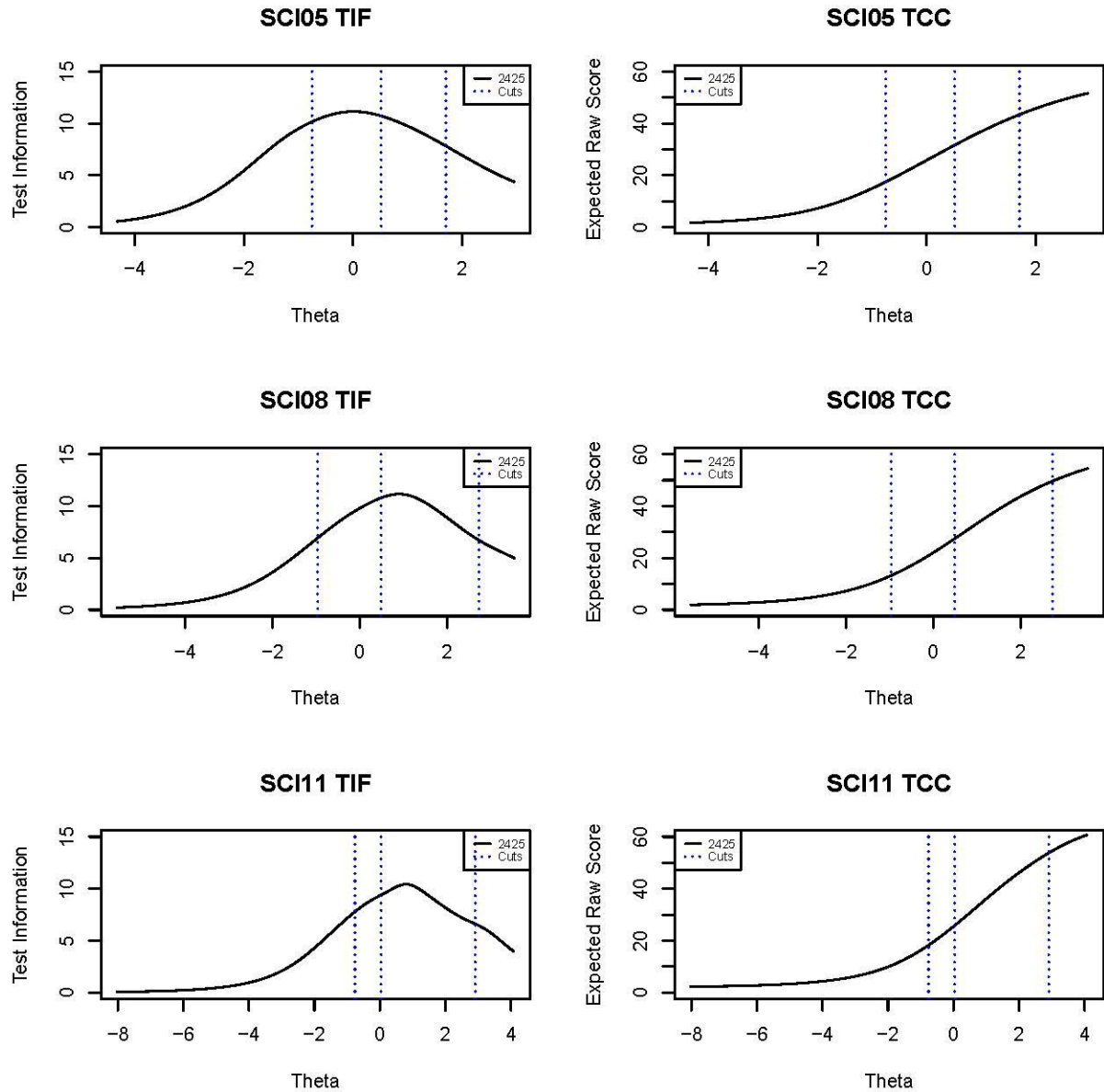
IREF	a	b	c	d0	d1	d2
636837	1.24502	0.88903	0.20065			
636843	0.82189	1.19908	0.23931			
711815	0.6411	0.45125	0.21615			
711825	1.05272	0.6885	0.13813			
741231	0.77176	-0.11035	0.08947			
779991	0.71876	0.74161	0.25045			
780001	0.75441	1.71203	0.25317			
710124	0.29583	1.62795	0			
713686	0.59822	0.11082	0			
713695	0.73355	1.1572	0			
741227	0.51411	3.08932	0			
750828	0.82649	1.64379	0			
631375	0.53477	0.28839	0	1.14789	-1.14789	0
631831	0.75139	-0.25616	0	1.43286	-1.43286	0
635564	0.67645	0.75564	0	1.79265	-1.79265	0
636830	0.58204	-0.2841	0	1.11864	-1.11864	0
636852	0.48058	0.62908	0	1.14652	-1.14652	0
637630	0.82024	0.28781	0	0.85137	-0.85137	0
639667	0.81322	0.08684	0	0.7298	-0.7298	0
641096	0.35837	-0.04348	0	2.42557	-2.42557	0
697271	0.59365	1.85517	0	1.75458	-1.75458	0
697312	0.21229	3.52501	0	3.17578	-3.17578	0
697316	0.88244	0.159	0	0.81557	-0.81557	0
707172	0.51382	0.9717	0	2.13148	-2.13148	0
711819	0.76588	0.51149	0	0.93123	-0.93123	0
711823	0.27788	3.45677	0	3.06823	-3.06823	0
713386	0.67185	0.37076	0	0.99105	-0.99105	0
713388	0.3141	0.88684	0	2.23302	-2.23302	0
741223	0.50664	0.52453	0	1.06758	-1.06758	0
741225	0.89123	-0.6481	0	0.62649	-0.62649	0
741317	0.46459	1.07971	0	1.89346	-1.89346	0
758864	0.37045	2.17733	0	1.99225	-1.99225	0
763145	0.42437	0.43832	0	2.0433	-2.0433	0
763149	0.42027	1.16028	0	1.6029	-1.6029	0
779995	0.67177	0.88442	0	1.08003	-1.08003	0
780003	0.57682	0.95361	0	1.2511	-1.2511	0
780122	0.398	3.08755	0	2.04427	-2.04427	0
799732	0.74962	0.00184	0	0.85063	-0.85063	0
800020	0.31993	2.47266	0	2.14856	-2.14856	0
850712	0.52004	0.12594	0	1.35656	-1.35656	0
664145	0.93651	1.62419	0	1.75601	0.61427	-0.52505
666113	1.11228	1.6174	0	1.61593	0.60996	-0.35253
787842	0.86613	2.55089	0	2.21609	0.89655	-0.57936

Table G-3. IRT Parameters for Operational Items on the BIE Grade 11 Science Assessment

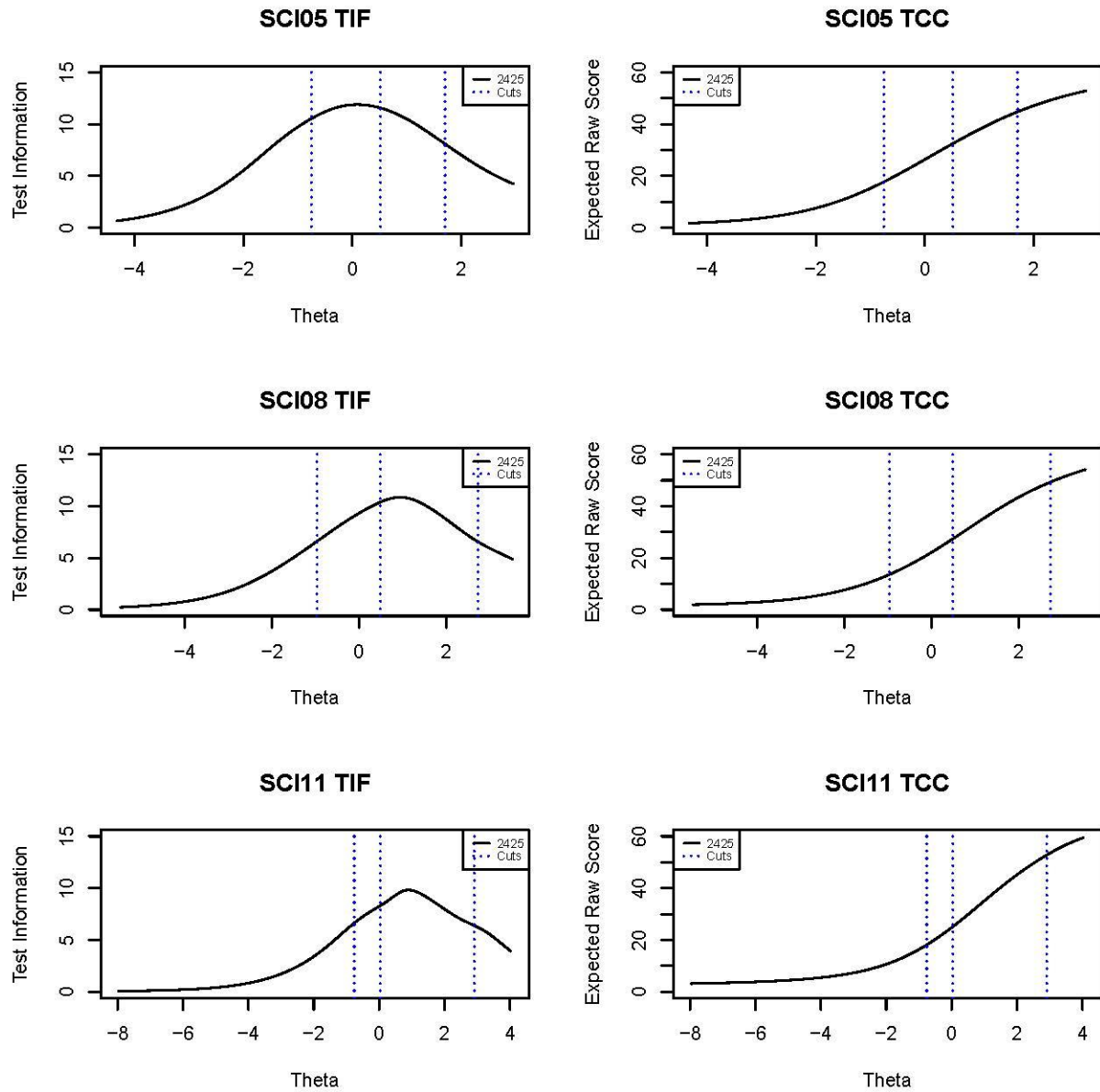
IREF	a	b	c	d0	d1	d2
637964	0.77636	1.07849	0.16874			
637994	0.70552	1.87717	0.27908			
638025	1.35714	3.12685	0.33508			
639319	0.88077	0.0358	0.11812			
642634	0.58688	1.64355	0.2296			
643598	0.80955	0.10797	0.26409			
709172	1.66626	0.72491	0.32706			
746281	0.73138	2.62864	0.20961			
707594	0.85505	1.24509	0			
707606	0.41482	1.70112	0			
709184	0.68279	1.73125	0			
746259	0.27151	-0.11743	0			
624850	0.33908	0.1728	0	2.29196	-2.29196	0
627081	0.63793	0.93026	0	1.13172	-1.13172	0
631404	0.60635	1.3411	0	1.19485	-1.19485	0
634611	0.11677	0.27816	0	4.29543	-4.29543	0
635856	0.38331	-0.59926	0	1.89634	-1.89634	0
637979	0.32127	1.26504	0	1.95506	-1.95506	0
638054	0.56303	0.14578	0	1.04978	-1.04978	0
639344	0.48912	-0.33432	0	1.39477	-1.39477	0
639346	0.57702	1.0448	0	1.11586	-1.11586	0
640122	0.58998	0.51542	0	1.16922	-1.16922	0
640447	0.41864	1.87538	0	2.15293	-2.15293	0
640580	0.588	0.84799	0	1.40344	-1.40344	0
641189	0.89683	-0.5705	0	0.77725	-0.77725	0
641194	0.08593	0.27816	0	5.2606	-5.2606	0
641199	0.74696	0.15667	0	0.6845	-0.6845	0
642454	0.89617	0.64585	0	0.836	-0.836	0
642552	0.63299	-0.80587	0	0.86637	-0.86637	0
705054	0.45085	0.72373	0	1.15238	-1.15238	0
706121	0.37378	1.52302	0	1.69456	-1.69456	0
706123	0.50857	0.84938	0	1.48219	-1.48219	0
707610	0.27459	2.03186	0	2.8365	-2.8365	0
707617	0.16661	0.27816	0	1.08179	-1.08179	0
709258	0.54755	0.18406	0	2.12582	-2.12582	0
709262	0.62682	1.96371	0	1.12534	-1.12534	0
717815	0.90645	-0.82914	0	0.65242	-0.65242	0
746267	0.61007	-0.13111	0	1.22359	-1.22359	0
746295	0.94342	0.09914	0	0.85605	-0.85605	0
752118	0.46448	0.28031	0	1.65913	-1.65913	0
755818	0.60318	-0.1724	0	1.09399	-1.09399	0
762012	0.61362	0.01422	0	1.06998	-1.06998	0
798143	0.3527	2.35956	0	1.88531	-1.88531	0
798455	0.26101	1.51122	0	2.15881	-2.15881	0
663619	0.77372	2.15365	0	2.67225	0.95451	-0.7654
733046	0.86662	2.05307	0	1.78199	0.56599	-0.49912
735361	0.94616	1.96296	0	2.2039	0.13724	-0.75723

APPENDIX H
TEST CHARACTERISTIC CURVES
AND
TEST INFORMATION FUNCTIONS

Scale Form 1



Scale Form 2



APPENDIX I
RAW TO SCALED SCORE LOOKUP TABLES

Table I-1. Raw to Scaled Score Look-up Table—Science Grade 5 (Scale Form 1)

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
0	-4.32534	500	1	1.34082	16.76021
1	-4.15431	502	1	1.22976	15.37203
2	-3.98328	504	1	1.12699	14.08735
3	-3.23962	513	1	0.76813	9.60162
4	-2.80423	519	1	0.61652	7.70647
5	-2.49351	522	1	0.53088	6.63602
6	-2.24939	525	1	0.47541	5.94266
7	-2.04629	528	1	0.43656	5.45694
8	-1.87075	530	1	0.40797	5.09960
9	-1.71480	532	1	0.38625	4.82813
10	-1.57338	534	1	0.36939	4.61742
11	-1.44304	536	1	0.35611	4.45136
12	-1.32135	537	1	0.34551	4.31894
13	-1.20652	538	1	0.33697	4.21219
14	-1.09720	540	1	0.33001	4.12514
15	-0.99238	541	1	0.32426	4.05324
16	-0.89124	542	1	0.31944	3.99306
17	-0.79314	543	1	0.31536	3.94205
18	-0.69756	545	2	0.31188	3.89845
19	-0.60406	546	2	0.30889	3.86111
20	-0.51228	547	2	0.30635	3.82936
21	-0.42190	548	2	0.30423	3.80285
22	-0.33263	549	2	0.30251	3.78143
23	-0.24422	551	2	0.30120	3.76504
24	-0.15643	552	2	0.30029	3.75364
25	-0.06905	553	2	0.29977	3.74715
26	0.01812	554	2	0.29964	3.74545
27	0.10529	555	2	0.29987	3.74838
28	0.19262	556	2	0.30046	3.75578
29	0.28030	557	2	0.30140	3.76749
30	0.36849	558	2	0.30267	3.78336
31	0.45737	559	2	0.30427	3.80332
32	0.54711	560	3	0.30619	3.82733
33	0.63789	562	3	0.30843	3.85543
34	0.72990	563	3	0.31102	3.88773
35	0.82334	564	3	0.31395	3.92440
36	0.91844	565	3	0.31726	3.96572
37	1.01543	566	3	0.32096	4.01203
38	1.11457	567	3	0.32510	4.06378
39	1.21617	569	3	0.32972	4.12153
40	1.32057	570	3	0.33488	4.18603
41	1.42816	571	3	0.34065	4.25817
42	1.53939	573	3	0.34713	4.33912
43	1.65478	573	3	0.35443	4.43032
44	1.77496	576	4	0.36268	4.53353
45	1.90067	577	4	0.37207	4.65087
46	2.03279	579	4	0.38279	4.78483
47	2.17239	581	4	0.39506	4.93820
48	2.32073	583	4	0.40912	5.11396
49	2.47935	585	4	0.42521	5.31509
50	2.65007	587	4	0.44356	5.54453
51	2.83510	589	4	0.46443	5.80540
52	2.95466	590	4	0.47837	5.97960

continued

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
53	2.95466	590	4	0.47837	5.97960
54	2.95466	590	4	0.47837	5.97960
55	2.95466	590	4	0.47837	5.97960
56	2.95466	590	4	0.47837	5.97960
57	2.95466	590	4	0.47837	5.97960
58	2.95466	590	4	0.47837	5.97960
59	2.95466	590	4	0.47837	5.97960
60	2.95466	590	4	0.47837	5.97960
61	2.95466	590	4	0.47837	5.97960
62	2.95466	590	4	0.47837	5.97960
63	2.95466	590	4	0.47837	5.97960
64	2.95466	590	4	0.47837	5.97960

Table I-2. Raw to Scaled Score Look-up Table—Science Grade 5 (Scale Form 2)

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
0	-4.32534	500	1	1.23407	15.42584
1	-4.19121	501	1	1.15551	14.44384
2	-4.05707	503	1	1.08196	13.52444
3	-3.33094	512	1	0.76100	9.51245
4	-2.89148	517	1	0.61999	7.74986
5	-2.57213	521	1	0.53770	6.72129
6	-2.31878	525	1	0.48283	6.03538
7	-2.10708	527	1	0.44338	5.54229
8	-1.92394	530	1	0.41369	5.17118
9	-1.76149	532	1	0.39069	4.88361
10	-1.61463	533	1	0.37251	4.65641
11	-1.47984	535	1	0.35796	4.47444
12	-1.35458	537	1	0.34616	4.32705
13	-1.23701	538	1	0.33651	4.20640
14	-1.12569	539	1	0.32852	4.10648
15	-1.01953	541	1	0.32182	4.02271
16	-0.91769	542	1	0.31613	3.95158
17	-0.81945	543	1	0.31124	3.89054
18	-0.72427	545	2	0.30702	3.83777
19	-0.63166	546	2	0.30336	3.79206
20	-0.54122	547	2	0.30021	3.75265
21	-0.45261	548	2	0.29753	3.71909
22	-0.36552	549	2	0.29529	3.69108
23	-0.27965	550	2	0.29347	3.66839
24	-0.19476	551	2	0.29207	3.65082
25	-0.11062	552	2	0.29105	3.63814
26	-0.02701	553	2	0.29041	3.63010
27	0.05627	554	2	0.29012	3.62644
28	0.13941	555	2	0.29016	3.62694
29	0.22259	556	2	0.29051	3.63140
30	0.30598	557	2	0.29118	3.63969
31	0.38977	558	2	0.29214	3.65179
32	0.47412	559	2	0.29342	3.66777
33	0.55922	561	3	0.29502	3.68780
34	0.64527	562	3	0.29697	3.71216
35	0.73248	563	3	0.29930	3.74122
36	0.82106	564	3	0.30203	3.77543
37	0.91128	565	3	0.30523	3.81533
38	1.00341	566	3	0.30892	3.86149

continued

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
39	1.09775	567	3	0.31317	3.91457
40	1.19464	568	3	0.31802	3.97527
41	1.29449	570	3	0.32355	4.04442
42	1.39772	571	3	0.32983	4.12293
43	1.50483	572	3	0.33695	4.21188
44	1.61642	573	3	0.34501	4.31261
45	1.73315	575	4	0.35414	4.42671
46	1.85582	577	4	0.36449	4.55615
47	1.98537	578	4	0.37626	4.70330
48	2.12294	580	4	0.38968	4.87094
49	2.26988	582	4	0.40497	5.06216
50	2.42784	584	4	0.42242	5.28020
51	2.59881	586	4	0.44226	5.52822
52	2.78519	588	4	0.46475	5.80941
53	2.95466	590	4	0.48579	6.07239
54	2.95466	590	4	0.48579	6.07239
55	2.95466	590	4	0.48579	6.07239
56	2.95466	590	4	0.48579	6.07239
57	2.95466	590	4	0.48579	6.07239
58	2.95466	590	4	0.48579	6.07239
59	2.95466	590	4	0.48579	6.07239
60	2.95466	590	4	0.48579	6.07239
61	2.95466	590	4	0.48579	6.07239
62	2.95466	590	4	0.48579	6.07239
63	2.95466	590	4	0.48579	6.07239
64	2.95466	590	4	0.48579	6.07239

Table I-3. Raw to Scaled Score Look-up Table—Science Grade 8 (Scale Form 1)

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
0	-5.56012	800	1	2.11505	21.15055
1	-5.50756	800	1	2.07661	20.76606
2	-5.45501	801	1	2.03873	20.38728
3	-3.91021	816	1	1.14655	11.46550
4	-3.17721	823	1	0.85023	8.50232
5	-2.70059	828	1	0.69669	6.96688
6	-2.34787	832	1	0.60263	6.02631
7	-2.06678	834	1	0.53944	5.39435
8	-1.83168	837	1	0.49431	4.94314
9	-1.62823	839	1	0.46063	4.60629
10	-1.44773	841	1	0.43459	4.34590
11	-1.28451	842	1	0.41388	4.13883
12	-1.13472	844	1	0.39703	3.97029
13	-0.99564	844	1	0.38305	3.83052
14	-0.86526	846	2	0.37129	3.71285
15	-0.74206	848	2	0.36127	3.61265
16	-0.62486	849	2	0.35266	3.52658
17	-0.51272	850	2	0.34521	3.45214
18	-0.40487	851	2	0.33874	3.38741
19	-0.30069	852	2	0.33308	3.33081
20	-0.19964	853	2	0.32810	3.28103
21	-0.10129	854	2	0.32370	3.23697
22	-0.00524	855	2	0.31977	3.19767
23	0.08884	856	2	0.31624	3.16237
24	0.18124	857	2	0.31305	3.13048

continued

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
25	0.27221	858	2	0.31016	3.10165
26	0.36200	859	2	0.30757	3.07572
27	0.45083	859	2	0.30527	3.05274
28	0.53890	860	3	0.30329	3.03293
29	0.62644	861	3	0.30166	3.01664
30	0.71363	862	3	0.30043	3.00429
31	0.80071	863	3	0.29963	2.99630
32	0.88787	864	3	0.29931	2.99307
33	0.97535	865	3	0.29949	2.99491
34	1.06338	866	3	0.30020	3.00205
35	1.15219	867	3	0.30146	3.01459
36	1.24203	868	3	0.30326	3.03257
37	1.33314	868	3	0.30559	3.05592
38	1.42579	869	3	0.30846	3.08458
39	1.52025	870	3	0.31185	3.11849
40	1.61681	871	3	0.31577	3.15768
41	1.71578	872	3	0.32023	3.20231
42	1.81749	873	3	0.32527	3.25269
43	1.92232	874	3	0.33093	3.30928
44	2.03070	875	3	0.33726	3.37265
45	2.14308	877	3	0.34433	3.44331
46	2.25999	878	3	0.35216	3.52160
47	2.38199	879	3	0.36074	3.60744
48	2.50972	880	3	0.37002	3.70018
49	2.64381	881	3	0.37986	3.79861
50	2.78498	883	4	0.39013	3.90135
51	2.93400	884	4	0.40078	4.00778
52	3.09180	886	4	0.41196	4.11961
53	3.25966	888	4	0.42428	4.24275
54	3.43944	889	4	0.43888	4.38885
55	3.53988	890	4	0.44809	4.48093
56	3.53988	890	4	0.44809	4.48093
57	3.53988	890	4	0.44809	4.48093
58	3.53988	890	4	0.44809	4.48093
59	3.53988	890	4	0.44809	4.48093
60	3.53988	890	4	0.44809	4.48093
61	3.53988	890	4	0.44809	4.48093
62	3.53988	890	4	0.44809	4.48093
63	3.53988	890	4	0.44809	4.48093
64	3.53988	890	4	0.44809	4.48093

Table I-4. Raw to Scaled Score Look-up Table—Science Grade 8 (Scale Form 2)

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
0	-5.56012	800	1	2.09083	20.90832
1	-5.50309	800	1	2.04594	20.45941
2	-5.44605	801	1	2.00186	20.01862
3	-3.94522	816	1	1.10204	11.02042
4	-3.24445	823	1	0.82762	8.27623
5	-2.78310	827	1	0.68865	6.88649
6	-2.43590	831	1	0.60363	6.03626
7	-2.15486	834	1	0.54602	5.46025
8	-1.91667	836	1	0.50436	5.04362
9	-1.70836	838	1	0.47278	4.72778
10	-1.52202	840	1	0.44794	4.47944

continued

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
11	-1.35249	842	1	0.42783	4.27827
12	-1.19624	843	1	0.41113	4.11131
13	-1.05074	844	1	0.39700	3.96999
14	-0.91411	846	2	0.38486	3.84859
15	-0.78492	847	2	0.37432	3.74316
16	-0.66204	848	2	0.36509	3.65091
17	-0.54454	850	2	0.35698	3.56980
18	-0.43169	851	2	0.34982	3.49821
19	-0.32283	852	2	0.34348	3.43481
20	-0.21743	853	2	0.33784	3.37842
21	-0.11503	854	2	0.33280	3.32798
22	-0.01523	855	2	0.32826	3.28258
23	0.08233	856	2	0.32414	3.24144
24	0.17795	857	2	0.32040	3.20400
25	0.27191	858	2	0.31699	3.16988
26	0.36446	859	2	0.31390	3.13898
27	0.45585	859	2	0.31114	3.11139
28	0.54629	861	3	0.30874	3.08740
29	0.63602	861	3	0.30674	3.06739
30	0.72525	862	3	0.30518	3.05185
31	0.81422	863	3	0.30412	3.04121
32	0.90316	864	3	0.30359	3.03587
33	0.99231	865	3	0.30361	3.03612
34	1.08191	866	3	0.30422	3.04215
35	1.17222	867	3	0.30540	3.05401
36	1.26349	868	3	0.30717	3.07167
37	1.35598	869	3	0.30950	3.09505
38	1.44998	870	3	0.31240	3.12403
39	1.54577	871	3	0.31586	3.15858
40	1.64366	872	3	0.31988	3.19876
41	1.74397	873	3	0.32448	3.24478
42	1.84708	874	3	0.32970	3.29704
43	1.95338	875	3	0.33561	3.35608
44	2.06334	876	3	0.34225	3.42255
45	2.17746	877	3	0.34970	3.49699
46	2.29632	878	3	0.35797	3.57970
47	2.42054	879	3	0.36705	3.67049
48	2.55081	881	3	0.37685	3.76853
49	2.68785	881	3	0.38725	3.87245
50	2.83245	883	4	0.39809	3.98086
51	2.98549	885	4	0.40936	4.09358
52	3.14807	887	4	0.42135	4.21350
53	3.32174	888	4	0.43486	4.34863
54	3.50878	889	4	0.45137	4.51366
55	3.53988	890	4	0.45439	4.54391
56	3.53988	890	4	0.45439	4.54391
57	3.53988	890	4	0.45439	4.54391
58	3.53988	890	4	0.45439	4.54391
59	3.53988	890	4	0.45439	4.54391
60	3.53988	890	4	0.45439	4.54391
61	3.53988	890	4	0.45439	4.54391
62	3.53988	890	4	0.45439	4.54391
63	3.53988	890	4	0.45439	4.54391
64	3.53988	890	4	0.45439	4.54391

Table I-5. Raw to Scaled Score Look-up Table—Science Grade 11 (Scale Form 1)

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
0	-8.02951	1100	1	3.95194	29.63952
1	-7.20235	1106	1	3.08631	23.14732
2	-6.37519	1112	1	2.37570	17.81772
3	-5.54803	1118	1	1.80012	13.50093
4	-4.27837	1128	1	1.13744	8.53082
5	-3.57885	1133	1	0.86924	6.51933
6	-3.10039	1136	1	0.72093	5.40700
7	-2.73741	1139	1	0.62663	4.69969
8	-2.44424	1141	1	0.56177	4.21327
9	-2.19715	1143	1	0.51487	3.86151
10	-1.98231	1145	1	0.47970	3.59774
11	-1.79108	1146	1	0.45254	3.39408
12	-1.61770	1148	1	0.43104	3.23280
13	-1.45825	1149	1	0.41362	3.10214
14	-1.30988	1150	1	0.39922	2.99411
15	-1.17053	1151	1	0.38710	2.90328
16	-1.03864	1152	1	0.37678	2.82587
17	-0.91298	1153	1	0.36791	2.75930
18	-0.79260	1153	1	0.36023	2.70176
19	-0.67670	1155	2	0.35359	2.65193
20	-0.56464	1155	2	0.34784	2.60878
21	-0.45587	1156	2	0.34286	2.57147
22	-0.34994	1157	2	0.33855	2.53916
23	-0.24644	1158	2	0.33480	2.51103
24	-0.14505	1159	2	0.33149	2.48618
25	-0.04549	1159	2	0.32849	2.46367
26	0.05249	1160	3	0.32568	2.44258
27	0.14910	1161	3	0.32295	2.42216
28	0.24452	1162	3	0.32026	2.40197
29	0.33894	1162	3	0.31761	2.38207
30	0.43254	1163	3	0.31508	2.36311
31	0.52550	1164	3	0.31283	2.34620
32	0.61804	1164	3	0.31103	2.33275
33	0.71038	1165	3	0.30988	2.32410
34	0.80273	1166	3	0.30950	2.32125
35	0.89534	1166	3	0.30995	2.32466
36	0.98843	1167	3	0.31123	2.33422
37	1.08219	1168	3	0.31325	2.34934
38	1.17681	1169	3	0.31589	2.36917
39	1.27244	1169	3	0.31903	2.39273
40	1.36922	1170	3	0.32256	2.41917
41	1.46729	1171	3	0.32638	2.44786
42	1.56678	1171	3	0.33045	2.47836
43	1.66782	1172	3	0.33473	2.51048
44	1.77056	1173	3	0.33922	2.54414
45	1.87515	1174	3	0.34391	2.57930
46	1.98177	1175	3	0.34878	2.61586
47	2.09061	1175	3	0.35382	2.65364
48	2.20187	1176	3	0.35897	2.69225
49	2.31577	1177	3	0.36416	2.73117
50	2.43259	1178	3	0.36931	2.76982
51	2.55261	1179	3	0.37436	2.80773
52	2.67621	1180	3	0.37933	2.84500

continued

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
53	2.80390	1180	3	0.38437	2.88278
54	2.93632	1182	4	0.38984	2.92380
55	3.07442	1183	4	0.39635	2.97264
56	3.21946	1184	4	0.40474	3.03552
57	3.37318	1185	4	0.41596	3.11972
58	3.53796	1186	4	0.43105	3.23290
59	3.71694	1188	4	0.45105	3.38289
60	3.91441	1189	4	0.47714	3.57852
61	4.10383	1190	4	0.50569	3.79267
62	4.10383	1190	4	0.50569	3.79267
63	4.10383	1190	4	0.50569	3.79267
64	4.10383	1190	4	0.50569	3.79267
65	4.10383	1190	4	0.50569	3.79267
66	4.10383	1190	4	0.50569	3.79267
67	4.10383	1190	4	0.50569	3.79267
68	4.10383	1190	4	0.50569	3.79267

Table I-6. Raw to Scaled Score Look-up Table—Science Grade 11 (Scale Form 2)

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
0	-8.02951	1100	1	3.84255	28.81915
1	-7.44858	1104	1	3.32974	24.97303
2	-6.86766	1108	1	2.84381	21.32860
3	-6.28674	1113	1	2.39402	17.95511
4	-5.70581	1117	1	1.98827	14.91205
5	-4.40833	1127	1	1.26456	9.48418
6	-3.65821	1132	1	0.96248	7.21859
7	-3.14032	1136	1	0.79869	5.99018
8	-2.74546	1139	1	0.69518	5.21385
9	-2.42553	1142	1	0.62313	4.67350
10	-2.15565	1144	1	0.56968	4.27257
11	-1.92139	1145	1	0.52826	3.96195
12	-1.71368	1147	1	0.49523	3.71424
13	-1.52638	1148	1	0.46838	3.51289
14	-1.35518	1150	1	0.44630	3.34721
15	-1.19688	1151	1	0.42798	3.20986
16	-1.04908	1152	1	0.41273	3.09547
17	-0.90991	1153	1	0.39999	2.99993
18	-0.77789	1153	1	0.38933	2.91999
19	-0.65186	1155	2	0.38039	2.85296
20	-0.53086	1156	2	0.37287	2.79653
21	-0.41413	1157	2	0.36649	2.74865
22	-0.30106	1157	2	0.36099	2.70743
23	-0.19115	1158	2	0.35614	2.67104
24	-0.08398	1159	2	0.35170	2.63778
25	0.02075	1159	2	0.34748	2.60608
26	0.12334	1161	3	0.34330	2.57473
27	0.22403	1161	3	0.33908	2.54309
28	0.32303	1162	3	0.33484	2.51126
29	0.42056	1163	3	0.33069	2.48021
30	0.51687	1164	3	0.32687	2.45150
31	0.61220	1164	3	0.32360	2.42697
32	0.70683	1165	3	0.32110	2.40828
33	0.80102	1166	3	0.31954	2.39656

continued

Raw Score	Theta	Scale Score	Performance Level	CSEM Theta	Scaled CSEM
34	0.89505	1166	3	0.31897	2.39224
35	0.98918	1167	3	0.31934	2.39508
36	1.08366	1168	3	0.32058	2.40432
37	1.17869	1169	3	0.32253	2.41896
38	1.27449	1169	3	0.32506	2.43795
39	1.37122	1170	3	0.32806	2.46044
40	1.46905	1171	3	0.33144	2.48578
41	1.56816	1171	3	0.33515	2.51361
42	1.66873	1172	3	0.33917	2.54376
43	1.77094	1173	3	0.34349	2.57618
44	1.87499	1174	3	0.34811	2.61083
45	1.98110	1175	3	0.35302	2.64763
46	2.08947	1175	3	0.35818	2.68633
47	2.20037	1176	3	0.36353	2.72649
48	2.31405	1177	3	0.36900	2.76751
49	2.43078	1178	3	0.37448	2.80863
50	2.55090	1179	3	0.37990	2.84926
51	2.67481	1180	3	0.38524	2.88932
52	2.80300	1180	3	0.39065	2.92984
53	2.93617	1182	4	0.39648	2.97359
54	3.07525	1183	4	0.40338	3.02535
55	3.22155	1184	4	0.41224	3.09178
56	3.37688	1185	4	0.42410	3.18073
57	3.54367	1186	4	0.44007	3.30052
58	3.72521	1188	4	0.46130	3.45975
59	3.92595	1189	4	0.48908	3.66812
60	4.10383	1190	4	0.51709	3.87819
61	4.10383	1190	4	0.51709	3.87819
62	4.10383	1190	4	0.51709	3.87819
63	4.10383	1190	4	0.51709	3.87819
64	4.10383	1190	4	0.51709	3.87819
65	4.10383	1190	4	0.51709	3.87819
66	4.10383	1190	4	0.51709	3.87819
67	4.10383	1190	4	0.51709	3.87819
68	4.10383	1190	4	0.51709	3.87819

APPENDIX J

DECISION ACCURACY AND CONSISTENCY RESULTS

*Calculations based on those students attempting 5 or more items on the given BIE assessment.
Statistic values are suppressed for those subjects/grades with fewer than 50 students. The
consistency statistics are in parentheses.*

Table J-1. Decision Accuracy for BIE Science Forms, as a Function of Grade and Performance Level*

Grade	Scale Form	DA	DC	Kappa	Performance Level							
		Overall	Overall		DA 1	DC 1	DA 2	DC 2	DA 3	DC 3	DA 4**	DC 4**
5	1	0.82	0.75	0.60	0.85	0.79	0.82	0.74	0.76	0.65	0.81	0.60
	2	0.83	0.76	0.60	0.87	0.80	0.81	0.75	0.80	0.66	0.78	0.56
8	1	0.83	0.75	0.54	0.74	0.61	0.84	0.80	0.86	0.73	0.92	0.71
	2	0.82	0.74	0.56	0.80	0.69	0.82	0.77	0.85	0.74	0.58	0.53
11	1	0.80	0.72	0.54	0.85	0.82	0.67	0.54	0.82	0.72		
	2	0.80	0.72	0.51	0.85	0.82	0.64	0.51	0.86	0.70		

* Calculations based on those students attempting 5 or more items on the given BIE assessment. Statistical values are suppressed for those subjects/grades with fewer than 50 students. The consistency statistics are in parentheses.

**Too Few BIE students scored at performance level 4 in 2025.

Table J-2. Decision Consistency for BIE Science Forms, as a Function of Grade and Cut Score*

Grade	Cut Scores Form	1				2				3**			
		DA	DC	FP	FN	DA	DC	FP	FN	DA	DC	FP	FN
5	1	0.89	0.85	0.06	0.05	0.94	0.92	0.03	0.02	0.99	0.98	0.01	0.00
	2	0.89	0.84	0.05	0.06	0.95	0.93	0.03	0.01	0.99	0.99	0.01	0.00
8	1	0.89	0.84	0.05	0.06	0.94	0.91	0.04	0.02	1.00	1.00	0.00	0.00
	2	0.87	0.82	0.05	0.07	0.95	0.92	0.03	0.02	1.00	1.00	0.00	0.00
11	1	0.87	0.81	0.08	0.05	0.93	0.90	0.04	0.03				
	2	0.86	0.80	0.09	0.05	0.94	0.91	0.04	0.02				

* Calculations based on those students attempting 5 or more items on the given BIE assessment. Statistical values are suppressed for those subjects/grades with fewer than 50 students.

**Too Few BIE students scored at performance level 4 (or above cut 3) in 2025.

APPENDIX K
BIE SCIENCE REPORTING BUSINESS REQUIREMENTS



BIE Science Reporting Business Requirements

1/3/2025

Woreen-Ann Bogle

cognia™

Date	Updated Content Description	Updated By
1/3/2025	<i>Initial Creation</i>	<i>W.Bogle</i>
1/24/2025	<i>Incorporate PgM edits</i>	<i>W.Bogle</i>

Glossary	
CBT	Computer Based Test
PBT	Paper Based Test
HS	High School
MC	Multiple Choice
SRB	Student Response Booklet
EL	English Learner
OE	Open Ended also called Open Response items
FT	Field Test
BIE	Bureau of Indian Education

Table of Contents

Overview	5
Points of Contact	5
Document References	5
General Information	5
Assessments	5
Reporting Cycles	5
Receivables	6
Deliverables	6
Reports Produced	6
Data File Deliverables	7
Pre-Test Administration	7
Forms	8
Item Types	9
Reporting Categories	9
Post Test Administration	10
Demographic Clean-up	10
Student Data Processing	10
Test Data Processing	10
Scan Paper Delivery and Data Denotation	11
Scoring Data	11
Student Participation and Exclusions	12
Test Attempt Rules	12
Not Tested Reasons	12
Student Participation Status	13
Student Participation Summary	14
Calculations	14
Aggregations	14
Scoring Method	14
Performance Levels	14
Data Suppression Rules	14
Specific Reporting Rules	15
LENS Online Reporting	15
Student Results Data File	15
eMetric Summary Data File	16
Student Results Labels	16
Student Reports	16

Student Results Rosters	17
Shipping Product Code Summary	17
Reporting Products	17
Appendix	18

Overview

Each year in the Spring, Cognia administers a Science assessment for BIE. The Science assessment utilizes items from Cognia's Secure Science Item Bank (SSIB). The online tests are administered in eMetric's iTester platform. The test is also administered on paper as needed. The test is administered across the country for students in grades 5, 8 and 11.

Points of Contact

Title	Name	Contact Email
IT Project Manager	Sarah McCain	Sarah.McCain@cognia.org
Program Manager	Mara Allaire	Mara.Allaire@cognia.org

Document References

- Student Results data file Layout
 - Student Results data file for both BIE and eMetric.
- Data Processing Specifications
 - Documents details around student and test data processing and staging of data for reporting purposes.

Change Log

Year to Year Changes

1. Human Reader accommodation bubble added to the scannable.
2. Added Off Grade tester to participation status.

General Information

Assessments

1. Assessments were administered to students beginning March 10, 2025, and ending on April 18, 2025.
2. Students are tested online (CBT) and on paper (PBT).
3. Tests are administered in grades 5, 8 and 11.
4. The reporting iCore Contract code is 800450.

Reporting Cycles

There are 2 reporting cycles. Each cycle is described below.

1. The first reporting cycle is prior to the test administration. In this cycle, the outbound rosters are produced. These rosters are produced using the Pre-ID file that contains paper testers. The roster is printed and shipped along with test materials for paper testers.
2. The final reporting cycle will include all reporting deliverables,

Receivables

Data Files received will pass all validation rules and formats based on the layout and specification documents.

1. Cognia was provided a data file by the client for the purpose of Pre-ID prior to the test administration window. This data was uploaded to the iTester platform in LightHouse™. Schools were provided a window to choose mode of testing and accommodations. The data was pulled to generate student identification labels to be affixed to answer documents. This data is also used to create the Outbound rosters which accompany the student identification labels.
2. Cognia will not receive a demographic file for reporting purposes.

Deliverables

Reports Produced

1. Pre-test Administration

- Outbound Roster
 - Lists each student marked for paper testing in the Pre-ID extract from the iTester platform.
 - Demographic fields from the Pre-ID file are reported on the roster.
 - A Roster is produced for each school included in the Pre-ID file.

2. Final Reporting

- Student Results Labels
 - Contains the student's earned scaled score and overall performance level.
 - Printed and shipped to schools.
 - Only Tested students will receive a results label.
 - See Report Specific Rules section for more information.
- Student Report
 - Contains the student's earned scaled score and overall performance level.
 - Contains the performance level earned for each reporting category.
 - Contains a comparison of the student's performance to the student's school, district and to all participating students.
 - Printed and shipped to schools.
 - PDFs provided to eMetric to be made available in the Download Hub of Data Interaction (DI).
 - The PDFs will be accessible in the DI Download Hub at the organization level.

- Only Tested students will receive a student report.
- See Report Specific Rules section for more information.
- School Results Roster
 - Contains a list of all students in the school.
 - Printed and shipped to schools.
 - The student's participation status is included in the report.
 - If the student receives a scale score and achievement level, the results are presented on the results roster.
 - See Report Specific Rules section for more information.

Data File Deliverables

3. Pre-test Administration

- No data file deliverables were produced in the pre-test administration period.

4. Final Reporting

- Student Results Data File
 - Contains the results for each participating student.
 - Participation Status for each student based on test attemptedness rules below.
 - Follows the Student Results data file layout.
 - Delivered to BIE and eMetric via the Cognia SFTP site.
 - eMetric Summary Data File
 - Data file used to QC the summary data in DI.
 - Summaries done for specified demographics.
 - eMetric Metadata file for Student PDFs
 - Data File used to describe the online set up and PDF files being handed off to eMetric to be loaded into the DI Download Hub.
 - See Specific Reporting Rules section for specifications.

Pre-Test Administration

This section describes the data preparation for student records pre-test administration:

- The Pre-ID data file is used to provide answer booklet labels for students in the Pre-ID data file.
 - A total record count will be provided with the final label data to iCore Distribution.
 - Each student label has a unique Barcode associated with a Student ID
 - One student label will be printed for each booklet being administered.
 - Counts are given for accommodations needing materials shipped to the school.
- The Pre-ID data is used to produce the Outbound Rosters that accompany the answer booklet labels.

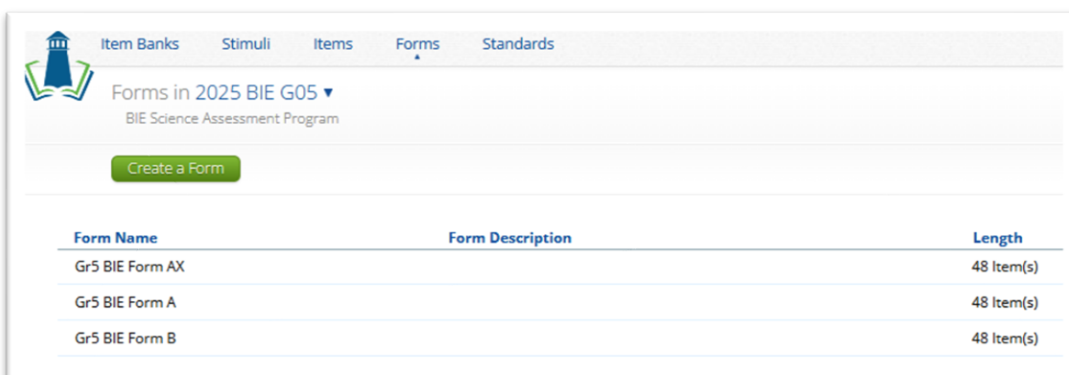
Forms

The test uses items from Cognia's Secure Science Item Bank (SSIB). The test is composed of Operational and Field Test items. The Operational items count toward the students' overall score and are included in subscore calculations.

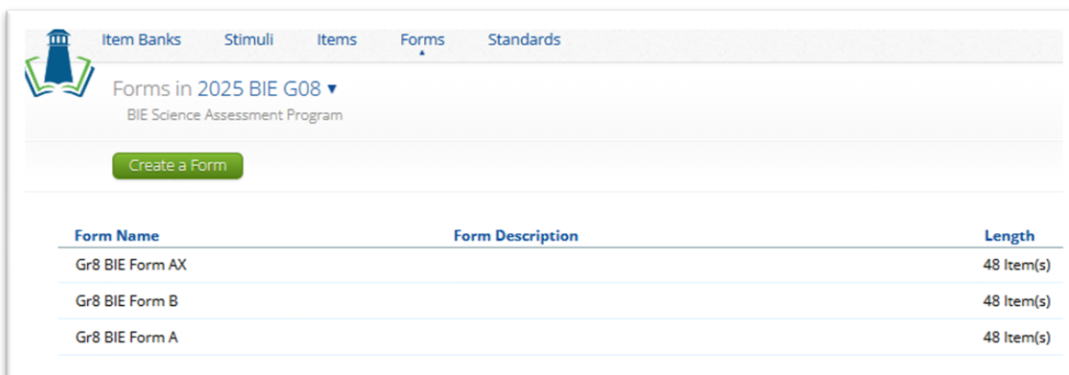
1. There are 2 main Core Forms: Core Form A and Core Form B. The form AX is the accommodated form used for paper replacement, and Large Print. Form AX is used for English TTS.
2. The Braille form is a reuse of the spring 2023 form AX.
3. English Paper accommodated forms (represented by Form AX which is the paper replacement of the online form) are provided for students with the following accommodations:
 - a. Large Print
 - b. Braille
4. Paper accommodated forms contain replacement items for the TEIs on the online Form A test form that becomes Form AX.
5. The online accommodated form has the TTS-accommodation and is CBT Form AX.
6. Tests are provided only in English.

Test Mode	Grade	Forms Administered	Additional Accommodated Form
Online	05	A+B	Yes (Text-to-Speech Form AX (TTS))
	08	A+B	Yes (Text-to-Speech Form AX (TTS))
	11	A+B	Yes (Text-to-Speech Form AX (TTS))
Paper	05, 08, 11	Form AX	Used for large print 2025 Braille = 2023 AX form

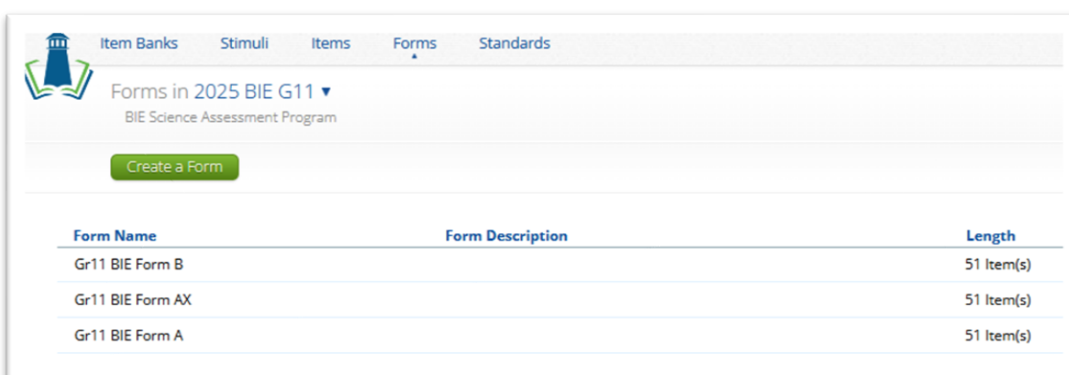
Online Test Forms



Form Name	Form Description	Length
Gr5 BIE Form AX		48 Item(s)
Gr5 BIE Form A		48 Item(s)
Gr5 BIE Form B		48 Item(s)



Form Name	Form Description	Length
Gr8 BIE Form AX		48 Item(s)
Gr8 BIE Form B		48 Item(s)
Gr8 BIE Form A		48 Item(s)



Form Name	Form Description	Length
Gr11 BIE Form B		51 Item(s)
Gr11 BIE Form AX		51 Item(s)
Gr11 BIE Form A		51 Item(s)

Item Types

The following item types are administered in the test:

Item Type	Definition	Valid Point Values	Scoring Method and Scoring Rules
MS-1	machine-scored item that may be multiple choice, multi-select, or TEI interaction	0,1	Machine scored; all or nothing scoring for the interaction, no partial credit
MS-2	machine-scored item with part a, part b; interactions may be any combination of multiple choice, multi-select, or TEI	0,1,2	Machine scored; part a and part b each worth one point; each part scored all or nothing (0,1); sum scoring for parts for total score of (0, 1, 2); each part scored independently
OE	hand-scored extended text interaction (traditional open response/constructed response item)	0,1,2,3,4	Hand scored; holistically scored; one rubric/one dimensional scoring

There is a Cluster/Passage which consists of two MS-1 items and two MS-2 items for a total of maximum 6 points.

Reporting Categories

- For Disciplinary Core Ideas (DCI) reporting categories, the subscore reporting will be based on three disciplinary domains.

- Physical Sciences (PS)
 - Life Sciences (LS)
 - Earth and Space Sciences (ESS)
2. The PrimaryContentStandard column in the NTS extract provides the coding for the DCI reporting categories.
 3. The Standards Column in the NTS extract indicates the DCI.
 4. DCI Reporting categories are defined as follows:

Code	DCI Reporting Category
LS	Life Science
PS	Physical Science
ES	Earth and Space Science

5. Psychometrics provides the performance indicators for student's performance on each reporting category. Values:
 - 1=Below Standard
 - 2=At/Near Standard
 - 3=Above Standard

Post Test Administration

Demographic Clean-up

Demographic cleanup is done in the iTester Portal for online testers. Paper testers are linked back to the pre-ID file provided when a valid ID is provided, or a valid label is used on the answer booklet. Data Processing Specifications contains more details on identification of students and linking to the pre-ID file.

Student Data Processing

1. Student Names will have all periods, commas and apostrophes removed.
 - a. Middle Name is the first initial of the middle name or blank if not available.
 - b. Special characters will be set to blank.
2. All records will be suppressed from processing if Name fields, Student ID, and Test Items are all blank.

Test Data Processing

Duplicates may exist where there is more than one data record with the same Student ID, be the record online or paper.

1. Duplicate Test records with the same Student ID/Grade will be combined or otherwise suppressed. See *Data Processing Specifications* for resolution of duplicate tests.
 - a. If there is a duplicate where the student takes one session in one test instance and another session in another test instance the 2 sessions will be combined/merged to create one complete test.
 - i. If the schools differ between session 1 and session 2, the school from the last session taken will be used for reporting (if it can be determined by the session updated dates for online tests)
 - ii. The record will be flagged in the data file as being a merged record.
 - iii. If an online session is merged with a paper session, test mode flag is set to “both”.
2. Duplicate Cross Grade tests are identified as more than one test taken with two different grades from the same student.
 - a. Should the Student have no work in the off-grade book, or the book is void and work in the matching grade test, suppress the off-grade test.
 - b. If both books have responses, send a report to Program Management for research and resolution.
3. When merging sessions to create one test record for a student, if there are overlapping attempted sessions, the session with the most attempted items will be used in the merge.

Scan Paper Delivery and Data Denotation

Each Paper test is scanned and delivered immediately to the Reporting Data Processing team. At the time of receipt, Data Processing will perform procedures to accurately identify inaccuracies in the data. The data will be formatted as specified in the Scan Delivery Layout Format

1. All discrepancies with the Scan File be resolved accordingly.
2. Any Student Response Booklet where VOID is bubbled and there is at least one item is attempted will be researched via Webdesk system. See *Data Processing Specifications* for resolution of Void bubbles.

Scoring Data

Scoring division will provide Data Processing with the open response scores for all tests.

1. Every score record will contain valid scores for all items.
 - a. A validation of score values will be performed.
 - b. If a score value is found to be invalid, resolution will be attempted with the Scoring Division.
2. Each score record will be associated with a Booklet ID or a Test ID.
 - a. If a score record is received without an associated Test or Booklet ID, resolution will be attempted with the Scoring Division.
3. All unresolved scoring records will be included in a report to the Scoring Division, for research and resolution.
4. The following values will be received
 - from Scoring : B=Blank
 - U=Unreadable with code number 51
 - F=Non-English with code number 53
 - W=Wrong Location with code number 52
 - O=Off



Topic with code number 54

5. Score values of U and W will be blanked out and reported with a null/blank value.
6. Score values of B, F and O will be given a score of 0 for analysis purposes.
7. Both operational and field test open response hand scored items will be scored.

Student Participation and Exclusions

Test Attempt Rules

Test Attempted indicates that a student has answered at least five (5) operational items on the test.

1. If a session is voided, any items attempted will be blanked out and will not count toward test attemptedness.
2. A valid attempt to an open response item is determined if the score is a non-blank (not B=blank or NULL) score.
3. Only field test items can have a null score meaning that the item was not scored.
4. A student is classified into 2 possible attempt groups of Attempt Status:
 - a. Attempt Status 0 is assigned to the test if the student did not provide a valid response to at least 5 operational items.
 - b. Attempt Status 1 is assigned to the test if the student provided a valid attempt to at least 5 operational items on the test.

Not Tested Reasons

The below not tested reasons or Test Report Codes are the valid options on the scannables and in the online testing platform. School Administrators set the codes as appropriate. Test invalidations will be provided to Cognia by BIE, if appropriate. Numbers correspond to the numbers in the testing platform. The numbers differ on the scannable.

- Valid Not Tested Reasons are:
 - Withdrew Before Test Completion (01)
 - Non-allowed Modification (02)
 - Medical Emergency (04)
 - Other non-completion (06)
 - Test Irregularities (07)
 - Absent (08)
 - Participating in MSAA (10)
 - Void or Do Not Report (bubbled on the scannable or DNR indicated in the online testing platform)



Student Participation Status

All student tests will reflect the Student Participation based on the test attempt status and any “Not Tested Reason” indicated.

1. Participation status is determined using both the test “Attempt Status” value and the “Not Tested Reason”.
 - a. If Attempt Status is 0 (less than 5 operational items with valid attempts), and the test has a “Not Tested Reason”, the “Not Tested Reason” is reported, otherwise, the test is reported as “Did Not Reach Minimum Attempt”.
 - b. If Attempt Status is 1 (the test has at least 5 operational items with valid attempts),
 - i. The student is classified as Tested.
 - ii. If the student has a Not Tested Reason, the Not Tested Reason is ignored.
2. Off grade testers (both up and down) will be reported as follows:
 - a. They will receive a student report and label if they meet the attemptedness rule.
 - b. Student Grade will be set to the Tested Grade.
 - c. They will NOT be included in any aggregations regardless of attemptedness. However, aggregations will be reported on their student report for their school, district, and BIE.
 - d. They will be reported in DI (Dynamic Reporting Platform).
 - e. They will be included in the results data file to BIE.
 - f. They will be included on the printed Roster.
 - g. Data Processing will assign a participation status using their student IDs.
 - h. If they meet attemptedness their participation status is set to N (off grade testers). If they do not meet attemptedness their participation status will be B (Did not Meet Minimum Attempt) or other not tested reason marked on the test record.
3. Regardless of the test attempt status, if a student is on the test invalidation list from the client, their test will be marked as Invalidated. The test is suppressed from the reporting dataset and is not part of any reporting deliverable or analyses.
4. Only “Tested” students, that is, students who meet attemptedness and are not off grade testers will be included in analyses.
5. The client will indicate in the exception list any tests that should be invalidated.
6. The following hierarchy is followed if more than one not tested reason is marked and the student’s attempt status is 0. The hierarchy is from highest to lowest priority:

Invalidated
test
Void/DNR test
Participating
in MSAA
Medical
Emergency
Absent
Withdrew Before Test
Completion Test
Irregularities
Non-Allowed
Modification
Other Non-
Completion

Student Participation Summary

The participation summary table below defines the participation codes:

Participation Status	Participation Code	Included in State File	Included in Aggregations	Receive a Student Report and Label	Included on Roster and DI
Tested	Z	Yes	Yes	Yes	Yes
Did Not Reach Minimum Attempt	B	Yes	No	No	Yes
Withdrew Before Test Completion	C	Yes	No	No	Yes
Non-Allowed Modification	D	Yes	No	No	Yes
Medical Emergency	F	Yes	No	No	Yes
Other Non-Completion	H	Yes	No	No	Yes
Test Irregularities	I	Yes	No	No	Yes
Absent	J	Yes	No	No	Yes
Participating in MSAA	M	Yes	No	No	Yes
Off grade testers (and meet attemptedness)	N	Yes	No	Yes	Yes
Void Test	K	Yes	No	No	Yes
Invalidated Test	L	Yes	No	No	Yes

Calculations

Aggregations

Aggregation inclusion rules are summarized in the Student Participation Summary table above.

Scoring Method

1. Overall raw score is calculated by summing the student's score on each operational item.
2. NTS, Cognia's item bank, indicates which items are operational and count toward the student's overall raw score.

Performance Levels

1. Psychometrics will provide the scale to be used for scaling the BIE tests.
2. Scale scores and achievement levels will be assigned based on the students' overall raw score and applying the scale score lookup from Psychometrics.

Data Suppression Rules

See Report Specific Rules section for data suppression rules.

Specific Reporting Rules

LENS Online Reporting

Cognia's online reporting system LENS is not used for BIE. Online reporting is done in Data Interaction.

General reporting rule: If the student's name is not available the name is reported as Name Not Provided. This applies to all deliverables.

Student Results Data File

- Naming Convention of the data files: BIE2425sciStudentResults.csv for the State Results data file.
 - Item level data will be provided according to the BIE2425sciItemLevelDataFileLayout.xlsx. The naming convention for the item level file is BIE2425sciStudentItemLevelResults.csv.
- If a student's test was merged to create one test, then the mergedtest flag is set to 1, otherwise it is set to 0.
- If the mergedtest flag is set to 1 and the student tested at 2 or more different locations, the student is reported back to where the last session was taken.
- A student appears in the file once with only one test.
- NumAttempted is the number of operational items in the test that met the item attemptedness rules described above.
- All grades are reported in the same file.
- Participation file is not handed off to eMetric.
- eMetric will receive the BIE2425sciStudentResults.csv after standards validation for loading DI. eMetric will use the 'Included in DI reporting' column in the file layout to determine what data to load to DI.

eMetric Metadata file for the Individual Student Report PDFs

- The column headings for the file are: ProgramName, ReportName, Year, Org_Num, PDF_name
- The file is a csv file
- The naming convention for the file is BIE2425sci_PDFmetadata.csv
- The file is posted to the ftp site for eMetric to access
- Org_Num=Districtcode-Schoolcode
- Year=2025
- ProgramName=Science (General Education)

- ReportName=Individual Student Report
- Web file naming convention: BIE2425sciStudentReport_<districtcode||schoolcode>.pdf

eMetric Summary Data File

- Data file used to compare summary results to aggregations in DI. Used for quality assurance checks by eMetric.
- The data file follows the BIE2425scieMetricSummaryDataTransfer layout.

Student Results Labels

- Student results labels are only produced for students who Tested (participation status='Z').
- The label contains the scale score and achievement level the student earned.
- The Test Date on the label is Spring 2025.
- There are 10 labels per page.
- Labels are produced by school and tested grade. Within each file the students in a school are sorted by student grade, lastname, firstname, mi, NASIS ID.
- The labels are shipped to the school.

Student Reports

- Student reports are only produced for students who Tested (participation status='Z').
- The student report contains the scale score and achievement level the student earned.
- The report is printed in color.
- 2 copies of the report are printed and shipped to the school.
- The reports are printed duplex on 8 ½ x 11 sheets of paper.
- Averages and percentages on the report are rounded to the nearest whole number.
- Percentages are formatted with a % sign.
- Orange color used on the report is used to indicate the student's earned scale score.
- The minimum N-size that must be reached for aggregations to be reported is 10. Any result including less than 10 students will be suppressed from the reports.
- Print the achievement level descriptor with examples on the second page corresponding to the student's earned achievement level.
- If the earned achievement level is not Advanced, print the achievement level descriptor without examples, for the achievement level directly above the earned achievement level. This appears adjacent to the earned achievement level descriptor and is surrounded by a box. For example, the student earns Nearing Proficiency, the achievement level descriptor for Nearing Proficiency with examples is printed and the achievement level descriptor for Proficient is printed in a text box.
- Web versions do not include the slip sheets.
- Web reports are run by school.
- Within the school web files the students are sorted by student grade and lastname,firstname.
- Web file naming convention: BIE2425SciStudentReport_<districtcode||schoolcode>.pdf



Student Results Rosters

- The roster lists all students in the school in a specific tested grade.
- The report is printed duplex.
- Student names are formatted as Lastname, Firstname (title case).
- If a student is not tested, the scale score and achievement level will be '-' on the roster.
- Up to 14 students are listed on the front page.
- A roster may have multiple pages if there are more than 14 students in the school and tested grade.
- Students are listed alphabetically by Lastname, Firstname (title case).
- One copy of the report is printed.
- The report is printed landscape and simplex on 8 ½ x 11 paper in color.
- There is alternate row shading on the roster.
- If a student is Tested, the achievement level is highlighted in a shade of blue that matches the shade on the student report for that achievement level.
- The scale score and achievement levels are printed in bold font.

Shipping Product Code Summary

Reporting Products

Reports are being shipped to schools. Reports are packed by grade. Below is the definition of the 4 reporting products.

Contract Code: [800450] Description BIE Science 2025 Admin ID 1	Report Type	Report For	Grade(s)	Report Subtype	Content Code	Qty
Student Report Parent Copy	07	1	05,08,11	02	00	1
Student Report-School Copy	07	1	05,08,11	01	00	1
Student Results Label	07	1	05,08,11	03	00	1
Student Results Roster	07	1	05,08,11	12	00	1

Appendix

The following schools received a waiver to not test Science in Spring 2025:

DistrictCode	School Code	School Name
TCS	D56S02	Beatrice Rafferty School/Sipayik Elementary
TCS	D78S21	Bogue Chitto Elementary School
TCS	D53F13	Bug-O-Nay-Ge-Shig School
TCS	D52S04	Cherokee Elementary School
TCS	D52S03	Cherokee High School
TCS	D10P15	Chief Leschi School
TCS	D78S23	Choctaw Central High School
TCS	D78S24	Choctaw Central Middle School
TCS	D78S22	Conehatta Elementary School
TCS	D61J03	Duckwater Shoshone Elementary
TCS	D53F15	Fond du Lac Ojibwe School
TCS	D60F07	Hannahville Indian School
TCS	D57S02	Indian Island School
TCS	D55S02	Indian Township School
TCS	D09B02	Jones Academy
TCS	D10P17	Lummi High School
TCS	D10P14	Lummi Tribal School System
TCS	D11A13	Mandaree Day School
TCS	D51F02	Meskwaki Settlement School
TCS	D53F18	Nay-Ah-Shing School
TCS	D25M14	Ohkay Owingeh Community School
TCS	D78S25	Pearl River Elementary School
TCS	D10P02	Quileute Tribal School
TCS	D78S13	Red Water Elementary School
TCS	D25M32	Santa Fe Indian School
TCS	D78S14	Standing Pine Elementary School
TCS	D78S15	Tucker Elementary School
TCS	D11A14	Twin Buttes Day School
TCS	D10P13	Wa He Lut Indian School
TCS	D11A15	White Shield School

APPENDIX L
SCORE REPORT INTERPRETATION QUICK GUIDE


BIE Science Assessment – Report Interpretation Quick Guide



The Bureau of Indian Education (BIE) Science Assessment is based on Next Generation Science Standards (NGSS), and is administered to students in grades 5, 8, and 11. The standards focus on important disciplinary core ideas, scientific and engineering practices, and crosscutting concepts that apply across scientific disciplines. The assessment provides information regarding each student’s progress toward achievement of essential knowledge and skills that will help them explain and make sense of phenomena in the world around them, solve problems, and apply their scientific literacy to understand that scientific dilemmas they may face as adults.

This document outlines how to interpret the information on the individual student report. The report provides general information about the BIE Science Assessment, resources for parents/guardians, and the student’s results, including the student’s scale score and achievement level. The report also provides information on how the student performed compared to other BIE students.

Sample Individual Student Report (page 1)



BIE Science Assessment
Individual Student Report
Spring 2024

Dear Parent or Guardian:

We are excited that your student participated in the high-quality Science Assessment aligned to the Next Generation Science Standards (NGSS) as required for all Bureau Funded Schools as part of *Every Student Succeeds Act (ESSA)*. This is your student’s performance report based on the Spring Science Assessment results.

As your student grows and learns during the school year, the Bureau of Indian Education (BIE) will administer the science assessment annually for grades 5, 8, and 11 to measure how well your student is performing on the Next Generation Science Standards. It will provide a better understanding of what your student knows and is able to do.

These assessment results will allow you to gain critical insights to your student’s progress, support your understanding of his/her strengths and areas identified as needing improvement. Additionally, the information will support instruction, identify resources, and provide enrichment activities to meet the needs of all learners.

We thank you for your continued support and cooperation.

The Bureau of Indian Education

A

Student Name: LASTNAME177, FIRST177
NASIS ID: D00177
Tested Grade: 05
Student Grade: 05
District: Demonstration District A
School: Demonstration School 2

B

RESOURCES

U.S. Department of the Interior/Bureau of Indian Education: To learn more about the Bureau of Indian Education, please navigate to <https://www.bie.edu> or scan the QR code.

Assessment & Accountability: Integrated content and assessments ensure broad access to accurate and timely student and school-level information which is used for educational decision making. For more information on BIE Assessments and Accountability, please navigate to <https://www.bie.edu/landing-page/assessments-and-accountability> or scan the QR code.

Standards: The BIE Science assessment is aligned to the Next Generation Science Standards. To read the standards, navigate to <https://www.nextgenscience.org/search-standards> or scan the QR code.

C

Achievement Level Descriptors

500	Novice	544	Nearing Proficiency	560	Proficient	574	Advanced	590
Students demonstrate evidence of emerging understanding and use of college and career readiness knowledge, skills, and abilities.		Students demonstrate evidence of partial understanding and use of college and career readiness knowledge, skills, and abilities.		Students demonstrate evidence of satisfactory understanding and use of college and career readiness knowledge, skills, and abilities.		Students demonstrate evidence of thorough understanding and use of college and career readiness knowledge, skills, and abilities.		

©2024 Cognia, Inc. All rights reserved.

A This section provides assessment information and a note to parents/guardians around the use of the report information in schools.

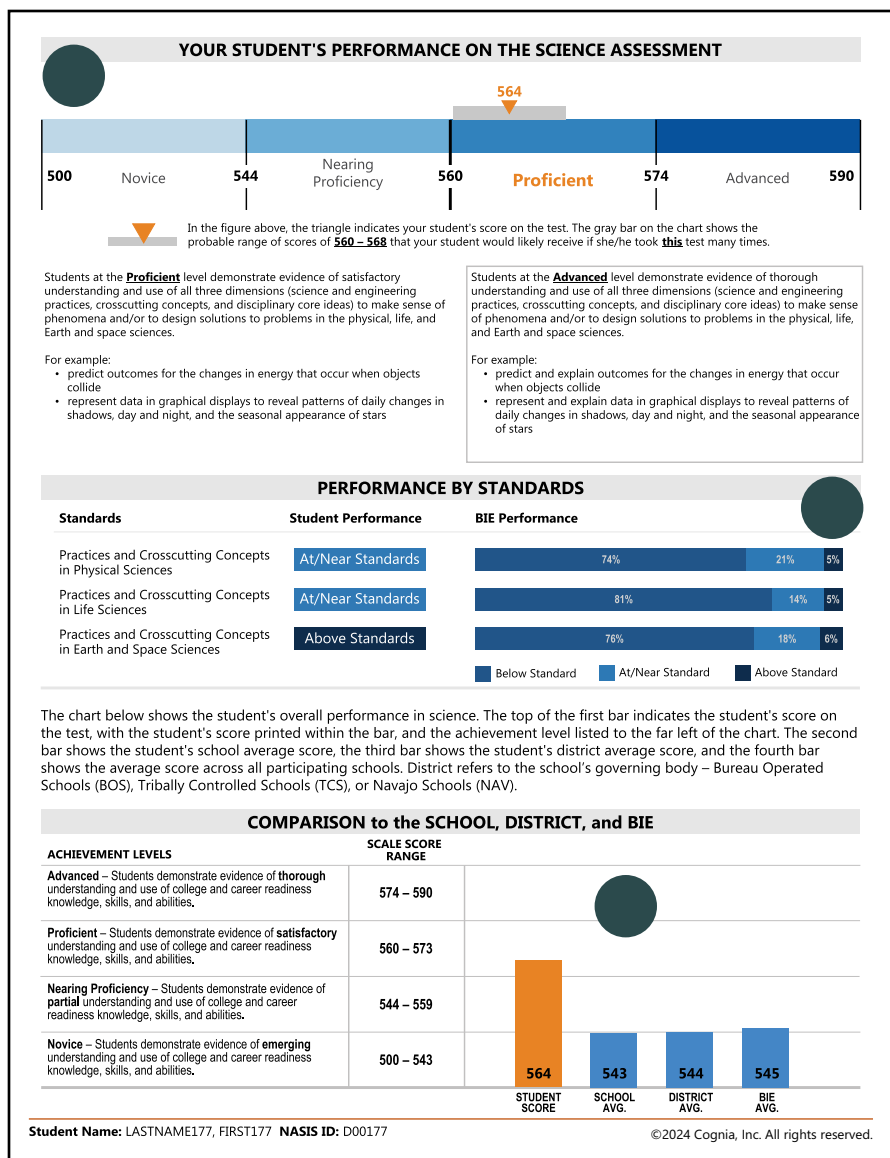
B This section provides links to additional resources for the BIE Science Assessment, including standards, general BIE information, and specific assessment information. Links can be accessed by typing the link into a browser or scanning the QR code with a smart phone/tablet.


C This section describes the four possible achievement levels a student can reach – Novice, Nearing Proficiency, Proficient, and Advanced.

BIE Science Assessment – Report Interpretation Quick Guide


While page one outlines general assessment and report information, page two provides specific information on the student's performance, both individually and in relation to their peers.

Sample Individual Student Report (page 2)



 This section describes your student's performance on the overall assessment. Their scale score is provided on the chart, indicating their achievement level. Additional information regarding the achievement level and the skills demonstrated at that level are outlined below. To the right, information is included on what skills the student would need to develop to reach the next achievement level.

● This section breaks down the student performance by each science discipline. The student performance is displayed along with a bar graph of overall BIE performance.

 This section displays the student's performance in relation to peers in his or her school, district, and across the BIE. District, as used on this report, refers to the school's governing body – Bureau Operated Schools, Tribally Controlled Schools, and Navajo Schools.

For additional information regarding the student report, contact your school test coordinator (STC). School inquiries should be directed to the BIE.

APPENDIX M
CUMULATIVE SCALED-SCORE DISTRIBUTIONS

Figure M-1. Cumulative Scale Score Distributions—Grade 5 Science
BIE Science Cumulative Scale Score Distributions

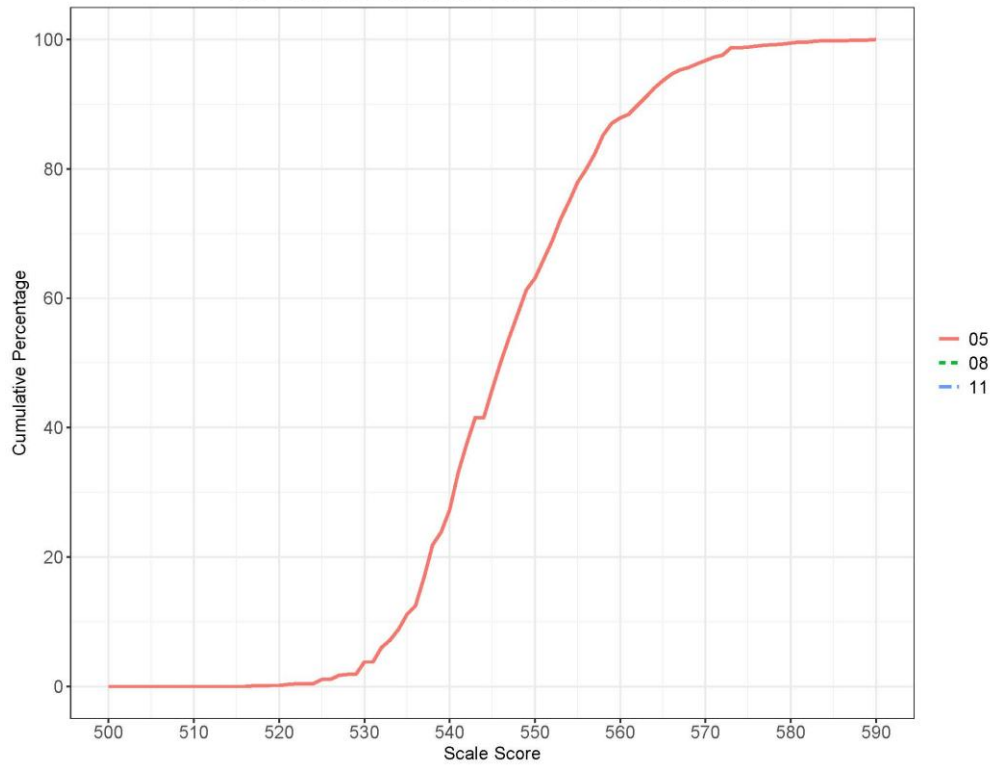


Figure M-2. Cumulative Scale Score Distributions—Grade 8 Science

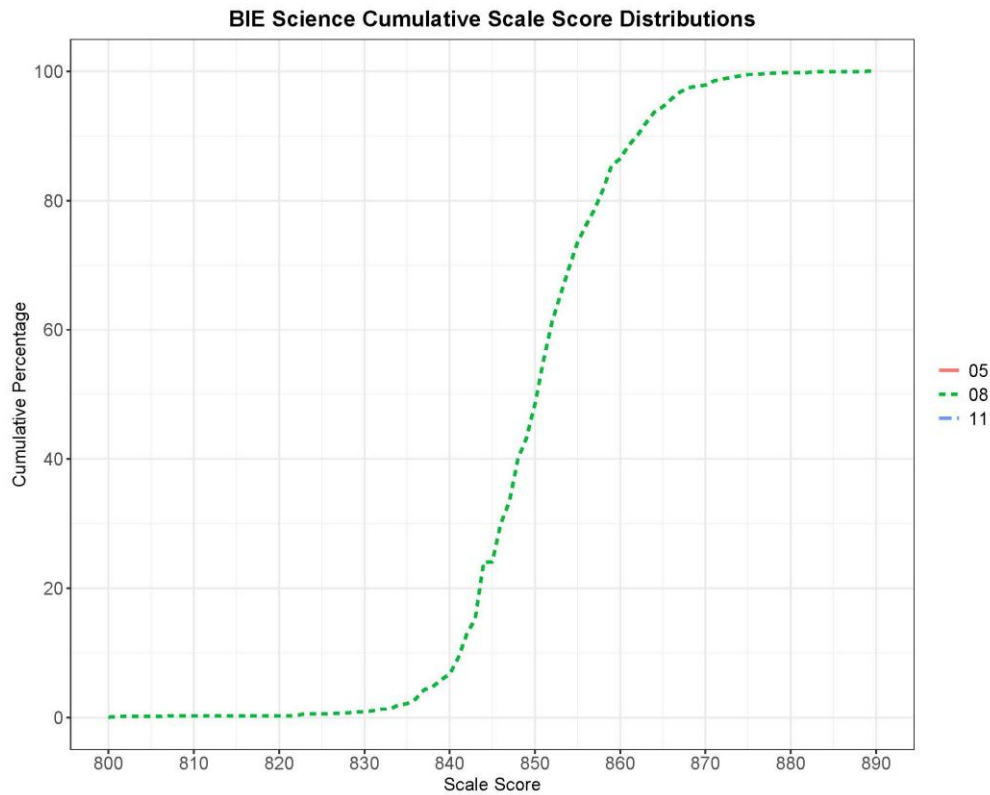
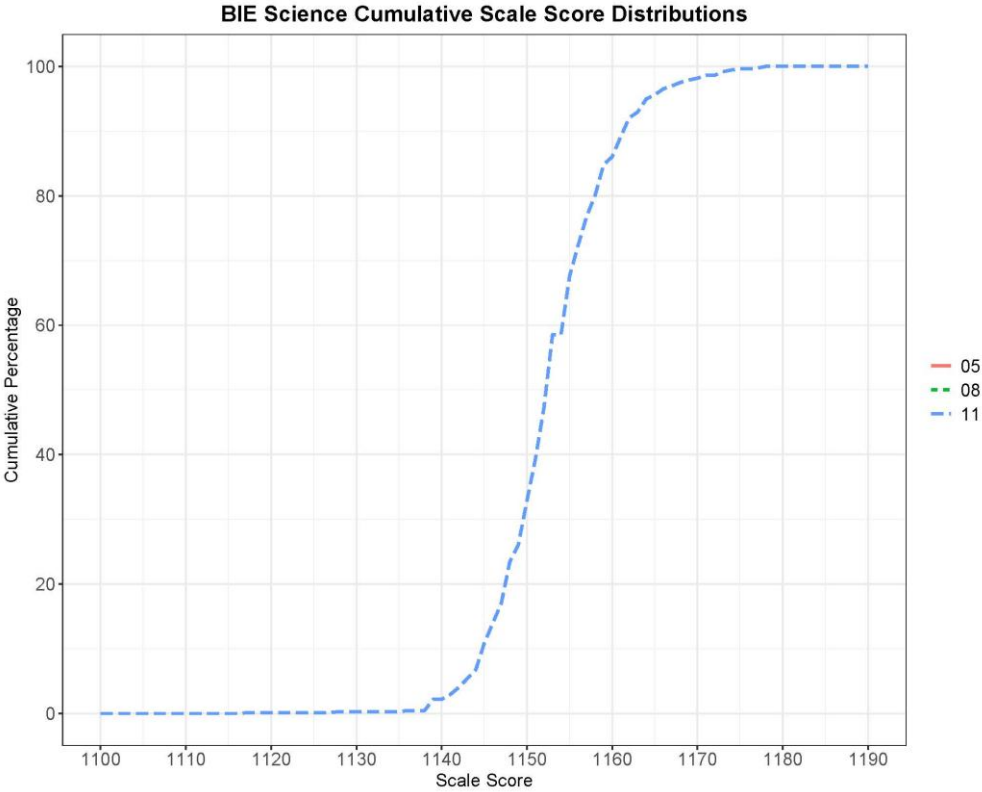


Figure M-3. Cumulative Scale Score Distributions—Grade 11 Science



APPENDIX N

SCALED SCORE DESCRIPTIVE STATISTICS

Calculations based on those students attempting 5 or more items on the given BIE assessment. Statistic values are suppressed for those subjects/grades with fewer than 50 students.

Table N-1. Scaled Score Descriptive Statistics for BIE Science Grade 5, as a Function of Subgroup*

Subgroup	Number of Students	Mean	Median	SD	Skewness	Kurtosis
Overall	1,928	547.61	546.00	10.90	0.53	0.33
Female	961	547.62	547.00	10.34	0.56	0.52
Male	967	547.60	546.00	11.43	0.51	0.15
Currently receiving LEP services	321	547.84	548.00	9.54	0.21	0.16
Not receiving LEP services	394	548.42	548.00	11.48	0.42	0.45
Special Ed	351	542.77	541.00	9.19	0.76	1.74
Non-Special Ed	1,573	548.69	548.00	10.96	0.47	0.21
Economically Disadvantaged Students	1,840	547.58	546.00	10.88	0.54	0.37
Gifted Students	81	553.00	551.00	13.54	0.29	-0.61
Non-gifted Students	103	546.29	545.00	10.02	0.22	-0.74
Title1 Students	1,870	547.55	546.00	10.87	0.51	0.27
Non-title1 Students	15	547.93	548.00	8.91	-0.03	-1.59

*Calculations based on those students attempting 5 or more items on the given BIE assessment. Statistical values are suppressed for those subjects/grades with fewer than 50 students.

Table N-2. Scaled Score Descriptive Statistics for BIE Science Grade 8, as a Function of Subgroup*

Subgroup	Number of Students	Mean	Median	SD	Skewness	Kurtosis
Overall	1,677	851.05	851.00	8.72	-0.14	2.94
Female	825	850.76	851.00	8.36	-0.47	3.54
Male	852	851.33	851.00	9.06	0.09	2.42
Currently receiving LEP services	202	850.31	850.00	6.52	0.13	0.09
Not receiving LEP services	405	852.62	852.00	8.71	0.42	0.38
Special Ed	299	846.99	847.00	7.96	-0.60	6.28
Non-Special Ed	1,374	851.90	851.00	8.59	-0.13	2.64
Economically Disadvantaged Students	1,625	851.11	851.00	8.75	-0.19	2.94
Gifted Students	116	858.55	859.00	9.62	0.00	-0.30
Non-gifted Students	109	850.11	850.00	7.91	0.17	0.76
Title1 Students	1,623	850.99	851.00	8.66	-0.19	3.06
Non-title1 Students	13	849.15	851.00	13.05	0.36	0.39

*Calculations based on those students attempting 5 or more items on the given BIE assessment. Statistical values are suppressed for those subjects/grades with fewer than 50 students.

Table N-3. Scaled Score Descriptive Statistics for BIE Science Grade 11, as a Function of Subgroup*

Subgroup	Number of Students	Mean	Median	SD	Skewness	Kurtosis
Overall	1,095	1153.29	1153.00	6.78	0.28	1.61
Female	542	1152.91	1153.00	6.33	0.36	1.21
Male	553	1153.67	1153.00	7.18	0.19	1.79
Currently receiving LEP services	108	1151.71	1152.00	6.18	-0.08	0.11
Not receiving LEP services	226	1153.56	1153.00	7.20	-0.33	3.19
Special Ed	152	1150.53	1151.00	6.02	-0.37	8.24
Non-Special Ed	925	1153.80	1153.00	6.79	0.32	0.97
Economically Disadvantaged Students	1,052	1153.30	1153.00	6.78	0.26	1.65
Gifted Students	63	1155.59	1155.00	5.84	0.81	1.13
Non-gifted Students	44	1152.86	1153.00	7.55	0.79	1.25
Title1 Students	1,042	1153.27	1153.00	6.74	0.24	1.64
Non-title1 Students	6	1153.00	1151.50	13.04	0.93	0.65

*Calculations based on those students attempting 5 or more items on the given BIE assessment. Statistical values are suppressed for those subjects/grades with fewer than 50 students.