

New Jersey

Student Learning Assessment—Science
(NJSLA–S)

TECHNICAL REPORT
Grades 5, 8, and 11
2023

June 2024
PTM XXXX.XX



State of New Jersey
Department of Education

Copyright © 2024 by New Jersey Department of Education
All rights reserved

State Board of Education

Kathy Goldenberg President	Burlington
Andrew J. Mulvihill Vice President	Sussex
Arcelio Aponte	Middlesex
Mary G. Bennett	Essex
Mary Beth Berry	Hunterdon
Elaine Bobrove	Camden
Fatimah Burnam-Watkins	Union
Ronald K. Butcher	Gloucester
Jack Fornaro	Warren
Nedd James Johnson	Salem
Jeanette Peña	Hudson
Joseph Ricca Jr.	Morris
Sylvia Sylvia-Cioffi	Monmouth

Mr. Kevin Dehmer, Acting Commissioner of Education
Secretary, State Board of Education

It is a policy of the New Jersey State Board of Education and the State Department of Education that no person, on the basis of race, creed, national origin, age, sex, handicap, or marital status, shall be subjected to discrimination in employment or be excluded from or denied benefits in any activity, program, or service for which the department has responsibility. The department will comply with all state and federal laws and regulations concerning nondiscrimination.

Table of Contents

PART 1: INTRODUCTION	1
1.1 Purpose of the Assessment.....	1
1.2 Description of the Assessment.....	2
1.2.1 Content Domains and Scientific Practices	2
1.2.2 Crosscutting Concepts	8
1.2.3 Types of Scores	8
1.3 Organizational Support	10
PART 2: TEST DEVELOPMENT	11
2.1 Test Specifications.....	11
2.1.1 Test Blueprints	11
2.1.2 Unit Design	12
2.1.3 Item Types	13
2.2 Item Development Processes	14
2.2.1 Item Writing.....	15
2.2.2 Content Specialist Review	15
2.2.3 Editorial Review	16
2.2.4 NJ Science Advisory Committee Content Review	16
2.2.5 Bias and Sensitivity Committee Review	17
2.2.6 Field Test.....	17
2.2.7 Statistical Review	17
2.2.8 Second Bias and Sensitivity Review	18
2.2.9 Ready for Operational Testing	18
2.3 Test Construction Process	18
2.3.1 Test Construction—First Draft	18
2.3.2 Test Construction Content and Psychometric Review	20
2.3.3 Test Construction NJDOE Review.....	20
2.4 2023 NJSLA–S Test Construction.....	20
2.4.1 Grade 5 Test Construction	21
2.4.2 Grade 8 Test Construction	23
2.4.3 Grade 11 Test Construction	26
2.5 2023 NJSLA–S State of the Item Bank.....	28
PART 3: TEST ADMINISTRATION	29
3.1 District Test Coordinator Training.....	29
3.2 Test Security and Administration Procedures	30
3.2.1 Computer-Based Testing.....	30
3.2.2 Paper-Based Testing.....	31

3.3 Test Irregularities and Breaches	32
3.4 Test Accessibility Features and Accommodations	34
3.4.1 Accessibility Features	35
3.4.2 Accommodations.....	35
PART 4: SCORING.....	38
4.1 Machine-Scored Items	38
4.1.1 Adjudication	38
4.2 Handscored Items.....	39
4.2.1 Selecting Handscoring Staff	39
4.2.2 Operational Range Finding	39
4.2.3 Field Test Range Finding	40
4.2.4 Developing Scoring Guides	40
4.2.5 Team Leader Training and Duties	40
4.2.6 Scorer Training and Qualifying.....	40
4.2.7 Monitoring Scorer Performance	42
4.2.8 Automatic Rescores	44
4.3 Quality Control	44
4.3.1 QC Sample	44
4.3.2 Key Information Sheets	45
4.3.3 Aggregate Data	45
PART 5: STANDARD SETTING	46
PART 6: ITEM AND TEST STATISTICS.....	47
6.1 Classical Test Theory Statistics.....	47
6.1.1 Item Difficulty and Discrimination Descriptive Statistics	47
6.1.2 Speededness.....	54
6.1.3 Operational DIF Analysis.....	55
6.2 Item Response Theory	60
6.2.1 Unidimensionality	61
6.2.2 Partial Credit Model Fit Statistics.....	64
6.2.3 Local Independence.....	79
6.2.4 Item Characteristic Curves—CR Items.....	79
6.3 Student Test Performance.....	85
6.3.1 Scale Score Distribution by Form.....	85
6.3.2 Scale Score Distributions by Demographic Group	85
6.3.3 Subscore Proficiency Classification	86

PART 7: EQUATING AND SCALING	87
7.1 Summary of Equating and Scaling Procedures	87
7.1.1 Rounding Rules	89
7.2 Accommodative Form Equivalence	90
7.2.1 Special Equating	90
7.3 Subscore Performance Levels.....	91
PART 8: RELIABILITY.....	92
8.1 Classical Test Theory Reliability Estimates	92
8.1.1 Reliability and Measurement Error	92
8.1.2 Raw Score Internal Consistency.....	93
8.2 Item Response Theory Reliability	97
8.2.1 Test Information Functions	97
8.2.2 Conditional Standard Error of Measurement	100
8.2.3 Item Maps	102
8.3 Reliability of Performance Classifications	106
8.3.1 Conditional Standard Error of Measurement at Each Cut Score	106
8.3.2 Classification Consistency Indices.....	106
8.4 Reliability of Subscore Performance Classifications	107
8.5 Rater Reliability	109
PART 9: VALIDITY	110
9.1 Evidence Based on Test Content.....	110
9.1.1 Alignment Study	111
9.2 Evidence Based on Response Processes	112
9.2.1 Cognitive Lab Study	113
9.3 Evidence Based on Internal Structure.....	114
9.3.1 Intercorrelations	114
9.3.2 Other Internal Structure Evidence.....	114
9.3.3 Confirmatory Factor Analysis for 2023 NJSLA–S.....	114
9.4 Evidence Based on Relationships to Other Variables	119
9.5 Evidence Based on the Consequences of Testing	121
9.6 Other Validity Evidence.....	121

9.7 Summary	123
9.7.1 Student Performance Level Classifications: Overall Scale Score.....	123
9.7.2 Student Performance Level Classifications: Domains and Practices Subscores	124
9.7.3 Future NJSLA–S Validity Studies	125
PART 10: REPORTING.....	126
10.1 Individual Student Report	126
10.2 Student Label	129
10.3 Student Roster	129
10.4 School Summary and District Summary of Schools	131
10.5 School and District Performance Level Summary Reports.....	134
REFERENCES	137
APPENDIX A: GLOSSARY OF ABBREVIATIONS	142
APPENDIX B: NEW JERSEY SCIENCE ADVISORY AND BIAS AND SENSITIVITY COMMITTEES– DISTRICT AND COUNTY REPRESENTATION.....	144
APPENDIX C: STATISTICAL REVIEW REFERENCE SHEET	147
APPENDIX D: 2019 NJSLA–S STANDARD SETTING: EXECUTIVE SUMMARY	149
APPENDIX E: NJSLA–S PERFORMANCE-LEVEL DESCRIPTORS	154
E.1 Policy PLDs.....	154
E.2 Threshold PLDs	155
E.2.1 Grade 5 Threshold PLDs.....	155
E.2.2 Grade 8 Threshold PLDs.....	165
E.2.3 Grade 11 Threshold PLDs.....	175
E.3 Reporting PLDs	187
E.3.1 Reporting PLDs–Level 1	187
E.3.2 Reporting PLDs–Level 2	187
E.3.3 Reporting PLDs–Level 3	187
E.3.4 Reporting PLDs–Level 4	187
APPENDIX F: DETAILED TEST MAPS	188
APPENDIX G: SCALE SCORE CUMULATIVE FREQUENCY DISTRIBUTIONS.....	196
APPENDIX H: ITEM PARAMETERS AND MODEL FIT TABLES	202

APPENDIX I: RAW SCORE-TO-SCALE-SCORE-CONVERSION TABLES	208
APPENDIX J: RAW SCORE-TO-THETA SUBSCORE TABLES	214
APPENDIX K: SUBSCORE PROFICIENCY CLASSIFICATIONS	232
APPENDIX L: EXECUTIVE SUMMARY OF THE NJSLA–S ALIGNMENT EVALUATION STUDY	238
APPENDIX M: EXECUTIVE SUMMARY OF EVALUATION OF THE COGNITIVE PROCESS STUDY	242
APPENDIX N: OBSERVED <i>P-VALUES</i> FOR THE FIT AND UNDERFIT SUBGROUPS OF STUDENTS	246
APPENDIX O: PARAMETER ESTIMATES FROM THE CONFIRMATORY FACTOR ANALYSES FOR THE 2023 NJSLA–S TESTS	254

Tables and Figures

Table 1.2.1: Earth and Space Science DCIs	3
Table 1.2.2: Life Science DCIs	4
Table 1.2.3: Physical Science DCIs	5
Table 1.2.4: SEP Consolidation	5
Table 1.2.5: Investigating Practices	6
Table 1.2.6: Sensemaking Practices	7
Table 1.2.7: Critiquing Practices	7
Table 1.2.8: Crosscutting Concepts	8
Table 1.2.9: NJSLA–S Scale Score Ranges	9
Table 2.1.1: Test Blueprints	12
Figure 2.1.1. Sample Grade 5 Unit	13
Table 2.1.2: NJSLA–S Item Types	14
Table 2.3.1: Summary of NJSLA–S Test Construction Statistical Constraints	20
Table 2.4.1: Points Available by Domain and Practice	21
Table 2.4.2: 2023 NJSLA–S Grade 5 Item and Point Totals by Reporting Category	21
Table 2.4.3: 2023 NJSLA–S Grade 5 DCIs	22
Table 2.4.4: 2023 NJSLA–S Grade 5 SEPs	22
Table 2.4.5: 2023 NJSLA–S Grade 5 CCCs	22
Table 2.4.6: 2023 NJSLA–S Grade 5 Test Construction Statistics	23
Table 2.4.7: 2023 NJSLA–S Grade 5 Test Construction DIF Classifications	23
Table 2.4.8: 2023 NJSLA–S Grade 8 Item and Point Totals by Reporting Category	24
Table 2.4.9: 2023 NJSLA–S Grade 8 DCIs	24
Table 2.4.10: 2023 NJSLA–S Grade 8 SEPs	24
Table 2.4.11: 2023 NJSLA–S Grade 8 CCCs	25
Table 2.4.12: 2023 NJSLA–S Grade 8 Test Construction Statistics	25
Table 2.4.13: 2023 NJSLA–S Grade 8 Test Construction DIF Classifications	25

Table 2.4.14: 2023 NJSLA–S Grade 11 Item and Point Totals by Reporting Category	26
Table 2.4.15: 2023 NJSLA–S Grade 11 DCIs	26
Table 2.4.16: 2023 NJSLA–S Grade 11 SEPs	27
Table 2.4.17: 2023 NJSLA–S Grade 11 CCCs	27
Table 2.4.18: 2023 NJSLA–S Grade 11 Test Construction Statistics	28
Table 2.4.19: 2023 NJSLA–S Grade 11 Test Construction DIF Classifications	28
Figure 3.1.1. Slide 2 from the 2023 DTC	30
Table 3.1.1: NJSLA–S 2023 Grades 5, 8, and 11 Science Testing Window	30
Table 3.2.1: CBT School Test Coordinator Checklist	31
Table 3.2.2: PBT School Test Coordinator Checklist	32
Table 4.2.1: Scoring Personnel by Grade	42
Table 4.2.2: Automatic Rescore Results	44
Table 6.1.1: Grade 5 Item Difficulty (<i>p-value</i>) Distribution and Summary Statistics	49
Table 6.1.2: Grade 5 Item Discrimination Distribution and Summary Statistics	50
Table 6.1.3: Grade 8 Item Difficulty (<i>p-value</i>) Distribution and Summary Statistics	51
Table 6.1.4: Grade 8 Item Discrimination Distribution and Summary Statistics	52
Table 6.1.5: Grade 11 Item Difficulty (<i>p-value</i>) Distribution and Summary Statistics	53
Table 6.1.6: Grade 11 Item Discrimination Distribution and Summary Statistics	54
Table 6.1.7: Operational Testing Schedule—Items and Time Allocations	55
Table 6.1.8: Percentage of Students Omitting the Last TE Item in Each Operational Unit	55
Table 6.1.9: Differential Item Functioning Evaluation Criteria	56
Table 6.1.10: Grade 5 DIF Classification by Item Type	57
Table 6.1.11: Grade 8 DIF Classification by Item Type	58
Table 6.1.12: Grade 11 DIF Classification by Item Type	59
Table 6.2.1: Correlation Matrix for Domains	62
Table 6.2.2: Correlation Matrix for Practices	62

Figure 6.2.1. Grade 5 Scree Plot	63
Figure 6.2.2. Grade 8 Scree Plot	63
Figure 6.2.3. Grade 11 Scree Plot	64
Table 6.2.4: Summary of Item Infit and Outfit Statistics	65
Table 6.2.5: Summary of Rasch Discrimination Statistics	65
Table 6.2.6: Summary of Rasch Lower Asymptote Statistics	66
Table 6.2.7: Summary of Person Infit Statistics by Demographic Group	68
Table 6.2.8: Summary of Person Outfit Statistics by Demographic Group	70
Table 6.2.9: Summary of Person Infit Statistics by Form	72
Table 6.2.10: Summary of Person Outfit Statistics by Form.....	73
Figure 6.2.4. Grade 5 Person Infit and Outfit Distributions	74
Figure 6.2.5. Grade 8 Person Infit and Outfit Distributions	74
Figure 6.2.6. Grade 11 Person Infit and Outfit Distributions	75
Table 6.2.11: Distribution of Person Fit for Grades 5, 8, and 11	76
Table 6.2.12: <i>p-values</i> of Items Flagged by Delta Method.....	78
Table 6.2.13: Summary of Yen’s Q3 Statistics.....	79
Table 6.2.14: Constructed-Response Point Distribution Percentages	80
Figure 6.2.7. ICC Plot for Grade 5 Constructed-Response Item 1	80
Figure 6.2.8. ICC Plot for Grade 5 Constructed-Response Item 2	81
Figure 6.2.9. ICC Plot for Grade 5 Constructed-Response Item 3	81
Figure 6.2.10. ICC Plot for Grade 8 Constructed-Response Item 1	82
Figure 6.2.11. ICC Plot for Grade 8 Constructed-Response Item 2	82
Figure 6.2.12. ICC Plot for Grade 8 Constructed-Response Item 3	83
Figure 6.2.13. ICC Plot for Grade 11 Constructed-Response Item 1	83
Figure 6.2.14. ICC Plot for Grade 11 Constructed-Response Item 2	84
Figure 6.2.15. ICC Plot for Grade 11 Constructed-Response Item 3	84

Table 6.3.1: Descriptive Statistics of Students’ Test Performance by Form	85
Table 7.1.1: Scale Score Ranges for Proficiency Levels by Grade	88
Table 7.1.2: Slope and Intercept of Theta-to-Scale Score Transformations and Performance- Level Cut Scores by Grade.....	89
Table 8.1.1: Coefficient Alpha and SEM by Form	94
Table 8.1.2: Coefficient Alpha and SEM by Reporting Category.....	95
Table 8.1.3: Coefficient Alpha and SEM by Demographic Group	96
Table 8.1.4: Coefficient Alpha and SEM by Item Type	97
Figure 8.2.1. Grade 5 Test Information Function.....	99
Figure 8.2.2. Grade 8 Test Information Function.....	99
Figure 8.2.3. Grade 11 Test Information Function.....	100
Figure 8.2.4. Grade 5 Conditional Standard Error of Measurement	101
Figure 8.2.5. Grade 8 Conditional Standard Error of Measurement	101
Figure 8.2.6. Grade 11 Conditional Standard Error of Measurement	102
Figure 8.2.7. Grade 5 Item Difficulty and Student Ability Distributions	103
Figure 8.2.8. Grade 8 Item Difficulty and Student Ability Distributions	104
Figure 8.2.9. Grade 11 Item Difficulty and Student Ability Distributions	105
Table 8.3.1: Cut Scores with Conditional Standard Error of Measurement.....	106
Table 8.3.2: Performance Level Classification Consistency.....	107
Table 8.4.1: Subscore Performance Classification Consistency and Conditional Standard Error of Measurement	108
Table 8.5.1: Inter-rater Agreement Rate of Constructed-Response Items.....	109
Table 9.2.1: Range PLD Alignment by DCI, SEP, and Grade Level	113
Figure 9.3.1. Domain Subscore Structure	115
Figure 9.3.2. Practice Subscore Structure	116
Table 9.3.1. Model Fit Indices for the Domain Model	117
Table 9.3.2. Correlations between the Latent Subscores Implied by the Domain Model	117

Table 9.3.3. Number of Items with Highest rpb_{sub} on Assigned Domain	118
Table 9.3.4. Number of items with Highest rpb_{sub} on Assigned Practice	118
Table 9.4.1: Grade 5 Intercorrelations by Content Area	119
Table 9.4.2: Grade 8 Intercorrelations by Content Area	120
Table 9.4.3: Grade 11 Intercorrelations by Content Area	120
Figure 10.1.1. Sample Individual Student Report–Page 1	127
Figure 10.1.2. Sample Individual Student Report–Page 2	128
Figure 10.2.1. Sample Student Label	129
Figure 10.3.1. Sample Student Roster	130
Figure 10.4.1. Sample School Performance Level Summary Report–Domains and Practices	132
Figure 10.4.2. Sample District Performance Level Summary Report–Domains and Practices	133
Figure 10.5.1. Sample School Performance Level Summary Report	135
Figure 10.5.2. Sample District Performance Level Summary Report	136
Table A.1: Glossary of NJSLA–S Abbreviations	142
Table B.1: Grade 5 NJSAC District and County Representation.....	144
Table B.2: Grade 8 NJSAC District and County Representation	145
Table B.3: Grade 11 NJSAC District and County Representation	146
Table B.4: NJBSC District and County Representation	146
Table ES-1: Final Recommendations from Standard-Setting Panelists.....	151
Figure ES-1. Percentages of Students Classified at Each Level after Round 3	151
Table ES-2: Responses to Key Evaluation Questions.....	152
Table ES-3: Summary of Reasonableness Ratings and Comments.....	153
Table F.1: Grade 5 Test Map–Metadata and Item Statistics.....	188
Table F.2: Grade 8 Test Map–Metadata and Item Statistics.....	190
Table F.3: Grade 11 Test Map–Metadata and Item Statistics.....	193
Table G.1: Grade 5–Scale Score Cumulative Frequency Distribution	196

Table G.2: Grade 8–Scale Score Cumulative Frequency Distribution	198
Table G.3: Grade 11–Scale Score Cumulative Frequency Distribution	200
Table H.1: Grade 5–IRT Item Parameters and Fit Statistics.....	202
Table H.2: Grade 8–IRT Item Parameters and Fit Statistics.....	204
Table H.3: Grade 11–IRT Item Parameters and Fit Statistics.....	206
Table I.1: Grade 5–Operational	208
Table I.2: Grade 8–Operational	210
Table I.3: Grade 11–Operational	212
Table J.1: Grade 5 Earth and Space Science Score Table.....	214
Table J.2: Grade 5 Life Science Score Table	215
Table J.3: Grade 5 Physical Science Score Table	216
Table J.4: Grade 5 Sensemaking Score Table.....	217
Table J.5: Grade 5 Critiquing Score Table.....	218
Table J.6: Grade 5 Investigating Score Table	219
Table J.7: Grade 8 Earth and Space Science Score Table	220
Table J.8: Grade 8 Life Science Score Table	221
Table J.9: Grade 8 Physical Science Score Table	222
Table J.10: Grade 8 Sensemaking Score Table	223
Table J.11: Grade 8 Critiquing Score Table	224
Table J.12: Grade 8 Investigating Score Table	225
Table J.13: Grade 11 Earth and Space Science Score Table.....	226
Table J.14: Grade 11 Life Science Score Table	227
Table J.15: Grade 11 Physical Science Score Table	228
Table J.16: Grade 11 Sensemaking Score Table.....	229
Table J.17: Grade 11 Critiquing Score Table	230
Table J.18: Grade 11 Investigating Score Table.....	231

Table K.1: Grade 5 Content Disaggregated Subscore Proficiency Classifications	232
Table K.2: Grade 5 Practice Disaggregated Subscore Proficiency Classifications	233
Table K.3: Grade 8 Content Disaggregated Subscore Proficiency Classifications	234
Table K.4: Grade 8 Practice Disaggregated Subscore Proficiency Classifications	235
Table K.5: Grade 11 Content Disaggregated Subscore Proficiency Classifications	236
Table K.6: Grade 11 Practice Disaggregated Subscore Proficiency Classifications	237
Table N.1: Grade 5 Fit and Underfit <i>p-values</i>	246
Table N.2: Grade 8 Fit and Underfit <i>p-values</i>	248
Table N.3: Grade 11 Fit and Underfit <i>p-values</i>	251
Table O.1: Grade 5 Domain Model Parameter Estimates	254
Table O.2: Grade 8 Domain Model Parameter Estimates	256
Table O.3: Grade 11 Domain Model Parameter Estimates	258

PART 1: INTRODUCTION

This Technical Report provides information about the technical characteristics of the 2023 administration of the New Jersey Student Learning Assessment–Science (NJSLA–S) to fifth-, eighth-, and eleventh-grade students. The NJSLA–S is administered under the direction of the New Jersey Department of Education (NJDOE). This report provides extensive detail about the development and operation of NJSLA–S and is intended for use by those who evaluate tests, interpret scores, or use test results for making educational decisions. The documentation in this report is based on the measurement procedures stated in the Standards for Educational and Psychological Testing (American Educational Research Association [AERA], American Psychological Association [APA], & National Council on Measurement in Education [NCME], 2014), hereafter referred to as the “Standards.”

NJSLA–S is an integrated program of testing, accountability, and curricular and instructional support. The test itself is one part of a complex network intended to help schools focus their energies on improving student learning. As such, it can only be evaluated properly within this full context. Detailed descriptions of the NJSLA–S 2023 test development, administration, scoring, and reporting are provided in Parts 2, 3, 4, and 10 of this document. Psychometric discussions of item and test statistics, equating and scaling, reliability, and validity can be found in Parts 6, 7, 8, and 9.

Data for the analyses presented in this Technical Report were collected from the NJSLA–S spring administration from May 1, 2023, through May 26, 2023.

- The standard setting discussed in Part 5 of this report is based on a standard setting study conducted in 2019. More details about the 2019 standard setting study can be found in the 2019 NJSLA–S technical report (NJDOE, 2019).
- Analyses in Parts 6 (Item and Test Statistics) and 8 (Reliability) of this report are based on test results from the entire state population of fifth-, eighth-, and eleventh-grade students.

1.1 Purpose of the Assessment

The 1965 Elementary and Secondary Education Act (ESEA), as reauthorized by the 2015 Every Student Succeeds Act (ESSA, 2015) contained requirements for each state to assess science at least once during grades 3–5, grades 6–9, and grades 10–12. The NJSLA–S measures student proficiency annually in grades 5, 8, and 11 with regard to the New Jersey Student Learning Standards for Science, adopted in 2014 for implementation by the start of the 2016–17 school year for grades 6–12 and by the start of the 2017–18 school year for grades K–5. These science standards are based upon the National Research Council’s (NRC; 2012) Framework for K–12 Science Education, which identifies the science knowledge and skills that all K–12 students should know, and the Next Generation Science Standards (NGSS; NGSS Lead States, 2013), developed collaboratively by stakeholders across 25 states. The emphasis in instruction and assessment is on learning and understanding core principles and theories.

The New Jersey Student Learning Assessments are part of an ongoing system of activities that provide evidence related to student learning. The data from the NJSLA–S, students’ daily interactions with teachers, and their performance on teacher– and district–developed assessments, combine to provide a complete picture of student achievement in science. Schools and local education agencies (LEAs) should use the results to identify strengths and weaknesses in their educational programs. The results may also be used, along with other indicators of student progress, to identify those students who may need instructional support to address any identified knowledge or skill gaps.

1.2 Description of the Assessment

The NJSLA–S assesses students in grades 5, 8, and 11 on their understanding and explanations of scientific phenomena and scenarios. The 2018–19 school year marked the first administration of the NJSLA–S; the spring 2019 operational administration was the assessment’s baseline year, and 2023 was the third year of administration. The assessment was not administered in 2020 and 2021 due to the COVID pandemic.

The NJSLA–S comprises two parts—the performance-based assessment (PBA) and the machine scorable assessment (MSA). The PBA contains one open-ended, constructed-response item and between two and four technology-enhanced items (TEI). The MSA contains a mixture of TEI and multiple-choice items.

Furthermore, the tests cover a range of material. To accomplish the necessary scope, each test item requires students to address multiple underlying variables, with items representing an interaction of disciplinary core ideas (DCIs—within the domains of Physical, Life, and Earth and Space Science), science and engineering practices (SEPs—Investigating, Sensemaking, or Critiquing), and crosscutting concepts (CCC). Every test item counts toward the students’ performance in exactly one reported domain and one reported practice. (Each item is also aligned to a CCC, and the CCC concepts and the knowledge, skills, and abilities (KSAs) associated with them contribute to the overall scale score; however, there is no specific reported CCC performance indicator for the NJSLA–S.)

1.2.1 Content Domains and Scientific Practices

The NJSLA–S is a unidimensional test designed to assess the New Jersey Student Learning Standards for Science (NJSLS–S). The robust standards have been subdivided into six distinct sub-categories for test construction and reporting purposes. The six foundational sub-categories are equally divided between three science content domain categories (Earth and Space, Life, and Physical) and three scientific practice categories (Sensemaking, Critiquing, and Investigating).

Science content domains. Disciplinary core ideas can be classified into three major science content domains: Earth and Space Science, Life Science, and Physical Science. The NJSLA–S is designed to measure student performance in each of the three science content domains. The test development processes focus on balancing each science content domain equally.

Furthermore, within each content domain, each DCI is balanced. (See the Framework for further information.)

1. *Earth and Space Science*. The *Framework* (NRC, 2012) states that “Earth and space sciences (ESS) investigate processes that operate on Earth and also address its place in the solar system” (p. 169). Table 1.2.1 shows the three ESS DCIs and the topics delineated within each.

Table 1.2.1: Earth and Space Science DCIs

DCI Topic Description
ESS1: Earth’s Place in the Universe
ESS1.A: The universe and its stars
ESS1.B: Earth and the solar system
ESS1.C: The history of planet Earth
ESS2: Earth’s Systems
ESS2.A: Earth materials and systems
ESS2.B: Plate tectonics and large-scale system interactions
ESS2.C: The roles of water in Earth’s surface processes
ESS2.D: Weather and climate
ESS2.E: Biogeology
ESS3: Earth and Human Activity
ESS3.A: Natural Resources
ESS3.B: Natural Hazards
ESS3.C: Human Impacts on Earth Systems

2. *Life Science*. The *Framework* (NRC, 2012) for the life sciences (LS) “focus on patterns, processes, and relationships of living organisms” (p. 139). Table 1.2.2 presents the four LS DCIs and their underlying topics.

Table 1.2.2: Life Science DCIs

DCI Topic Description
LS1: From Molecules to Organisms: Structures and Processes
LS1.A: Structure and function
LS1.B: Growth and development of organisms
LS1.C: Organization for matter and energy flow in organisms
LS1.D: Information processing
LS2: Ecosystems: Interactions, Energy, and Dynamics
LS2.A: Interdependent relationships in ecosystems
LS2.B: Cycles of matter and energy transfer in ecosystems
LS2.C: Ecosystem dynamics, functioning, and resilience
LS2.D: Social interactions and group behavior
LS3: Heredity: Inheritance and Variation of Traits
LS3.A: Inheritance of traits
LS3.B: Variation of traits
LS4 Biological Evolution: Unity and Diversity
LS4.A: Evidence of common ancestry and diversity
LS4.B: Natural selection
LS4.C: Adaptation
LS4.D: Biodiversity and humans

3. *Physical Science*. According to the *Framework* (NRC, 2012) the goal of learning physical science (PS) “is to help students see that there are mechanisms of cause and effect in all systems and processes that can be understood through a common set of physical chemical principles” (p. 103). Table 1.2.3 illustrates the four PS DCIs along with the associated detailed topics for each.

Table 1.2.3: Physical Science DCIs

DCI Topic Description
PS1: Matter and its Interactions
Structure and matter
Chemical reactions
PS2: Motion and Stability: Force and Interactions
Force and motion
Types of interactions
Stability and instability in physical systems
PS3: Energy
Definitions of energy
Conservation of energy and energy transfer
Relationship between energy and forces
Energy in chemical processes and everyday life
PS4: Waves and their Applications in Technologies for Information Transfer
Wave properties
Electromagnetic radiation
Information technologies and instrumentation

Scientific practices. The Framework (2012) contains eight different Scientific and Engineering Practices (SEPs). One of the goals of the SEPs is to help “students understand how scientific knowledge develops; such direct involvement gives them an appreciation of the wide range of approaches that are used to investigate, model, and explain the world” (p. 42). Within the context of the NJSLA–S, the SEPs are consolidated into three categories of scientific practices: Investigating, Sensemaking, and Critiquing. Table 1.2.4, adapted from the work of McNeill, et al. (2015), shows how the eight Framework SEPs were consolidated for the purposes of the NJSLA–S.

Table 1.2.4: SEP Consolidation

SEP	Grouping
Asking Questions and Defining Problems (AQDP)	Investigating
Planning and carrying out investigations (PACI)	Investigating
Using mathematics and computational thinking (UMCT)	Investigating
Analyzing and interpreting data (AID)	Sensemaking
Constructing explanations and designing solutions (CEDS)	Sensemaking
Developing and using models (DUM)	Sensemaking
Engaging in argument from evidence (EAE)	Critiquing
Obtaining evaluating and communicating information (OECI)	Critiquing

1. *Investigating*. Investigating Practices (McNeill et al., 2015) involve asking questions, conducting investigations, and using mathematical skills to probe naturally occurring phenomena. Table 1.2.5 delineates the *Framework* definition of each of the Investigating Practices.

Table 1.2.5: Investigating Practices

SEP	NRC Framework
Asking questions and defining problems (AQDP)	Students at any grade level should be able to ask questions of each other about the texts they read, the features of the phenomena they observe, and the conclusions they draw from their models or scientific investigations. For engineering, they should ask questions to define the problem to be solved and to elicit ideas that lead to the constraints and specifications for its solution. (p. 56)
Planning and carrying out investigations (PACI)	Students should have opportunities to plan and carry out several different kinds of investigations during their K–12 years. At all levels, they should engage in investigations that range from those structured by the teacher—in order to expose an issue or question that they would be unlikely to explore on their own (e.g., measuring specific properties of materials)—to those that emerge from students’ own questions. (p. 61)
Using mathematics and computational thinking (UMCT)	Although there are differences in how mathematics and computational thinking are applied in science and in engineering, mathematics often brings these two fields together by enabling engineers to apply the mathematical form of scientific theories and by enabling scientists to use powerful information technologies designed by engineers. Both kinds of professionals can thereby accomplish investigations and analyses and build complex models, which might otherwise be out of the question. (p. 65)

2. *Sensemaking*. Sensemaking Practices (McNeill et al., 2015) are conceptualized as analyzing the data that is produced from an investigation and developing models and explanations that can explain naturally occurring phenomena. Table 1.2.6 illustrates the *Framework* definition of each of the Sensemaking Practices.

Table 1.2.6: Sensemaking Practices

SEP	NRC Framework
Developing and using models (DUM)	Modeling can begin in the earliest grades, with students’ models progressing from concrete “pictures” and/or physical scale models (e.g., a toy car) to more abstract representations of relevant relationships in later grades, such as a diagram representing forces on a particular object in a system. (p. 58)
Analyzing and interpreting data (AID)	Once collected, data must be presented in a form that can reveal any patterns and relationships and that allows results to be communicated to others. Because raw data as such have little meaning, a major practice of scientists is to organize and interpret data through tabulating, graphing, or statistical analysis. Such analysis can bring out the meaning of data—and their relevance—so that they may be used as evidence. (p. 61)
Constructing explanations and designing solutions (CEDS)	Asking students to demonstrate their own understanding of the implications of a scientific idea by developing their own explanations of phenomena, whether based on observations they have made or models they have developed, engages them in an essential part of the process by which conceptual change can occur. (p. 68)

3. *Critiquing*. Critiquing Practices (McNeill et al., 2015) are conceptualized as the ability of students to evaluate information, engage in argument, and communicate whether the models, explanations, or interpretations are adequate representations of naturally occurring phenomena. Table 1.2.7 shows the *Framework* definition of each of the Critiquing Practices.

Table 1.2.7: Critiquing Practices

SEP	NRC Framework
Engaging in argument from evidence (EAE)	The study of science and engineering should produce a sense of the process of argument necessary for advancing and defending a new idea or an explanation of a phenomenon and the norms for conducting such arguments. In that spirit, students should argue for the explanations they construct, defend their interpretations of the associated data, and advocate for the designs they propose. (p. 73)
Obtaining, evaluating and communicating information (OEI)	Any education in science and engineering needs to develop students’ ability to read and produce domain-specific text. As such, every science or engineering lesson is in part a language lesson, particularly reading and producing the genres of texts that are intrinsic to science and engineering. (p. 76)

1.2.2 Crosscutting Concepts

The *Framework* (2012) contains seven different Crosscutting Concepts (CCCs). They were selected to help “students with an organizational framework for connecting knowledge from the various disciplines into a coherent and scientifically based view of the world” (p. 83). Due to reporting constraints, the CCCs are the lowest priority of the three dimensions described in the Framework. However, because each item is aligned to a CCC, the CCC concepts and the knowledge, skills, and abilities associated with them are still being assessed by the NJSLA–S and contribute to the overall NJSLA–S scale score. Table 1.2.8 shows the CCCs being measured by the NJSLA–S.

Table 1.2.8: Crosscutting Concepts

CCC	NRC Framework (p. 84)
Patterns	Observed patterns of forms and events guide organization and classification, and they prompt questions about relationships and the factors that influence them.
Cause and Effect	Events have causes, sometimes simple, sometimes multifaceted. A major activity of science is investigating and explaining causal relationships and the mechanisms by which they are mediated. Such mechanisms can then be tested across given contexts and used to predict and explain events in new contexts.
Scale, Proportion, and Quantity	In considering phenomena, it is critical to recognize what is relevant at different measures of size, time, and energy and to recognize how changes in scale, proportion, or quantity affect a system’s structure or performance.
Systems and System Models	Defining the system under study—specifying its boundaries and making explicit a model of that system—provides tools for understanding and testing ideas that are applicable throughout science and engineering.
Energy and Matter	Tracking fluxes of energy and matter into, out of, and within systems helps one understand the systems’ possibilities and limitations.
Structure and Function	The way in which an object or living thing is shaped and its substructure determine many of its properties and functions.
Stability and Change	For natural and built systems alike, conditions of stability and determinants of rates of change or evolution of a system are critical elements of study.

1.2.3 Types of Scores

Student performance on the NJSLA–S is described using scale scores and performance levels. Each grade level has its own grade-specific scale that represents a composite score of student performance on the three NJSLA–S dimensions (DCIs, SEPs, and CCCs). Student performance is classified into four grade-specific performance levels based on the NJSLA–S Performance-Level Descriptors (PLDs). Both the scale score and the performance levels are described below.

- **Scale Scores.** The NJSLA–S reports scale scores to indicate a student’s performance. A scale score is a conversion of the raw score (the total number of points a student earned on the test), using a predetermined mathematical algorithm, to permit legitimate and meaningful comparisons. As such, they provide the best generalized information about overall performance. The total scores in science are reported as scale scores with a range of 100 to 300.
- **Performance Levels.** One of the primary purposes of the NJSLA–S is to identify areas of curricular strength and weakness by examining the extent to which students meet the established performance expectations in science. Based on test results, a student’s performance is categorized as being at one of four performance levels, each of which is defined by a student’s scale score and used to report overall student performance on the NJSLA–S. Grade-appropriate Performance-Level Descriptors (PLDs) translate these performance levels into words. They describe the KSAs students should have at each performance level, Level 1 through Level 4. Each performance level is associated with a range of scale scores, as indicated in Table 1.2.9:

Table 1.2.9: NJSLA–S Scale Score Ranges

Grade	Level 1	Level 2	Level 3	Level 4
5	100–149	150–199	200–242	243–300
8	100–149	150–199	200–230	231–300
11	100–157	158–199	200–249	250–300

Students performing at Level 3 and Level 4 are considered proficient and above; they demonstrate appropriate or exemplary understanding of the DCIs and SEPs. Students performing at Level 1 and Level 2 are considered below the state minimum proficiency level. They demonstrate minimal or partial understanding of the DCIs and SEPs. Students at this performance level may need additional instructional support, which could be individual or programmatic intervention.

Student performance is also classified as “Below,” “Near/Met,” or “Above” expectations in each of the three content domains (Earth and Space, Life, and Physical Science) and the three scientific practices (Investigating, Sensemaking, Critiquing). These subscore performance classifications are primarily meant to provide teachers, schools, and administrators with feedback as to the specific KSAs that their students displayed on the NJSLA–S. Individual students and their parents and teachers receive student-level data on these subscores.

1.3 Organizational Support

The New Jersey Department of Education’s Office of Assessments coordinates the development and implementation of the NJSLA–S. In addition to planning, scheduling, and directing all NJSLA–S activities, the staff is extensively involved in numerous test-design, item and statistical review, security, quality-assurance, and analytical procedures. Measurement Incorporated (MI), the primary contractor for the NJSLA–S at grades 5, 8, and 11, is responsible for all aspects of the testing program, including activities such as program management, development of tests, publishing documents for test administration, handscoring constructed-response items, and psychometric support (including standard setting). Pearson, the sub-contractor for NJSLA–S, provides item banking; test registration, administration, and digital delivery; and reporting. MI and Pearson work closely together under the direction of the Office of Assessments to ensure ancillary materials and administrative procedures closely match those of the NJSLA–Math and NJSLA–ELA assessments.

PART 2: TEST DEVELOPMENT

The NJSLA–S is aligned to the New Jersey Student Learning Standards for Science (NJSLS–S), adopted in 2014, which in turn are based upon the National Research Council’s *Framework for K–12 Science Education* and the Next Generation Science Standards (NGSS).

The Test Design and Development chapter within the Standards (2014) outlines a series of five primary phases of the test development process: (1) test specifications; (2) item development and review; (3) assembling and evaluating test forms; (4) development of procedures and materials for test administration and scoring; and (5) test revisions (p. 83). The following sections in Part 2 detail the NJSLA–S test specifications, item development processes, and both the test construction processes and their results in 2023. The development of procedures and materials for test administration and scoring is covered in Parts 2 and 3.

2.1 Test Specifications

According to the *Standards*, “[t]he term *test specifications* is sometimes limited to description of the content and format of the test. In the *Standards*, test specifications are defined more broadly to also include documentation of the purpose and intended uses of the test, as well as detailed decisions about content, format, test length, psychometric characteristics of the items and test, delivery mode, administration, scoring, and score reporting” (p. 76).

The NJSLA–S was developed to measure the knowledge, skills, and abilities (KSAs) identified in the NJSLS–S in grades 5, 8, and 11. The test is designed to provide reporting information for student ability levels at the holistic level and each of the three science content domains (Earth and Space, Life, and Physical) and the three scientific practices (Investigating, Sensemaking, and Critiquing). The test specifications call for a balanced test design that prioritizes each science content domain and each DCI, each scientific practice, and each SEP, as well as all seven CCCs. (Please refer to Section 1.2 of this document for an explanation of the DCIs, SEPS, and CCCs.) The detailed information recommended in the Standards is presented in the sections that follow.

2.1.1 Test Blueprints

Table 2.1.1 depicts the test blueprint—the numbers of items comprising each part of the test—for all grades. Note that each multiple-choice (MC) item is worth one point; each technology-enhanced (TE) item is worth one point; each constructed-response (CR) item is worth three or four points. Each constructed-response item is scored using an item-specific rubric. The table summarizes the number of items on the operational NJSLA–S for each of the six reporting categories as well as for both the Performance-Based Assessment (PBA) and Machine-Scorable Assessment (MSA) components. An explanation of the PBA and MSA components is provided in the following section.

Table 2.1.1: Test Blueprints

Domain	Practice	Grade 5 PBA	Grade 5 MSA	Grade 8 PBA	Grade 8 MSA	Grade 11 PBA	Grade 11 MSA
PS	Investigating <i>AQDP, PACI, UMCT</i>	1–2	3–5	1–2	4–7	1–2	4–8
PS	Sensemaking <i>DUM, AID, CEDS</i>	1–2	3–5	1–2	4–7	1–2	4–8
PS	Critiquing <i>EAE, OECI</i>	1–2	3–5	1–2	4–7	1–2	4–8
PS	Total Items	3–5	11–13	3–5	14–18	3–5	15–21
LS	Investigating <i>AQDP, PACI, UMCT</i>	1–2	3–5	1–2	4–7	1–2	4–8
LS	Sensemaking <i>DUM, AID, CEDS</i>	1–2	3–5	1–2	4–7	1–2	4–8
LS	Critiquing <i>EAE, OECI</i>	1–2	3–5	1–2	4–7	1–2	4–8
LS	Total Items	3–5	11–13	3–5	14–18	3–5	15–21
ESS	Investigating <i>AQDP, PACI, UMCT</i>	1–2	3–5	1–2	4–7	1–2	4–8
ESS	Sensemaking <i>DUM, AID, CEDS</i>	1–2	3–5	1–2	4–7	1–2	4–8
ESS	Critiquing <i>EAE, OECI</i>	1–2	3–5	1–2	4–7	1–2	4–8
ESS	Total Items	3–5	11–13	3–5	14–18	3–5	15–21

2.1.2 Unit Design

The NJSLA–S consists of four units—three operational and one field test. The units are numbered 1–4, and the field test unit placement varies from year to year. Each unit contains a machine-scorable (MSA) and a performance-based (PBA) component; a balance of Earth and Space, Life, and Physical Science items; a balance of Investigating, Sensemaking, and Critiquing Practice items; a prescribed proportion of MC, TE, and CR item types; and psychometric constraints that are discussed in Section 2.4 of this technical report.

Each MSA and PBA component of a unit is linked to naturally occurring phenomena that provide the impetus for scenarios. The students are provided with the scenario and subsequently presented with two to five items that measure their mastery of the NJSLA–S. All items attached to a phenomenon-based scenario are independent—that is, for example, if a PBA section contains four total items, a student’s response to one of the four items will not impact that student’s ability to correctly answer any of the other three. Figure 2.1.1 illustrates the composition of a sample grade 5 unit.

MSA: 4 stimuli, 3 items each				PBA: 1 stimulus, 5 items
<u>Stim. 1</u> 3 TE 3 points	<u>Stim. 2</u> 2 TE, 1 MC 3 points	<u>Stim. 3</u> 2 TE, 1 MC 3 points	<u>Stim. 4</u> 2 TE, 1 MC 3 points	<ul style="list-style-type: none"> • 4 one-point TEs • 1 four-point CR
Total # items, MSA: 12 Total points, MSA: 12				
Total # items, Unit: 17 Total points, Unit: 20				

Figure 2.1.1. Sample Grade 5 Unit

Machine-scorable assessment (MSA). The MSA component of the NJSLA–S is defined as the portion of the assessment that is scored by a computer. Each cluster of MSA items contains a context-dependent stimulus that presents the students with a naturally occurring phenomenon. Depending on the grade level, each unit contains anywhere from four to seven stimuli, and each stimulus is associated with three to six items. MSA items can be either multiple-choice (MC) or technology-enhanced (TE) items, but within each unit no more than 50% of the MSA items can be MC items.

Performance-based assessment (PBA). The PBA component of the NJSLA–S is defined as that portion of the test which requires students to display KSAs to a greater degree of cognitive depth, the degree to which the student displayed depth of knowledge and expertise; it is based on more complex phenomena than the MSA section. The PBA components (one per unit) contain one stimulus, each of which can accommodate two to four TE items and one constructed-response (CR) item. In 2023, NJDOE required that the PBA section contain 7 to 8 total points, with three or four of those points coming from the CR item.

2.1.3 Item Types

Table 2.1.2 describes each NJSLA–S item type. Three main types of items comprise the NJSLA–S: multiple-choice (MC), technology-enhanced (TE), and constructed-response (CR).

- MC items all have a key (A, B, C, or D) associated with them, and students are asked to select the best of the four options. MC items are scored dichotomously, 0/1.
- TE items require students to interact with more complex methods of answering the items. Examples of TE item interactions include drop-down choice; hot spot; text entry; drag and drop; multiple selection; and ordering. TE items are scored dichotomously, 0/1.
- CR items are open-ended questions designed to elicit a student response to a range of KSAs that are challenging to measure with traditional MC or TE items. All CR items are rubric-dependent and scored by a human reader.

Table 2.1.2: NJSLA–S Item Types

Item Type	Description
MC: Multiple Choice	Select one response from four possible options (A, B, C, D).
TE: Multiple Selection	Select two or more answer options.
TE: Drop-Down Choice	Select from a drop-down menu embedded in the prompt.
TE: Ordering	Drag text or image-based options into a particular order.
TE: Drag and Drop	Place one or more text or graphic choices into blank spots within a sentence, table, or diagram.
TE: Matching in a Table	Check a box in the table to match the row to the column.
TE: Text Entry	Type a brief constrained response to the question.
TE: Bar Graph	Drag each bar to the correct length on the graph.
TE: Hot Spot	Select one or more regions on a graphic or image to identify an answer.
TE: Hot Text	Select one or more sentences within a paragraph of text.
CR: Constructed-Response	Type an extended open-ended response to the prompt.

2.2 Item Development Processes

NJSLA–S item development was conducted by MI and Pearson with oversight from NJDOE staff and the New Jersey Science Advisory Committee (NJSAC). The item development process is rigorous and involves item writers, content specialists, editors, graphic artists, programmers, scoring experts, and psychometricians. The resulting products are phenomenon-based scenarios (PBSs) and items that are aligned to the NJSLS–S and the NJSLA–S reporting categories. The PBSs and their items are all housed in Pearson’s Assessment Banking for Building and Interoperability (ABBI) item banking system. ABBI is specifically designed to handle next-generation online, interactive, and accessible content. The steps in the item development process are detailed in the sections below. It warrants emphasis that between the NJSAC and the New Jersey Bias and Sensitivity Committee (NJBSC), New Jersey educators and administrators were intimately and actively involved in the item development process; each item that appears on the NJSLA–S was reviewed and approved multiple times. The principles of universal design were incorporated into the development of NJSLA–S phenomenon-based stimuli and their items. There are seven elements of assessments designed to meet the expectations of universal design (Thompson, et al., 2002). The seven elements are listed below. All seven elements are incorporated into each step within the item writing process; however, there are specific steps where elements are emphasized and reviewed more extensively by experts.

1. Inclusive assessment population
2. Precisely defined constructs
3. Accessible, non-biased items
4. Amenable to accommodations
5. Simple, clear, and intuitive instructions and procedures
6. Maximum readability and comprehensibility
7. Maximum legibility

2.2.1 Item Writing

The item development process begins with the training of item writers on the specifications of NJSLS–S item development. Per the principles of universal design, item writers are trained on how to write PBSs and items that clearly communicate the task at hand for the students while also carefully maintaining alignment to the construct the NJSLS–S is intending to measure.

Once the item writers start item development, they initially identify naturally occurring phenomena that are pertinent for assessing the NJSLS–S. Next, the item writers research and develop a scenario that contains specific examples of how a phenomenon manifests itself in nature. (Priority is given to scenarios that are specifically relevant to New Jersey, such as native species of plants and animals, weather patterns, geological features, etc.)

Item writers then begin writing clusters of items related to the phenomenon-based scenario. Each item is aligned to a single scientific content domain and DCI, a scientific practice and SEP, and a CCC. To measure as many KSAs as possible with a single item cluster, item writers are instructed to vary the SEPs and CCCs within each cluster of items. An item type is typically assigned according to the item type’s effectiveness and efficiency in measuring the targeted KSAs. To best align the test to the NJSLS–S blueprint, item writers are instructed to use no more than 50% MC items in each cluster of items. All items are also aligned to one of Webb’s (1997; 2002) Depth-of-Knowledge (DOK) classifications.

Once a phenomenon-based scenario has a diverse cluster of four to ten items, it enters the item writing peer review process. Two different item writers review the scientific justification for the phenomenon and scenario, the alignment of the items to the NJSLS–S, the readability and appropriateness of the content, and any other conceptual understandings inherent to either the scenario or item cluster. The item writers functioning as peer reviewers iteratively rework the scenario with the original item writer until they all reach agreement.

2.2.2 Content Specialist Review

Up to three content specialists review each PBS. The first content specialist review focuses on reviewing references and evaluating the science, scope, and structure of the PBS. If major revisions are needed, then the PBS is sent back to the initial item writer; if the revisions are minor, then the PBS is moved onto the second stage of the content specialist review process.

The second content specialist review focuses on universal design element 2: precisely defined constructs. The content specialist ensures the correct alignment of the PBS and all its associated items to:

- NJSLS–S
- DCI
- SEP
- CCC
- Content Domain Reporting Category
- Scientific Practices Reporting Category

If revisions are suggested, then the first content specialist and the second content specialist discuss the revisions with the item writer. If all parties agree, then the PBS is revised. If resolution is needed, then a third content specialist settles any disputes.

As a final step in the content specialist review process, the third content specialist is also charged with verifying that all the science in the PBS is accurate, that each item is answerable based on the information presented in the PBS, that all answer keys are correct, and that the alignment is in accordance with the NJSLS–S. During this step, universal design elements 5 and 6 are thoroughly reviewed to confirm that the PBS and its items have clear student instructions, that its readability is appropriate, and that it strictly adheres to the New Jersey Science Style Guidelines. Upon the final content review, the PBS is sent to editorial for its review.

2.2.3 Editorial Review

Two editors review each PBS. Their focus is on verifying that universal design elements 5, 6, and 7 are respected. The editors are charged with verifying the readability of the PBS (i.e., the PBS is easy to read and not unnecessarily complex) and checking for grammatical, spelling, and careless errors in the text. They also review each graphic or table for legibility (e.g., graphics have proper legends). Other editorial tasks include ensuring the direction lines and other components within the PBS all adhere to the New Jersey Style Guidelines. Once the PBS has passed both editorial reviews, then it is ready for review by the New Jersey Science Advisory Committee (NJSAC).

2.2.4 NJ Science Advisory Committee Content Review

All items on the NJSLS–S are reviewed by the New Jersey educators who compose the New Jersey Science Advisory Committee (NJSAC). In 2023, the NJSAC comprised a diverse group of New Jersey science educators representing 19 of the 21 New Jersey counties. The districts each NJSAC member represents and the counties they come from are presented in Appendix B.

The NJSAC is the final authority on universal design principle 2: precisely defined constructs. For the 2023 administration, committee members ensured that each item was aligned to the vision set forth in the NJSLS–S, which includes properly aligning each item to a DCI, SEP, and CCC and confirming that the PBS’s content was accurate. They also reviewed the PBS and its items in accordance with universal design principles 5 and 6 by confirming that the items had grade-appropriate vocabulary, that the reading level was appropriate, and that item instructions were simple and clear.

The NJSAC took an active role in editing the content of the items during their item reviews. They collectively interacted with each other, NJDOE, and the content specialists to make suggestions and offer solutions to improve the quality of item development and the NJSLS–S test. The NJSAC item reviews were held both in-person at locations approved by NJDOE, and in secure, online platforms. The PBSs and items were all reviewed in ABBI.

2.2.5 Bias and Sensitivity Committee Review

If an item passes the NJSAC’s content review, it proceeds to review by the New Jersey Bias and Sensitivity Committee (NJBSC). This step in the item development processes is where extra emphasis is placed on universal design elements 1, 3, and 4. The NJBSC makes sure that all students have the opportunity to show what they know regardless of their background or the test form they took. They ensure that each item is free from bias and meets the industry guidelines for fairness and sensitivity (ETS, 2015). As described in Standard 3.3 (AERA, APA, NCME, 2014), this step helps guard against the introduction of construct-irrelevant language, images, or situations that might either offend or be more familiar to one group of New Jersey students than another.

Of the ten NJBSC members, nine taught special education status students, seven specialized in teaching students designated as English learners, and five were bilingual. Collectively, they had over 100 years of teaching experience. As with NJSAC content reviews, the NJBSC reviews were conducted in-person and in ABBI; the NJBSC actively worked with each other, NJDOE, and the content specialists to limit test bias. The NJBSC’s district and county representation is presented in Appendix B.

2.2.6 Field Test

Once an item has passed both reviews from the NJSAC and the NJBSC, it is eligible for placement onto one of that year’s field test units. The purpose of field testing is to gather data to evaluate whether an item is performing as it was intended. The field test items are placed into different field test units. Each grade has at least 10 field test units, and there may be as many as 18 field test units. The units are placed into the operational test form in designated positions that rotate from year to year. Each unit is reviewed by content specialists and NJDOE to ensure that none of the field test items cue answers to the operational test items. The field test units are spiraled at the student level, which ensures that the students who take any of the field test units are a demographically representative sample of New Jersey students. A minimum of 4,000 students respond to each NJSLA–S field test item so that the samples are large enough that the resulting item statistics that are presented at the NJSLA–S Statistical Reviews are stable.

2.2.7 Statistical Review

The NJSAC reviews a battery of statistics for all field test items at the NJSLA–S statistical review. MI’s psychometric staff leads the statistical review and either trains or re-trains all NJSAC members on how to interpret the item statistics so that they can make effective evaluative judgments as to the usefulness of the item. Each committee member is presented the NJSLA–S Statistical Review Reference Sheet that provides them with quick access to definitions of the statistics and the optimal range of values. The NJSAC decides whether the item should be “Accepted,” “Rejected,” or “Revised and Re-Field Tested.” MI’s lead content specialists and an NJSAC committee member simultaneously log the decisions made by the committee, including whether an item is to be revised and how to best improve the item.

MI’s psychometric staff emphasizes to the NJSAC that feedback from statistical review is used to refine future item development in an effort to constantly improve the quality of NJSLA–S stimuli and items. The NJSLA–S Statistical Review Reference Sheet given to panelists is presented in Appendix C.

2.2.8 Second Bias and Sensitivity Review

As a crucial part of statistical review, the NJBSC reviews all items flagged for being possibly biased against groups of New Jersey students. Groups of students include Male/Female, White/Black, White/Hispanic, and White/Asian. The NJBSC members are trained by MI staff prior to reviewing the items on how to interpret the statistics they will see, which include differential item functioning (DIF) statistics and the percentage of each group of students that selected each answer option. DIF is described in Section 2.3.1.1.

2.2.9 Ready for Operational Testing

Once an item has passed both statistical review and the second bias and sensitivity review, it is then eligible to be placed onto an operational test form, and its status in ABBI is updated accordingly.

2.3 Test Construction Process

The NJSLA–S test construction process ensures that the operational test forms balance the specifications set forth in the test blueprint, along with other psychometric constraints. Each form is built to measure students across the whole spectrum of ability levels and to foster valid interpretations of test scores in adherence to the standards for test design and development put forth in the Standards (AERA, APA, NCME, 2014). The steps and constraints associated with constructing the NJSLA–S operational tests are detailed in the following sections. An evaluation of the results of the test construction process is presented in Section 2.4.

2.3.1 Test Construction—First Draft

The first step in the NJSLA–S test construction process involves MI’s content staff manually selecting approved items that best match the NJSLA–S test blueprint and statistical constraints. The process of selecting items is contingent upon the state of the item bank at each grade level. If specific content constraints are challenging to fulfill given the types of items present within the item bank, then those content constraints are given priority in the initial selection of items. Next, items are selected iteratively based on which content constraints need to be fulfilled while simultaneously balancing the various statistical constraints. Detailed descriptions of the statistical constraints are presented in Section 2.3.1.1.

2.3.1.1 Test construction statistical constraints. To ensure that the NJSLA–S operational test form is reliable and fosters valid interpretations, the following statistical constraints are used by MI’s content staff during the test construction process. The primary goal is to balance the content and statistical constraints for the test as a whole; when possible, each unit is designed based on the same statistical constraints. Table 2.3.1 provides a summary of the NJSLA–S test construction constraints.

Item difficulty. Each test form is constructed to a specific difficulty level. The most important decision made from the NJSLA–S is at the Level 3 cut score because it is the place on the scale associated with whether students are classified as proficient. To maximize the reliability of those decisions, the average item difficulty parameter of the test form should be as close to the Level 3 cut score as possible.

Item discrimination. Item discrimination refers to the ability of the item to discriminate between students with different abilities. A poorly discriminating item could indicate ineffective measurement of the NJSLA–S scale and reduces test form reliability. Under classical test theory, item discrimination is measured via the item-total correlation, which can range from –1.0 to 1.0; items with item-total correlations that are below .2 are only selected for placement on the operational test form if no other viable options are available. Items with negative discrimination are not selected.

IRT model fit. The NJSLA–S uses an item response theory (IRT) model called the partial credit model (PCM; Masters, 1982) to estimate student ability levels. The PCM makes certain assumptions that, if violated, could impact the validity of interpretations made from NJSLA–S test scores. Statistical constraints based on PCM model fit statistics include infit, outfit, Rasch discrimination, and lower asymptote, which are discussed in detail in Section 6.2.2 of this report. During test construction, the mean item infit, outfit, and Rasch discrimination statistics are all constrained to be as close to 1.0 as possible. If an individual item has an infit or outfit statistic outside of the acceptable range of 0.7 to 1.3 or a Rasch discrimination statistic outside of the acceptable range of 0.5 to 1.5, it is only used if no other viable options are available. The lower asymptote statistic is constrained to be as close to zero as possible; any item whose lower asymptote is greater than 0.1 is flagged and only used if necessary.

Time on items. The NJSLA–S is not designed to be a speeded test; consequently, almost all students should be able to finish it within the allotted time. Items are selected to minimize the median time spent on the test. If the median time spent on items is greater than the total test time for a test unit minus 30 minutes, then items that are taking students too long are replaced by items that take less time, unless no other options are available.

Differential Item Functioning. Differential Item Functioning (DIF) exists when different groups of students have different probabilities of getting an item correct, after controlling for their ability levels. NJSLA–S comparison groups include Male/Female, White/Black, White/Hispanic, and White/Asian. If any item favors one group over another based on the ETS Mantel-Haenszel (Dorans & Holland, 1993; Zieky, 1993) and Penfield (2007) DIF classification methods, that item is classified as demonstrating either “B” or “C” level DIF. All items classified as either “B” or “C” are reviewed by the New Jersey Bias and Sensitivity Committee during the statistical review process. If they deem an item biased, then it is ineligible for placement on the operational

NJSLA–S regardless of DIF classification. A small number of “B” items can be used to maintain the test blueprint, whereas “C” items are not used on the operational NJSLA–S.

Table 2.3.1: Summary of NJSLA–S Test Construction Statistical Constraints

Statistical Constraint	Description
Item Difficulty	Average Rasch B is as close as possible to the Level 3 theta cut score.
Item Discrimination	Item-total correlations are greater than 0.2.
IRT Model Fit	<ul style="list-style-type: none"> • Item Infit and Outfit statistics range from 0.7 to 1.3 and average 1.0. • Item Discrimination statistics range from 0.5 to 1.5 and average 1.0. • Item Lower Asymptote statistics < 0.1 and average as close to 0.0 as possible.
Time on Items	Total median time on operational items < (total operational test time – 30 minutes).
DIF	<ul style="list-style-type: none"> • “B” items are only used if necessary. • “C” items are not used.

2.3.2 Test Construction Content and Psychometric Review

After MI’s content staff finishes the first draft of the operational test forms, content specialists at each grade level check the forms to ensure that no items cue each other or have content that is too similar. The content and psychometric review is an iterative process between content specialists and psychometricians. If, during the review, psychometricians identify items that better meet the statistical constraints and other psychometric properties of the NJSLA–S, the candidate items are replaced. The content and psychometric review then resumes until the test matches the content and statistical criteria of the NJSLA–S. It should be noted that certain candidate items have only undergone field testing at this stage, which means the item statistics are based on smaller field-test samples. However, the psychometric analyses presented in Part 6 of this report rely on the full operational test data from the current year. The content and psychometric review then resumes until the test matches the content and statistical constraints of the NJSLA–S.

2.3.3 Test Construction NJDOE Review

All NJSLA–S test forms are reviewed and approved by NJDOE. Once content and psychometrics have agreed upon the operational test forms, they are sent to NJDOE for approval. After NJDOE approves the test forms they are released for final editorial review and publishing.

2.4 2023 NJSLA–S Test Construction

Overall, the test construction process achieved forms that matched the balance required by the test blueprint in Section 2.1.1. All grade levels had 7- or 8-point PBA sections representing each of the three content domains. Moreover, Table 2.4.1 shows that at each grade level, the science content domains were sufficiently balanced across the scientific practices reporting categories. The largest offset among content domains was in the Sensemaking category of the grade 11

test, where there were 11 points aligned to Life Science content domain; 7 points aligned to the Earth Science domain; and 9 points aligned to the Physical domain.

Table 2.4.1: Points Available by Domain and Practice

Grade	Practice	Earth	Life	Physical
5	Investigating	5	4	7
5	Sensemaking	5	8	4
5	Critiquing	11	8	8
8	Investigating	8	7	7
8	Sensemaking	7	9	10
8	Critiquing	7	10	7
11	Investigating	8	7	11
11	Sensemaking	7	11	9
11	Critiquing	9	8	8

2.4.1 Grade 5 Test Construction

For grade 5, out of 60 total score points, the three content domains were well balanced, ranging from 19 to 21 points each, as illustrated in Table 2.4.2. Each content domain had one PBA section devoted to it. The scientific practices were less balanced than content domains for grade 5, with only 16 out of 60 points being allocated to the Investigating reporting category. Despite being less than ideal, the 16 points were still enough to produce reliable measures of student Investigating abilities. Other content considerations that were met included: MC items only made up 14 points of the total test score (less than 50%), each unit contained a CR item, and all eight SEPs and all seven CCCs were represented by multiple points on the test. All 11 of the major DCI clusters were represented by multiple points. Table 2.4.2 details the item and point totals for each of the six reporting categories. Tables 2.4.3 through 2.4.5 show the distributions of DCIs, SEPs, and CCCs.

Table 2.4.2: 2023 NJSLA–S Grade 5 Item and Point Totals by Reporting Category

Domains/Practices	MC Items	TE Items	CR Items	Items	Points
Earth and Space	5	12	1	18	21
Life	2	14	1	17	20
Physical	7	8	1	16	19
Total–Domains	14	34	3	51	60
Investigating	7	9	0	16	16
Sensemaking	1	16	0	17	17
Critiquing	6	9	3	18	27
Total–Practices	14	34	3	51	60

Table 2.4.3: 2023 NJSLA–S Grade 5 DCIs

DCI	Items	Points
ESS1	3	3
ESS2	15	18
LS1	3	3
LS2	7	7
LS3	4	7
LS4	3	3
PS1	4	7
PS2	4	4
PS3	5	5
PS4	3	3

Table 2.4.4: 2023 NJSLA–S Grade 5 SEPs

SEP	Items	Points
AQDP	8	8
PACI	5	5
UMCT	3	3
DUM	4	4
AID	7	7
CEDS	6	6
EAE	12	18
OECI	6	9

Table 2.4.5: 2023 NJSLA–S Grade 5 CCCs

CCC	Items	Points
C & E	11	14
E & M	2	2
Patterns	11	14
S & SM	13	16
S, P, & Q	5	5
SC	3	3
SF	6	6

The statistical constraints for the 2023 Grade 5 NJSLA–S operational test form were met. Three items had item-total correlations below the 0.20 threshold indicating a low discriminating item (see Section 6.1). However, as shown in Table 2.4.6, the mean item-total correlation was higher than the mean target for the form (0.35). Additionally, each of the model-fit statistics averaged close to their target values. The mean and median test times of 97.599 and 78.567 minutes, respectively, were well below the 105-minute threshold, and out of 51 DIF classifications for each of the four group comparisons (i.e., Male/Female, White/Black, White/Hispanic and

White/Asian), there were zero values categorized as “C” and only five values categorized as “B” for all the DIF analyses performed. All “B” DIF items were approved for operational test use by the NJBSC as described in Section 2.3.1.1. Tables 2.4.6 and 2.4.7 summarize the test construction and DIF statistics.

Table 2.4.6: 2023 NJSLA–S Grade 5 Test Construction Statistics

Statistic	Average	Target	Flags
Rasch B	0.458	0.904	N/A
IT Correlation	0.432	> 0.35	3
Infit	1.064	1.00	7
Outfit	1.143	1.00	17
PCM Discrim.	0.888	1.00	8
Lower Asymptote	0.048	0.00	4
Median Time (min)	78.567	< 105	N/A

Table 2.4.7: 2023 NJSLA–S Grade 5 Test Construction DIF Classifications

Groups	A	B	C
Male/Female	49	2	0
White/Black	49	2	0
White/Hispanic	50	1	0
White/Asian	51	0	0

2.4.2 Grade 8 Test Construction

At grade 8, the content domains were well balanced. Out of 72 total score points, the three content domains ranged from 22 to 26 points each, as illustrated in Table 2.4.8. Each content domain had one PBA section devoted to it. The scientific practices were also well balanced, with less than 4 points difference among the reporting categories: Investigating 22 points, Critiquing 24 points, and Sensemaking 26 points. Other content considerations that were met included: MC items only made up 18 points (less than 50%) of the total test score; each unit contained a CR item, and all eight SEPs and all seven CCCs were represented by multiple points on the test. Similarly, all 11 major DCI clusters were represented by at least four items. Table 2.4.8 details the item and point totals for each of the six reporting categories; Tables 2.4.9 through 2.4.11 show the distributions of DCIs, SEPs, and CCCs for grade 8.

Table 2.4.8: 2023 NJSLA–S Grade 8 Item and Point Totals by Reporting Category

Domains/Practices	MC Items	TE Items	CR Items	Items	Points
Earth and Space	5	14	1	20	22
Life	8	15	1	24	26
Physical	5	15	1	21	24
Total–Domains	18	44	3	65	72
Investigating	7	12	1	20	22
Sensemaking	6	16	1	23	26
Critiquing	5	16	1	22	24
Total–Practices	18	44	3	65	72

Table 2.4.9: 2023 NJSLA–S Grade 8 DCIs

DCI	Items	Points
ESS1	10	10
ESS2	5	7
ESS3	5	5
LS1	6	6
LS2	7	9
LS3	7	7
LS4	4	4
PS1	6	6
PS2	7	10
PS3	3	3
PS4	5	5

Table 2.4.10: 2023 NJSLA–S Grade 8 SEPs

SEP	Items	Points
AQDP	9	9
PACI	6	8
UMCT	5	5
DUM	5	5
AID	11	11
CEDS	7	10
EAE	14	14
OECI	8	10

Table 2.4.11: 2023 NJSLA–S Grade 8 CCCs

CCC	Items	Points
C & E	20	25
E & M	3	3
Patterns	24	24
S & SM	4	4
S, P, & Q	6	6
SC	3	5
SF	5	5

The statistical constraints for the 2023 Grade 8 NJSLA–S operational test form were met. Three grade 8 items were flagged for having item-total correlations below the 0.20 threshold indicating a low discriminating item (see Section 6.1), including one with the lowest value of 0.142. However, the other two had values above 0.15, and as shown in Table 2.4.12, the average item-total correlation was close to the mean target value for the form (0.35). The infit, outfit, and PCM discrimination model fit statistics each averaged close to their ideal values of 1.00. The mean and median test time of 97.877 and 81.300 minutes, respectively, were well below the 105-minute threshold, and out of 65 DIF classifications for each of the four group comparisons (i.e., Male/Female, White/Black, White/Hispanic, and White/Asian), there were zero values categorized as “C” and eight values categorized as “B.” All “B” DIF items were approved for operational test use by the NJBSC as described in Section 2.3.1.1. Tables 2.4.12 and 2.4.13 summarize the test construction and DIF statistics.

Table 2.4.12: 2023 NJSLA–S Grade 8 Test Construction Statistics

Statistic	Average	Target	Flags
Rasch B	0.318	0.416	N/A
IT Correlation	0.368	> 0.35	3
Infit	1.011	1.00	0
Outfit	1.034	1.00	3
PCM Discrim.	0.990	1.00	0
Lower Asymptote	0.019	0.00	1
Median Time (min)	81.300	< 105	N/A

Table 2.4.13: 2023 NJSLA–S Grade 8 Test Construction DIF Classifications

Groups	A	B	C
Male/Female	62	3	0
White/Black	62	3	0
White/Hispanic	63	2	0
White/Asian	65	0	0

2.4.3 Grade 11 Test Construction

The grade 11 content domains were well balanced. Out of 78 total score points, the three content domains ranged from 24 to 28 points each. Each content domain had one PBA section. The scientific practices were also well balanced, with only 2 points difference among the reporting categories: Critiquing 25 points, Investigating 26 points, and Sensemaking 27 points. Other content considerations that were met included: MC items only made up 31 points (less than 50%) of the total test score; each unit contained a CR item, and all eight SEPs and all eleven DCIs were represented by multiple points on the test. Each of the seven CCCs was adequately represented. Table 2.4.14 details the item and point totals for each of the six reporting categories; Tables 2.4.15 through 2.4.17 show the distributions of DCIs, SEPs, and CCCs for grade 11.

Table 2.4.14: 2023 NJSLA–S Grade 11 Item and Point Totals by Reporting Category

Domains/Practices	MC Items	TE Items	CR Items	Items	Points
Earth and Space	8	12	1	21	24
Life	10	12	1	23	26
Physical	13	12	1	26	28
Total–Domains	31	36	3	70	78
Investigating	13	10	1	24	26
Sensemaking	10	13	1	24	27
Critiquing	8	13	1	22	25
Total–Practices	31	36	3	70	78

Table 2.4.15: 2023 NJSLA–S Grade 11 DCIs

DCI	Items	Points
ESS1	7	7
ESS2	6	6
ESS3	8	11
LS1	5	5
LS2	14	17
LS3	1	1
LS4	3	3
PS1	9	9
PS2	1	1
PS3	12	14
PS4	4	4

Table 2.4.16: 2023 NJSLA–S Grade 11 SEPs

SEP	Items	Points
AQDP	7	7
PACI	8	8
UMCT	9	11
DUM	5	5
AID	15	15
CEDS	4	7
EAE	16	19
OECI	6	6

Table 2.4.17: 2023 NJSLA–S Grade 11 CCCs

CCC	Items	Points
C & E	10	13
E & M	6	6
Patterns	11	11
S & SM	22	27
S, P, & Q	10	10
SC	5	5
SF	6	6

The 2023 Grade 11 NJSLA–S operational test form construction saw 15 items flagged for outfit, 10 items flagged for PCM discrimination, and nine items flagged for lower asymptote. Nevertheless, the flagged items had an outfit, PCM discrimination, or lower asymptote values near the target thresholds. Additionally, the values of outfit on the flagged items indicated that their inclusion would not distort or degrade the measures (Engelhard & Wang, 2021; Linacre, 2002; 2016), and the average value of infit or outfit were close to their ideal values of 1.00. Two grade 11 items were flagged for having item-total correlations below the 0.20 threshold indicating a low discriminating item (see Section 6.1), with values of 0.180 and 0.147. However, as shown in Table 2.4.18, the average item-total correlation was close to the mean target for the form (0.35). The mean and median test times were 107.121 and 86.783 minutes, respectively, which were well below the 150-minute target. Of 70 DIF classifications for each of the four group comparisons (i.e., Male/Female, White/Black, White/Hispanic, and White/Asian), there were zero values categorized as “C” and only one value categorized as “B.” All “B” DIF items were approved for operational test use by the NJBSC as described in Section 2.3.1.1. Tables 2.4.18 and 2.4.19 summarize the test construction and DIF statistics for grade 11.

Table 2.4.18: 2023 NJSLA–S Grade 11 Test Construction Statistics

Statistic	Average	Target	Flags
Rasch B	0.455	0.475	N/A
IT Correlation	0.393	> 0.35	2
Infit	1.057	1.00	1
Outfit	1.122	1.00	15
PCM Discrim.	0.870	1.00	10
Lower Asymptote	0.044	0.00	9
Median Time (min)	86.783	< 150	N/A

Table 2.4.19: 2023 NJSLA–S Grade 11 Test Construction DIF Classifications

Groups	A	B	C
Male/Female	69	1	0
White/Black	70	0	0
White/Hispanic	70	0	0
White/Asian	70	0	0

2.5 2023 NJSLA–S State of the Item Bank

Upon the completion of the 2023 test construction process, MI’s psychometricians analyzed the item bank and facilitated a discussion of the results with content specialists and NJDOE staff. The goal of the discussion was to guide future item development so that it could support valid test score interpretations. The item bank analysis looked at how many items were developed, how many survived the field test and statistical review processes, and how many items were available for creating the 2024 NJSLA–S. Item counts were disaggregated by item type, content domain, scientific practice, DCI, SEP, and CCC. Content areas where the bank had been severely depleted were discussed to determine why they had been problematic and how the next round of item development could improve upon the results.

PART 3: TEST ADMINISTRATION

Standard 6.1 (AERA, NCME, APA, 2014) requires that “[t]est administrators should follow carefully the standardized procedures for administration and scoring specified by the test developer” (p. 114). The test developer is responsible for providing “appropriate training, documentation, and oversight so that the individuals who administer or score the test(s) are proficient in the appropriate test administration or scoring procedures and understand the importance of adhering to the directions provided by the test developer” (p. 114). The following sections detail the myriad processes, procedures, and trainings that were undertaken to properly administer the NJSLA–S.

3.1 District Test Coordinator Training

District Test Coordinators (DTCs) were trained in proper test administration procedures during the annual NJSLA District Test Coordinator Training. In turn, they were “responsible for ensuring that all district and school personnel involved in the administration of New Jersey state assessment programs have been trained” (see Figure 3.1.1; NJDOE, 2024). Information about the NJSLA–S administration is in the Test Coordinator Manual (TCM). That information is not fully replicated here, but the following elements are specific topics that the DTCs were trained on and are also of importance to this technical report:

- Scheduling and testing site requirements
- NJSLA–S participation requirements
- Accessibility features and accommodations available for use on the NJSLA–S
- Materials and tools that would be shipped to schools prior to administration
- Student registration and placement procedures
- Protocols for securely handling materials
- Post-testing responsibilities
- Links and contact information related to the NJSLA–S

The NJSLA TCM can be read in full at the [NJSLA–S website](#) under **Documents and Downloads**.

Your Contribution and Impact



- › Turn-key training is a vital component to ensuring that students are supported through the assessment process and data is secure and accurate.
- › District Test Coordinators (DTCs) are responsible for ensuring that **all district and school personnel** involved in the administration of New Jersey state assessment programs have been trained.
- › State Assessment Coordinators are available to support districts in ensuring the statewide assessment program is implemented with fidelity.
- › **Thank you** for your tireless efforts and leadership in supporting New Jersey's students!



Spring 2023 District Test and Technology Coordinator Training

Figure 3.1.1. Slide 2 from the 2023 DTC

Table 3.1.1 shows the NJSLA–S 2023 testing window dates as well as testing time. Testing times do not include the extra time needed for administrative tasks such as logging students into their testing sessions or reading them directions.

Table 3.1.1: NJSLA–S 2023 Grades 5, 8, and 11 Science Testing Window

Grade	CBT	PBT	Testing Time
5	5/1/23–5/26/23	5/1/23–5/26/23	45 minutes per unit
8	5/1/23–5/26/23	5/1/23–5/26/23	45 minutes per unit
11	5/1/23–5/26/23	5/1/23–5/26/23	60 minutes per unit

3.2 Test Security and Administration Procedures

This section provides information regarding the NJSLA–S test administration procedures. Descriptions of both the computer-based test (CBT) and paper-based test (PBT) procedures are detailed below. For a complete description of all test administration activities, refer to the NJSLA TCM.

3.2.1 Computer-Based Testing

The NJSLA–S CBT forms are delivered via Pearson’s test delivery system, TestNav. TestNav is a secure browser that restricts students’ actions so that they are unable to access or interact with other applications that are outside of the online test materials. Likewise, the student login process is secure; for every test session, Test Administrators (TAs) provide students with testing tickets that include their unique login and password information. If a student needs to exit the test prior to its completion, the TAs can, to ensure test security, lock a test section for the student to access when they return.

Each School Test Coordinator (STC) is provided with a checklist of tasks that they are required to complete during CBT (see Table 3.2.1). The STCs and TA use PearsonAccess^{next} to manage each test session; they can monitor the progress of each of their students and lock and unlock units. PearsonAccess^{next} is a next-generation web-based platform that allows end-to-end monitoring of test administrations for the TAs. Students are only assigned one unit at a time in a prescribed order. STCs and TAs are also charged with assisting with technical issues if they arise. The TCM provides them with a list of typical CBT issues and gives procedures for addressing them. The District Test Coordinator (DTC) and STC are strongly advised to monitor testing and ensure security procedures. Furthermore, they must ensure that TAs provide students with the correct accommodations and accessibility features. After the completion of each unit, STCs collect test materials from the TAs, which include scratch paper, accommodated test materials, and paper copies of the periodic table. Finally, at the end of each day, all NJSLA–S materials must be returned to a secure storage area. Table 3.2.1 shows the checklist of CBT-related tasks that the STCs are charged with completing. For a complete discussion of these procedures, please refer to the TCM.

Table 3.2.1: CBT School Test Coordinator Checklist

Tasks	TCM Section(s)
Ensure that TAs have a computer or tablet available.	Section 3.5
Distribute test materials to TAs.	Section 3.9
Manage test sessions in PearsonAccess ^{next} .	Section 4.1.2
Monitor each testing room to ensure that test administration and security protocols are followed and that required administration information is being documented and collected. Be available during testing to answer questions from TAs.	Section 4.1.4
Investigate all testing irregularities and security breaches, and follow New Jersey policy for reporting these incidents.	Section 2.2
Ensure that TAs provide applicable students with their approved testing accommodations and pre-identified accessibility features.	Section 4.1.4
Schedule and supervise make-up testing.	Sections 2.4.2 and 4.1.5
Create make-up test sessions in PearsonAccess ^{next} .	Section 4.1.5
Respond to all technology-related issues.	Section 4.1.3
Collect materials from TAs.	Section 4.1.5
Ensure that all units are locked after testing on each testing day.	Section 4.1.2

3.2.2 Paper-Based Testing

The following section describes the responsibilities of the DTC and STC during PBT administration. Like the CBT administration, the DTC and STC are required to complete a checklist of tasks (see Table 3.2.2). The tasks are similar to the CBT checklist, except that they are specific to the PBT administration. For instance, the PBT checklist requires STCs to follow protocols for damaged test materials such as test booklets or answer documents. For a complete discussion of these procedures, please refer to the TCM.

Table 3.2.2: PBT School Test Coordinator Checklist

Tasks	TCM Section(s)
Distribute test materials to TAs.	Section 3.10
Monitor each testing room to ensure that test administration and security protocols are followed and that required administration information is being documented and collected. Be available during testing to answer questions from TAs.	Section 4.2.2
Investigate all testing irregularities and security breaches, and follow New Jersey policy for reporting these incidents.	Section 2.2
Ensure that TAs provide applicable students with their approved testing accommodations and pre-identified accessibility features.	Section 4.2.2
Schedule and supervise make-up testing.	Sections 2.4.2 and 4.2.4
Follow the protocol for contaminated or damaged test materials, and refer to New Jersey policy for reporting these incidents.	Section 4.2.3
Collect materials from TAs, and ensure that all test booklets and answer documents have a student name or student ID label.	Section 4.2.4

3.3 Test Irregularities and Breaches

If test security is compromised, the validity of the inferences made from test scores can be affected. Thus, any action that compromises test security is prohibited. These actions are classified as testing irregularities or security breaches. A more complete discussion of test irregularities and breaches can be found in the NJSLA TCM.

Examples of test irregularities and breaches include, but are not limited to:

- **Test Administration Irregularities**
 - Student reviewing or working on the wrong unit of the test; if the student **completes** the wrong unit of a test, the DTC must **immediately** contact the appropriate State Assessment Program Coordinator for directions.
- **Electronic Devices Irregularities**
 - Using a cell phone or other prohibited electronic device (e.g., smartphone, iPod®, smartwatch, personal scanner, eReader) while secure test materials are still distributed, while students are testing, after a student turns in his or her test materials, or during a break.
 - Exception: Test Coordinators, Technology Coordinators, Test Administrators, and proctors are permitted to use cell phones in the testing environment **only** in cases of emergencies or when timely administration assistance is needed. Districts may set additional restrictions on allowable devices as needed.

- Exception: Certain electronic devices may be allowed for medical or audiological purposes during testing. For specific information, refer to the *NJSLA & NJGPA Accessibility Features and Accommodations Manual* at the [New Jersey Assessments Resource Center](#) under **Educator Resources > Test Administration Resources > Accessibility Features and Accommodations Resources > Manuals > NJSLA & NJGPA Accessibility Features and Accommodations.**
- **Test Supervision Irregularities**
 - Coaching students during testing, including giving students verbal or nonverbal cues, hints, suggestions, or paraphrasing or defining any part of the test
 - Engaging in activities (e.g., grading papers, reading a book, newspaper, or magazine) that prevent proper student supervision at all times while secure test materials are still distributed or while students are testing
 - Leaving students unattended without a Test Administrator for any period of time while secure test materials are still distributed or while students are testing. (Proctors must be supervised by a Test Administrator at all times.)
 - Deviating from testing time procedures
 - Allowing cheating of any kind
 - Providing unauthorized persons with access to secure materials
 - Unlocking a test in PearsonAccess^{next} during non-testing times without NJDOE approval
 - Failing to provide a student with a documented accommodation or providing a student with an accommodation that is not documented and therefore is not appropriate
 - Allowing students to test before or after the test administration window without NJDOE approval
- **Test Materials Irregularities and Breaches**
 - Losing a student testing ticket
 - Losing a student test booklet or answer document
 - Losing tactile graphics booklets
 - Leaving test materials unattended or failing to keep test materials secure at all times
 - Reading or viewing tests before, during, or after testing
 - Exception: Administration of a Human Reader/Signer accessibility feature or accommodation which requires a Test Administrator to access the tests
 - Copying or reproducing (e.g., taking a picture of) any part of the test or any secure test materials or online test forms
 - Revealing or discussing test items with anyone, including students and school staff, through verbal exchange, email, social media, or any other form of communication
 - Removing secure test materials from the school building or removing them from locked storage for any purpose other than administering the test

- **Testing Environment Irregularities**

- Failing to follow administration directions exactly as specified in the *Test Administrator Manual* (TAM) (An electronic version of the manual can be viewed at the [NJ Assessments Resource Center](#), located under **Educator Resources > Test Administration Resources > Test Administrator Manuals** as well as on the [NJSLA-S website](#).)
- Displaying any resource (e.g., poster, model, display, teaching aid) that defines, explains, or illustrates terminology or concepts, or otherwise provides unauthorized assistance during testing
- Allowing preventable disruptions such as talking, making noises, or excessive student movement around the classroom
- Allowing unauthorized visitors in the testing environment
 - Unauthorized Visitors: Visitors, including parents/guardians, school board members, reporters, and school staff not authorized to serve as Test Administrators or proctors, are prohibited from entering the testing environment.
 - Authorized Visitors: Observation visits by the principal, monitors from the NJDOE Office of Assessment, monitors from the district, and NJDOE-authorized observers are allowed as long as these individuals do not disturb the testing process.

Protocols are established to report and document any testing irregularity or security breach. All Test Administrators are trained to ensure the proper protocols are implemented. First, both the School and District Test Coordinators must be immediately notified. The DTC is then charged with immediately contacting their NJSLS–S State Contact. The DTC may require the STC to complete the New Jersey Testing Irregularity or Security Breach Form available at the [New Jersey Assessments Resource Center](#) under **Educator Resources > Test Administration Resources > Forms > NJSLS/NJGPA Testing Irregularity and Security Breach Form** to properly document the event. Finally, more information or investigation may be requested by either the DTC or the NJSLS–S State Contact.

3.4 Test Accessibility Features and Accommodations

Standard 3.9 states that “[t]est developers and/or test users are responsible for developing and providing test accommodations, when appropriate and feasible, to remove construct-irrelevant barriers that otherwise would interfere with examinees’ ability to demonstrate their standing on the target constructs” (p. 67). Federal and state regulations require that all students—including those classified as English learners (EL) and those with disabilities—be included in the statewide assessment program and assessed annually. The Every Student Succeeds Act of 2015 (ESSA) mandates that all states must test science one time each in three different grade bands: 3–5, 6–8, and 9–12. The NJSLA Test Coordinator Manual states:

Students who are full-time home-schooled or full-time at a private or parochial school are not eligible to take any statewide assessment. Students with disabilities who attend an approved private school for the disabled and whose tuition is not the financial responsibility of the district are also not eligible to take any statewide assessment. (p. 13)

To ensure that the diverse population of students taking the NJSLA–S is tested under appropriate conditions and to adhere to the principles of universal design (Thompson et al., 2002), NJDOE has adopted test accommodations and accessibility features that may be used when testing special populations of students. The content of the test remains the same, but administration procedures, setting, and answer modes may be adapted. Students requiring accommodations may be tested in a separate location from general education students.

The *NJSLA and NJGPA Accessibility Features and Accommodations Manual (AF&A Manual)* is available online at the [New Jersey Assessments Resource Center](#) under **Educator Resources > Test Administration Resources > Accessibility Features and Accommodations (AF&A) Resources > Manuals > NJSLA & NJGPA Accessibility Features And Accommodations, 11th Edition**. It contains detailed information about each accessibility feature and accommodation. Schools must refer to the *AF&A Manual* for full information about identifying and administering accessibility features and accommodations.

3.4.1 Accessibility Features

The purpose of accessibility features is to ensure that a diverse population of students is being tested fairly and that construct-irrelevant factors are not unduly impacting their test scores. According to the NJSLA and NJGPA AF&A Manual (2022) accessibility features are defined as “tools or preferences that are either built into the testing platform or provided externally by Test Administrators” (p. 54). All students have access to accessibility features. However, for some accessibility features to be available for students during testing, an administrator must have identified the student as needing the accessibility feature prior to testing. It is essential that students using accessibility features get to practice with them prior to operational testing. Thus, NJSLA–S practice tests that contain the accessibility features are available throughout the year at the [NJSLA–S website](#).

3.4.1.1 Text-to-Speech. The most used NJSLA–S accessibility feature is Text-to-Speech (TTS). Prior to testing, an administrator activates the TTS accessibility feature for individual students. When the selected student gets placed into a testing session, their form automatically defaults to the designated TTS form. During testing the student can select the TTS player, and the test will be read aloud to them via the TTS software embedded within TestNav. Students using the TTS accessibility feature must be wearing headphones. The items on the TTS form all contain the same phenomenon-based scenarios, item stems, and response options as are presented to the students taking the traditional CBT form. All final TTS forms are verified by NJDOE to ensure that the TTS functionality is working correctly.

3.4.2 Accommodations

The role of accommodations is to minimize the impact of a student’s disabilities or English language proficiency level on his or her assessment performance. The NJSLA and NJGPA AF&A Manual (2022) defines an accommodation as “an assessment practice or procedure that changes the presentation, response, setting, and/or time and scheduling of assessments” (p. 64). Accommodations are only available to students who have an Individualized Education Program (IEP), a Section 504 plan, or an English learner (EL) plan.

Different accommodations are necessary depending on whether the test was administered using a CBT or PBT format. Per NJDOE policy, all students who received PBT versions of the

NJSLA–S had appropriate accommodations. A comprehensive explanation of each NJSLA–S accommodation is presented in the NJSLA and NJGPA AF&A Manual. The NJSLA–S CBT accommodations include:

- Assistive Technology–Screen Reader
- Assistive Technology–Non-Screen Reader
- American Sign Language (ASL) Text-to-Speech (TTS)
- Human Reader
- Spanish
- Spanish Text-to-Speech
- Spanish Human Reader

PBT accommodations are received as kits, and they include:

- Braille
- Large Print
- Read-Aloud
- Spanish
- Spanish Large Print
- Spanish Read Aloud
- Tactile Graphics

3.4.2.1 Accommodated test form development. The Standards (AERA, APA, NCME, 2014) state that “an appropriate accommodation is one that responds to specific individual characteristics but does so in a way that does not change the construct the test is measuring or the meaning of the scores” (p. 67). Each of the accommodated test forms requires specific processes to ensure they are addressing the needs of their intended users. After NJDOE approval, the accommodated test forms were sent to various subcontractors so that they could adapt the items to Spanish, braille, and American Sign Language (ASL). The adaptation processes for those forms are presented in Sections 3.4.2.1.1 through 3.4.2.1.3. The Paper-Based Test (PBT) form adaptation process is presented in Section 3.4.2.1.4. Following adaptation, NJDOE verifies each accommodated test form.

3.4.2.1.1 Spanish. All Spanish accommodations were made by Teneo Linguistics Company (TLC). TLC received the NJDOE-approved tests and created the translations within ABBI. Once the items were translated, a committee of New Jersey Spanish teachers reviewed the items online, with TLC representatives in attendance. Edits were made during the review, and then the final versions of the online forms were verified by NJDOE. The translation that was created for the online version was then used to create the paper version of the Spanish tests.

3.4.2.1.2 Braille. All braille accommodations were created by the National Braille Press (NBP). NBP received the downloaded paper versions of the operational test forms. NBP provided MI with feedback about any items that were unable to be brailled. Once the tests were brailled, external reviewers received the draft braille versions and reviewed for any issues a student might have taking the braille tests. For the 2023 NJSLA–S, all items were able to be brailled.

3.4.2.1.3 American Sign Language. All ASL accommodations were created by the ADS Group in Plymouth, MN. They provided ASL video production with two ASL content specialist translators and one ASL proofer. Their video production engineer provided studio editing. Additionally, they provided proofing/QC services as well as closed captioning. Once NJDOE approved the operational test forms, the ADS group created the videos in American Sign Language for each item. These items were verified by external expert reviewers under the guidance of MI.

3.4.2.1.4 Paper-Based Test.

The conversion of the NJSLSA–S CBT into PBT form was undertaken by MI’s Editorial Department. Most PBT items were the same as their CBT counterparts. However, some aspects needed adaptation. The following bullets represent the major changes that took place with the stimuli and items during the adaptation processes:

- All artwork was converted from color to grayscale.
- Video items were converted to still images. This was accomplished by MI’s Editorial staff working in conjunction with content specialists to select specific frames from the video that effectively conveyed its essence. In some cases, the captured images were redrawn to ensure that no essential information was lost in the adaptation process.
- TE items were converted to PBT format via multiple methods depending on the TE item type.

3.4.2.2 Accommodated test form equivalence. Occasionally during the accommodated test form conversion process, an item is deemed unable to be accommodated. This can occur for a multitude of reasons—some items do not translate well from English to Spanish, while others are challenging to braille, for example. The procedures for calculating the separate scale score tables, if needed, are detailed in Part 7: Equating and Scaling. In 2023, all items were deemed adequately accommodated by external reviewers, content specialists, and NJDOE.

PART 4: SCORING

It is the responsibility of the test developer to establish scoring procedures (AERA, APA, NCME, 2014). Standard 6.8 states that “[a] scoring protocol should be established, which may be as simple as an answer key for multiple-choice questions” (p. 118). For constructed-response items, the procedures outlined by the Standards require that test developers provide “scoring training materials, scoring rubrics, and examples of test takers’ responses at each score level” (p. 118). The procedures for both the machine-scoring and handscoring of NJSLA–S student responses are described in the following sections.

4.1 Machine-Scored Items

All multiple-choice (MC) and technology-enhanced (TE) items are machine-scored. Each item has a key (correct answer) associated with it, which has been supplied and verified by content specialists and approved by NJDOE prior to test administration. All student responses are machine-scored based on these prior approved keys. Prior to the administration, Pearson’s Customer Data Quality (CDQ) team creates multiple sets of mock test responses for each test form. These responses are scored and processed just as the real tests will be during the administration. The CDQ team verifies that the student responses were accurately captured from the test and that they were scored accurately. Verification steps include comparing responses to the possible ranges of responses to the item, comparing raw overall scores and subscores for entire tests to the maximum values, validating ID unique item numbers (UINs) against the test map, and flagging inconsistent student records for investigation. After the administration, the same checks are made on the data files containing real student tests before they are transferred to MI for psychometric analysis and the adjudication process.

4.1.1 Adjudication

Adjudication involves the careful review of all student responses to an item, ensuring that its key was applied correctly and that no possible correct answer has been overlooked in the many prior key checks. All machine-scored items are adjudicated by MI’s psychometric department. During adjudication, the psychometric team analyzes the student response patterns for each item. The response patterns are simple for items with limited possible options; for instance, an MC item only has 5 possible student responses (A, B, C, D, or blank). However, some TE items can have hundreds of different student responses. The student response data are used to produce one file for each operational item. This file contains each unique response option, the point-value associated with it (i.e., 0, 1, or 2), the total number and the percentage of students selecting each response, and the item-total correlation associated with each response option that was selected more than 100 times. The item means and item-total correlations are also calculated at the item level, and items are flagged for aberrant behavior across all these metrics. Details of the flagging criteria are presented in Part 6 of this document. Upon completion, the files are securely transferred to each grade level’s lead content specialists for review.

The role of the content specialists during the adjudication process is to use the information housed in the adjudication files to identify any possible miskeys. They are instructed to first check items that were flagged for having low item means and item-total correlations because those statistics could indicate that the item is not performing as intended. Next, they look at combinations of student responses that are keyed as receiving “0” points but have item-total

correlations above 0. That combination of response-level data could also be an indication of a possible student response that deserves credit for a correct response, but that has been keyed as incorrect. Finally, through a sorting process, the content specialists can relatively quickly review all other combinations of student responses. If there are any miskeys, key changes are submitted to NJDOE, and upon approval, Pearson incorporates them into the scoring algorithm. These steps are essential to ensure both the reliability of student test scores and their valid interpretations.

4.2 Handscored Items

All NJSLA–S CR items are scored by human scorers according to the procedures outlined in the sections that follow.

4.2.1 Selecting Handscoring Staff

MI's recruiting team first recruits qualified scorers who have experience scoring New Jersey Science assessments. To supplement this core pool, MI's recruiting team contacts other scorers in MI's database who have experience successfully scoring other large-scale assessments. Returning staff are selected based on experience and performance, as well as attendance, punctuality, and cooperation with work procedures and MI policies. MI maintains evaluations and performance data for all staff who work on each scoring project in order to determine employment eligibility for future projects. For new scorers, the recruiting team reviews applications—including prospective scorers' resumes, references, proof of degree, and recognition of scorer requirements—before offering employment. All our scorers have a minimum of a four-year college degree, and many are current or former educators.

In selecting Team Leaders, MI management staff and scoring directors review the files of all returning staff. They look for people who are experienced Team Leaders with a record of good performance on previous projects and consider scorers who have been recommended for promotion to the Team Leader position.

MI requires that all handscoring staff (Scoring Directors, Team Leaders, scorers, and clerical staff) sign a confidentiality/nondisclosure agreement before receiving training or accessing secure project materials. The employment agreement indicates that participants may not reveal information about the test, the scoring criteria, or the scoring methods to any person.

4.2.2 Operational Range Finding

Range-finding meetings are conducted to establish “true” scores from a representative sample of papers (i.e., responses). One hundred sample papers per task are chosen from the available field-test papers. At the beginning of the range-finding meeting, the scoring rubrics of the items are discussed and refined by the committee. The sample responses brought to the range-finding meetings are selected from a broad range of New Jersey LEAs in order to ensure that the sample is representative of overall student performance. To maximize the probability that papers eligible for the highest score points are included in the sample, special efforts are made by MI management and scoring staff to include high-performing responses. The range-finding committees consist of NJDOE content specialists, New Jersey teacher representatives, and MI management personnel, as well as the Scoring Director responsible for each content area.

4.2.3 Field Test Range Finding

Prior to field-test scoring, content committees consisting of NJDOE personnel, New Jersey teacher representatives, and MI leadership personnel meet virtually to determine “true” scores for 30 selected papers representing each of the score points for each item to be tested. Field-test scoring guides and training sets are developed using the papers scored at the range finding. Time is spent determining whether any changes need to be made to the scoring rubrics associated with the items being reviewed before any field-test scoring takes place.

4.2.4 Developing Scoring Guides

After the range finding meetings, training materials are developed consisting of an anchor set (examples of responses for each score point) and training/qualifying sets (practice papers) for each task using the responses scored at range finding. Anchor sets usually consist of two or more annotated examples of each score point, arranged in score point order. To maximize consistency, the same anchor sets are used each year for items administered in multiple administrations. Anchor sets include annotations that explain how the scoring criteria are applied to each response’s specific features and why the response merits a particular score. These annotations connect to highlighted sections of the student response in training lessons, drawing scorers’ attention to the critical training pieces to elucidate the precise scoring rationale and to help scorers define the lines between score points. Training/qualifying sets consist of clearly anchored papers in random score point order. These sets are constructed using responses from the Operational Range Finding, with the scores assigned by the range-finding committee for each response.

4.2.5 Team Leader Training and Duties

After the anchor, training, and qualifying papers have been identified and finalized, the Scoring Director conducts Team Leader training for each task. This process typically takes up to four days depending on the content. Procedures are similar to those for training scorers (described in more detail below) but are more comprehensive, dealing with identification of non-scorable responses, unusual approaches to a prompt, alert situation responses (e.g., child-in-danger), and other duties performed only by leadership. Team Leaders assist in training scorers by serving as a resource when scorers are training.

During scoring, Team Leaders respond to questions, read behind scorers’ scored responses, and counsel scorers having difficulty with the criteria. Team Leaders also monitor the scoring patterns of each scorer throughout the project, conduct retraining as necessary through responses to scorer questions and reading behind scorers, perform second readings, and maintain a professional working environment.

4.2.6 Scorer Training and Qualifying

All scorers are trained using the rubrics, anchor papers, training papers, and qualifying papers selected during the range-finding meetings and approved by the NJDOE. MI’s Virtual Scoring Center™ (VSC™) includes an online training interface that presents rubrics, anchor sets, and training/qualifying sets. VSC™ is used for all training and qualifying, whether site-based or remote. VSC™ provides for effortless and timely communication with scoring leadership throughout training and allows scorers to efficiently navigate the training materials.

Recruited staff must maintain rigorous adherence to established training methodologies to ensure the quality and credibility of our scoring. MI enforces strict attendance during training. Scorers are trained as a group to maintain consistency and are trained on all relevant training materials. Scorers have access to all training materials during live scoring. The same training protocol is followed for both site-based and remote scorers.

After scorers have signed contracts and nondisclosure forms and have been provided with an introduction to the project, training begins. Scorer training and Team Leader training follow the same format. Scorers and Team Leaders are introduced to the constructed-response task and the anchor set. This process includes modeling how to identify the essential information in anchor responses to establish a consistent scoring vocabulary. Any nuances in interpreting and applying the scoring rubric are also highlighted at this stage.

Scoring personnel log in to Measurement Incorporated Remote Access (MIRA) to review the rubric and anchor responses. MIRA includes all online training modules, is the portal to the VSC™ interface, and is the data repository of all scoring reports that are used for scorer monitoring. Here, Team Leaders and scorers assign scores to a practice/qualifying set of responses. They are reminded to compare each practice response to comparable anchor responses to ensure accuracy and consistency in scoring the practice responses. MI trains scoring personnel to reference those student responses as representative of the rubric. The rubric is a tool, but the anchor responses represent how the rubric is applied. After Team Leaders and scorers score practice responses, they are provided with the correct scores. The same process is followed for all subsequent practice/qualifying sets.

Scorers must demonstrate their ability to score accurately by attaining 70% perfect agreement and 100% adjacent agreement (within one point) percentage on two of the qualifying sets before they read packets of operational student responses. Any scorer unable to meet the standards set by the NJDOE is dismissed.

Training is carefully orchestrated so that scorers understand how to apply the rubric in scoring the papers, learn how to reference the scoring guide, develop the flexibility needed to deal with a variety of responses, and retain the consistency needed to score all papers accurately. In addition to completing all of the initial training and qualifying, scorers are trained in the use of the VSC™ handscoring system, “flagging” of unusual responses for Team Leader review, and other procedures necessary for the conduct of a smooth project.

Levels of staffing for scoring the 2023 NJSLA–S are presented in Table 4.2.1. Specifically, Table 4.2.1 shows the number of scorers, Team Leaders, and Scoring Directors at each grade level who participated in scoring.

Table 4.2.1: Scoring Personnel by Grade

Grade	Scorers	Team Leaders	Scoring Director
5	89	3	1
8	147	3	1
11	70	3	1

4.2.7 Monitoring Scorer Performance

In addition to thorough and consistent training, reliable scoring depends upon careful evaluation of scorer performance to support a continuous loop of feedback among the scorers, Team Leaders, Scoring Directors, and Scoring Monitors. Scoring Directors offer direct leadership and guidance to Team Leaders as they monitor individual scorer performance. Scoring Directors also furnish scorers with general guidance and clarify appropriate application of the training materials, while Team Leaders provide direct supervision, which allows for a higher degree of scrutiny of scorer performance, individual attention, and opportunities for immediate intervention or correction if required.

Real-time reports that provide both daily and cumulative (project-to-date) data are used to monitor and evaluate scoring performance. Scoring Monitors and Scoring Directors review these reports daily. As they review these data, they can identify any issues evident in scores being generated and address them with Team Leaders and individual scorers when necessary. These reports are described in more detail below.

The quality of MI's handscoring program is maintained through ongoing monitoring by experienced scoring leadership. Scoring Directors and Team Leaders are skilled in detecting scoring trends and remediating any issues that arise. Scorers who are unable to meet accuracy and productivity standards after feedback and retraining will not be allowed to continue scoring. When this occurs, MI can reset any scores assigned by a dismissed scorer and have the responses immediately rescored.

MI's handscoring process incorporates ongoing checks for and controls against scorer error. Specifically, MI implements the following quality-assurance procedures:

- **Validity checks.** MI's VSC™ scoring system randomly seeds validity responses among operational responses during scoring. A small set of validity responses are selected and approved by Scoring Monitors and Scoring Directors. The "true" scores for these responses are entered into a validity database. Validity responses are indistinguishable from operational responses. Scorer accuracy and drift are evaluated using validity results. The validity responses are dispersed evenly across all of an item's score point levels, and they are selected based on how well they represent typical examples of each score point. Readers are encouraged to send responses that are difficult to score to their team leader; thus, those types of papers are not selected as validity responses.

- **Blind double reads.** For each item, a minimum of 10% of responses are randomly selected to receive blind double reads. Scorer agreement is used to evaluate the reliability of scoring across all scorers.
- **Daily systematic review of handscoring reports.** Scoring Directors monitor and evaluate scorers' performance daily using an array of handscoring reports, described below. MI provides any retraining necessary to ensure scorer accuracy. Retraining strategies are implemented under the direction of the Scoring Monitors in conjunction with Scoring Directors and Team Leaders.
- **Targeted read-behinds.** Team Leaders conduct targeted read-behinds for scorers who have been identified, based on Validity performance, or based on other performance data, as targets for close monitoring. When conducting targeted read-behinds, Team Leaders pay careful attention to the particular score points with which individual scorers have difficulty. This information is obtained by reviewing the results of validity and score point distribution reports. Team Leaders provide feedback by discussing incorrectly scored responses with the individual scorer and continue to monitor to ensure the scorer has understood and applied the feedback appropriately.
- **Score verifications.** MI implements a series of automated score verifications to ensure the accuracy of scores. For example, a blank check is then conducted, which resets scores when a condition code of "blank" is assigned to a response that has one or more characters in the response string (e.g., a response comprising spaces or tabs). In this case, only after three independent scorers have assigned a condition code of "blank" to a response that appears blank but includes characters in the response string is the score recorded. A similar check is run when a score or condition code other than "blank" is assigned to a response that includes no characters in the response string. Automatic resetting of double-scored responses occurs when two scorers assign non-adjacent scores, mismatched condition codes, or a combination of a condition code and a numeric score, thus providing an additional score verification. In addition to automatically resetting and rescoring these responses, the scorer information is captured in a report and reviewed by Scoring Directors, as one of many tools used to determine retraining needs.

VSC™ provides an appropriate infrastructure for facilitating our extensive quality-assurance procedures. Through VSC™, handscoring leadership can review scorer performance, conduct read-behinds, provide feedback and respond to questions, deliver retraining and/or recalibration responses on demand and at regularly scheduled intervals, and prevent scorers from scoring additional live responses if they require additional monitoring.

Scorers are dismissed when, in the opinion of the appropriate Scoring Monitor and/or Scoring Director, they have been counseled, retrained, and given a reasonable opportunity to improve and continue to perform below an acceptable standard for accuracy or production. In the case of the former, all scores assigned by a scorer during a given timeframe can be identified and reset, and the responses can be released back into the scoring pool for immediate rescoring.

4.2.8 Automatic Rescores

As shown in Section 8.5, the raters are not in perfect agreement 100% of the time. Thus, to ensure that no student is unjustly penalized because a rater may have been a little too stringent, rescoring is conducted automatically for any student who scores one raw score point below the proficient cut score. MI reviews student responses to constructed response items and verifies the original scores or makes changes where warranted. A score is never lowered during the automatic rescoring process even if it was deemed to be too high. LEAs do not need to request rescoring. Table 4.2.2 provides automatic rescoring results for all three grade levels. All open-ended/constructed-response item types were scored by a single rater.

Table 4.2.2: Automatic Rescore Results

Grade	Eligible for Automatic Rescore	Number of Changes	Percentage Changed (of those Eligible)
5	1,813	259	14.3
8	1,573	358	22.8
11	1,616	279	17.3

4.3 Quality Control

To confirm that the processing of student tests and test scores was done correctly, Measurement Incorporated conducted quality control checks in addition to supporting NJDOE doing its own quality control checks. To produce the score reports, Pearson started with a large data file called the Student Record File (SRF) that includes all the information that will be shown on the reports. MI began by verifying the information in this file.

MI checked the student demographic information against what the districts had entered into the test registration system. MI also verified the results of the machine-scored items. This was straightforward for multiple choice items but more complicated for technology-enhanced items that could have multiple correct answers. For open-ended items, MI compared the scores to its handscoring system.

After Pearson produced the final score reports, MI used the previously verified SRF to verify the data shown on a large sample of these reports.

4.3.1 QC Sample

For NJDOE's part of QC, MI selected a sample of several hundred student tests for them to manually review. NJDOE staff compared the scores for these students to the final SRF and score reports produced by Pearson. The sample included all test forms, as well as students with a wide variety of values for demographic variables such as gender, ethnicity/race, English learner status, and disability status. These selected students were provided specifically for validation of ISRs and Rosters. After individual test scores were verified, NJDOE used them to calculate aggregate figures, such as average scale scores, that are shown in data files and score reports. The following section details the processes NJDOE used to ensure that student test scores were accurate.

4.3.2 Key Information Sheets

To help organize NJDOE's QC process for the student-level data, MI produced a Key Information Sheet (KIS) for each student's test. The KIS is a spreadsheet that is used to keep track of all the information from a test, as a helpful aid for the QC process. The KIS is pre-populated with student information from a test (such as name and accommodations), the key for each machine-scored item, and a spot to record the points earned for each item.

First, MI verified the student information on the KIS against the student information in the test registration system. Next, MI scored the student's responses to each selected-response item against the key and recorded the score on the KIS. For open-ended items, MI exported scores from the handscoring system to record on the spreadsheet. Formulas in the KIS automatically tallied the student's overall points, as well as the points in each domain and practice. Any discrepancies between these totals and the preliminary data file from Pearson required scrutiny of the points earned for each item. The KIS helped to narrow down the problem to a particular domain, practice, and unit. MI provided NJDOE with the verified KIS on August 7, 2023.

After external review and NJDOE's approval of the scale score tables, Pearson imported the scale score tables and produced a final set of data files and score reports. NJDOE staff used the KISs prepared by MI to determine each student's scale score, overall performance level, and performance level for each subscore. Then NJDOE compared this information to student-level score reports. This stage provided NJDOE with confidence that each piece of student-level information on these reports was accurately derived from the original sources of test data.

4.3.3 Aggregate Data

Certain numbers shown on the Individual Student Report and School Student Roster are aggregated figures, such as averages of scores at the school, district, or state level; and the percentages of students achieving each overall performance level. In addition, other score reports and the summary data file only show aggregated data. NJDOE verified all of these values by calculating the same figures from the raw data in the SRF. This step was not necessarily limited to the schools in the QC sample.

PART 5: STANDARD SETTING

Cizek and Bunch (2007) define standard setting as “the process of establishing one or more cut scores on examinations” (p. 5). Cut scores divide a distribution of test scores into two or more categories. The purpose of conducting a standard setting is to assist the users of test scores in making valid interpretations. Standard 5.21 states that “[w]hen proposed score interpretations involve one or more cut scores, the rationale and procedures used for establishing cut scores should be documented clearly” (p. 107). The 2019 NJSLA–S Technical Report details the processes, procedures, and analyses used to accomplish the 2019 NJSLA–S Standard Setting. The executive summary from the 2019 NJSLA–S Standard Setting Report is presented in Appendix D of this report.

The 2019 NJSLA–S Standard Setting was externally reviewed by NJTAC member Stephen Koffler (Koffler, 2019). He evaluated the process based on the Standards (2014) and the framework established by Kane (2001). Koffler focused on three major sources of validity evidence: procedural, internal, and external. Overall, he concluded that “the NJSLA–S Standard Setting Study was sound, followed best practice and met the professional standards for performing a Standard Setting Study and recommending valid and defensible cut scores.” (p. iv).

PART 6: ITEM AND TEST STATISTICS

Standard 5.0 states that “[t]est scores should be derived in a way that supports the interpretations of test scores for proposed uses of tests. Test developers and users should document evidence of fairness, reliability, and validity of test scores for their proposed uses” (p. 102). The NJSLA–S was designed to support inferences based on the classification of students into four performance levels, as has been described throughout this technical report. The interpretations of the performance-level classifications are dependent upon the test performing as intended. As was described in Section 2.3, the NJSLA–S was constructed using a combination of classical test theory (CTT) statistics, item response theory (IRT) statistics, and the content constraints. The following sections detail how well the 2023 NJSLA–S performed based on those CTT and IRT statistics, along with other criteria. Detailed test maps containing item metadata, various statistics, and Range PLD alignment are presented in Appendix F of this report. The final section in this part presents disaggregated descriptive statistics of scale scores and subscore proficiency classifications.

The data for these and all subsequent analyses were verified by Pearson’s Customer Data Quality (CDQ) team prior to delivery to MI. Responses from students who did not attempt to take the test or who had their test scores voided were removed from the data prior to analyses. NJDOE requires a student to attempt at least one item in at least two different operational test units to obtain a scale score. Student responses were voided for cheating, security breaches, or other reasons.

6.1 Classical Test Theory Statistics

For each administration, a set of statistics based on CTT was generated and reviewed for item calibrations and scaling. The statistics can be grouped into measures of four psychometric concepts:

- Item Difficulty
- Item Discrimination
- Speededness
- Differential Item Functioning

These statistics were calculated for every operational item; each statistic provides some key information about the quality of each item from an empirical perspective. If any of the four statistics suggested an item was negatively impacting the reliability or validity of test score interpretations, a recommendation was made to NJDOE to remove the item from operational use. Descriptions of each type of item statistic appear in the following sections. Please note that one MC item was dropped from the Grade 11 test due to content concerns. As a result, it was excluded from the 2023 psychometric analyses presented in this technical report.

6.1.1 Item Difficulty and Discrimination Descriptive Statistics

Monitoring item difficulty is essential for ensuring that the test is reliable and will foster valid test score interpretations. If items tend to be too challenging or too easy for a population of test takers, then the reliability and validity of test score interpretations will suffer. In CTT, the item difficulty of a dichotomous item is assessed via the *p-value*, which is defined as the proportion of students who answered an item correctly. *P-values* can range from 0.00 to 1.00;

an item with a high *p-value* is easier to answer correctly, whereas an item with a low *p-value* is more challenging. Dichotomous items with *p-values* either below .25 or above .90 were flagged for review during the adjudication process described in Section 4.1.1. For polytomous items, such as the 0–4-point CR items, item difficulty is expressed as an item mean. The polytomous item flagging criteria involves converting the item mean to a proportion by dividing it by the maximum points possible on the item (i.e., making it an adjusted item mean or a *p-value*). Polytomous items are then flagged if their converted *p-value* falls outside of the .25 to .90 range. It should be noted that the flagging criteria only provide a general guideline, and some productive items have *p-values* outside of the .25 to .90 range.

Item discrimination is also important to monitor. If items are unable to discriminate between students with different ability levels, then both the reliability and the validity of test score interpretations can suffer. In CTT, the item discrimination is expressed as the correlation between item scores and the total score of the remaining items on the test, the latter being a proxy for overall student ability. The item-total correlation (denoted by *rpb* in this technical report) can range from –1.00 to 1.00. Items with discrimination values below 0.2 are flagged for review during the adjudication process. Items with item-total correlations that are below zero (i.e., negative) are considered for removal from the test because they could harm both the reliability and the validity of test score interpretations.

For NJSLA–S items, Tables 6.1.1, 6.1.3 and 6.1.5 summarize the item difficulties in terms of *p-values* for grades 5, 8, and 11, respectively; Tables 6.1.2, 6.1.4 and 6.1.6 show the item discrimination (*rpb*) summaries for grades 5, 8, and 11, respectively. Several intervals of item difficulties or discriminations were created for computing the frequency distributions. In these tables, the descriptive statistics and frequency distributions for each item type are disaggregated by content domain and scientific practice.

Overall, the average item difficulties and discriminations appear to be productive for measuring students in New Jersey. At each grade level, the average TE items tended to be slightly more challenging and more discriminating than MC items. The CR items were, as expected, more discriminating than the MC and TE items.

At grade 5, most of the MC and TE items had item difficulties between .25 and .75, indicating an average item difficulty level on the scale. The grade 5 CR item in physical science tended to be slightly more challenging than the CR items in Earth and Space Science and Life Science. At grade 8, only four MC items and one TE item had *p-values* at or above .50, indicating most of the items were at the harder end of the scale. The grade 8 CR item in Earth and Space Science tended to be more challenging than the CR items in Life Science and Physical Science. At grade 11, zero MC or TE items had *p-values* above .75, indicating that there were no items on the easier end of the scale. Additionally for Grade 11, the CR item in Physical Science tended to be more challenging than the CR items in Earth and Space Science and Life Science.

The average item-total correlations (*rpb*) for CR items were .61, .61, and .60 for grades 5, 8, and 11, respectively, indicating good item discrimination. Also, the frequency distributions of item-total correlations for MC and TE items appear to be productive for discriminating between high- and low-achieving students. At grade 5, only one MC item and two TE items had item-total correlations below .20, while eight MC and twenty-one TE items had item-total correlations

above .40. At grade 8, two MC items and one TE item had item-total correlations below .20, while two MC items and twenty TE items had discriminations above .40. At grade 11, one MC item and zero TE items had item-total correlations below .20, while nine MC and twenty-one TE items had discriminations above .40.

Table 6.1.1: Grade 5 Item Difficulty (*p-value*) Distribution and Summary Statistics

Item Type	Domain/ Practice	N of Items	Distribution of Item Difficulty (<i>p-value</i>)					Descriptive Statistics		
			[0,.25)	[.25,.5)	[.5,.75)	[.75,.9)	[.9,1]	Mean	S.D.	Median
MC	NJSLA-S	14	0	6	8	0	0	.53	.10	.53
	Earth and Space	5	0	1	4	0	0	.61	.10	.65
	Life	2	0	1	1	0	0	.50	.01	.50
	Physical	7	0	4	3	0	0	.49	.09	.44
	Critiquing	6	0	3	3	0	0	.54	.12	.54
	Investigating	7	0	3	4	0	0	.53	.11	.51
	Sensemaking	1	0	0	1	0	0	.55	N/A	.55
TE	NJSLA-S	34	2	28	4	0	0	.39	.10	.39
	Earth and Space	12	0	9	3	0	0	.40	.12	.40
	Life	14	1	12	1	0	0	.40	.10	.40
	Physical	8	1	7	0	0	0	.34	.07	.35
	Critiquing	9	1	8	0	0	0	.44	.07	.32
	Investigating	9	1	8	0	0	0	.32	.08	.34
	Sensemaking	16	0	12	4	0	0	.45	.10	.43
CR	NJSLA-S	3	1	2	0	0	0	.29	.12	.28
	Earth and Space	1	0	1	0	0	0	.28	N/A	.28
	Life	1	0	1	0	0	0	.41	N/A	.41
	Physical	1	1	0	0	0	0	.18	N/A	.17
	Critiquing	3	1	2	0	0	0	.29	.12	.28
	Investigating	0	0	0	0	0	0	N/A	N/A	N/A
	Sensemaking	0	0	0	0	0	0	N/A	N/A	N/A

Table 6.1.2: Grade 5 Item Discrimination Distribution and Summary Statistics

Item Type	Domain/ Practice	N of Items	Distribution of Item Discrimination (<i>rpb</i>)					Descriptive Statistics		
			[0, .2)	[.2, .3)	[.3, .4)	[.4, .5)	[.5, 1]	Mean	S.D.	Median
MC	NJSLA-S	14	1	2	3	6	2	.40	.12	.44
	Earth and Space	5	0	1	2	1	1	.40	.11	.37
	Life	2	0	0	1	0	1	.48	.11	.48
	Physical	7	1	1	0	5	0	.38	.14	.44
	Critiquing	6	0	1	1	2	2	.44	.12	.47
	Investigating	7	1	1	2	3	0	.36	.13	.40
	Sensemaking	1	0	0	0	1	0	.44	N/A	.44
TE	NJSLA-S	34	2	6	5	7	14	.43	.16	.45
	Earth and Space	12	1	2	1	3	5	.41	.15	.45
	Life	14	0	2	3	2	7	.49	.15	.52
	Physical	8	1	2	1	2	2	.36	.16	.41
	Critiquing	9	1	2	2	2	2	.37	.16	.40
	Investigating	9	0	4	0	3	2	.39	.15	.42
	Sensemaking	16	1	0	3	2	10	.49	.15	.52
CR	NJSLA-S	3	0	0	0	0	3	.61	.05	.63
	Earth and Space	1	0	0	0	0	1	.63	N/A	.63
	Life	1	0	0	0	0	1	.65	N/A	.65
	Physical	1	0	0	0	0	1	.56	N/A	.56
	Critiquing	3	0	0	0	0	3	.61	.05	.63
	Investigating	0	0	0	0	0	0	N/A	N/A	N/A
	Sensemaking	0	0	0	0	0	0	N/A	N/A	N/A

Table 6.1.3: Grade 8 Item Difficulty (*p-value*) Distribution and Summary Statistics

Item Type	Domain/ Practice	N of Items	Distribution of Item Difficulty (<i>p-value</i>)					Descriptive Statistics		
			[0,.25)	[.25,.5)	[.5,.75)	[.75,.9)	[.9,1]	Mean	S.D.	Median
MC	NJSLA-S	18	2	12	4	0	0	.39	.11	.38
	Earth and Space	5	0	3	2	0	0	.44	.10	.41
	Life	8	1	7	0	0	0	.34	.06	.34
	Physical	5	1	2	2	0	0	.42	.17	.39
	Critiquing	5	0	4	1	0	0	.42	.09	.40
	Investigating	7	2	4	1	0	0	.33	.11	.32
	Sensemaking	6	0	4	2	0	0	.43	.12	.40
TE	NJSLA-S	44	16	27	1	0	0	.30	.11	.29
	Earth and Space	14	3	11	0	0	0	.33	.10	.35
	Life	15	8	7	0	0	0	.28	.11	.25
	Physical	15	5	9	1	0	0	.30	.12	.27
	Critiquing	16	4	11	1	0	0	.33	.13	.32
	Investigating	12	7	5	0	0	0	.27	.11	.24
	Sensemaking	16	5	11	0	0	0	.30	.10	.29
CR	NJSLA-S	3	1	2	0	0	0	.28	.16	.33
	Earth and Space	1	1	0	0	0	0	.10	N/A	.10
	Life	1	0	1	0	0	0	.41	N/A	.41
	Physical	1	0	1	0	0	0	.33	N/A	.33
	Critiquing	1	0	1	0	0	0	.41	N/A	.41
	Investigating	1	1	0	0	0	0	.10	N/A	.10
	Sensemaking	1	0	1	0	0	0	.33	N/A	.33

Table 6.1.4: Grade 8 Item Discrimination Distribution and Summary Statistics

Item Type	Domain/ Practice	N of Items	Distribution of Item Discrimination (<i>rpb</i>)					Descriptive Statistics		
			[0, .2)	[.2, .3)	[.3, .4)	[.4, .5)	[.5,1]	Mean	S.D.	Median
MC	NJSLA-S	18	2	4	10	2	0	.33	.10	.34
	Earth and Space	5	0	1	3	1	0	.36	.08	.34
	Life	8	1	3	3	1	0	.31	.10	.32
	Physical	5	1	0	4	0	0	.32	.12	.38
	Critiquing	5	0	2	3	0	0	.34	.05	.35
	Investigating	7	2	1	3	1	0	.30	.15	.34
	Sensemaking	6	0	1	4	1	0	.35	.05	.35
TE	NJSLA-S	44	1	11	12	16	4	.38	.11	.39
	Earth and Space	14	0	3	5	3	3	.40	.12	.37
	Life	15	1	5	1	8	0	.36	.12	.40
	Physical	15	0	3	6	5	1	.39	.09	.40
	Critiquing	16	0	3	3	8	2	.40	.10	.42
	Investigating	12	0	3	5	3	1	.36	.11	.35
	Sensemaking	16	1	5	4	5	1	.37	.12	.37
CR	NJSLA-S	3	0	0	0	0	3	.61	.09	.60
	Earth and Space	1	0	0	0	0	1	.53	N/A	.53
	Life	1	0	0	0	0	1	.60	N/A	.60
	Physical	1	0	0	0	0	1	.70	N/A	.71
	Critiquing	1	0	0	0	0	1	.60	N/A	.60
	Investigating	1	0	0	0	0	1	.53	N/A	.53
	Sensemaking	1	0	0	0	0	1	.70	N/A	.71

Table 6.1.5: Grade 11 Item Difficulty (*p-value*) Distribution and Summary Statistics

Item Type	Domain/ Practice	N of Items	Distribution of Item Difficulty (<i>p-value</i>)					Descriptive Statistics		
			[0,.25)	[.25,.5)	[.5,.75)	[.75,.9)	[.9,1]	Mean	S.D.	Median
MC	NJSLA-S	30	1	20	9	0	0	.43	.10	.42
	Earth and Space	8	0	7	1	0	0	.42	.08	.42
	Life	10	0	9	1	0	0	.38	.09	.39
	Physical	12	1	4	7	0	0	.47	.12	.53
	Critiquing	7	0	7	0	0	0	.37	.08	.37
	Investigating	13	1	9	3	0	0	.40	.11	.41
	Sensemaking	10	0	4	6	0	0	.49	.08	.52
TE	NJSLA-S	36	13	18	5	0	0	.34	.15	.30
	Earth and Space	11	3	5	3	0	0	.41	.19	.31
	Life	12	6	4	2	0	0	.30	.16	.26
	Physical	13	4	9	0	0	0	.31	.09	.31
	Critiquing	14	6	7	1	0	0	.32	.11	.28
	Investigating	9	2	6	1	0	0	.34	.16	.30
	Sensemaking	13	5	5	3	0	0	.35	.20	.31
CR	NJSLA-S	3	1	1	1	0	0	.35	.17	.33
	Earth and Space	1	0	0	1	0	0	.53	N/A	.53
	Life	1	0	1	0	0	0	.33	N/A	.33
	Physical	1	1	0	0	0	0	.19	N/A	.19
	Critiquing	1	0	0	1	0	0	.53	N/A	.53
	Investigating	1	1	0	0	0	0	.19	N/A	.19
	Sensemaking	1	0	1	0	0	0	.33	N/A	.33

Table 6.1.6: Grade 11 Item Discrimination Distribution and Summary Statistics

Item Type	Domain/ Practice	N of Items	Distribution of Item Discrimination (<i>rpb</i>)					Descriptive Statistics		
			[0, .2)	[.2, .3)	[.3, .4)	[.4, .5)	[.5,1]	Mean	S.D.	Median
MC	NJSLA–S	30	1	16	4	6	3	.33	.10	.28
	Earth and Space	8	0	5	1	1	1	.32	.10	.28
	Life	10	0	5	3	1	1	.34	.10	.32
	Physical	12	1	6	0	4	1	.33	.12	.28
	Critiquing	7	0	6	0	1	0	.28	.07	.27
	Investigating	13	1	5	4	3	0	.32	.09	.33
	Sensemaking	10	0	5	0	2	3	.37	.13	.36
TE	NJSLA–S	36	0	6	9	15	6	.41	.10	.44
	Earth and Space	11	0	1	2	8	0	.41	.08	.45
	Life	12	0	1	4	3	4	.43	.11	.45
	Physical	13	0	4	3	4	2	.39	.12	.37
	Critiquing	14	0	2	4	4	4	.43	.10	.45
	Investigating	9	0	2	1	6	0	.39	.10	.44
	Sensemaking	13	0	2	4	5	2	.41	.11	.43
CR	NJSLA–S	3	0	0	0	0	3	.60	.01	.61
	Earth and Space	1	0	0	0	0	1	.61	N/A	.61
	Life	1	0	0	0	0	1	.59	N/A	.59
	Physical	1	0	0	0	0	1	.61	N/A	.61
	Critiquing	1	0	0	0	0	1	.61	N/A	.61
	Investigating	1	0	0	0	0	1	.61	N/A	.61
	Sensemaking	1	0	0	0	0	1	.59	N/A	.59

6.1.2 Speededness

The consequence(s) of time limits on examinees’ scores is called speededness (Swineford, 1949). A traditional measure of speededness is the number of items that are not attempted by students. Logically, in each separately timed subsection of a test, it can be assumed that a student may have run out of time if the student did not attempt the last item. The percentage of students omitting an item provides information about speededness, although it must be kept in mind that students can omit an item for reasons other than speededness (for example, choosing to not put effort into answering a constructed-response item). Thus, if the percentage of omits is low, that implies that there is little speededness. Conversely, if the percentage of omits is high, speededness, as well as other factors, may be the cause.

The NJSLA–S was not designed to be a speeded test, but rather a power test. That is, all students are expected to have ample time to finish all items and prompts. NJSLA–S assessments were administered during a testing window with a specified amount of time per unit by grade. Students were assumed to have enough time to complete the test. The numbers of items and item types composing each operational test unit for each grade level, along with the testing time, are detailed in Table 6.1.7. Additionally, Table 6.1.8 presents the percentage of students

omitting the last TE item in each test section. Overall, the small percentages of students shown in the table indicated that each grade level test did not show speededness.

Table 6.1.7: Operational Testing Schedule—Items and Time Allocations

Grade	Unit	Items	Time in Minutes
5	1	6 MC, 10 TE, 1 CR	45
5	2	4 MC, 12 TE, 1 CR	45
5	3	4 MC, 12 TE, 1 CR	45
8	1	8 MC, 12 TE, 1 CR	45
8	2	6 MC, 15 TE, 1 CR	45
8	3	4 MC, 17 TE, 1 CR	45
11	1	12 MC, 10 TE, 1 CR	60
11	2	8 MC, 15 TE, 1 CR	60
11	3	10 MC, 12 TE, 1 CR	60

Table 6.1.8: Percentage of Students Omitting the Last TE Item in Each Operational Unit

Grade	Unit	Location	% Student
5	1	17	4.0
5	2	17	2.2
5	3	16	0.5
8	1	21	4.2
8	2	21	1.8
8	3	21	1.6
11	1	23	3.4
11	2	23	0.6
11	3	23	4.0

6.1.3 Operational DIF Analysis

The *Standards* define Differential Item Functioning (DIF) as “when different groups of test takers with similar overall ability, or similar status on an appropriate criterion, have, on average, systematically different responses to a particular item” (p. 16). If items perform differently for sub-groups of students after controlling ability, the test might disadvantage some groups of students over others.

By convention, the two groups of test takers involved in DIF analyses are referred to as the focal and reference groups. Different methods are used for DIF detection depending on whether the item is dichotomous or polytomous. For dichotomous items, DIF was identified using the Mantel-Haenszel (MH) procedure (Mantel & Haenszel, 1959). It is considered effective and efficient (Clauser & Mazor, 1998; Hills, 1989). For the NJSLA–S, under the MH procedure, a statistical significance test (MH Chi-square test) of DIF and an evaluation of effect sizes in DIF measures (MH D-DIF statistic) were performed in conjunction with the ETS DIF classification

system (Dorans & Holland, 1993). The letters A, B, and C are used to denote DIF categories in the ETS DIF classification system with A-level indicating a negligible degree of DIF, B-level indicating slight to moderate DIF, and C-level indicating large DIF. Items classified as C-level DIF require a careful review for possible biases. For polytomous items, DIF was identified using the Liu-Agresti (LA) procedure (Liu & Agresti, 1996; Penfield & Algina, 2003, 2006). The LA estimator of the cumulative common odds ratio for DIF detection is a generalization of the MH procedure. This allows the ETS DIF categorization system to be applied to DIF studies of polytomous items. Table 6.1.9 exhibits the DIF evaluation criteria for dichotomous and polytomous items. The effect size in DIF measures under the MH procedure is denoted by MH D-DIF; that under the LA procedure is denoted by Log(LA).

Table 6.1.9: Differential Item Functioning Evaluation Criteria

DIF Category	Dichotomous Items	Polytomous Items
A (Negligible)	Nonsignificant MH Chi-square test ($p \geq .05$) or $ MH\ D-DIF < 1.0$	Nonsignificant LA Chi-square test ($p \geq .05$) or $ Log(LA) < 0.43$
B (Slight to moderate)	Significant MH Chi-square test ($p < .05$) and $1.0 \leq MH\ D-DIF < 1.5$	Significant LA Chi-square test ($p < .05$) and $0.43 \leq Log(LA) < 0.63$
C (Moderate to high)	Significant MH Chi-square test ($p < .05$) and $ MH\ D-DIF \geq 1.5$	Significant LA Chi-square test ($p < .05$) and $ Log(LA) \geq 0.63$

*Log indicates the logarithm function.

NJSLA–S DIF detection analyses for the field test items only focused on four major comparisons of students: Male/Female, White/Black, White/Hispanic, and White/Asian. For the operational assessment, four other comparisons were made: non-English learner (EL-No)/English learner (EL-Yes), students with disabilities (SWD-Yes)/ students without disabilities (SWD-No), Not economically disadvantaged (EconDis-No)/economically disadvantaged (EconDis-Yes), and TTS/CBT test takers due to the large numbers of students taking the TTS forms. The traditional CBT test takers were the reference group, whereas the TTS test takers were the focal group.

Table 6.1.10, Table 6.1.11, and Table 6.1.12 show the DIF classifications for all eight comparison groups for grade 5, 8, and 11, respectively. The results of the operational DIF analysis were positive except for a small number of items classified as “C,” including two TE items for the Male/Female comparison and one TE item for the EL-No/EL-Yes comparison at grade 8. Additionally, one MC item, two TE items, and one CR item for the EL-No/EL-Yes comparison at grade 11 were classified as “C.” For all other comparisons, zero items across all grade levels were classified as “C.” Moreover, each grade level, comparison group, and item type contained minimal classifications of “B” items. All items were classified as “A” for CBT/TTS DIF at grades 5, 8, and 11. The “C” DIF items will be reinvestigated when more test data become available.

Table 6.1.10: Grade 5 DIF Classification by Item Type

Grade	Group	Item Type	A	B	C
5	Male/Female	MC	14	0	0
5	Male/Female	TE	32	2	0
5	Male/Female	CR	3	0	0
5	Male/Female	Total	49	2	0
5	White/Black	MC	14	0	0
5	White/Black	TE	33	1	0
5	White/Black	CR	3	0	0
5	White/Black	Total	50	1	0
5	White/Hispanic	MC	14	0	0
5	White/Hispanic	TE	34	0	0
5	White/Hispanic	CR	3	0	0
5	White/Hispanic	Total	51	0	0
5	White/Asian	MC	14	0	0
5	White/Asian	TE	34	0	0
5	White/Asian	CR	3	0	0
5	White/Asian	Total	51	0	0
5	EL-No/EL-Yes	MC	14	0	0
5	EL-No/EL-Yes	TE	32	0	0
5	EL-No/EL-Yes	CR	3	0	0
5	EL-No/EL-Yes	Total	49	2	0
5	SWD-No/SWD-Yes	MC	14	0	0
5	SWD-No/SWD-Yes	TE	34	0	0
5	SWD-No/SWD-Yes	CR	3	0	0
5	SWD-No/SWD-Yes	Total	51	0	0
5	EconDis-No/EconDis-Yes	MC	14	0	0
5	EconDis-No/EconDis-Yes	TE	34	0	0
5	EconDis-No/EconDis-Yes	CR	3	0	0
5	EconDis-No/EconDis-Yes	Total	51	0	0
5	CBT/TTS	MC	14	0	0
5	CBT/TTS	TE	34	0	0
5	CBT/TTS	CR	3	0	0
5	CBT/TTS	Total	51	0	0

Table 6.1.11: Grade 8 DIF Classification by Item Type

Grade	Group	Item Type	A	B	C
8	Male/Female	MC	18	0	0
8	Male/Female	TE	42	0	2
8	Male/Female	CR	3	0	0
8	Male/Female	Total	63	0	2
8	White/Black	MC	18	0	0
8	White/Black	TE	43	1	0
8	White/Black	CR	2	1	0
8	White/Black	Total	63	2	0
8	White/Hispanic	MC	18	0	0
8	White/Hispanic	TE	42	2	0
8	White/Hispanic	CR	3	0	0
8	White/Hispanic	Total	63	2	0
8	White/Asian	MC	18	0	0
8	White/Asian	TE	44	0	0
8	White/Asian	CR	3	0	0
8	White/Asian	Total	65	0	0
8	EL-No/EL-Yes	MC	18	0	0
8	EL-No/EL-Yes	TE	42	1	1
8	EL-No/EL-Yes	CR	2	1	0
8	EL-No/EL-Yes	Total	62	2	1
8	SWD-No/SWD-Yes	MC	18	0	0
8	SWD-No/SWD-Yes	TE	44	0	0
8	SWD-No/SWD-Yes	CR	3	0	0
8	SWD-No/SWD-Yes	Total	65	0	0
8	EconDis-No/EconDis-Yes	MC	18	0	0
8	EconDis-No/EconDis-Yes	TE	44	0	0
8	EconDis-No/EconDis-Yes	CR	3	0	0
8	EconDis-No/EconDis-Yes	Total	65	0	0
8	CBT/TTS	MC	18	0	0
8	CBT/TTS	TE	44	0	0
8	CBT/TTS	CR	3	0	0
8	CBT/TTS	Total	65	0	0

Table 6.1.12: Grade 11 DIF Classification by Item Type

Grade	Group	Item Type	A	B	C
11	Male/Female	MC	30	0	0
11	Male/Female	TE	36	0	0
11	Male/Female	CR	1	2	0
11	Male/Female	Total	67	2	0
11	White/Black	MC	30	0	0
11	White/Black	TE	36	0	0
11	White/Black	CR	3	0	0
11	White/Black	Total	69	0	0
11	White/Hispanic	MC	30	0	0
11	White/Hispanic	TE	36	0	0
11	White/Hispanic	CR	3	0	0
11	White/Hispanic	Total	69	0	0
11	White/Asian	MC	30	0	0
11	White/Asian	TE	36	0	0
11	White/Asian	CR	3	0	0
11	White/Asian	Total	69	0	0
11	EL-No/EL-Yes	MC	28	1	1
11	EL-No/EL-Yes	TE	33	1	2
11	EL-No/EL-Yes	CR	2	0	1
11	EL-No/EL-Yes	Total	63	2	4
11	SWD-No/SWD-Yes	MC	30	0	0
11	SWD-No/SWD-Yes	TE	36	0	0
11	SWD-No/SWD-Yes	CR	3	0	0
11	SWD-No/SWD-Yes	Total	69	0	0
11	EconDis-No/EconDis-Yes	MC	30	0	0
11	EconDis-No/EconDis-Yes	TE	36	0	0
11	EconDis-No/EconDis-Yes	CR	3	0	0
11	EconDis-No/EconDis-Yes	Total	69	0	0
11	CBT/TTS	MC	30	0	0
11	CBT/TTS	TE	36	0	0
11	CBT/TTS	CR	3	0	0
11	CBT/TTS	Total	69	0	0

6.2 Item Response Theory

The grade-specific NJSLA–S student ability estimates and subsequent scale scores are calibrated via item response theory (IRT) statistical processes. Section 6.2 of this report explains how IRT is used in the context of the NJSLA–S. First, the concept of IRT is explained. Then, the specific IRT model used for the NJSLA–S is described in conjunction with the assumptions underlying the model. The remainder of Section 6.2 presents evaluations of how well the assumptions of IRT are met.

IRT is conceptualized as a family of mathematical models that provide an equation for the relationship between the probability of a student response to a test item and student latent ability on the construct of interest (Hambleton & Swaminathan, 1985). While latent traits (e.g., anxiety, intelligence, or mastery of the NJSLA–S) are not directly observable, student responses to items are directly observable. Within the context of the NJSLA–S, the latent trait theoretically being measured by the items is student understanding of the New Jersey science curriculum: the NJSLA–S. The directly observable behaviors resulting from that latent trait are the responses of students to those items.

IRT addresses many of the limitations of classical test theory (CTT), such as sample and test dependency, and can improve both the construction and uses of tests (Hambleton & van der Linden, 1982). Hence, IRT can enhance the validity of the inferences made from test scores. Under IRT, item parameters (e.g., item difficulty) are independent of the students who took the test. Similarly, student ability estimates are independent of the test items. Moreover, the test information function (TIF; see Section 8.2 for a more detailed explanation) allows for test construction to be targeted to specific places on the student ability spectrum where decisions are most important to maximize the test’s ability to reliably classify examinees. The increased power of IRT in comparison to CTT requires that certain assumptions be met. When the assumptions of IRT are met, data can be used for psychometric analyses such as equating.

Comprising dichotomous and polytomous items with varying score points (e.g., 0–4-point CR items), the NJSLA–S was constructed to meet the assumptions of a specific IRT model: the Rasch-based partial credit model (PCM; Masters, 1982). The Rasch family of IRT models is a special case of IRT models. That is, Rasch models assume that items discriminate equally and that guessing on items is minimal (Hambleton & Swaminathan, 1985). The PCM is a flexible Rasch-based model that can be used with both dichotomous and polytomous item response data (Ostini & Nering, 2010).

For each polytomous item, there are some ordered levels of performance and an associated number of steps required to move from one level to the next. Statistically, under the PCM, the probability, $\pi_{ij}(x)$, of student j obtaining an item score x with $x = 0, 1, \dots, m$ on polytomous item i can be written as follows:

$$\pi_{ij}(x; \theta_j; \beta_i, \tau_{in}) = \frac{\exp[\sum_{n=0}^x \theta_j - (\beta_i + \tau_{in})]}{\sum_{k=0}^m \exp[\sum_{n=0}^k \theta_j - (\beta_i + \tau_{in})]} \quad \text{Equation 6.1}$$

where θ_j is the student proficiency score, β_i denotes the difficulty or location parameter for item i and τ_{in} with $n = 0, 1, \dots, m$ denotes the threshold or step parameters. For model identification, it is defined that $\tau_{i0} = 0$, $\sum_{n=0}^m \tau_{in} = 0$ and $\exp[\sum_{n=0}^0 \theta_j - (\beta_i + \tau_{in})] = 1$

Accordingly, the predicated probability of a correct response (i.e., item score $x = 1$) to a dichotomous item is given by the following:

$$P(x_{ij} = 1; \theta_j, \beta_i) = \frac{\exp(\theta_j - \beta_i)}{1 + \exp(\theta_j - \beta_i)}, \quad \text{Equation 6.2}$$

where $P(x_{ij} = 1; \theta_j, \beta_i)$ is the probability of student j with a proficiency score θ_j to obtain a correct response to item i and β_i denotes the item difficulty parameter for item i .

Assessing the IRT model fit (i.e., how well the NJSLA–S data meet the assumptions of the PCM) is imperative before using the PCM to analyze NJSLA–S data. If the NJSLA–S data do not meet the assumptions made by the PCM, then any results obtained by using the PCM would be suspect. However, if the NJSLA–S does meet the assumptions required by the PCM, then equating (as presented in Part 7 of this technical report) can be performed under the PCM to place item parameter and ability estimates on a common scale. This allows meaningful, grade-specific comparisons across forms and is important for ensuring the equivalence of the test across years.

The main assumptions of the PCM as they apply to the NJSLA–S are that the test is unidimensional, the items discriminate relatively equally, guessing on items is minimal, each individual item is independent of the others, and the resulting item parameter estimates are invariant regardless of who answered the items. Each of these five IRT assumptions will be explained in detail in the sections below as they relate to the PCM. Also, the PCM item category characteristic functions are graphically presented to show the relationships between student ability estimates and the probability of achieving a specific score point on the 0–3- or 0–4-point CR items. Overall, the results of the 2023 NJSLA–S indicate that the assumptions of the PCM were adequately met.

6.2.1 Unidimensionality

Unidimensionality was checked via multiple methods. First, the intercorrelations among the subscores were evaluated. High correlations would indicate strong linear relationships among the subscore variables, providing evidence of unidimensionality. Second, the eigenvalues of the principal components analysis (PCA) were evaluated. A dominant first eigenvalue, in comparison to the other eigenvalues, is evidence of unidimensionality. Overall, there is ample evidence that the NJSLA–S is a unidimensional test and that the PCM assumption of unidimensionality has been met.

6.2.1.1 Intercorrelations. Tables 6.2.1 and 6.2.2 show the Pearson product-moment correlations among the domains and practices, respectively. High correlations would be evidence of a unidimensional test. Generally, more items in a cluster (i.e., a domain or a practice) will lead to a higher correlation between that cluster and the total test score.

At each grade level, all domains and practices correlated with the total NJSLA–S test score at .90 or above. The lowest correlation among any clusters was .77. The intercorrelations among subscores indicate that the NJSLA–S is a unidimensional test.

Table 6.2.1: Correlation Matrix for Domains

Grade	Domain	NJSLA–S	Earth and Space	Life	Physical
5	Earth and Space	0.93	1.00	-	-
	Life	0.95	0.82	1.00	-
	Physical	0.91	0.77	0.80	1.00
8	Earth and Space	0.93	1.00	-	-
	Life	0.93	0.79	1.00	-
	Physical	0.94	0.81	0.80	1.00
11	Earth and Space	0.93	1.00	-	-
	Life	0.94	0.81	1.00	-
	Physical	0.94	0.80	0.82	1.00

Table 6.2.2: Correlation Matrix for Practices

Grade	Practice	NJSLA–S	Investigation	Sensemaking	Critiquing
5	Investigating	0.90	1.00	-	-
	Sensemaking	0.94	0.80	1.00	-
	Critiquing	0.96	0.79	0.85	1.00
8	Investigating	0.91	1.00	-	-
	Sensemaking	0.94	0.79	1.00	-
	Critiquing	0.94	0.79	0.83	1.00
11	Investigating	0.93	1.00	-	-
	Sensemaking	0.94	0.81	1.00	-
	Critiquing	0.94	0.81	0.83	1.00

6.2.1.2 Principal Component Analysis. Principal Components Analysis (PCA) is a data reduction technique that attempts to account for the variance in measures by converting them into uncorrelated principal components (Brown, 2006). The resulting principal components can be ordered according to the eigenvalues (i.e., the magnitudes of variance accounted for) from the largest to the smallest. The first principal component accounts for as much measured variance as possible, and each succeeding factor does the same until there are as many principal components as original variables (Gorsuch, 1983). Then, a scree plot displays the eigenvalues on the Y-axis and the number (i.e., the order) of principal components on the X-axis. Gorsuch (1983) noted that this method of interpretation works well when sample sizes are large, and the factors are well-defined. The scree plots are interpreted by finding the place on the plot where the slope leveled off. The principal components to the left of that point on the plot are deemed practically significant.

Figures 6.2.1 through 6.2.3 show the scree plots for grades 5, 8, and 11, respectively. As exhibited in these plots, the second most prominent eigenvalue for each grade level is below 2, whereas the most prominent eigenvalues range from approximately 11–13. Each grade’s scree plot shows that only one major dimension is contributing to the variability in student responses to items. The results of each grade’s PCA provide further evidence of the unidimensionality of the NJSLA–S.

Grade 5 Scree Plot

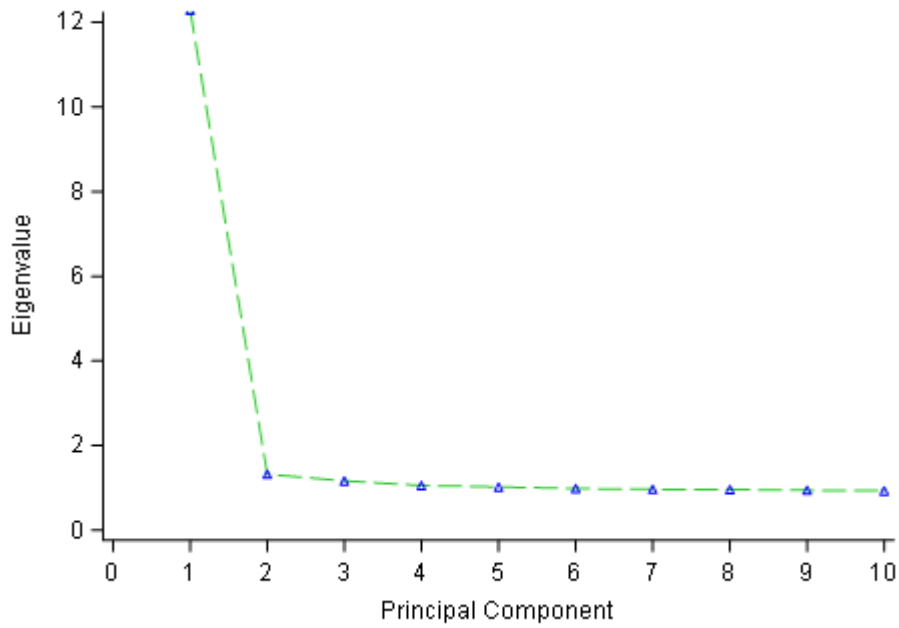


Figure 6.2.1. Grade 5 Scree Plot

Grade 8 Scree Plot

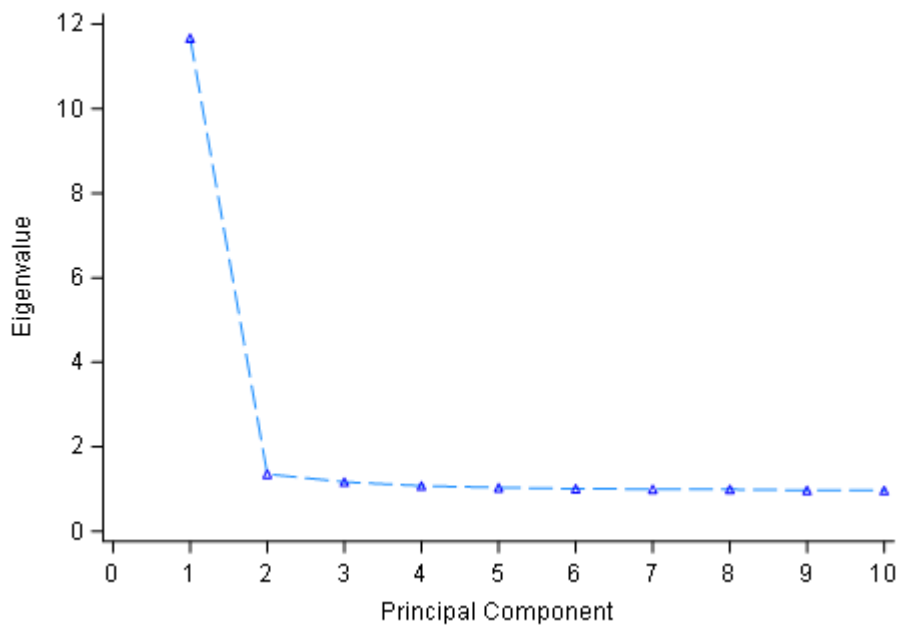


Figure 6.2.2. Grade 8 Scree Plot

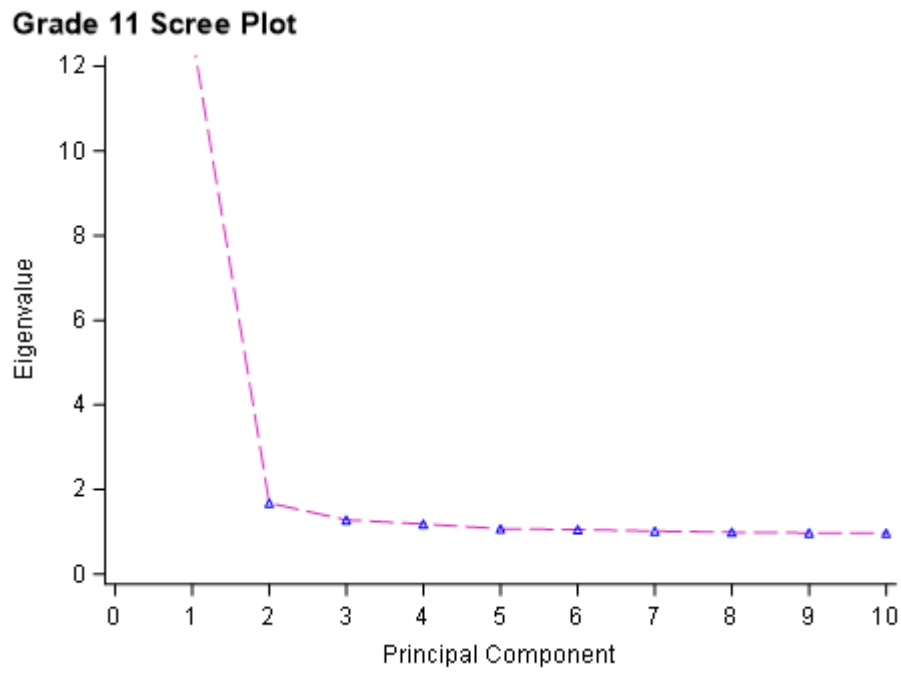


Figure 6.2.3. Grade 11 Scree Plot

6.2.2 Partial Credit Model Fit Statistics

Hambleton, Swaminathan, and Rogers (1991) noted that “[a] poorly fitting IRT model will not yield invariant item and ability parameters” (p. 53), which diminishes the beneficial properties inherent to IRT. The PCM model fit was assessed at the item level via Rasch-based item infit and outfit, discrimination, and guessing statistics. At the person level, model fit was evaluated using Rasch-based person infit and outfit statistics. These statistics were calculated using the 2023 NJSLA–S test data via Winsteps 3.92.1 (Linacre, 2016) Appendix H of this technical report provides the resulting item fit statistics. Overall, there is ample evidence that items fit the assumptions of the PCM for all grades. In particular, the grades 8 and 11 item performance were remarkable, with a small (or zero) percent of items flagged for each of the four model fit categories.

6.2.2.1 Item infit and outfit. Rasch infit and outfit statistics range from zero to infinity, with 1.0 representing ideal model fit (i.e., no misfit). For the NJSLA–S, items were flagged for having infit or outfit statistics outside of the 0.7 to 1.3 range (Wright and Linacre, 1994). Infit statistics are influenced by unexpected responses from students on items that have item difficulties near their ability level (Wright and Masters, 1982). Conversely, outfit statistics are heavily influenced by unexpected student responses to items that are either relatively easy or relatively hard for the student based on the student’s ability level.

Table 6.2.4 provides a summary of item infit and outfit statistics at each grade level. Three grade 5, zero grade 8, and zero grade 11 items were flagged for problematic infit statistics. Slightly more items were flagged based on the outfit statistics, with seven, three, and zero items being flagged for grades 5, 8, and 11, respectively. However, problematic outfit statistics are less of a threat to the validity of test score interpretations than problematic infit statistics. Thus,

while there is clearly room for improving the item outfit, the infit and outfit statistics provide reasonable evidence that the assumptions of the PCM have been met.

Table 6.2.4: Summary of Item Infit and Outfit Statistics

Grade	Fit Statistic	Mean	Min	Max	Outside 0.7 to 1.3	% Flagged
5	Infit	1.00	0.71	1.39	3 out of 51	6.0
5	Outfit	1.03	0.62	1.76	7 out of 51	14.0
8	Infit	1.00	0.80	1.23	0 out of 65	0.0
8	Outfit	1.02	0.71	1.44	3 out of 65	5.0
11	Infit	1.00	0.75	1.18	0 out of 69	0.0
11	Outfit	1.01	0.70	1.29	0 out of 69	0.0

6.2.2.2 Rasch discrimination. The PCM assumes that all items discriminate equally. Practically, items never discriminate equally, but if they are within reasonable thresholds then the assumption will be met. The PCM does not model item discrimination, nor does it adjust item difficulty or person ability estimates based on item discrimination. However, Winsteps provides an index of item discrimination that approximates the discrimination parameter from the 2PL model (Linacre, 2016). Rasch discrimination statistics are centered at 1.0, which indicates that the item is discriminating exactly as expected by the PCM. Items are flagged when their discrimination statistics fall outside of the range of 0.5 to 1.5.

Table 6.2.5 provides a summary of Rasch discrimination statistics at each grade level. The Rasch discrimination values were good across each grade. There were eight (16%), zero (0%), and three (4%) items flagged for having values outside the 0.5 to 1.5 threshold for grades 5, 8, and 11, respectively. While five of these items were also flagged for item outfit statistics in grade 5, the Rasch discrimination analysis results provide evidence that the PCM assumptions have been met.

Table 6.2.5: Summary of Rasch Discrimination Statistics

Grade	Fit Statistic	Mean	Min	Max	Outside 0.5 to 1.5	% Flagged
5	Discrimination	1.01	0.04	1.59	8 out of 51	16.0
8	Discrimination	1.00	0.61	1.45	0 out of 60	0.0
11	Discrimination	1.00	0.49	1.60	3 out of 69	4.0

6.2.2.3 Rasch lower asymptote. The PCM assumes that there is minimal guessing on the test items. Practically, however, students guess, and sometimes they guess correctly. Thus, as with the assumption of equal discrimination, the assumption of minimal guessing is met if item guessing statistics remain within a reasonable threshold. The PCM models guessing as misfit (i.e., infit and outfit) and does not adjust item difficulty or person ability estimates based on guessing. However, Winsteps provides an index that approximates a guessing parameter in the form of lower asymptote statistics (Linacre, 2016). Rasch lower asymptote statistics are ideally 0.0, which indicates that an item is displaying little to no guessing. Items are flagged when their lower asymptote statistics fall outside of the range of 0.0 to 0.1.

Table 6.2.6 provides a summary of the lower asymptote statistics at each grade level. Each grade level saw only a few items flagged for having a lower asymptote value outside of the .1 threshold (four items in grade 5, zero items in grade 8, and three items in grade 11). Four of these items were also flagged for at least one other statistic (i.e., infit, outfit, or discrimination). Unsurprisingly, these items had low item-total correlations. Nevertheless, the Rasch lower asymptote statistics provided evidence that the PCM assumptions have been satisfied as few items displayed lower asymptote values outside the acceptable threshold.

Table 6.2.6: Summary of Rasch Lower Asymptote Statistics

Grade	Fit Statistic	Mean	Min	Max	Greater Than .1	% Flagged
5	Lower Asymptote	.03	.00	.21	4 out of 51	8.0
8	Lower Asymptote	.02	.00	.08	0 out of 60	0.0
11	Lower Asymptote	.03	.00	.18	3 out of 69	4.0

6.2.2.4 Rasch person infit and outfit. PCM person fit statistics are useful for evaluating whether student response patterns are reasonable. Reasonableness includes not only response patterns that are improbable, but those that are too probable. Multiple factors can cause distortions in the expected patterns of test scores, including:

- Carelessness—examinees miss items that they should have answered correctly.
- Cheating—examinees receive information to correctly answer items that they would have normally missed.
- Guessing—examinees correctly answer items without knowing the correct answer.
- Creative responses—examinees misinterpret the item.
- Test administration errors.

Two PCM person-fit statistics were used: infit and outfit. Person infit is more influenced by responses to items that are targeted at the person’s ability level; outfit is more influenced by responses to items that are relatively easy or hard for a student (Wright & Masters, 1982). Ideally, both statistics would be close to 1.0. For the NJSLA–S, values larger than 1.3 indicate model underfit, while values smaller than .7 indicate model overfit.

Tables 6.2.7 and 6.2.8 show, respectively, the person infit and outfit descriptive statistics by demographic variables. For NJSLA–S, person fit statistics were evaluated based on the following demographics: gender, ethnicity, English learner (EL) status, economically disadvantaged (EconDis) status, and students with disabilities (SWD) status. Tables 6.2.9 and 6.2.10 breakdown, respectively, the person infit and outfit descriptive statistics by test forms including CBT, PBT, TTS, Spanish, Spanish TTS, and Human Reader forms. Figures 6.2.4 through 6.2.6 exhibit grade level distributions of both the person infit and outfit statistics for all students.

At the overall level across all combinations of grade and demographic variables, as shown in Table 6.2.7, 9.45% of grade 5 students, 2.36% of grade 8 students, and 4.74% of grade 11 students were flagged for person infit statistics. Asian students tended to have the highest percentage of students flagged for person infit among the ethnic groups investigated at their grade level. As shown in Table 6.2.9, the grade 11 HR forms flagged 15.38% of students for person infit. However, it should be noted that there were only 26 students taking the HR form.

Overall, there were relatively more students flagged for aberrant person outfit statistics than for person infit statistics. As shown in Table 6.2.8, the English learners and students with disabilities tended to have a higher percentage of students that were flagged for person outfit statistics relative to other demographic groups at their grade level. Within those groups, the students that were flagged also tended to be lower performing. Additionally, they were more likely to have taken accommodated forms, which themselves had higher percentages of students flagged for person outfit than did the CBT forms.

As stated earlier, aberrant person outfit statistics are less of a threat to the validity of the test score inferences than are aberrant person infit statistics. It is likely that the reason for the large percentages of person-outfit flags is that while these students tended to be lower performing, there were some items that they were able to unexpectedly answer correctly. Moreover, because the students that were flagged were so low performing, it is unlikely that the misfit was having any meaningful impact on the reliability of the student proficiency classification. However, a deeper investigation into the person outfit statistics for English learners, students with disabilities, and the accommodated forms was conducted to ensure there were no concerns. The results of this investigation are summarized in Section 6.2.2.5.

Table 6.2.7: Summary of Person Infit Statistics by Demographic Group

Grade	Group	N	Mean Scale Score	Mean Person Infit	Person Infit Min	Person Infit Max	N Flagged	% Flagged	Flagged Mean Scale Score
5	NJSLA-S	96,392	166.01	1.03	0.57	3.56	9,113	9.45	181.79
5	Male	49,082	167.53	1.03	0.57	3.56	4,488	9.14	187.16
5	Female	47,299	164.42	1.04	0.57	3.56	4,624	9.78	176.57
5	Am. Indian	169	168.79	1.04	0.72	2.00	16	9.47	185.94
5	Asian	10,765	202.00	1.03	0.59	2.63	1,266	11.76	205.35
5	Black	14,028	143.47	1.04	0.58	2.75	1,112	7.93	166.29
5	Hispanic	31,700	147.85	1.04	0.59	3.56	2,594	8.18	165.80
5	Pacific Islander	179	169.62	1.04	0.71	2.38	13	7.26	157.00
5	White	36,375	178.71	1.03	0.57	3.56	3,770	10.36	188.64
5	EL – Yes	9,158	125.24	1.04	0.64	2.59	463	5.06	146.99
5	EL – No	87,234	170.29	1.03	0.57	3.56	8,650	9.92	183.65
5	EconDis – Yes	36,109	143.93	1.04	0.59	2.95	2,797	7.75	163.61
5	EconDis – No	60,283	179.23	1.03	0.57	3.56	6,316	10.48	189.84
5	SWD – Yes	20,003	143.22	1.04	0.62	3.54	1,430	7.15	167.30
5	SWD – No	76,389	171.97	1.03	0.57	3.56	7,683	10.06	184.49
8	NJSLA-S	101,478	162.05	0.99	0.66	2.42	2,390	2.36	178.85
8	Male	49,201	161.42	0.99	0.66	2.22	1,131	2.30	175.15
8	Female	52,212	162.63	0.99	0.66	2.42	1,255	2.40	182.13
8	Am. Indian	154	159.33	0.99	0.74	1.51	4	2.60	181.25
8	Asian	10,718	192.98	1.00	0.66	2.22	392	3.66	202.32
8	Black	14,998	143.73	1.00	0.66	2.21	285	1.90	161.32
8	Hispanic	32,921	147.92	1.00	0.68	2.22	687	2.09	162.08
8	Pacific Islander	206	170.60	0.99	0.73	1.47	4	1.94	170.00
8	White	39,768	171.73	0.98	0.66	2.42	947	2.38	186.00
8	EL – Yes	7,151	131.23	1.01	0.71	2.22	115	1.61	133.01

Grade	Group	N	Mean Scale Score	Mean Person Infit	Person Infit Min	Person Infit Max	N Flagged	% Flagged	Flagged Mean Scale Score
8	EL – No	94,327	164.39	0.99	0.66	2.42	2,275	2.41	181.17
8	EconDis – Yes	35,709	145.52	1.00	0.66	2.22	707	1.98	160.14
8	EconDis – No	65,769	171.03	0.99	0.66	2.42	1,683	2.56	186.71
8	SWD – Yes	20,520	144.56	1.00	0.67	2.28	337	1.64	163.56
8	SWD – No	80,958	166.49	0.99	0.66	2.42	2,053	2.54	181.36
11	NJSLA–S	94,023	170.87	1.01	0.60	2.45	4,452	4.74	179.94
11	Male	45,924	171.62	1.01	0.63	2.37	2,184	4.76	171.67
11	Female	47,959	170.08	1.01	0.60	2.45	2,259	4.71	187.81
11	Am. Indian	141	163.55	1.03	0.75	1.50	8	5.67	193.13
11	Asian	10,003	211.01	1.02	0.65	2.26	551	5.51	205.02
11	Black	12,731	148.02	1.02	0.63	2.40	581	4.56	157.39
11	Hispanic	28,687	151.90	1.01	0.60	2.45	1,241	4.33	160.02
11	Pacific Islander	313	174.57	1.01	0.72	1.89	14	4.47	179.79
11	White	40,005	181.20	1.01	0.63	2.44	1,954	4.88	191.66
11	EL – Yes	5,290	126.87	1.02	0.65	2.28	161	3.04	126.33
11	EL – No	88,733	173.49	1.01	0.60	2.45	4,291	4.84	181.95
11	EconDis – Yes	28,095	150.55	1.01	0.63	2.45	1,197	4.26	158.21
11	EconDis – No	65,928	179.53	1.01	0.60	2.44	3,255	4.94	187.93
11	SWD – Yes	18,600	148.90	1.02	0.63	2.45	827	4.45	163.58
11	SWD – No	75,423	176.29	1.01	0.60	2.37	3,625	4.81	183.67

Table 6.2.8: Summary of Person Outfit Statistics by Demographic Group

Grade	Group	N	Mean Scale Score	Mean Person Outfit	Person Outfit Min	Person Outfit Max	N Flagged	% Flagged	Flagged Mean Scale Score
5	NJSLA-S	96,392	166.01	1.03	0.17	6.81	8,838	9.17	133.79
5	Male	49,082	167.53	1.02	0.17	6.81	4,569	9.31	136.02
5	Female	47,299	164.42	1.03	0.31	3.19	4,269	9.03	131.42
5	Am. Indian	169	168.79	1.05	0.70	1.84	18	10.65	116.78
5	Asian	10,765	202.00	0.98	0.31	3.42	676	6.28	214.30
5	Black	14,028	143.47	1.06	0.36	3.19	1,867	13.31	114.70
5	Hispanic	31,700	147.85	1.06	0.31	6.81	3,793	11.97	116.81
5	Pacific Islander	179	169.62	1.02	0.63	1.93	13	7.26	165.69
5	White	36,375	178.71	1.00	0.17	4.44	2,245	6.17	151.60
5	EL-Yes	9,158	125.24	1.12	0.31	3.42	1,826	19.94	109.00
5	EL-No	87,234	170.29	1.02	0.17	6.81	7,012	8.04	140.25
5	EconDis-Yes	36,109	143.93	1.06	0.31	3.72	4,648	12.87	115.29
5	EconDis-No	60,283	179.23	1.00	0.17	6.81	4,190	6.95	154.33
5	SWD-Yes	20,003	143.22	1.07	0.31	6.81	2,934	14.67	114.73
5	SWD-No	76,389	171.97	1.01	0.17	3.72	5,904	7.73	143.27
8	NJSLA-S	101,478	162.05	1.02	0.31	3.06	6,562	6.47	126.40
8	Male	49,201	161.42	1.01	0.34	2.70	2,997	6.09	126.46
8	Female	52,212	162.63	1.02	0.31	3.06	3,563	6.82	126.35
8	Am. Indian	154	159.33	1.03	0.71	1.40	8	5.19	115.50
8	Asian	10,718	192.98	1.01	0.54	2.00	327	3.05	173.02
8	Black	14,998	143.73	1.04	0.34	3.06	1,542	10.28	119.39
8	Hispanic	32,921	147.92	1.03	0.31	2.70	2,856	8.68	121.70
8	Pacific Islander	206	170.60	1.00	0.65	1.64	7	3.40	121.14
8	White	39,768	171.73	1.00	0.39	2.55	1,659	4.17	131.85
8	EL-Yes	7,151	131.23	1.06	0.39	2.51	952	13.31	118.43

Grade	Group	N	Mean Scale Score	Mean Person Outfit	Person Outfit Min	Person Outfit Max	N Flagged	% Flagged	Flagged Mean Scale Score
8	EL–No	94,327	164.39	1.01	0.31	3.06	5,610	5.95	127.76
8	EconDis–Yes	35,709	145.52	1.03	0.34	2.57	3,354	9.39	120.96
8	EconDis–No	65,769	171.03	1.01	0.31	3.06	3,208	4.88	132.10
8	SWD–Yes	20,520	144.56	1.04	0.34	3.06	2,209	10.77	119.58
8	SWD–No	80,958	166.49	1.01	0.31	2.70	4,353	5.38	129.87
11	NJSLA–S	94,023	170.87	1.01	0.36	3.60	5,621	5.98	129.86
11	Male	45,924	171.62	1.01	0.48	3.60	2,423	5.28	129.71
11	Female	47,959	170.08	1.02	0.36	3.00	3,194	6.66	129.90
11	Am. Indian	141	163.55	1.02	0.74	1.54	11	7.80	143.00
11	Asian	10,003	211.01	1.00	0.51	2.27	359	3.59	198.47
11	Black	12,731	148.02	1.04	0.36	3.00	1,164	9.14	115.88
11	Hispanic	28,687	151.90	1.02	0.53	2.61	2,201	7.67	117.96
11	Pacific Islander	313	174.57	1.00	0.66	1.56	14	4.47	121.86
11	White	40,005	181.20	1.00	0.51	3.60	1,771	4.43	139.06
11	EL–Yes	5,290	126.87	1.06	0.57	2.31	661	12.50	111.36
11	EL–No	88,733	173.49	1.01	0.36	3.60	4,960	5.59	132.33
11	EconDis–Yes	28,095	150.55	1.02	0.36	3.60	2,285	8.13	116.62
11	EconDis–No	65,928	179.53	1.01	0.48	3.00	3,336	5.06	138.94
11	SWD–Yes	18,600	148.90	1.04	0.51	3.60	1,813	9.75	117.46
11	SWD–No	75,423	176.29	1.00	0.36	3.00	3,808	5.05	135.78

Table 6.2.9: Summary of Person Infit Statistics by Form

Grade	Form	N	Mean Scale Score	Mean Person Infit	Person Infit Min	Person Infit Max	N Flagged	%Flagged	Flagged Mean Scale Score
5	CBT	75,574	172.30	1.03	0.57	3.56	7,589	10.04	185.01
5	PBT	143	129.87	1.04	0.70	1.65	6	4.20	169.83
5	TTS	17,960	145.83	1.04	0.61	3.54	1,385	7.71	167.50
5	SP	1,473	125.19	1.04	0.70	2.35	65	4.41	145.71
5	SP TTS	983	126.12	1.04	0.70	2.04	42	4.27	149.33
5	HR	201	133.24	1.06	0.72	1.81	20	9.95	153.05
8	CBT	84,298	165.86	0.99	0.66	2.42	2,061	2.44	182.31
8	PBT	72	141.71	1.03	0.71	1.45	2	2.78	156.50
8	TTS	14,433	145.07	1.00	0.66	2.19	282	1.95	161.11
8	SP	1,876	134.29	0.99	0.72	1.63	22	1.17	135.59
8	SP TTS	729	134.46	1.00	0.74	2.22	23	3.16	129.74
8	HR	47	128.49	1.02	0.81	1.22	0	0.00	-
11	CBT	84,251	173.44	1.01	0.60	2.44	4,060	4.82	182.10
11	PBT	242	140.43	1.04	0.71	1.52	9	3.72	158.56
11	TTS	7,732	153.92	1.01	0.64	2.45	316	4.09	165.09
11	SP	1,474	127.59	1.02	0.69	2.28	50	3.39	121.54
11	SP TTS	273	128.11	1.01	0.71	1.68	11	4.03	117.09
11	HR	26	128.42	1.05	0.78	1.47	4	15.38	118.75

Note. CBT: Computer-Based Test; PBT: Paper-Based Test; TTS: Text-to-Speech; SP: Spanish; SP TTT: Spanish Text-to-Speech; HR: Human-Reader

Table 6.2.10: Summary of Person Outfit Statistics by Form

Grade	Form	N	Mean Scale Score	Mean Person Outfit	Person Outfit Min	Person Outfit Max	N Flagged	% Flagged	Flagged Mean Scale Score
5	CBT	75,574	172.30	1.01	0.31	4.44	5,850	7.74	143.42
5	PBT	143	129.87	1.11	0.65	2.22	26	18.18	119.46
5	TTS	17,960	145.83	1.06	0.17	6.81	2,438	13.57	116.01
5	SP	1,473	125.19	1.12	0.50	2.70	293	19.89	108.89
5	SP TTS	983	126.12	1.12	0.38	2.35	185	18.82	109.80
5	HR	201	133.24	1.08	0.44	1.85	30	14.93	116.73
8	CBT	84,298	165.86	1.01	0.31	2.70	4,768	5.66	128.75
8	PBT	72	141.71	1.22	0.71	2.05	22	30.56	127.73
8	TTS	14,433	145.07	1.04	0.38	3.06	1,532	10.61	120.16
8	SP	1,876	134.29	1.01	0.53	2.00	154	8.21	120.04
8	SP TTS	729	134.46	1.01	0.53	1.93	68	9.33	117.82
8	HR	47	128.49	1.08	0.65	1.63	9	19.15	117.78
11	CBT	84,251	173.44	1.01	0.36	3.00	4,749	5.64	132.23
11	PBT	242	140.43	1.10	0.66	2.27	40	16.53	111.35
11	TTS	7,732	153.92	1.02	0.51	3.60	613	7.93	118.78
11	SP	1,474	127.59	1.05	0.57	2.27	181	12.28	113.03
11	SP TTS	273	128.11	1.04	0.59	1.79	31	11.36	114.06
11	HR	26	128.42	1.02	0.73	1.42	2	7.69	100.00

Note. CBT: Computer-Based Test; PBT: Paper-Based Test; TTS: Text-to-Speech; SP: Spanish; SP TTT: Spanish Text-to-Speech; HR: Human-Reader

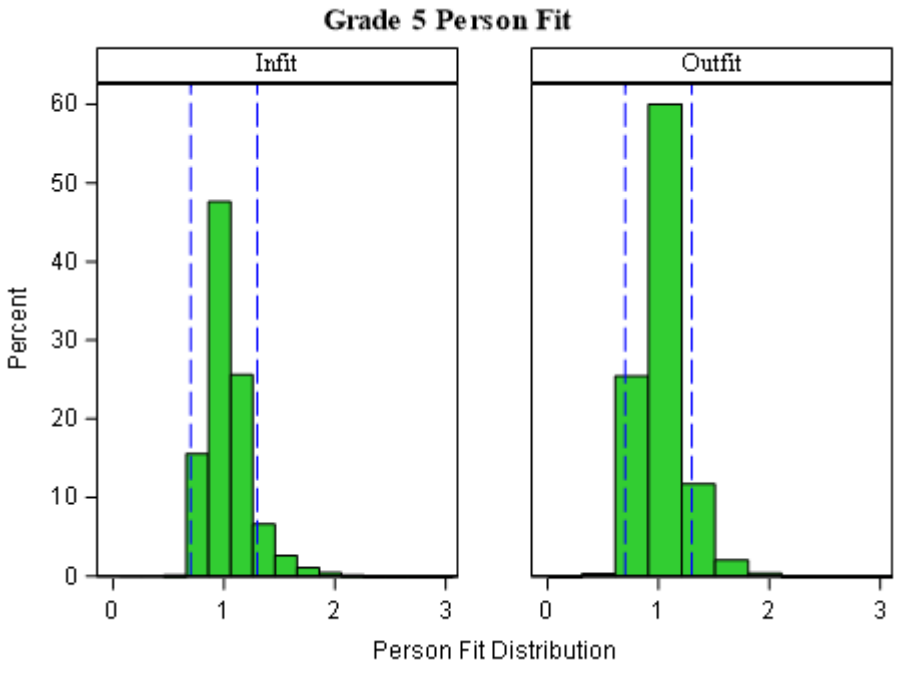


Figure 6.2.4. Grade 5 Person Infit and Outfit Distributions

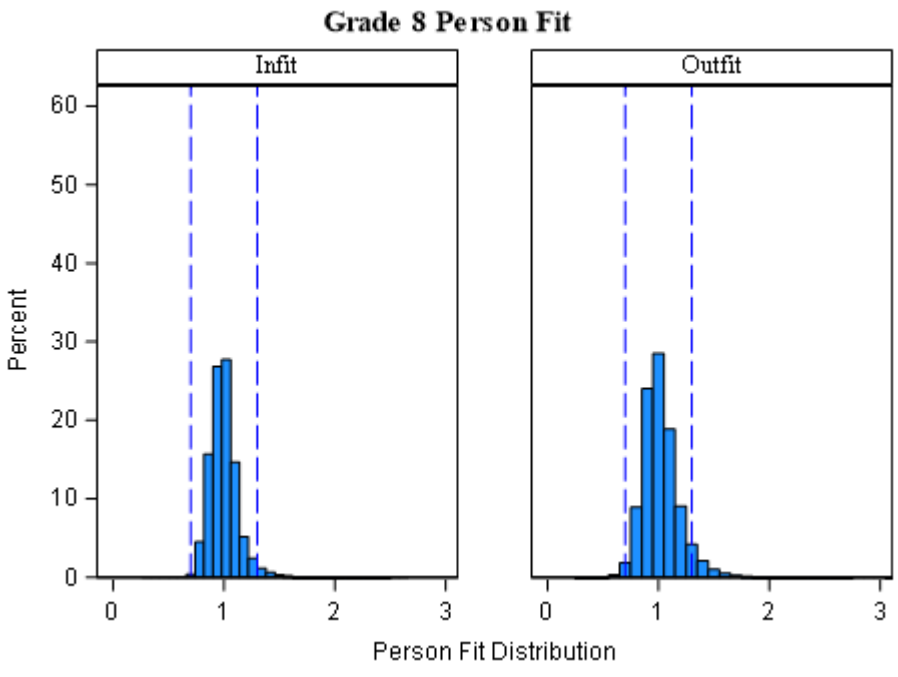


Figure 6.2.5. Grade 8 Person Infit and Outfit Distributions

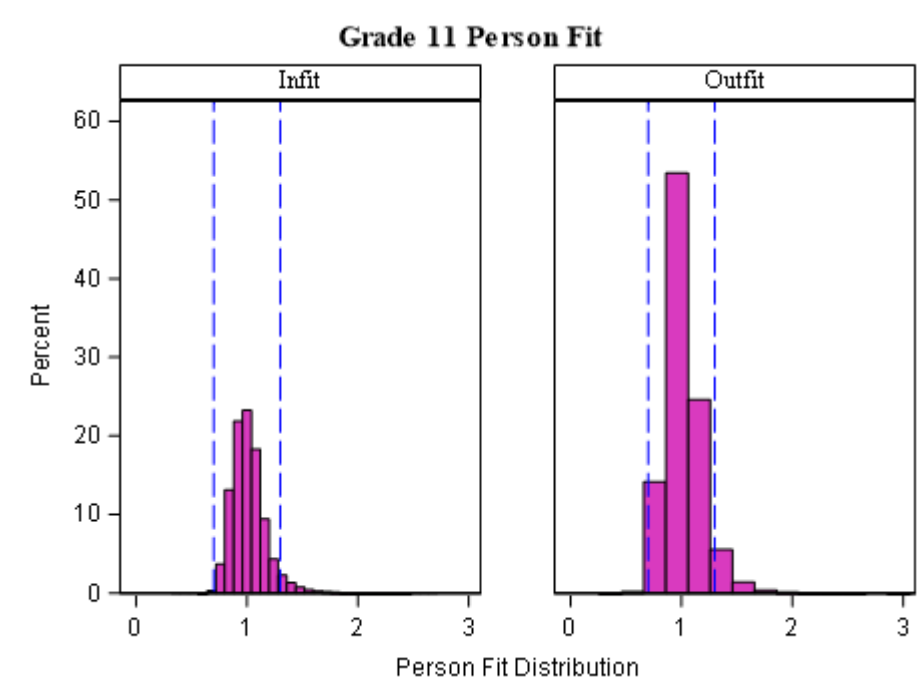


Figure 6.2.6. Grade 11 Person Infit and Outfit Distributions

6.2.2.5 Further investigation of person outfit statistics. Within the context of the Rasch model, person infit and outfit statistics are based on model residuals and are interpreted in terms of fit and misfit. It is essential that the model fit residuals are neither larger (indicating underfit) nor smaller (indicating overfit) than expected. In brief, overfit implies that the item scores were too predictable while underfit implies item scores were too unpredictable. Notably, underfit is more troublesome than overfit because an excess of expected scores does not necessarily imply invalidity in the measurement process (Engelhard & Wind, 2018; Linacre, 2016). However, an excess of unexpected scores could suggest there are some unaccounted-for factors affecting item scores. Therefore, this investigation focused on further examining the underfitting students in the student cohorts that showed a high frequency of outfit flags. Table 6.2.11 presents the distributions of total fit and total misfit, which further breakdowns into overfit and underfit by student cohort for each grade.

Table 6.2.11: Distribution of Person Fit for Grades 5, 8, and 11

Grade	Subgroup	EL (%)	EcoDisad (%)	SWD (%)	ACC (%)
5	Total	9158 (100%)	36109 (100%)	20003 (100%)	20760 (100%)
	Fit	7332 (80.1 %)	31461 (87.1%)	17069 (85.3%)	17788 (85.7%)
	Total Misfit	1826 (19.9%)	4648 (12.9%)	2934 (14.7%)	2972 (14.3%)
	Overfit	135 (1.5%)	364 (1.0%)	255 (1.3%)	254 (1.2%)
	Underfit	1691 (18.5%)	4284 (11.9%)	2679 (13.4%)	2718 (13.1%)
8	Total	7151 (100%)	35709 (100%)	20520 (100%)	17157 (100%)
	Fit	6199 (86.7%)	32355 (90.6%)	18311 (89.2%)	15372 (89.6%)
	Total Misfit	952 (13.3%)	3354 (9.4%)	2209 (10.8%)	1785 (10.4%)
	Overfit	123 (1.7%)	487 (1.4%)	291 (1.4%)	252 (1.5%)
	Underfit	829 (11.6%)	2867 (8.0%)	1918 (9.3%)	1533 (8.9%)
11	Total	5290 (100%)	28095 (100%)	18600 (100%)	9747 (100%)
	Fit	4629 (87.5%)	25810 (91.9%)	16787 (90.3%)	8880 (91.1%)
	Total Misfit	661 (12.5%)	2285 (8.1%)	1813 (9.7%)	867 (8.9%)
	Overfit	54 (1.0%)	244 (0.9%)	155 (0.8%)	89 (0.9%)
	Underfit	607 (11.5%)	2041 (7.3%)	1658 (8.9%)	778 (8.0%)

Note. EL = English learner; EcoDisad = Economically disadvantaged; SWD = Students with disabilities; ACC = students taking an accommodated test form

While the PCM provides summary statistics to identify overall person misfit (i.e., infit and outfit), other methods are required to identify which items were contributing to the unexpected student performance. To that end, this study used the item difficulty parameters to identify which items showed unexpected performance for the group of students showing underfit. Specifically, the delta plot method (Angoff & Ford, 1973; discussed in Section 7.1) was employed. To conduct the delta method, the *p-values* of fitting students within a cohort were compared to the *p-values* of underfitting students within a cohort. This comparison was made between the students exhibiting no misfit and the underfitting students for each cohort (i.e., Students with Disabilities, English Learners, Economically Disadvantaged, and students taking an accommodated test form) for all three grades. The observed *p-values* for the students exhibiting no misfit and students exhibiting underfit are provided in Appendix N.

Table 6.2.12 exhibits the items flagged resulting from a single round of delta analysis. There were three, four, and seven items flagged for grades 5, 8, and 11, respectively. For an item flagged for a student cohort, the associated *p-values* for the group of students showing no misfit (P_FT) and the *p-values* for the underfit subgroup (P_UF) are presented in Table 6.2.12. While there was some overlap in items flagged across cohorts, not all items were flagged for all cohorts. When an item was not flagged for a specific cohort, no P_FT or P_UF values were provided. As shown in Table 6.2.12, when the flagged TE items had Rasch B values larger than 1.2, the *p-values* associated with the underfit subgroup were higher than those associated with the fit subgroup. This pattern was observed for three TE items in Grade 5, two TE items in Grade 8, and three TE items in Grade 11. In contrast, on the flagged CR items, the underfit subgroup had lower *p-values* than the fit subgroup.

Based on the study results, the higher rate of outfit flags seen in these student cohorts likely stems from the fact that these, on average, lower performing (see Table 6.2.8) students unexpectedly answered the challenging TE items correctly and did not perform as expected on the CR items. As previously mentioned, aberrant person infit statistics pose a greater threat to the validity of inferences than aberrant person outfit statistics. Furthermore, considering that the students flagged for outfit statistics were predominantly low performing, it is improbable that the presence of misfit significantly impacted the reliability of student proficiency classification.

Table 6.2.12: *p-values* of Items Flagged by Delta Method

Flagged Item						EL	EcoDisad		SWD		ACC		
Grade	Item UIN	DCI	SEP	Item Type	Rasch B	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
5	2205M011_01	ESS	INV	TE	1.401	0.15	0.32	0.17	0.33	0.19	0.33	0.18	0.32
5	1905B009_05	LS	INV	TE	1.579	0.12	0.25	-	-	-	-	-	-
5	2205M012_04	PS	CRI	TE	1.587	0.18	0.35	0.19	0.37	0.17	0.36	0.18	0.36
8	2208M028_07	LS	CRI	TE	1.242	0.10	0.23	-	-	-	-	-	-
8	1908M003_07	LS	SEN	TE	1.413	0.07	0.26	0.10	0.27	0.11	0.28	0.10	0.28
8	2008M001_01	PS	SEN	TE	0.354	-	-	0.17	0.04	-	-	-	-
8	1908B000_11	PS	SEN	CR	-0.027	0.12	0.02	0.22	0.03	0.23	0.03	0.21	0.03
11	1911B009_03A	LS	SEN	TE	0.479	-	-	-	-	-	-	0.23	0.06
11	2111M004_02	LS	SEN	TE	2.137	0.04	0.25	0.09	0.25	0.09	0.24	0.08	0.26
11	2211M003_04	LS	SEN	TE	2.320	0.02	0.16	-	-	-	-	-	-
11	2011M010_05	PS	CRI	TE	1.358	-	-	-	-	0.19	0.30	-	-
11	2211B006_09	ESS	CRI	CR	-0.242	-	-	0.48	0.10	0.43	0.10	0.45	0.11
11	1911B009_07A	LS	SEN	CR	0.683	0.10	0.02	0.25	0.07	-	-	0.23	0.06
11	2211B000_12	PS	INV	CR	1.522	-	-	0.12	0.04	0.12	0.04	0.12	0.04

Note. EL = English learner; EcoDisad = Economically disadvantaged; SWD = Students with disabilities; ACC = students taking an accommodated test form; P_FT = average *p-value* of students showing no misfit; P_UF = average *p-value* of students showing underfit

6.2.3 Local Independence

The PCM assumes that student responses to items are independent of responses to other items. In other words, student performance on one item does not affect performance on the other items on the test. If the assumption of local independence is violated, then the validity of inferences made from test scores could be threatened, the reliability of the assessment could be overestimated, and item-total correlations could be inflated. The assumption of local independence was tested via calculations of Yen’s Q3 (Yen, 1984), which is an item residual correlation. The item residual (d_i) for item i at a student ability estimate θ_j is defined as follows:

$$d_i = X_i - E(X_i; \theta_j), \quad \text{Equation 6.3}$$

where X_i is an observed item score and $E(X_i; \theta_j)$ is the conditional expected item score under the IRT model of interest. The Q3 statistics for items i and k ($i \neq k$) are then computed as the Person correlation of d_i and d_k over all test takers.

Table 6.2.13 summarizes Yen’s Q3 statistics for the NJSLA–S test at each grade level. All pairwise combinations of items were checked, and they were flagged if their Q3 value was above .2 or below $-.2$ (Chen & Thissen, 1997). The results at all grades indicate that the assumption of local independence was met because very few combinations of items displayed Q3 values outside the acceptable threshold.

Table 6.2.13: Summary of Yen’s Q3 Statistics

Grade	Mean	Min	Max	Outside $-.2$ to $.2$	% Flagged
5	-.02	-.17	.34	3 out of 1,275	.24
8	-.01	-.08	.15	0 out of 2,080	.00
11	-.01	-.10	.27	1 out of 2,346	.04

6.2.4 Item Characteristic Curves—CR Items

Under IRT, the item characteristic curves (ICC; the item categorical response functions) for a CR item show the relationship between latent student ability (theta) and the probability of achieving a specific score point on that item. The ICCs for each of the hand-scored, constructed-response, 0–3- or 0–4-point items are presented in Figures 6.2.7 through 6.2.15 below. The vertical, dashed lines represent, from left to right, the Level 1/2, 2/3, and 3/4 cut scores on the theta scale. In addition, Table 6.2.14 shows the percentages of students receiving each score point for all nine CR items.

Table 6.2.14: Constructed-Response Point Distribution Percentages

Grade	Item	%0	%1	%2	%3	%4
5	CR Item 1	31.52	17.31	23.42	11.75	16.01
5	CR Item 2	63.87	16.23	9.18	7.31	3.42
5	CR Item 3	54.17	10.89	13.50	13.13	8.31
8	CR Item 1	39.48	19.69	17.42	14.57	8.85
8	CR Item 2	79.03	13.83	5.07	2.06	N/A
8	CR Item 3	34.69	25.61	22.37	17.32	N/A
11	CR Item 1	32.60	25.79	25.25	9.86	6.50
11	CR Item 2	60.93	25.11	10.90	3.06	N/A
11	CR Item 3	25.65	10.65	16.18	21.27	26.25

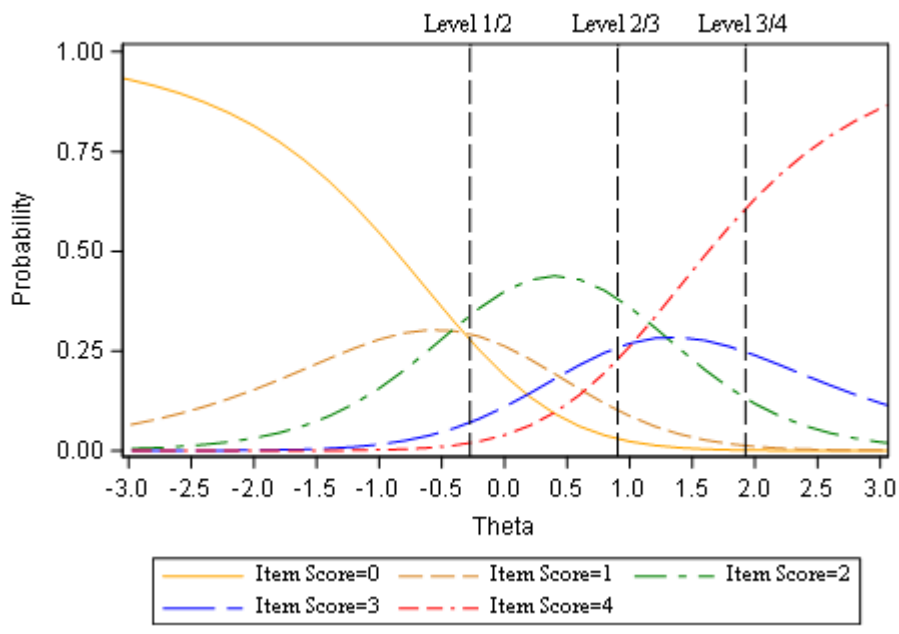


Figure 6.2.7. ICC Plot for Grade 5 Constructed-Response Item 1

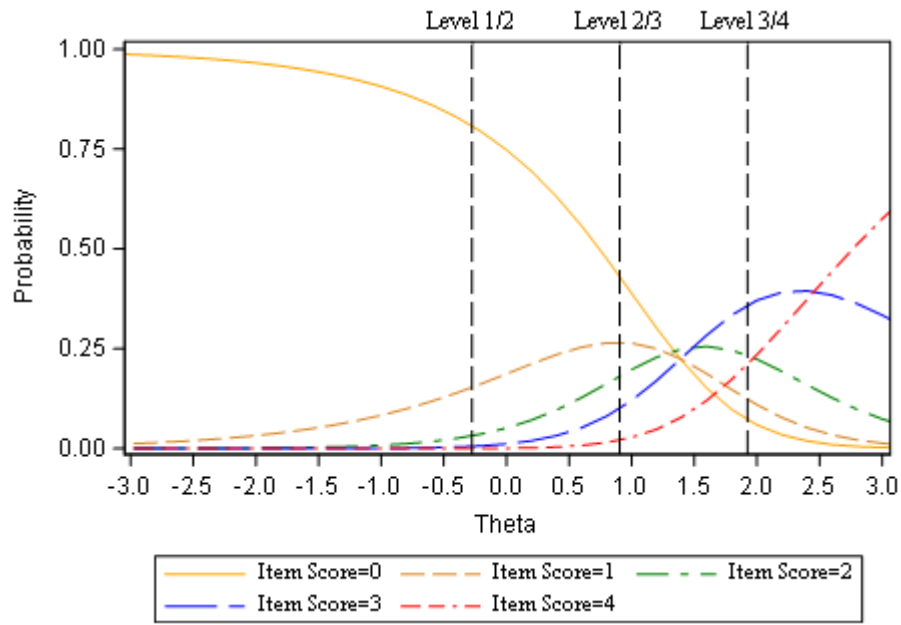


Figure 6.2.8. ICC Plot for Grade 5 Constructed-Response Item 2

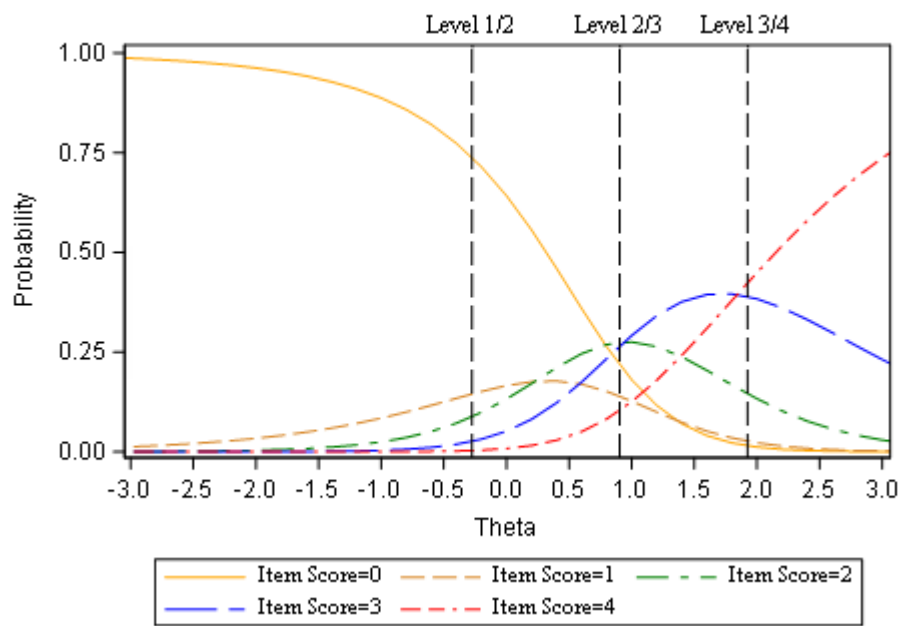


Figure 6.2.9. ICC Plot for Grade 5 Constructed-Response Item 3

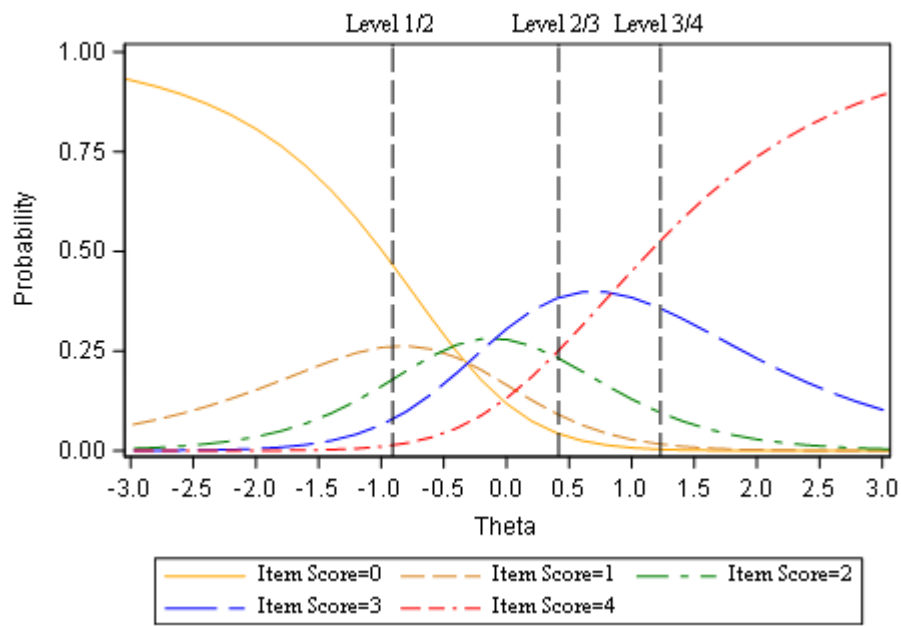


Figure 6.2.10. ICC Plot for Grade 8 Constructed-Response Item 1

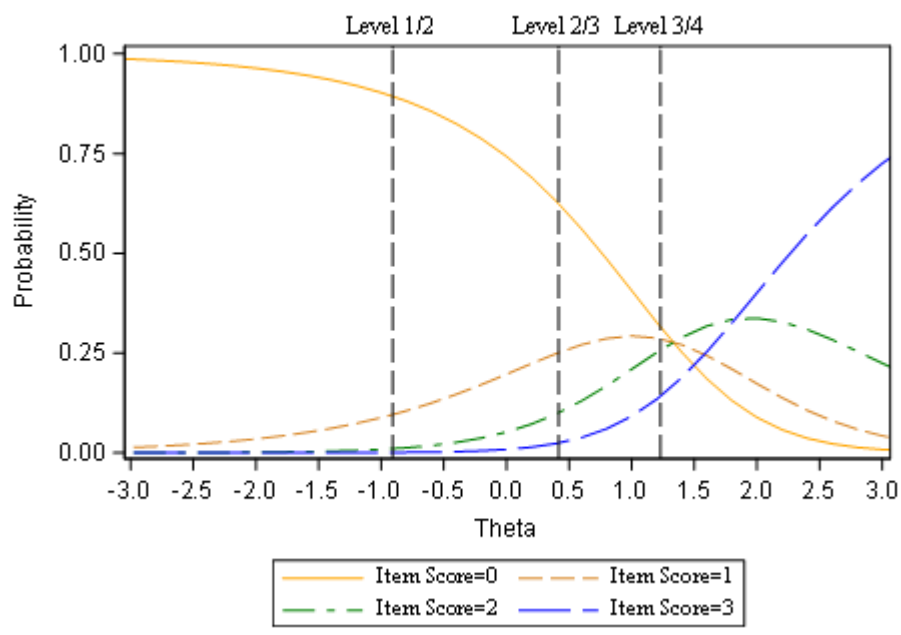


Figure 6.2.11. ICC Plot for Grade 8 Constructed-Response Item 2

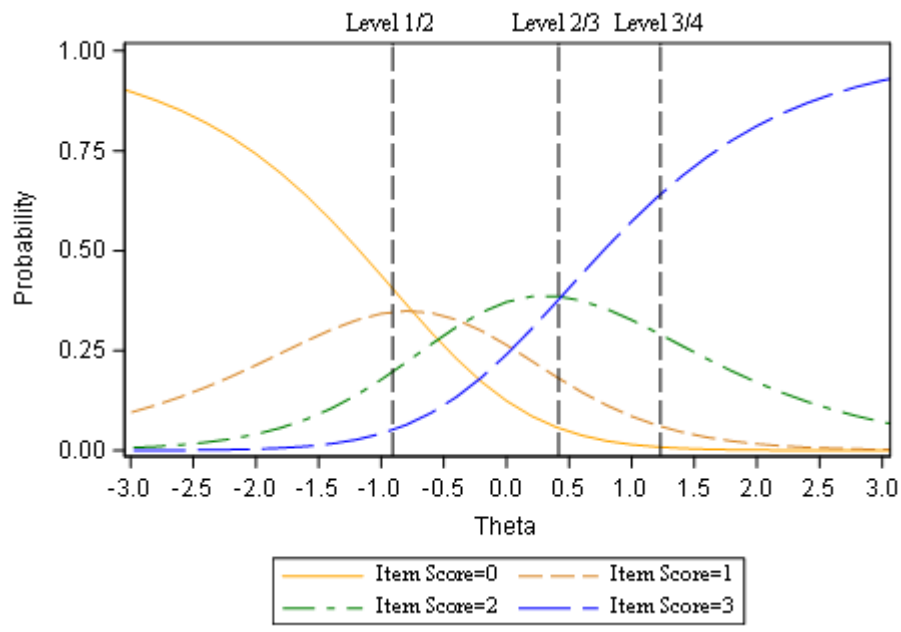


Figure 6.2.12. ICC Plot for Grade 8 Constructed-Response Item 3

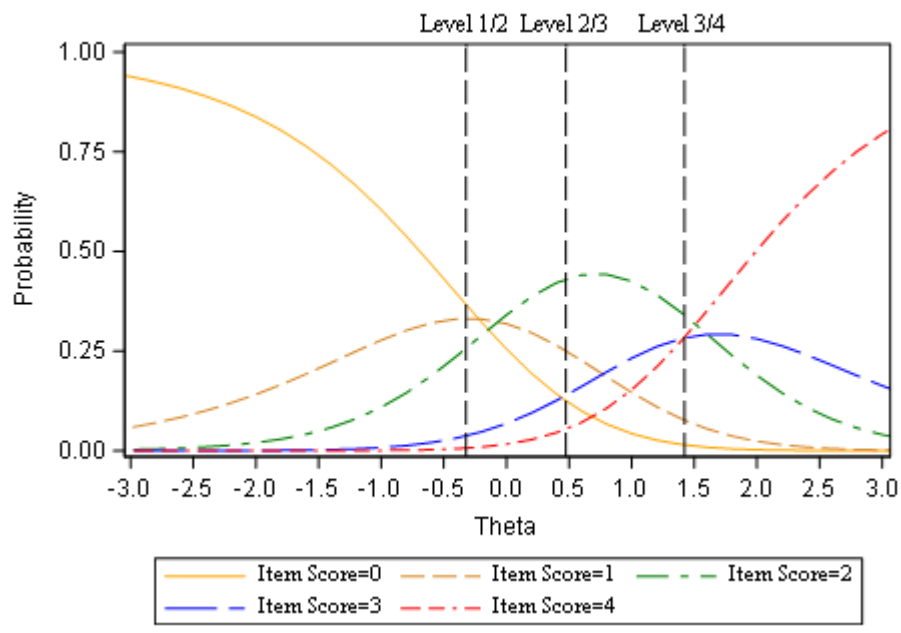


Figure 6.2.13. ICC Plot for Grade 11 Constructed-Response Item 1

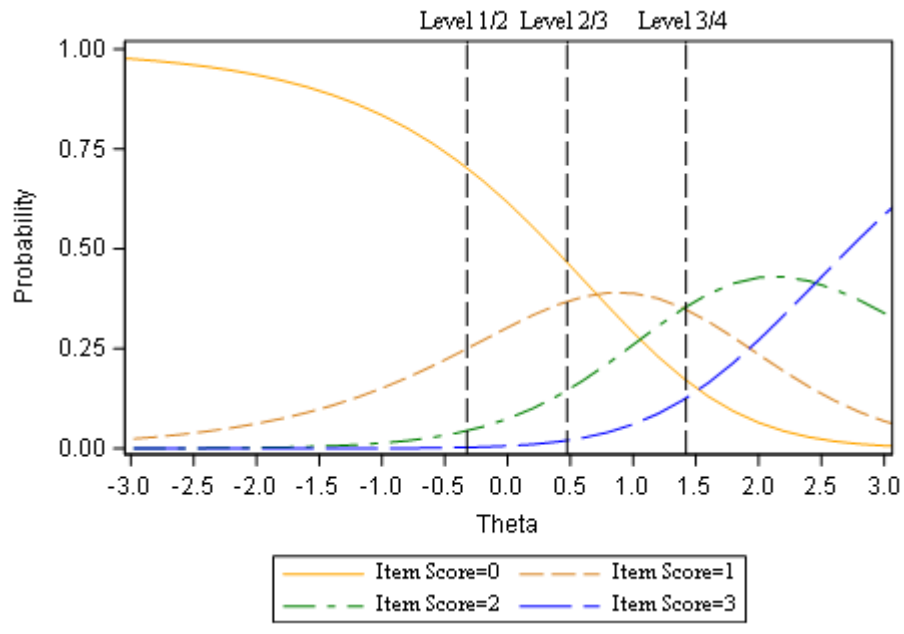


Figure 6.2.14. ICC Plot for Grade 11 Constructed-Response Item 2

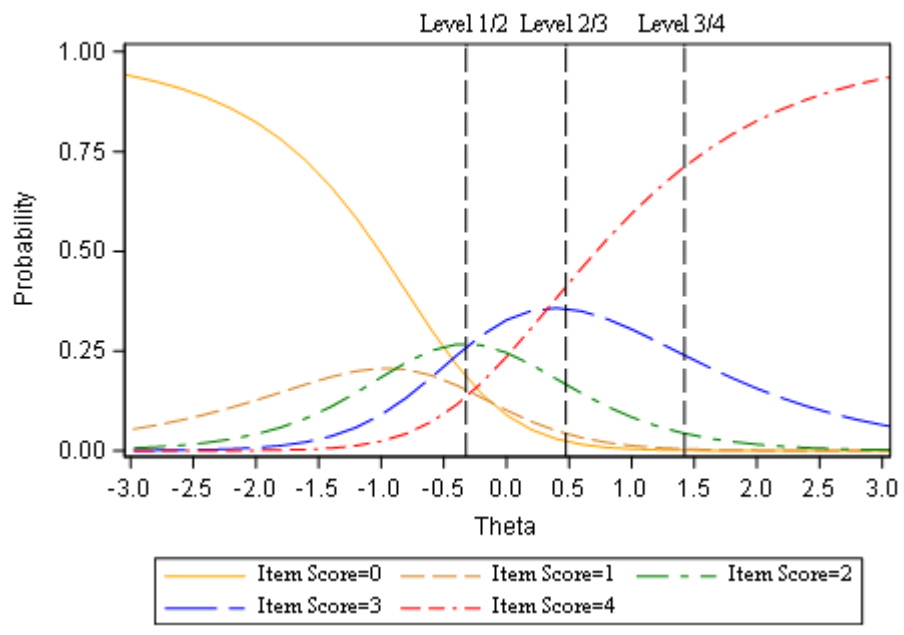


Figure 6.2.15. ICC Plot for Grade 11 Constructed-Response Item 3

6.3 Student Test Performance

Descriptive statistics for scale scores and performance-level distributions by form are presented in the following sections. For all the forms, scale scores have a range of 100 to 300. The Level 2/3 cut score is 200 at each grade level. Students who score at Level 3 or above are deemed proficient according to the results of the 2019 NJSLA–S Standard Setting. The Level 1/2 and 3/4 cut score ranges are more complex, and details regarding them can be found in Section 7.1 of this technical report. It should be noted that no scale score comparisons should be made across grade levels.

6.3.1 Scale Score Distribution by Form

Descriptive statistics for scale scores and percentage distributions of students' performance levels by form are summarized by grade in Table 6.3.1. The forms include CBT, PBT, TTS, SP, SP TTS, and HR.

Table 6.3.1: Descriptive Statistics of Students' Test Performance by Form

Grade	Form	N	Mean	SD	Min	Max	%L1	%L2	%L3	%L4
5	CBT	75,574	172.30	44.79	100	300	32.16	37.09	24.04	6.72
5	PBT	143	129.87	37.23	100	281	77.62	14.69	5.59	2.10
5	TTS	17,960	145.83	41.55	100	300	58.35	27.80	11.61	2.24
5	SP	1,473	125.19	27.16	100	238	80.86	16.97	2.17	0.00
5	SP TTS	983	126.12	26.97	100	234	79.86	18.31	1.83	0.00
5	HR	201	133.24	30.04	100	251	72.64	24.38	1.99	1.00
8	CBT	84,298	165.86	36.52	100	300	35.20	43.82	15.93	5.05
8	PBT	72	141.71	25.90	102	215	68.06	27.78	4.17	0.00
8	TTS	14,433	145.07	31.55	100	300	60.93	31.37	6.52	1.18
8	SP	1,876	134.29	20.64	100	227	77.03	22.44	0.53	0.00
8	SP TTS	729	134.46	20.87	100	215	75.31	24.01	0.69	0.00
8	HR	47	128.49	21.75	100	200	87.23	10.64	2.13	0.00
11	CBT	84,251	173.44	49.84	100	300	41.42	26.99	22.70	8.89
11	PBT	242	140.43	36.29	100	281	69.42	23.55	5.37	1.65
11	TTS	7,732	153.92	44.34	100	300	59.01	22.96	14.42	3.61
11	SP	1,474	127.59	24.65	100	246	87.25	11.53	1.22	0.00
11	SP TTS	273	128.11	25.32	100	211	83.52	14.29	2.20	0.00
11	HR	26	128.42	26.93	100	208	92.31	3.85	3.85	0.00

Note. CBT: Computer-Based Test; PBT: Paper-Based Test; TTS: Text-to-Speech; SP: Spanish; SP TTS: Spanish Text-to-Speech; HR: Human-Reader

6.3.2 Scale Score Distributions by Demographic Group

Descriptive statistics of scale scores and percentage distributions of students' test performance by demographic groups can be found on the [New Jersey Statewide Assessment Reports webpage](#). Scale score cumulative frequency distributions are attached as Appendix G of this technical report.

6.3.3 Subscore Proficiency Classification

There are no scale scores for the various subscores. As is explained in Section 7.3, student performance on the subscore categories was classified into three levels: Below, Near/Met, and Above Expectations. Appendix K presents the percentages of students who were placed in the three subscore proficiency classifications. The data are disaggregated by form type, gender, ethnicity, and other demographic variables for the content domains and practices at each grade level.

At grade 5, among the three content domains, Earth and Space Science saw the highest percentage (56.35%) of students classified as Below Expectation; among the three practices, Critiquing saw the highest percentage (58.54%) of students classified as Below Expectation. At grade 8, the highest percentages of students classified as Below Expectations were observed for Physical Science (66.00%) and Critiquing (66.54%). At grade 11, the Below Expectations percentages varied from 49.40% for Earth and Space Science to 52.27% for Physical Science, while the percentage of students classified as Below Expectation were relatively consistent across the practices with 52.03% for Sensemaking to 53.00% for Investigating. Overall, across content domains and practices at each grade level, there was a noticeable difference between the percentage of students classified as Below Expectation within the demographic groups. The exception to this was gender, where the percentage of students classified as Below Expectation was similar for both male and female students across grades.

PART 7: EQUATING AND SCALING

Standard 5.12 states that “A clear rationale and supporting evidence should be provided for any claim that scale scores earned on alternative forms of a test may be used interchangeably” (p. 105). Equating is the process that allows for the interchangeability of test scores from year to year and within year test forms (Kolen & Brennan, 2004).

7.1 Summary of Equating and Scaling Procedures

The NJSLA–S uses an internal anchor item equating design, in which an anchor item set is a subset of the operational items, as well as the partial credit model (PCM; Masters, 1982; discussed in Section 6.2 of this report) for maintaining the scale. In addition to students who took the regular NJSLA–S test, the equating samples include students who took the accommodated forms, as New Jersey Department of Education policy requires the same score tables be used for all accommodated test forms. The equating samples are demographically representative of the population of NJSLA–S test takers in 2023 in terms of demographic distributions of gender, ethnicity, and socioeconomic status. Before the equating samples were created, additional analyses were conducted to guarantee the appropriateness of items for use in generating student test scores. A preliminary item analysis was conducted on multiple-choice items to validate the keys. Item scores of all the multiple-choice items were determined to have been correct.

For the NJSLA–S, equating and item calibrations include three phases of psychometric analyses. A free (unconstrained) calibration was first conducted under the PCM using the equating sample for each grade. The free calibration run converged successfully for the 2023 NJSLA–S equating sample at each grade. The free-run item parameter estimates were employed for the second phase of equating analyses: anchor item stability evaluations.

The NJSLA–S assessment uses two methods for anchor item stability evaluation. The first is the displacement evaluation method, which investigates the deviation (i.e., displacement) in the IRT item difficulty parameter estimates (i.e., Rasch B) of the anchored parameter values in comparison to the free-run parameter estimates. An item is flagged using the 0.3-logit absolute difference criterion (Miller et al., 2004). The second anchor item stability evaluation method is the Delta Plot method (Angoff & Ford, 1973), which compares the item means (using *p-values* for dichotomous items and adjusted mean scores for polytomous items) obtained from the current year equating samples to those obtained from previous test administration(s). The item means are converted to Delta values, which in turn are used to compute a best-fitted line. For all the anchor items under investigation, their perpendicular distances (PD) in Delta scores to the best-fitted line are computed and evaluated. An item is flagged if it is more than two standard deviation units of the PDs away from the best-fitted line.

For the NJSLA–S, among the items flagged by both methods in a round of the anchor item evaluation process, the one with the largest absolute displacement value is dropped from the anchor set. The anchor item evaluation analysis is iteratively conducted until no items are flagged by both evaluation methods or the number of dropped items reaches 20% of the original set of anchor items. For the 2023 NJSLA–S equating, for both grades 5 and 8,

one anchor item was dropped from the respective anchor item sets. At grade 11, the anchor evaluation processes were terminated at round three and yielded two (8.6%) dropped items.

After the completion of anchor evaluations, Winsteps 3.92.1 was used to calibrate the Rasch values (i.e., the B-parameter estimates and the step parameter estimates) of all operational items to the base theta scale (i.e., the base scale). This was done by constraining the remaining anchor items to their Rasch values from the previous administrations or item bank that were already calibrated to the base scale. The results of the fixed Winsteps calibration run are used to develop the raw-to-theta-to-scale conversion tables for scoring. The development of scaling constants (i.e., intercept and slope) for converting theta scores to NJSLA–S scale scores is discussed in the following paragraph.

The NJSLA–S was scaled via a linear transformation that converted the IRT student ability estimates into scale scores. New Jersey has historically used a 100–300 scale for statewide assessments; in the past, with only three performance levels, scale scores of 200 and 250 represented proficient and advanced proficient performance, respectively. The NJSLA–S scaling procedure maintained the 100–300 scale; however, the scaling was slightly more complex due to the introduction of a third cut score (i.e., four performance levels). Policy decisions based on minimizing the number of students receiving the lowest obtainable scale score (LOSS) and the highest obtainable scale score (HOSS) necessitated that, at grades 5 and 8, the Level 1/2 and the Level 2/3 cut scores be anchored during the linear transformation and at grade 11, the Level 2/3 and Level 3/4 cut scores be anchored. The linear transformation is described in detail below.

At all grades, a scale score of 200 still represents the proficient cut point (i.e., Level 2/3 cut). Students who score below 200 are placed in either Level 1 or Level 2. They are classified as below proficient and display minimal or partial understanding of the NJSLA–S. Students who score 200 or above are classified as either Level 3 or Level 4. Their performance is deemed proficient, and it represents an appropriate or exemplary understanding of the NJSLA–S.

The scale score ranges are reflected in Table 7.1.1 below. The scale scores representing the cut score differentiating Level 1 from Level 2 and differentiating Level 3 from Level 4 vary depending on each grade. At grades 5 and 8 the Level 1/2 cut score was anchored at a scale score of 150, whereas at grade 11 the scale score cut was 158. The Level 3/4 cut score was anchored at 250 for grade 11, while it was 243 for grade 5 and 231 for grade 8.

Table 7.1.1: Scale Score Ranges for Proficiency Levels by Grade

Grade	Level 1	Level 2	Level 3	Level 4
5	100–149	150–199	200–242	243–300
8	100–149	150–199	200–230	231–300
11	100–157	158–199	200–249	250–300

To produce the scale score ranges above, linear transformations were applied to theta (θ) estimates and scale scores. The following formula, adapted from Kolen and Brennan (2004, p. 337), was used to obtain the slopes and intercepts for the transformation functions:

$$sc(y) = \left[\frac{sc(y_2) - sc(y_1)}{\theta_2 - \theta_1} \right] y + \left\{ sc(y_1) - \left[\frac{sc(y_2) - sc(y_1)}{\theta_2 - \theta_1} \right] \theta_1 \right\}, \quad \text{Equation 7.1}$$

where θ_1 and θ_2 are student ability estimates that correspond to the approved cut score points, and $sc(y_1)$ and $sc(y_2)$ are scale score points corresponding to θ_1 and θ_2 , respectively. The resulting slopes and intercepts of the linear transformations at each grade level are shown in Table 7.1.2.

Table 7.1.2: Slope and Intercept of Theta-to-Scale Score Transformations and Performance-Level Cut Scores by Grade

Grade	Level 1/2 Cut		Level 2/3 Cut		Level 3/4 Cut		Slope	Intercept
	Theta	Scale Score	Theta	Scale Score	Theta	Scale Score		
5	-0.2739	150	0.9035	200	1.9243	243	42.4639	161.6317
8	-0.9077	150	0.4156	200	1.2306	231	37.7800	184.2960
11	-0.3230	158	0.4751	200	1.4217	250	52.8189	174.9036

The following sections specify how these slopes and intercepts were used to generate the scale scores at each grade level. The complete raw-to-scale score conversion tables can be found in Appendix I.

7.1.1 Rounding Rules

NJDOE policy requires that scaled scores below 100 are rounded up to 100, and that scaled scores above 300 are rounded down to 300. Additional rules of adjustments to scale score tables required for scaling are as follows:

- If a raw score maps to an unrounded scale score that is greater than 199.499 and less than or equal to 200.000, it will serve as the proficient (Level 2/3) cut score. Otherwise, the highest raw score that maps to a scale score less than or equal to 199.499 will serve as the cut score. The selected cut score will be assigned a value of exactly 200.
- If a raw score maps to an unrounded scale score that is greater than 249.499 and less than or equal to 250.000 for Level 4 at grade 11, it will serve as the advanced (Level 3/4) cut score. Otherwise, the highest raw score that maps to a scale score less than or equal to 249.499 will serve as the cut score. The selected cut score will be assigned a value of exactly 250. The same rounding procedures apply to the cut scores for Levels 1/2 and 2/3 for grade 11 as well as Levels 1/2, 2/3, and 3/4 for grades 5 and 8.
- If two unrounded scale scores fall in the range of greater than 199.499 and less than 200.000 (i.e., the Level 2/3 cut scores for all the grades), the lower of these two scores would become the cut score. The same rounding procedures apply to the cut scores of Levels 2 and 4 for all the grades.
- When the implementation of the above rounding rules results in two raw scores mapping to an equivalent rounded scale score, the scale score associated with the higher of the two raw scores will be adjusted upward by one (1) scale score.

7.2 Accommodative Form Equivalence

NJDOE has traditionally used the same score tables for their accommodative forms as for their traditional operational test forms, a decision that is predicated on several assumptions. These were checked for all accommodative forms by either content experts versed in universal design or, in the case of the braille and Spanish forms, external reviewers.

First, it must be assumed that the latent trait measured by the accommodative forms is the same as the latent trait measured by the operational test forms. Given that the same items measuring the same skills and abilities were used across the tests, it seems reasonable to assume that changes to item format or item presentation would not greatly change the overall latent trait or construct measured by each assessment form. Moreover, all items were written based on the principles of universal design as described in Section 3.4.

A second assumption is that item parameters across the test forms within each content cluster are identical. This, of course, is a potentially tenuous assumption considering the different item formats across the test forms. However, NJDOE's policy requiring that the same score tables be used for all accommodative test forms rendered this assumption necessary. Thus, all the accommodative forms for the NJSLA–S were assumed to be equivalent. If an operational item is unable to be properly adapted to a specific accommodative form, then the assumption of equivalence is violated, and a special equating is required. For any given administration year, if this assumption is violated for any accommodative form(s), the special equating procedure is described in the following section. For the 2023 NJSLA–S administration, no special equating was needed.

7.2.1 Special Equating

In the event of errors during the test construction process that led to the removal of item(s) from the test, special equating was conducted to re-calculate score tables so that the students who received those forms were placed onto a scale equivalent to that underlying the other CBT forms. The following steps were taken to ensure the special equatings and CBT forms were on the same scale.

1. **Anchored item calibration.** The inequivalent items were removed prior to the special equating calibrations, and the item parameters and steps of the accommodated test items were fixed with the estimates resulting from the corresponding regular test items.
2. **Theta to the scale score metric transformation.** Because the theta values obtained from the anchored calibration and those obtained from the regular test score calibration are on the same metric, the transformation functions applied to the regular test scores could likewise be applied to the accommodated test scores.
3. **Raw-to-scale score tables for each special equating.** The rounding rules described in Section 7.1.1 were applied to the transformed scale scores, resulting in a separate raw-to-scale score table for each special equating that could be interpreted exactly the same as the other operational forms.

7.3 Subscore Performance Levels

The NJSLA–S assessments reports student performance in three content domains/disciplinary core ideas (DCI) including Earth and Space, Life, and Physical. The NJSLA–S also reports performance in three scientific and engineering practices (SEP) including Investigating, Sensemaking, and Critiquing. In each DCI and SEP, subscore performances are classified as “Below,” (Level 1) “Near/Met,” (Level 2), or “Above” (Level 3) expectations. The subscores for these six reporting categories are themselves described in Part 1 of this Technical Report. This section details the processes used to create the NJSLA–S subscore performance-level classifications.

For a given DCI or SEP at a grade level, the process for classifying NJSLA–S subscore performance first involved creating a subscore table. The subscore table was generated through a Winsteps fixed-parameter calibration run with the item parameter estimates of each item associated with the given DCI or SEP held constant (i.e., anchored) at the values obtained from operational item calibration results. A subscore table consisted of raw subscores, their associated thetas (θ), and the conditional standard errors of measurement (CSEM). The subscore performance level classifications were based on the extent to which the subscore theta values within the subscore score tables were statistically significantly above or below the overall scale’s Level 3 (proficient) theta cut score (denoted by θ^*). Based on the subscore table, the CSEM associated with θ^* denoted by CSEM* was estimated for subscore performance classification analyses. The “1.5 standard error rule” (Smarter Balanced Assessment Consortium, 2018) was then used to generate the subscore performance-level classifications as follows:

- A raw subscore is classified as “Above” if its associated θ is at or above ($\theta^* + 1.5 \text{ CSEM}^*$) units.
- A raw subscore is classified as “Below” if its associated θ is below ($\theta^* - 1.5 \text{ CSEM}^*$) units.
- A raw subscore is classified as “Near/Met” if its associated θ does not meet the definition of Above or Below.

The subscore score tables for each combination of grade and reporting category are presented in Appendix J.

PART 8: RELIABILITY

Test reliability refers to the consistency of test scores. Ultimately, valid interpretations of test scores are dependent upon those scores being reliable. *Standard 2.0* states that “[a]ppropriate evidence of reliability/precision should be provided for the interpretation for each intended score use” (p. 42). Examples of appropriate evidence include reliability coefficients, conditional standard errors of measurement (CSEM), test information functions, and decision consistency measures, amongst others. The following sections detail evidence supporting the reliability of the NJSLA–S test scores and subscores.

8.1 Classical Test Theory Reliability Estimates

This section describes the Classical Test Theory (CTT) reliability estimates calculated for the NJSLA–S. Section 8.1.1 describes the concept of reliability in the CTT framework, and Section 8.1.2 displays the reliability analysis results based on CTT.

8.1.1 Reliability and Measurement Error

Under the assumptions of CTT any observed measurement—such as a test score, X —is defined as a composite of true score, T , and its associated error:

$$X = T + \text{error} \quad \text{Equation 8.1}$$

Errors in measurement can result from any of a multitude of factors, including environmental factors (e.g., testing conditions) and examinee factors (e.g., fatigue, stress). CTT provides a means for this quantification of examinee inconsistency (i.e., measurement error). Student test scores are reliable when measurement error is minimized. Increasing reliability by minimizing measurement error is an important goal in the construction of any test.

Estimating the size of the measurement error associated with the true score is the key to estimating reliability. The definitions or assumptions in CTT lead to several important properties. For example, it can be demonstrated that observed score variance (σ_X^2) equals the sum of true score variance (σ_T^2) and error variance (σ_e^2) or mathematically,

$$\sigma_X^2 = \sigma_T^2 + \sigma_e^2 \quad \text{Equation 8.2}$$

The relationships among the variance terms (i.e., $\sigma_X^2, \sigma_T^2, \sigma_e^2$) are critical to a more thorough understanding of important CTT concepts, including reliability and the standard error of measurement. Under CTT, reliability (ρ_{X_1, X_2}) is defined as the correlation between observed scores (X_1, X_2) on parallel forms, which is equal to true score variance (σ_T^2) divided by observed score variance (σ_X^2):

$$\rho_{X_1 X_2} = \frac{\sigma_T^2}{\sigma_X^2} \quad \text{Equation 8.3}$$

With just a few algebraic steps, the CTT definition of the standard error of measurement (SEM, σ_e) can be shown as:

$$\sigma_e = \sigma_X \sqrt{1 - \rho_{X_1 X_2}} \quad \text{Equation 8.4}$$

Although the concepts of reliability and SEM are relatively straightforward, issues underlying the estimation of reliability are not. Reliability can be estimated via the correlation of scores on parallel forms or from test-retest data, or it can be estimated from a single test administration using any one of a variety of techniques (e.g., Brown, 1910; Cronbach, 1951; Kuder & Richardson, 1937).

For NJSLA–S, consistency of individual student performance was estimated using Cronbach’s (1951) coefficient alpha. Coefficient alpha is conceptualized as the proportion of total raw score variance that may be attributed to a student’s true score variance. Ideally, more score variance should be attributable to true test scores than to measurement errors. Coefficient alpha was estimated using the following formula:

$$\alpha = \frac{n}{n-1} \left[1 - \frac{\sum_{i=1}^n \sigma_i^2}{\sigma_X^2} \right], \quad \text{Equation 8.5}$$

where n is the number of items on the test, σ_i^2 is the item score variance of item i , and σ_X^2 is the variance of the observed total test score. Accordingly, SEMs were estimated and calculated using the following formula:

$$SEM = S_X \sqrt{1 - \alpha}, \quad \text{Equation 8.6}$$

where S_X is the standard deviation of observed total scores. For the NJSLA–S assessments, separate analyses were performed for each grade level. Scores from all item types were used in the computations.

8.1.2 Raw Score Internal Consistency

In order to accommodate the state’s diverse testing population, the NJSLA–S was delivered in multiple formats. The most used forms were the traditional online (CBT), the Text-to-Speech (TTS), the Spanish (SP), the paper-based test (PBT), and the Human Reader (HR). Reliability measures decrease when the students taking a given test form are more homogeneous in their test performance.

Table 8.1.1 displays the coefficient alpha and SEM for each form by grade. Overall, the reliability coefficients at each grade level indicate that students’ raw scores were reliable. The results at grade 5 stand out as particularly exceptional given that the grade 5 test is shorter than either the grade 8 or 11 tests. The grade 5 reliability coefficients ranged from .82 to .92. The most likely reason for the better results at grade 5, despite it being a shorter test, is that the grade 5 items were closer to the ability levels of the grade 5 students, thereby increasing the variance among test scores. At grade 8, where the distribution of test scores was heavily skewed toward the low end of the ability spectrum, reliability ranged from .73 to .92. The

relatively low-reliability measures for the Spanish, Spanish TTS, and Human Reader forms are due to those populations doing poorly on the test, which limits the amounts of variance in test scores. The grade 11 alpha coefficients ranged from .73 for the Spanish form to .93 for the CBT form. As shown in Table 8.1.1, the Grade 11 Spanish, Spanish TTS, and Human Reader form test takers did poorly on the test. This would result in less variance in test scores for these groups, which may explain the lower reliability estimates for these forms compared to the other test forms.

Table 8.1.1: Coefficient Alpha and SEM by Form

Grade	Form	N	Mean Raw Score	SD	Alpha	SEM
5	CBT	75,574	25.79	13.01	0.92	3.68
5	PBT	143	13.94	10.32	0.91	3.09
5	TTS	17,960	18.28	11.78	0.92	3.41
5	SP	1,473	12.48	7.24	0.83	3.00
5	SP TTS	983	12.73	7.18	0.82	3.03
5	HR	201	14.51	8.13	0.84	3.24
8	CBT	84,298	24.57	13.39	0.92	3.76
8	PBT	72	15.85	8.6	0.84	3.42
8	TTS	14,433	17.28	10.78	0.90	3.45
8	SP	1,876	13.34	6.14	0.73	3.18
8	SP TTS	729	13.38	6.18	0.73	3.19
8	HR	47	11.89	6.37	0.77	3.09
11	CBT	84,251	29.63	14.69	0.93	4.00
11	PBT	242	19.93	10.36	0.87	3.71
11	TTS	7,732	23.95	12.87	0.91	3.86
11	SP	1,474	16.36	6.66	0.73	3.48
11	SP TTS	273	16.51	6.82	0.74	3.47
11	HR	26	16.73	7.31	0.76	3.57

Note. CBT: Computer-Based Test; PBT: Paper-Based Test; TTS: Text-to-Speech; SP: Spanish; SP TTS: Spanish Text-to-Speech; HR: Human-Reader

Table 8.1.2 summarizes the coefficient alpha and SEM of raw scores of the six reporting categories by grade. In general, longer tests yield higher reliability coefficient estimates than shorter tests (Traub & Rowley, 1991). Thus, reporting categories such as Critiquing and Investigating, which had fewer items, tended to have lower reliability measures. For practice, the lowest subscore reliability of .75 was for Investigating at grades 5 and 8. For content domains, the lowest subscore reliability of .75 was for Physical Science at grade 5, which had the fewest items among the content domains for grade 5.

Table 8.1.2: Coefficient Alpha and SEM by Reporting Category

Grade	Reporting Category	Total # Items	#MC Items	#TE Items	#CR Items	Max Points	Alpha	SEM
5	Total	51	14	34	3	60	0.92	3.63
5	Earth and Space	18	5	12	1	21	0.80	2.19
5	Life	17	2	14	1	20	0.84	2.09
5	Physical	16	7	18	1	19	0.75	1.99
5	Sensemaking	17	1	16	0	17	0.85	1.74
5	Critiquing	18	6	9	3	16	0.80	2.65
5	Investigating	19	7	9	0	27	0.75	1.74
8	Total	65	44	18	3	72	0.92	3.72
8	Earth and Space	20	14	5	1	22	0.80	1.96
8	Life	24	15	8	1	26	0.79	2.24
8	Physical	21	15	5	1	24	0.80	2.23
8	Sensemaking	23	16	6	1	26	0.81	2.33
8	Critiquing	22	16	5	1	24	0.82	2.18
8	Investigating	20	12	7	1	22	0.75	1.90
11	Total	69	30	36	3	77	0.93	3.98
11	Earth and Space	20	8	11	1	23	0.77	2.37
11	Life	23	10	12	1	26	0.82	2.24
11	Physical	26	12	13	1	28	0.82	2.29
11	Sensemaking	24	10	13	1	27	0.83	2.31
11	Critiquing	22	7	14	1	25	0.79	2.42
11	Investigating	23	13	9	1	25	0.79	2.17

Table 8.1.3 shows the coefficient alpha and SEMs by demographic group. These calculations are based on the entire test. In general, the coefficient alphas are consistently high among the various demographic groups. At grade 5, the lowest value was .83, for English learner (EL) students, which is still very strong. At grade 8, the coefficient alphas hovered close to .90 except for the English learners ($\alpha_{EL-Yes} = .75$). The pattern for grade 11 was the same as for grade 8. The coefficient alpha values for all groups at grade 11 were above .89 except for the English learners ($\alpha_{EL-Yes} = .76$).

Table 8.1.3: Coefficient Alpha and SEM by Demographic Group

Grade	Group	N	Mean	SD	Alpha	SEM
5	NJSLA–S	96,392	24.01	13.14	0.92	3.63
5	Male	49,082	24.48	13.53	0.93	3.61
5	Female	47,299	23.52	12.70	0.92	3.64
5	Am. Indian	169	24.89	13.78	0.93	3.66
5	Asian	10,765	34.47	12.54	0.91	3.71
5	Black	14,028	17.55	10.73	0.90	3.38
5	Hispanic	31,700	18.75	11.07	0.90	3.45
5	Pacific Islander	179	25.01	13.20	0.92	3.67
5	White	36,375	27.65	12.46	0.91	3.71
5	EL–Yes	9,158	12.50	7.41	0.83	3.04
5	EL–No	87,234	25.22	13.02	0.92	3.66
5	EconDis–Yes	36,109	17.64	10.50	0.90	3.39
5	EconDis–No	60,283	27.83	13.08	0.92	3.71
5	SWD–Yes	20,003	17.56	11.66	0.92	3.35
5	SWD–No	76,389	25.70	12.98	0.92	3.68
8	NJSLA–S	101,478	23.23	13.25	0.92	3.72
8	Male	52,212	23.53	13.78	0.93	3.71
8	Female	49,201	22.90	12.66	0.91	3.72
8	Am. Indian	154	22.28	12.91	0.92	3.67
8	Asian	10,718	34.76	13.89	0.92	3.92
8	Black	14,998	16.74	9.98	0.88	3.42
8	Hispanic	32,921	18.07	10.38	0.89	3.50
8	Pacific Islander	206	26.08	12.94	0.91	3.83
8	White	39,768	26.61	12.96	0.91	3.82
8	EL–Yes	7,151	12.55	6.22	0.75	3.13
8	EL–No	94,327	24.04	13.30	0.92	3.75
8	EconDis–Yes	35,709	17.26	9.94	0.88	3.46
8	EconDis–No	65,769	26.47	13.69	0.92	3.81
8	SWD–Yes	20,520	17.16	11.08	0.90	3.42
8	SWD–No	80,958	24.77	13.32	0.92	3.77
11	NJSLA–S	94,023	28.88	14.62	0.93	3.98
11	Male	47,959	28.69	15.25	0.93	3.95
11	Female	45,924	29.06	13.93	0.92	4.00
11	Am. Indian	141	26.60	14.16	0.92	3.93
11	Asian	10,003	40.77	15.01	0.93	4.04
11	Black	12,731	22.19	11.50	0.89	3.83
11	Hispanic	28,687	23.29	11.88	0.89	3.86

Grade	Group	N	Mean	SD	Alpha	SEM
11	Pacific Islander	313	29.97	13.84	0.92	4.00
11	White	40,005	31.90	14.39	0.92	4.03
11	EL–Yes	5,290	16.18	7.10	0.76	3.48
11	EL–No	88,733	29.64	14.61	0.93	4.00
11	EconDis–Yes	28,095	22.91	11.74	0.89	3.85
11	EconDis–No	65,928	31.43	14.98	0.93	4.02
11	SWD–Yes	18,600	22.54	12.95	0.91	3.80
11	SWD–No	75,423	30.45	14.59	0.92	4.01

Table 8.1.4 displays coefficient alpha and SEM by the three main item types: multiple-choice (MC), technology-enhanced (TE), and constructed-response (CR). Those item types are more thoroughly described in Part 2 of this technical report. As would be expected, as the number of points associated with a specific item type increase, so does the corresponding coefficient alpha. More than half of the points available on each test were associated with TE item types; thus, it is not surprising that at each grade level, the TE items displayed alphas close to .9. The alphas associated with each grade level’s CR items were all close to .7, which is relatively strong given the limited number of points associated with them.

Table 8.1.4: Coefficient Alpha and SEM by Item Type

Grade	Item Type	Items	Points	Mean	S.D.	Alpha	SEM
5	MC	14	14	7.46	3.30	0.75	1.66
5	TE	34	34	13.11	7.72	0.90	2.48
5	CR	3	12	3.44	3.11	0.69	1.74
8	MC	18	18	6.98	3.52	0.71	1.89
8	TE	44	44	13.39	8.25	0.89	2.72
8	CR	3	3	2.86	2.52	0.68	1.43
11	MC	30	30	12.76	5.58	0.80	2.49
11	TE	36	36	12.12	7.32	0.89	2.43
11	CR	3	3	4.00	2.86	0.68	1.63

8.2 Item Response Theory Reliability

The reliability of the scale scores ascertained from the partial credit model (PCM; Masters, 1982) discussed in Section 6.2 was assessed in multiple ways. Test information functions (TIFs), conditional standard error measurements (CSEMs), and person-fit statistics were evaluated at each grade level. Overall, the 2023 NJSLA–S was reliable from the perspective of IRT and the PCM.

8.2.1 Test Information Functions

In IRT, the reliability of an assessment is conceptualized via the test information function (TIF, Hambleton & Swaminathan, 1985). Unlike coefficient alpha (Cronbach, 1951), the TIF is not

uniform across the entire range of test scores. Instead, the TIF can assess test reliability across the full range of scores. This is particularly important to a criterion-referenced test such as the NJSLA–S because it allows for the reliability of the assessment to be evaluated at the most important decision points (i.e., the Level 2/3 cut scores).

Psychometrically, under the IRT assumption of local independence, the TIF for a test is the summation of all the item information functions (IIF; Lord & Novick, 1968; Hambleton, 1989) as follows:

$$I(\theta) = \sum_{i=1}^N I_i(\theta) \quad \text{Equation 8.7}$$

where $I(\theta)$ is the amount of test information at an ability level of θ , $I_i(\theta)$ is the amount of information for item i at an ability level of θ , and N is the number of items on a test. It should be noted that the mathematical definition of the amount of item information depends on the IRT model employed. Under the partial credit model where the responses to item i are scored as the integers $0, 1, \dots, m_i$, the item information in item i is given by (Donoghue, 1994):

$$I_i(\theta) = \sum_{k=0}^{m_i} k^2 P_{ik}(\theta) - \left(\sum_{k=0}^{m_i} k P_{ik}(\theta) \right)^2 \quad \text{Equation 8.8}$$

where $P_{ik}(\theta)$ is the probability that an examinee of a given ability level θ will obtain a score of k on item i . With a few algebraic steps, the item information for a dichotomous item under the Rasch model is given by the following:

$$I_i(\theta) = P_i(\theta)(1 - P_i(\theta)) \quad \text{Equation 8.9}$$

Figures 8.2.1 to 8.2.3 illustrate, respectively, the TIFs for grades 5, 8, and 11 at person ability estimates ranging from -6 to $+6$. Within each figure, there are three vertical dash lines representing the test performance cut scores. More information at a specific ability level implies less measurement error. Ideally, the Level 2/3 cut score would occur at the peak of the information function where the most information and the least measurement error occur. Given the importance of making decisions at the Level 1/2 and 3/4 cut scores, the graph would also maintain ample information at those places along the scale.

The TIFs at each grade level were assessed primarily by whether they peaked close to the Level 2/3 cut score, and whether there was a precipitous drop in information at the Level 1/2 and 3/4 cut scores. At grade 5, the TIF peaked almost directly on the Level 2/3 cut scores, but there was a drop in information at the Level 3/4 cut. At grade 8, the TIF peaked close to the Level 2/3 cut scores, but there was a drop in information at the Level 1/2 cut. The grade 11 TIF peaked almost directly at the Level 2/3 cut score. However, there was a drop in information at the Level 3/4 cut for grade 11. Overall, the TIFs provide ample evidence that student ability estimates are reliable at the most important decision points. Nonetheless, both grades 5 and 8 would benefit from more information around the Level 1/2 cut score on future tests. In addition, grade 5 would benefit from more information around the level 3/4 cut score.

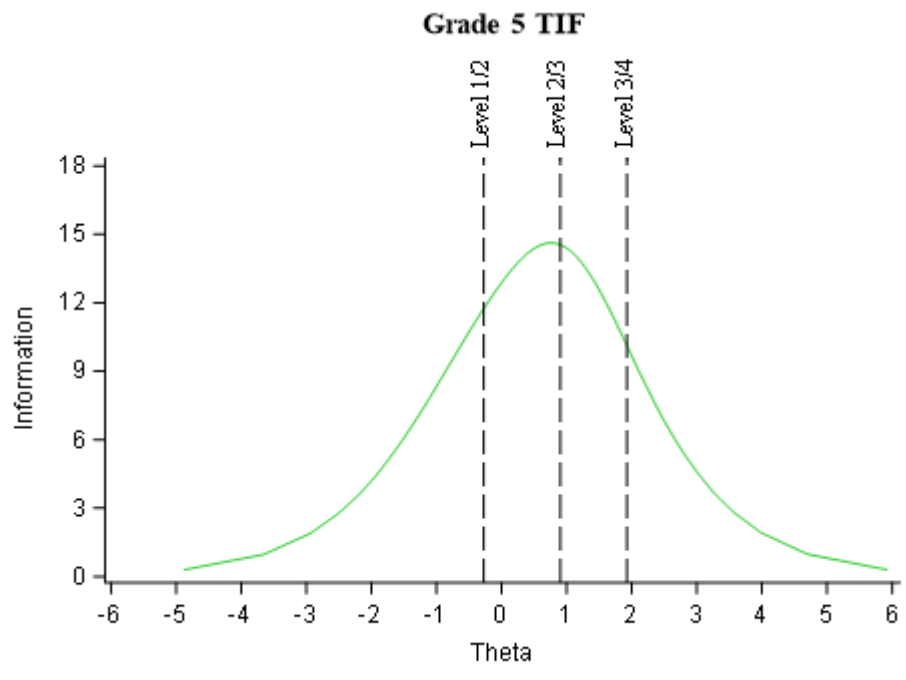


Figure 8.2.1. Grade 5 Test Information Function

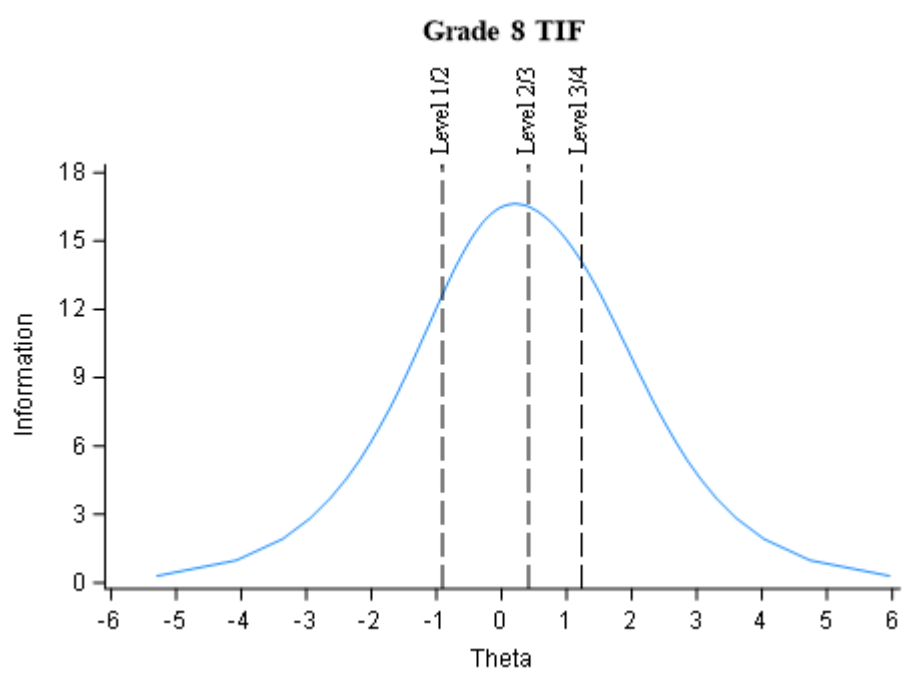


Figure 8.2.2. Grade 8 Test Information Function

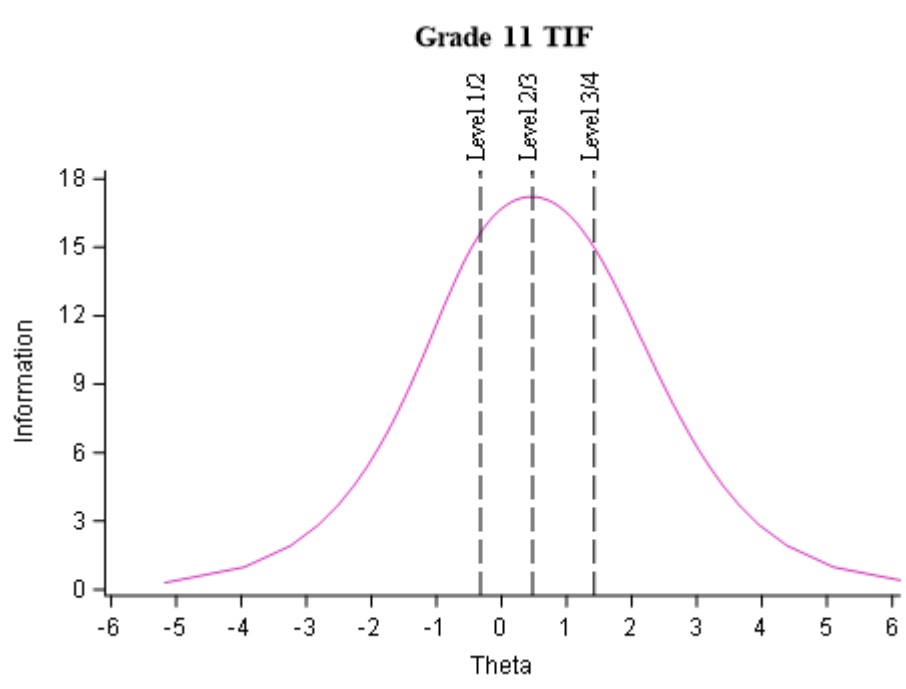


Figure 8.2.3. Grade 11 Test Information Function

8.2.2 Conditional Standard Error of Measurement

Under IRT, the conditional measurement error (CSEM) of an examinee’s estimated ability plays an important role in psychometric analyses. Mathematically, CSEMs are inversely related to the TIF, $I(\theta)$, and given by the following:

$$CSEM(\theta) = \frac{1}{\sqrt{I(\theta)}} \quad \text{Equation 8.10}$$

where $I(\theta)$ indicates the amount of test information (TIF) at an ability level of θ . TIF was discussed in Section 8.2.1 of this report. Given that the CSEMs are the inverse of the TIF, their interpretations are similar. If the amount of information at a given level of θ is large and hence the corresponding CSEM is small, it means an examinee whose true ability is at that level can be estimated with precision. That is, the estimates will be reasonably close to the true value. It should be noted that the TIF and CSEM do not depend on the distribution of examinees over the ability scale.

Figures 8.2.4 through 8.2.6 illustrate, respectively, the CSEMs for grades 5, 8, and 11 at ability estimates ranging from -6 to $+6$. As shown in these figures, the CSEMs around the three cut scores are around .25 for each grade level, indicating ability scores around the three cut scores are estimated with precision.

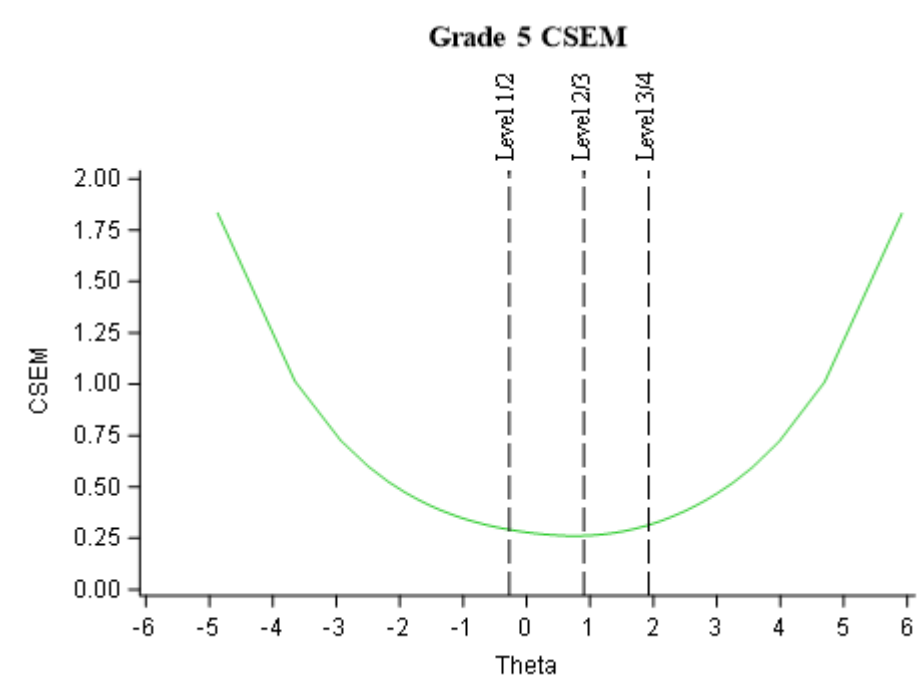


Figure 8.2.4. Grade 5 Conditional Standard Error of Measurement

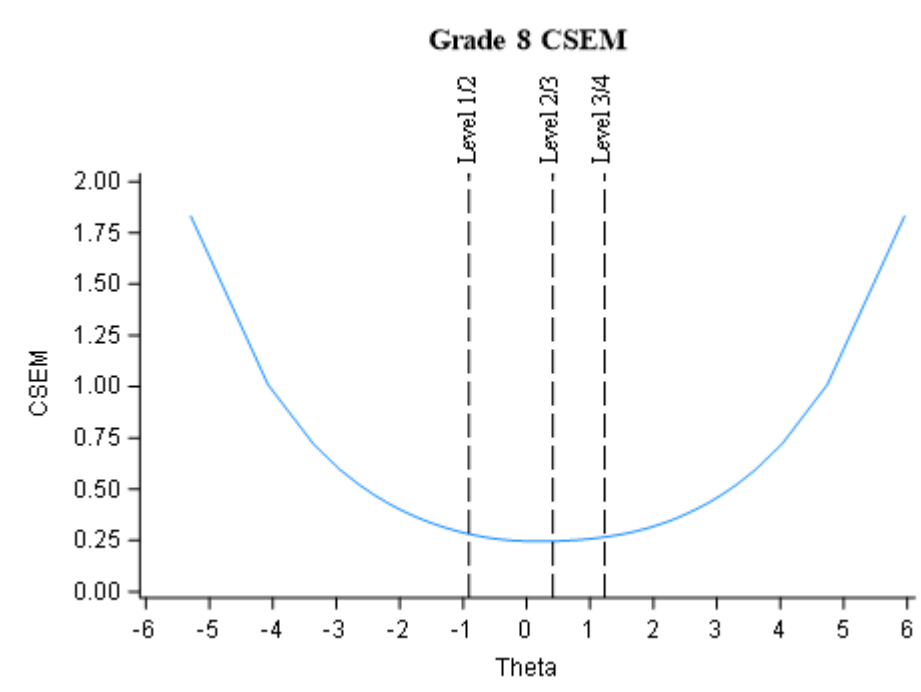


Figure 8.2.5. Grade 8 Conditional Standard Error of Measurement

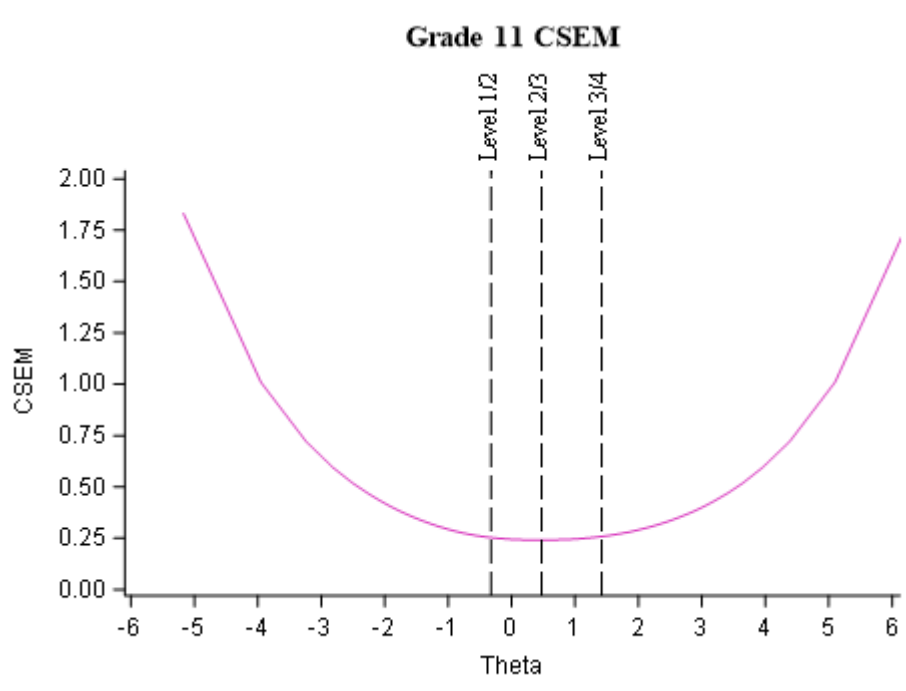


Figure 8.2.6. Grade 11 Conditional Standard Error of Measurement

8.2.3 Item Maps

An item map exhibits the distribution of person ability estimates, and the distribution of item difficulty parameter estimates along the latent scale (i.e., theta). Item maps are useful to compare the range and positions of the item difficulty distribution to those of the person ability measure distribution. Items that are targeted to the ability levels of the students taking the test will result in more reliable measures of student ability.

Figures 8.2.7 through 8.2.9 show the 2023 NJSLA–S item maps for grade levels 5, 8, and 11, respectively. Each item map figure is delineated into two panels, the top containing the item difficulty estimate distribution and the lower containing the ability (theta) distribution. As shown in the figures, at each grade level, NJSLA–S items were appropriately targeted to the student ability distribution. At grade 5, the item difficulty distributions peaked between the Level 1/2 and the Level 2/3 cut scores; the theta distributions peaked around the Level 1/2 cut score. At grade 5, there were few students above the Level 3/4 cut score and zero items along that part of the scale. The grade 8 item difficulty distribution was lacking items at the lower (easier) part of the scale in comparison to the student ability distribution. At grade 11 the item difficulty distribution peaked at the Level 2/3 cut and saw several items at the upper (harder) part of the scale, while the student ability distribution peaked near the Level 1/2 cut.

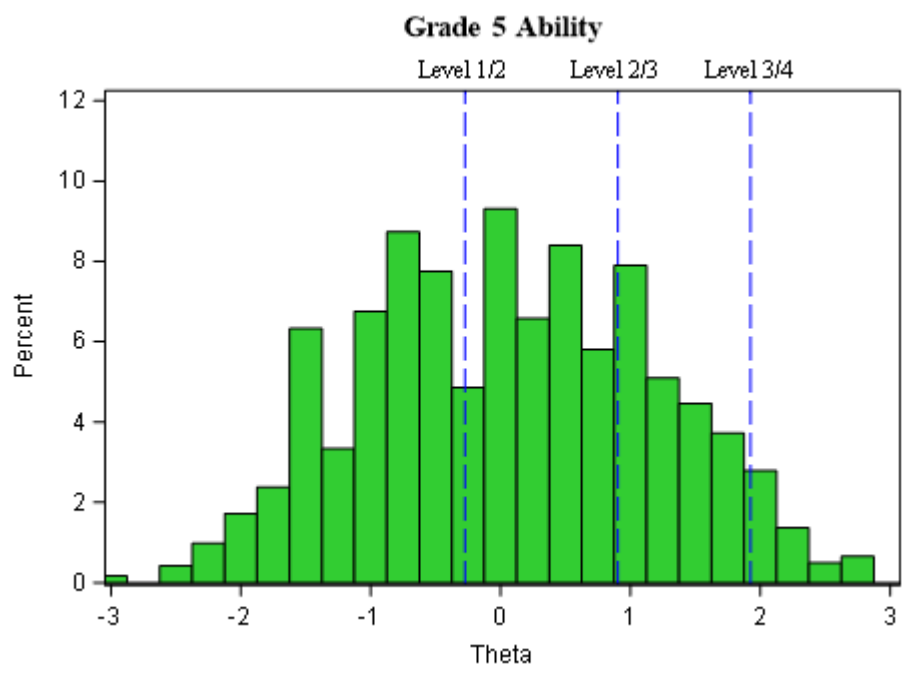
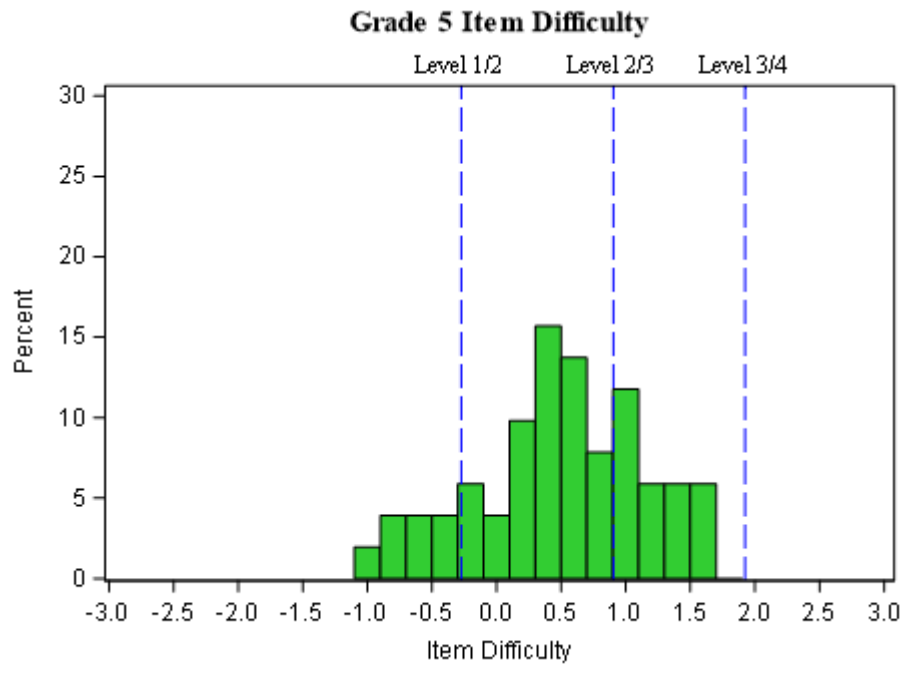


Figure 8.2.7. Grade 5 Item Difficulty and Student Ability Distributions

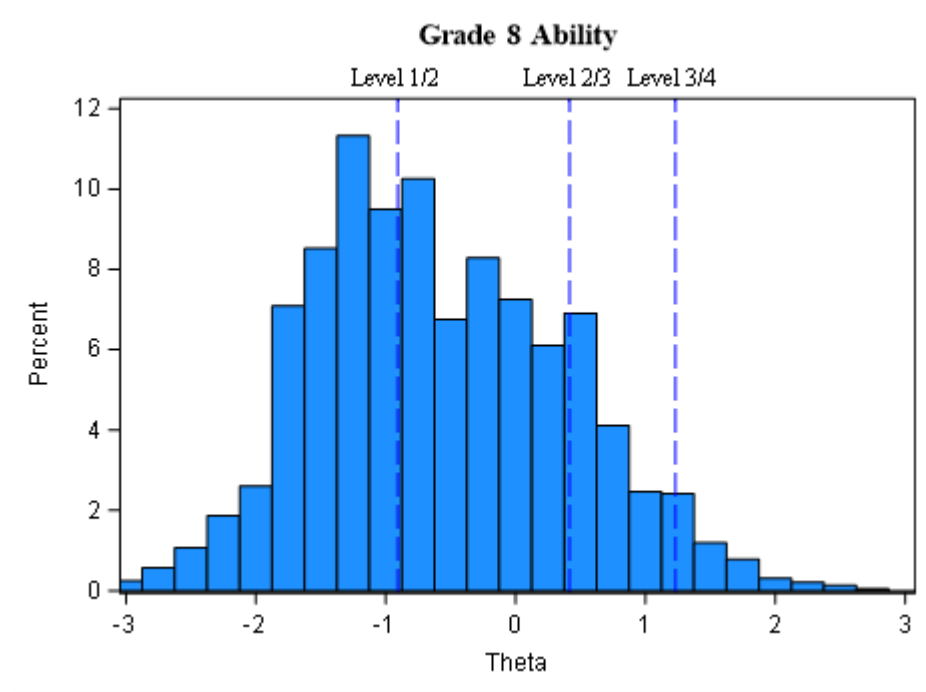
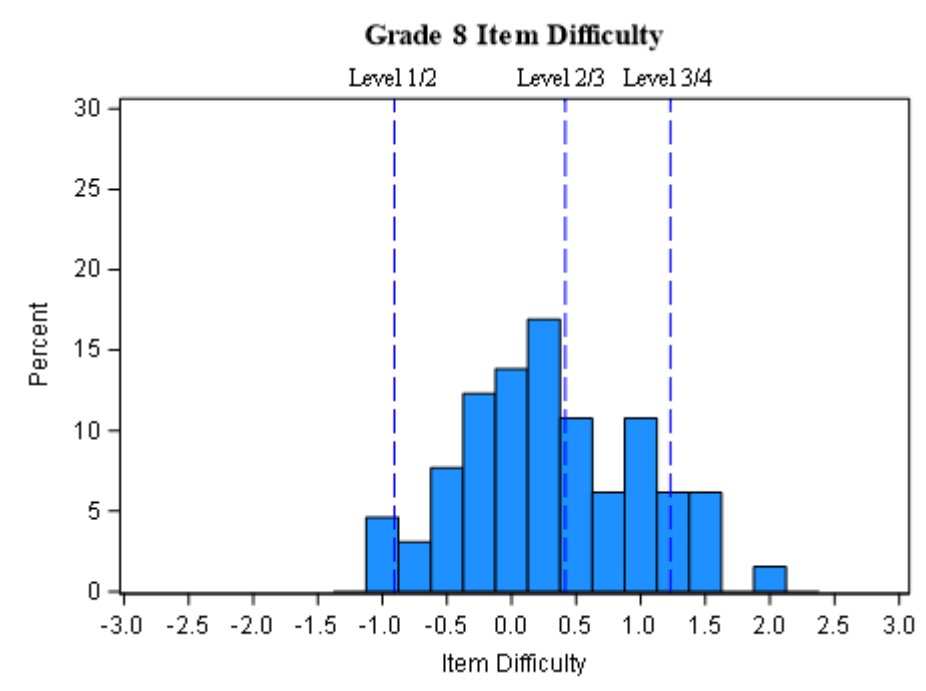


Figure 8.2.8. Grade 8 Item Difficulty and Student Ability Distributions

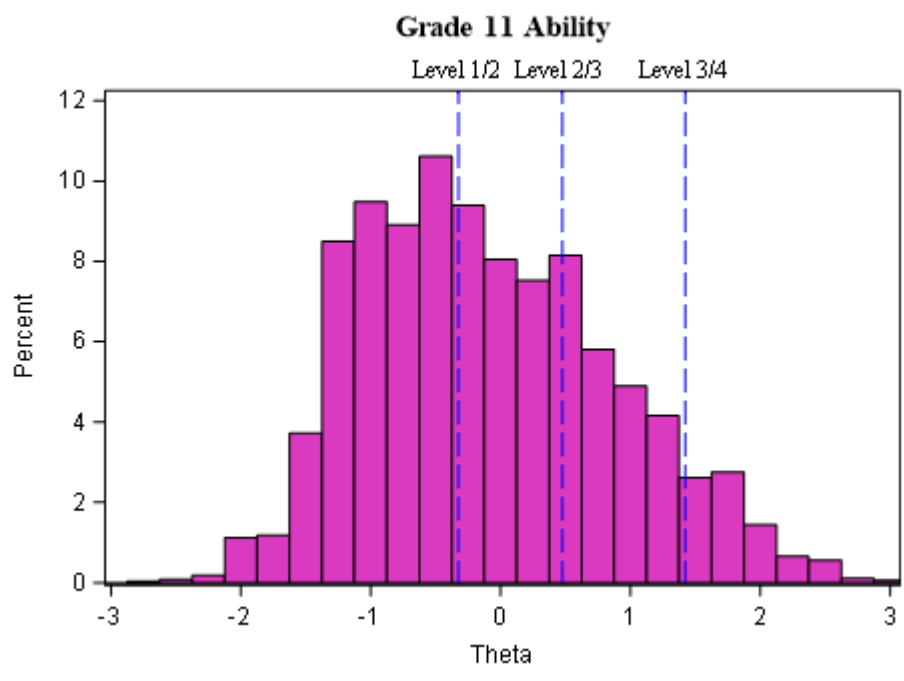
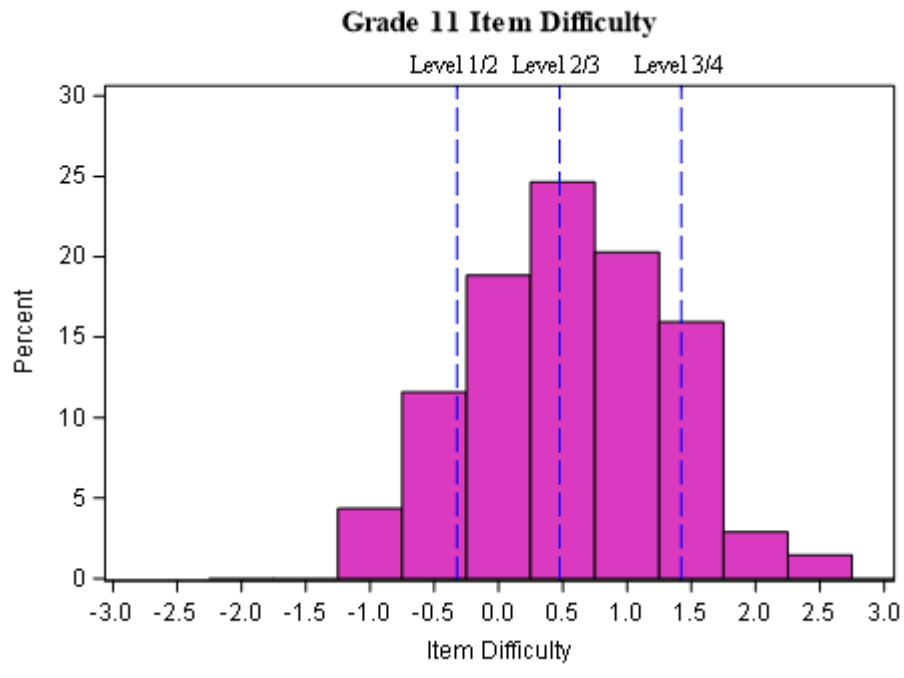


Figure 8.2.9. Grade 11 Item Difficulty and Student Ability Distributions

8.3 Reliability of Performance Classifications

The reliability of the performance-level classifications was evaluated via two methods. First, error bands were placed around each cut score using the CSEM. Next, the BB-CLASS (Brennan, 2004) program was used to calculate performance level classification consistency indices. The results of both methods indicate that the 2023 NJSLA–S performance-level classifications were reliable.

8.3.1 Conditional Standard Error of Measurement at Each Cut Score

As discussed in Section 8.2.2, the conditional standard error of measurements (CSEM) can be computed and evaluated along the theta (θ) scale. Also, the CSEM can be converted and placed on reported scales as needed and appropriate.

The 2023 NJSLA–S cut scores and the corresponding CSEM on the NJSLA–S scales are summarized in Table 8.3.1, and the CSEM tables for all raw and scale scores are presented in Appendix I. The values in Table 8.3.1 have been placed on the same scale as the scale score. At each grade, the cut score with the least amount of error is the level 2/3 cut score (200). At grade 5, the Level 3/4 cut score’s CSEM was slightly higher than at the Level 1/2, meaning that there was slightly less error in the scale score at 150 than at 243. At grades 8 and 11, the CSEM at the Level 1/2 cut is approximately the same as the CSEM for the Level 3/4 cut score. Table 8.3.1 also presents error bands that were placed around each of the cut scores to create upper and lower boundaries. The upper and lower bounds were defined by multiplying the cut score’s CSEM by two and either adding it to or subtracting it from the cut score. Overlap between the upper or lower bounds of a cut score and one of the other cut scores may indicate reliability issues among the performance level classifications. In 2023, no overlap between the upper or lower bound of a cut score and another cut score was found for any grade. The reliability of classification is investigated further with classification consistency indices discussed in Section 8.3.2.

Table 8.3.1: Cut Scores with Conditional Standard Error of Measurement

Grade	Level	Cut Scale Score	CSEM	Lower Bound	Upper Bound
5	Level 1/2	150	12.4	125.2	174.8
5	Level 2/3	200	11.1	177.8	222.2
5	Level 3/4	243	13.2	216.6	269.4
8	Level 1/2	150	10.7	128.6	171.4
8	Level 2/3	200	9.3	181.4	218.6
8	Level 3/4	231	10.0	211.0	251.0
11	Level 1/2	158	13.7	130.6	185.4
11	Level 2/3	200	12.7	174.6	225.4
11	Level 3/4	250	13.6	222.8	277.2

8.3.2 Classification Consistency Indices

A classification consistency index can be regarded as the percentage of examinees that would hypothetically be assigned to the same achievement level if the same test was administered a second time or an equivalent test was administered under the same conditions. Cohen’s

Kappa (Cohen, 1960, 1968) is a statistic that is often used to assess classification consistency. Coefficient Kappa (K) is given by:

$$K = \frac{P_o - P_c}{1 - P_c},$$

Equation 8.11

where P_o is the probability of a consistent classification and P_c is the probability of a consistent classification by chance. For the NJSLA–S, the classification consistency index for proficiency classifications is an estimate of how reliably the test classifies students into the performance categories (i.e., Levels 1–4).

Table 8.3.2 displays the results from BB-CLASS (Brennan, 2004) using the Livingston and Lewis (1995) consistency results. At each grade level, the classification consistency rates (P_o) ranged from .74 to .77. Thus, if the NJSLA–S had been administered a second time, approximately 75% of the students would have been classified at the exact same performance level. The most important decision is at the Level 2/3 cut score (200) because it demarcates the point along the scale where students are deemed proficient or not. The decision consistency at the Level 2/3 cut score or above was remarkable at .90 to .92, indicating 90% to 92% probability of being correctly classified as Level 3 or above. The overall NJSLA–S performance classification should be interpreted as consistent across grades.

Table 8.3.2: Performance Level Classification Consistency

Grade	Level 1/2 Cut	Level 2/3 Cut	Level 3/4 Cut	Kappa	P_o	P_o for Level 2/3 or above
5	150	200	243	.61	.74	.90
8	150	200	231	.64	.77	.92
11	158	200	250	.62	.74	.90

8.4 Reliability of Subscore Performance Classifications

The methodology used to create the subscore performance-level classifications was dependent on the CSEMs in the raw-to-theta subscore tables. Subscores associated with large CSEMs would indicate unreliable subscore performance-level classifications. The complete raw-to-theta subscore tables are presented in Appendix J.

Table 8.4.1 shows that the CSEMs associated with the subscore proficiency cut scores for each content domain and practice by grade were relatively small, indicating reliable subscore classifications. As presented in Table 8.4.1, the classification consistency rates (P_o) were above .70, given the short tests for content domains or practices at each grade level.

Table 8.4.1: Subscore Performance Classification Consistency and Conditional Standard Error of Measurement

Grade	Domain/Practice	Kappa	P_o	Level	Raw Subscore Cut	Theta	CSEM
5	Earth and Space	0.53	0.73	Near/Met	10	0.406	0.440
				Above	16	1.641	0.514
5	Life	0.56	0.74	Near/Met	9	0.307	0.456
				Above	15	1.630	0.518
5	Physical	0.49	0.71	Near/Met	7	0.268	0.500
				Above	13	1.608	0.480
5	Investigating	0.48	0.71	Near/Met	7	0.263	0.533
				Above	12	1.752	0.600
5	Sensemaking	0.57	0.75	Near/Met	8	0.204	0.501
				Above	14	1.950	0.647
5	Critiquing	0.53	0.75	Near/Met	11	0.447	0.386
				Above	19	1.591	0.401
8	Earth and Space	0.55	0.78	Near/Met	9	-0.065	0.461
				Above	15	1.186	0.469
8	Life	0.55	0.78	Near/Met	10	-0.103	0.412
				Above	17	1.082	0.429
8	Physical	0.55	0.79	Near/Met	10	-0.148	0.417
				Above	17	1.157	0.471
8	Investigating	0.50	0.75	Near/Met	7	-0.134	0.483
				Above	13	1.098	0.443
8	Sensemaking	0.55	0.78	Near/Met	11	-0.115	0.398
				Above	18	1.062	0.442
8	Critiquing	0.57	0.80	Near/Met	11	-0.080	0.420
				Above	18	1.276	0.489
11	Earth and Space	0.49	0.70	Near/Met	10	-0.102	0.423
				Above	17	1.322	0.502
11	Life	0.56	0.74	Near/Met	8	-0.123	0.440
				Above	16	1.265	0.418
11	Physical	0.56	0.74	Near/Met	10	-0.026	0.417
				Above	17	1.118	0.407
11	Investigating	0.53	0.73	Near/Met	9	-0.009	0.441
				Above	16	1.280	0.443
11	Sensemaking	0.56	0.74	Near/Met	11	0.013	0.418
				Above	18	1.244	0.437
11	Critiquing	0.53	0.72	Near/Met	9	-0.085	0.415
				Above	16	1.146	0.437

8.5 Rater Reliability

For constructed-response (CR) items, raters used item-specific scoring rubrics with a score range of 0 to 3 or 0 to 4, depending on the CR item. There were no half points assigned for any of the CR items. Only 10% of the constructed-response items were read by a second rater; the purpose of the second read was to investigate the consistency between raters. If the second read score was non-adjacent, then the scores for the response were erased and the paper was re-scored. Thus, all scores in the 10% of second reads were either perfect or adjacent agreement.

Table 8.5.1 shows, at the item level, the percentages of constructed-response items scored with exact or adjacent agreement and weighted Kappa. Weighted Kappa is a variation of Cohen's Kappa designed for ordinal variables. As shown in Table 8.5.1, the exact agreement rates ranged from 70.7% to 79.8% for grade 5, from 71.0% to 79.8% for grade 8, and from 68.5% to 81.6% for grade 11. While there was only one grade 8 CR item that showed weighted Kappa above .90, all the CR items in grade 5 and two of the CR items in grade 11 had weighted Kappas above .90. Overall, rater agreement on the NJSLA–S CR items was excellent.

Table 8.5.1: Inter-rater Agreement Rate of Constructed-Response Items

Grade	Item	% Raters in Exact Agreement	% Raters in Adjacent Agreement	Weighted Kappa
5	CR 1	70.7	29.3	0.93
5	CR 2	79.8	20.2	0.92
5	CR 3	77.6	22.4	0.94
8	CR 1	72.3	27.7	0.92
8	CR 2	79.8	20.2	0.80
8	CR 3	71.0	29.0	0.87
11	CR 1	71.3	28.7	0.90
11	CR 2	81.6	18.4	0.87
11	CR 3	68.5	31.5	0.91

PART 9: VALIDITY

The *Standards* state that “[v]alidity is a unitary concept. It is the degree to which all the accumulated evidence supports the intended interpretation of test scores for the proposed use” (AERA, APA, NCME, p. 14). If there is ample evidence to support reasonable interpretations and test uses, then they are considered to possess high validity (Kane, 2013). Conversely, interpretations and test uses that lack evidence possess low validity. Conceptually, Kane (2006) labeled the process of evaluating that evidence as validation. Test validation is an ongoing, ever-evolving process that extends through the duration of an assessment program. Every component within this technical report, from test development to score reporting, is evidence both for and against the valid interpretation and uses of test scores.

The *Standards* categorize validity evidence into five sections:

- evidence based on test content.
- evidence based on response processes.
- evidence based on internal structure.
- evidence based on relation to other variables.
- evidence based on the consequences of testing.

The following sections detail what evidence exists both for and against those five categories of validity evidence. Overall, the evidence suggests that the NJSLA–S fosters valid interpretations and uses of test scores as they pertain to the overall performance-level classifications of students.

9.1 Evidence Based on Test Content

Validity evidence based on test content refers to the relevance of the content of the test to the construct the test is purporting to measure. *Standard 1.11* states that:

[w]hen the rationale for test score interpretation for a given use rests in part on the appropriateness of test content, the procedures followed in specifying and generating content should be described and justified with reference to the intended population to be tested and the construct the test is intended to measure or the domain it is intended to represent. (AERA, APA, NCME, p. 26)

The content-related evidence of validity includes the extent to which the test items represent the specified content domains and cognitive dimensions. Adequacy of the content representation of the NJSLA–S is critical because the tests must provide an indication of student progress toward achieving the KSAs identified in the NJSL–S, and the tests must fulfill the requirements under ESSA (2015).

Adequate representation of the content domains defined in the NJSL–S is assured by using a test blueprint and a responsible test construction process as was described in Part 2. The NJSL–S is taken into consideration in the writing of all NJSLA–S items. In accordance with the test blueprint, the test construction process attempts to balance the six reporting categories and to ensure that the NJSLA–S contains an adequate representation of each content domain and scientific practice. Furthermore, all DCIs, SEPs, and CCCs are represented on the test.

Section 2.4 provides a summary of test construction in comparison to the goals established in the test blueprint.

The test content was well-balanced at the content domain level (i.e., Earth and Space, Life, and Physical Science). At each grade level, the content domains were all within five points of being perfectly balanced. The scientific practices (i.e., Investigating, Sensemaking, and Critiquing) were less balanced for grade 5 and within four points of being perfectly balanced for grades 8 and 11. At grade 5, the Critiquing practice was over-represented, accounting for 11 more points than the Investigating practice and 10 more points than the Sensemaking practice. At a more granular level all DCIs, SEPs, and CCCs were represented on each grade level's test. The relative balance of the DCIs, SEPs, and CCCs was less impressive with many categories being either over- or under-represented. Overall, the content domains and the range of DCIs, SEPs, and CCCs provide evidence that the test is adequately measuring the KSAs defined by the NJSLS–S. However, the relative lack of balance in the scientific practices and individual DCIs, SEPs, and CCCs provides evidence that the scale may be over-represented by certain components within the NJSLS–S, which could affect interpretations of test scores at both the overall and subscore level.

9.1.1 Alignment Study

In August of 2022, the NJDOE commissioned an independent evaluation of the alignment quality of the NJSLS–S administered at grades 5, 8, and 11. Evidence of alignment quality is critical to validity evaluation for standards-based assessments (Forte, 2017; Webb, 1997, 1999). Such evidence must draw upon an examination of how a test has been designed and developed, as well as instances of the test itself (Forte, 2013). As is the case for all validity evidence, evidence of alignment quality is necessary to support the interpretation and use of test scores. A well-aligned test is one that elicits a sample of student performance that is adequate to support inferences about student achievement in relation to the standards-based domains on which the test is based. To address the unique aspects of the three-dimensional nature of the NJSLS–S and the NJSLS–S items and test forms, the following alignment questions guided the evaluation: (1) To what extent do the blueprints support the consistent creation of test forms that reflect the standards and the score scale? (2) To what extent do the Performance-Level Descriptors (PLDs) reflect meaningful and appropriate score interpretations across the full range of the score scale? and (3) To what extent does the set of phenomena, tasks, and items reflect the blueprints and provide performance opportunities across the full range of the score scale?

The results of the study found that the blueprint development was well documented across all three grades (5, 8, and 11), and included a clear description of the review and revision process by stakeholders. Each blueprint met the criteria of strong evidence of alignment for Domain Concurrence, Balance of Representation, and Phenomena Design. The PLDs for all three grade levels were determined to have strong evidence of alignment with the NJSLS–S and were found to describe increasingly sophisticated and reasonable levels of performance for the concepts defined in the standards. All three test forms met the criteria for strong evidence of alignment with the intended DCI and were judged as strongly representing the multidimensionality of the standards with 100% of items aligning to the additional dimensions of the standards (SEP and CCC). Finally, all test forms met expectations for Domain Concurrence, Range of Knowledge, and Balance of Representation. Further, panelists evaluated the items on the form as being

cognitively challenging, though panelists noted that, while a range of cognitive challenge levels is present within the form, items tend to skew toward the higher levels of cognitive challenge, with less representation at the lower levels. The results of this alignment evaluation will be used to inform future item and assessment development activities. The Executive Summary of the alignment evaluation study is included in Appendix L.

9.2 Evidence Based on Response Processes

Standard 1.12 states that “[i]f the rationale for a test score interpretation for a given use depends on premises about the psychological processes or cognitive operations of test takers, then theoretical or empirical evidence in support of those premises should be provided” (AERA, APA, NCME, p. 26). Evidence based on response processes is complementary to evidence based on test content; it can come from several sources including response times, eye-tracking, think-aloud protocols, interviews, and/or focus groups. This complementary evidence is different from content evidence because its source is not content experts or teachers, but rather the actual student test takers. Padilla and Benitez (2014) noted that “validation studies aimed at obtaining evidence from response processes are scant” (p. 139). The NJSLA–S evidence based on judgment from the NJSAC, content specialists, and a cognitive lab study is described below.

The alignment of each item to the Range PLDs provides limited evidence of the cognitive processes theoretically being assessed by the NJSLA–S. As described in the 2019 NJSLA–S technical report (NJDOE, 2019), the Range PLDs were created in a collaborative effort by NJDOE, the NJSAC, content specialists, and psychometricians; they are based upon the NJSLS–S content standards. Note that the Range PLDs were not finalized until well after the completion of the item development process for the 2019 NJSLA–S.

The Range PLDs are the theoretical cognitive structure underlying all current NJSLA–S item and test development. They contain detailed descriptions of the knowledge, skills, and abilities (KSAs) that a student needs to display to be classified at a given performance level. Each item on the NJSLA–S was aligned to two Range PLDs: one based on the DCI, and one based on the SEP. Those alignments were verified by the NJSAC. The alignment of each item to the Range PLDs offers a theoretical link from the NJSLA–S’s underlying cognitive structure to the student responses, which provides limited validity evidence based on response processes. The detailed test maps presented in Appendix F display the Range PLD alignment for each item.

Table 9.2.1 shows the distributions of the performance levels associated with each item by grade level and by DCI and SEP. The DCI distribution of items at grade 5 and the DCI and SEP distributions at Grade 11 clustered at Levels 1 and 2, tapering off at Level 3. The grade 5 SEP distribution was clustered at Levels 2 and 3, as was the grade 8 DCI distribution. The grade 8 SEP distribution was more heavily centered at Level 2. These distributions largely correspond to the item difficulty distributions illustrated in Figures 8.2.4 through 8.2.6.

Table 9.2.1: Range PLD Alignment by DCI, SEP, and Grade Level

Grade	Domain/Practice	Level 1	Level 2	Level 3	Level 4
5	DCI	19	22	8	1
5	SEP	12	22	16	0
8	DCI	6	27	22	4
8	SEP	10	33	10	6
11	DCI	19	29	13	7
11	SEP	18	34	15	1

9.2.1 Cognitive Lab Study

To evaluate the degree to which the items and tasks on the NJSLA–S in grades 5, 8, and 11 elicit the intended response processes as represented in the NJSLS for Science, cognitive interviews with students were conducted. The purpose of this study was to gather evidence of the response process. Messick (1995) argued that the substantive validity of test scores relates to the theoretical underpinnings of the construct that is meant to be measured. In the case of statewide, standards-based, academic assessments, the construct that is meant to be measured derives from the set of standards in each content area and grade level. The validity evidence necessary to support score interpretation and use includes evidence regarding the alignment of test tasks to the standards in terms of breadth and depth (Webb, 1997; Forte, 2013, 2017) as well as consideration of whether the test tasks elicit the intended cognitive processes as students generate responses to the tasks (Thelk et al., 2006). Items must be developed to elicit those cognitive processes and examined to determine whether, in practice, students’ cognitive processing is influenced by variables other than the ones test designers are interested in measuring, which introduces construct irrelevant variance (Thelk et al., 2009). Two evaluation questions guided the evaluation: (1) To what extent do the tasks on the NJSLA–S tap the intended cognitive processes as represented in the NJSLS for Science? And (2) How do students interact with the task types within the NJSLA–S?

To answer the evaluation questions, cognitive laboratories (often referred to as cog-labs) were conducted with 12 students in each grade level across two New Jersey districts (one urban and one suburban) in November 2022. The cog labs used a think-aloud protocol, in which each student outlined his/her thinking as they worked to answer each item. The study included both a concurrent account of problem-solving, as well as a retrospective cognitive interview. Because the study was conducted in the fall, an off-grade approach was used to ensure participating students had the opportunity to learn the assessed standards. Students in grades 6, 9, and 12 participated in the study as they received instruction on the assessed standards during the previous school year. Two evaluators observed each student, audio-recorded the session, and independently coded their observations using a standard protocol. The Executive Summary of the cog-labs study is included in Appendix M.

9.3 Evidence Based on Internal Structure

According to the *Standards*, “[a]nalysis of the internal structure of a test can indicate the degree to which the relationships among test items and test components conform to the construct on which the proposed test score interpretations are based” (AERA, APA, NCME, p 16). The NJSLA–S was constructed as a unidimensional test. However, it also assesses student performance in several content clusters. It is important to study the pattern of relationships among the content clusters and testing methods. Therefore, this section addresses evidence based on responses and internal structure. Overall, the evidence supports the notion that the internal structure of the NJSLA–S is unidimensional and that its items are measuring the same construct. However, at the subscore level, results from a confirmatory factor analysis provided some evidence that the internal structure was not performing as intended.

9.3.1 Intercorrelations

One method for studying patterns of relationships to provide evidence supporting the inferences made from test scores is to evaluate the correlations between the total test score and its subscores. If the subscores are highly correlated, then that provides evidence that the test is unidimensional. Section 6.2.1.1 of this document summarizes correlation coefficients among test content domains and clusters by grade level. The intercorrelations of the NJSLA–S provide clear evidence that the NJSLA–S is unidimensional. Among the content domain subscores at all grade levels, the lowest correlation was .77 at grade 5 between Physical Science and Earth and Space Science. Among the scientific practices, the intercorrelation ranged from .79 to .81.

9.3.2 Other Internal Structure Evidence

Further evidence of the internal structure of the NJSLA–S was also presented via a principal component analysis (PCA). The PCA results are presented in Section 6.2.1.2. These scree plots show further evidence that the variability in the NJSLA–S test scores is due to a single dimension. No secondary factors at any grade level practically contributed to explaining the variation in the overall NJSLA–S test scores, while subtest scores could convey pedagogical information on a specific content domain or scientific practice.

Part 8 of this Technical Report provides ample evidence to support NJSLA–S reliability. Reliability is the extent to which items within a test measure aspects of a singular construct. Internal consistency reliability (for which evidence presented in Part 8) is one measure of reliability. The grade-level internal consistency reliability coefficients presented in Section 8.1 were strong, ranging from .92 to .93 for the CBT form. At the subscore level the reliability coefficients were relatively impressive, with the lowest estimates of .75 for both the grade 5 Physical Science and the grade 8 Investigating subscores.

9.3.3 Confirmatory Factor Analysis for 2023 NJSLA–S

To provide further evidence supporting the internal subscore structure of the NJSLA–S, confirmatory factor analyses (CFA) were conducted using the 2023 operational NJSLA–S test results. CFA is a powerful statistical technique used to verify proposed measurement models based on the underlying covariance structure of the data. With this technique, a measurement model is specified, the data are fit to the specified model, and then fit indices (and other

criteria) can be examined to determine how well the model fits the data (Brown, 2006; Kline, 2011).

Figure 9.3.1 presents the *a priori* test structure specified for the content domain CFA across all grades. Figure 9.3.2 presents the *a priori* test structure that was specified for the scientific practice CFA across all grades. Adjusted unweighted least squares with means and variances adjusted (ULSMV) was used to estimate these models. ULSMV is a robust estimation method that is appropriate to use when data are categorical or ordinal and there are potential concerns about non-normality in the latent variables (Beauducel & Herzberg, 2006; Muthén, 1993; Muthén et al., 1997).

To assess model fit, the model chi-square test was used as a test of exact model fit. However, it is well documented that the model chi-square test is sensitive to sample size (Cheung & Rensvold, 2002), so approximate fit indices were used to supplement the model chi-square. Three approximate fit indices were examined for each model: the comparative fit index (CFI; Bentler 1990), the root mean square error of approximation (RMSEA; Steiger, 1990), and the standardized root mean square residual (SRMR; Jöreskog & Sörbom, 1988). Following the recommendations of Brown and Cudeck (1993) and Hu and Bentler (1999), values of CFI greater than .90 and .95 were considered evidence of acceptable model fit and good model fit, respectively; values of RMSEA less than .08 and .05 were considered evidence of acceptable model fit and good model fit, respectively; and values of SRMR less than .08 were considered evidence of good model fit.

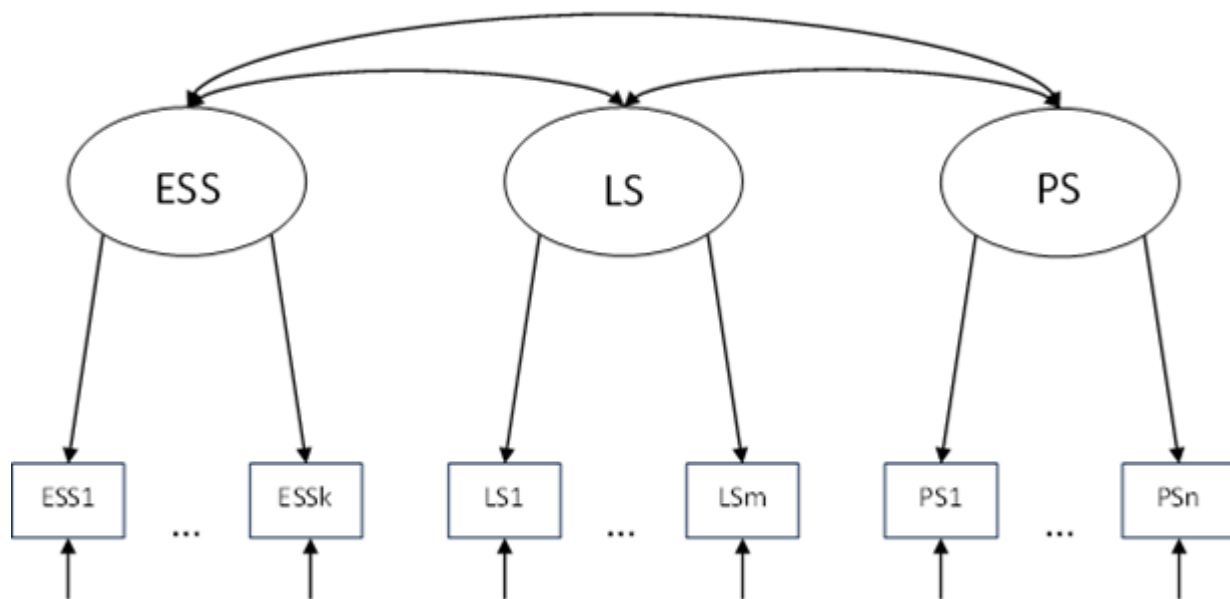


Figure 9.3.1. Domain Subscore Structure

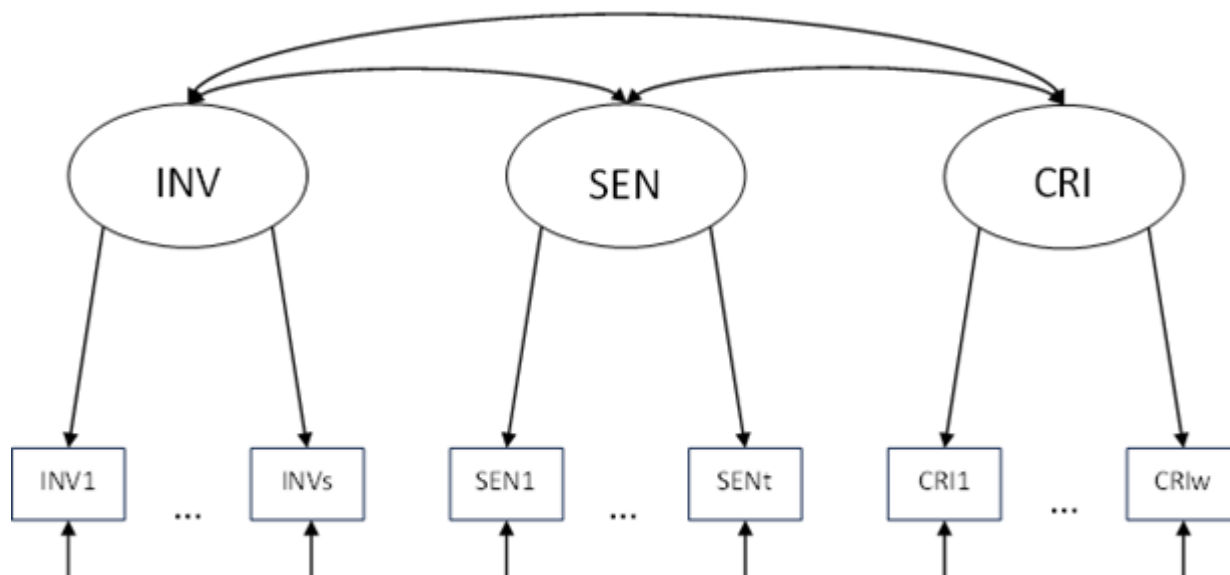


Figure 9.3.2. Practice Subscore Structure

There are a few considerations regarding the use of the model fit statistics in this study. First, since a robust estimation method was used, the scaled chi-square statistic is reported (Satorra & Bentler, 1994). The scaled chi-square statistic is simply the chi-square statistic modified by a scaling parameter to adjust for violations of normality. However, since the CFI and the RMSEA are functions of the chi-square statistic, they are calculated using this scaled chi-square statistic. Second, studies have shown that the CFI and RMSEA are impacted by the estimation method. Specifically, ULSMV model estimation has been shown to result in higher CFI values and lower RMSEA values when compared to maximum likelihood estimation methods (Garrido et al., 2016; Nye & Drasgow, 2011; Shi & Maydeu-Olivares, 2020; Xia & Yang, 2019). This can result in a lower likelihood of identifying a poorly fitting model. However, the SRMR has been shown to be consistent across estimation methods (Shi & Maydeu-Olivares, 2020), so it was given considerable weight when judging model fit.

9.3.3.1. Domain results. The domain model converged without issues for all three grades. The model fit indices for each grade are presented in Table 9.3.1. For all three grades, the chi-square test indicated that the differences between the observed- and model-predicted covariances were statistically significant. However, the approximate model fit indices indicated that the domain model was a good representation of the data. As shown in Table 9.3.2, the correlations between the latent subscores were very high due to the underlying unidimensionality of the NJSLA–S, but the acceptable model fit indices suggest that they were also empirically distinct from one another. The parameter estimates obtained from fitting the domain model are provided in Appendix O for all three grades.

Table 9.3.1. Model Fit Indices for the Domain Model

Grade	Chi-square	CFI	RMSEA	SRMR
5	$\chi^2_5 (1221) = 43140.93^*$	0.985	0.019	0.024
8	$\chi^2_8 (2012) = 40174.26^*$	0.985	0.014	0.020
11	$\chi^2_{11} (2346) = 74085.05^*$	0.975	0.018	0.026

Note. $N_5 = 96,392$; $N_8 = 101,478$; $N_{11} = 94,023$; * = $p < .001$

Table 9.3.2. Correlations between the Latent Subscores Implied by the Domain Model

Grade 5			Grade 8			Grade 11		
-	ESS	LS	-	ESS	LS	-	ESS	LS
ESS	-	-	ESS	-	-	ESS	-	-
LS	0.98	-	LS	0.98	-	LS	0.99	-
PS	0.97	0.99	PS	0.99	0.97	PS	0.99	0.98

9.3.3.2. Practice results. The practice model converged without issues for all three grades but returned nonpositive definite latent covariance matrices in each grade. A nonpositive definite latent covariance matrix suggests that some aspect of the model is not accurately reflecting the covariance structure observed in the data (i.e., the model is misspecified). For the NJSLA–S, this misspecification was determined to be caused by collinearity between the practice subscore categories. That is, the practice subscore categories were not empirically distinct enough to be treated as separate subscore categories. Since the practice model returned a nonpositive definite latent covariance matrix for all three grades, the parameter estimates should not be interpreted. Therefore, the parameter estimates obtained from the practice models are omitted from Appendix O.

9.3.3.3. Follow-up analysis. The item subscore correlations (rpb_{sub}) were examined to further examine the subscore categories. The item subscore correlations are analogous to item total correlations except the raw subscore rather than total test raw score is used for calculation. For each item, the item subscore correlation was calculated for all three domains (i.e., rpb_{ESS} , rpb_{LS} , and rpb_{PS}) and all three practices (i.e., rpb_{CRI} , rpb_{INV} , and rpb_{SEN}).

Table 9.3.3 and Table 9.3.4 show the number of items with the highest rpb_{sub} on their assigned domain and practice, respectively. For all three grades, at least 50% of the items showed the highest rpb_{sub} on their assigned domain. Conversely, only at least 50% of the Grade 5 Sensemaking items, the Grade 8 Critiquing and Sensemaking items, and the Grade 11 Sensemaking items showed the highest rpb_{sub} on their assigned practice.

Table 9.3.3. Number of Items with Highest rpb_{sub} on Assigned Domain

Grade	Domain	Item Count	Number (%) of items with highest rpb_{sub} on assigned Domain
5	ESS	18	18 (100.00%)
5	LS	17	17 (100.00%)
5	PS	16	16 (100.00%)
8	ESS	20	20 (100.00%)
8	LS	24	14 (58.33%)
8	PS	21	21 (100.00%)
11	ESS	20	20 (100.00%)
11	LS	23	15 (62.55%)
11	PS	26	26 (100.00%)

Note. ESS = Earth and Space Science; LS = Life Science; PS = Physical Science

Table 9.3.4. Number of items with Highest rpb_{sub} on Assigned Practice

Grade	Practice	Item Count	Number (%) of items with highest rpb_{sub} on assigned Practice
5	CRI	18	4 (22.22%)
5	INV	16	1 (6.25%)
5	SEN	17	10 (58.82%)
8	CRI	22	11 (50.00%)
8	INV	20	1 (5.00%)
8	SEN	23	13 (56.52%)
11	CRI	22	6 (27.27%)
11	INV	23	6 (26.09%)
11	SEN	24	14 (58.33%)

Note. CRI = Critiquing; INV = Investigating; SEN = Sensemaking

The results of the CFAs and the item subscore correlations support the structure of the NJSLA–S in terms of the content domains. Specifically, the results suggest that while the domain subscore categories are highly related, they are empirically distinct enough to provide some unique information in the form of subscores. However, the results of the CFAs and the item subscore correlations did not support the structure of the scientific practices. The presence of nonpositive definite latent covariance matrices and much commingling between the rpb_{sub} values suggest that the practice categories are not empirically distinct from one another. Explaining this finding requires further information from the item and test development processes of the NJSLA–S, which is presented below.

Each content domain: Earth and Space Science (ESS), Life Science (LS), and Physical Science (PS) is divided into a subset of fundamental, core ideas that are necessary for understanding a given science discipline. The core ideas are woven throughout the K–12 standards, providing themes

that build in complexity and allow for a more interdisciplinary approach to learning science. As a result, each DCI standard is a discrete scientific concept that can be directly assessed, and the cluster design of the NJSLA–S allows a DCI to be assessed multiple times. The discreteness of the definitions and the directness with which the content domains can be measured has resulted in distinct, well-defined subscore categories.

However, while all eight of the practices are defined in the standards, the SEP standards are not as discreetly defined as the DCI standards. This is likely because the DCIs are based on content knowledge and not application/skill-based standards. Overlap among the practices and perhaps, more importantly between reporting categories exists within the standards, which is most likely why collinearity was observed between the SEP subscores. Additionally, the skills characterized in the SEPs are inherently connected. There is no way for a student to critique (Critiquing) a dataset without analyzing the data (Sensemaking), and probably hypothesizing about it (Investigating) first. This inherent connectedness is likely another reason why the SEP subscores displayed collinearity.

9.4 Evidence Based on Relationships to Other Variables

Evidence based on relationships to other variables takes the form of relationships between test scores and other variables that are external to the test (AERA, APA, NCME, 2014). This evidence can come from investigating the relationships among tests that measure similar constructs, tests that measure different constructs, or other outcomes that a test purports to predict. NJDOE conducted an internal validity study that investigated the relationships among the NJSLA–S and other New Jersey large-scale, statewide subject scale scores (i.e., NJSLA–ELA and NJSLA–Math). The results indicate that the scientific KSAs the NJSLA–S is intended to measure comprise a construct distinct from other disciplines measured by the New Jersey statewide assessment program.

The results at grade 5 are displayed in Table 9.4.1. Students with valid scale scores in ELA, math, and science (NJSLA-ELA/M/S) in spring 2023 were included in the analysis. ELA consists of two major claims: Reading Complex Text and Writing. The scale scores for those two major claims were added to the matrix. The relationships among science, ELA, and math were consistent with expectations and showed correlations of .80 and .82, respectively. The correlation between science and ELA writing was .64, and the correlation between science and ELA reading was .82.

Table 9.4.1: Grade 5 Intercorrelations by Content Area

Content Area	N	Science	ELA	ELA-R	ELA-W	Math
Science	94,516	1.00	-	-	-	-
ELA	94,516	0.80	1.00	-	-	-
ELA Reading	94,516	0.82	0.95	1.00	-	-
ELA Writing	94,516	0.64	0.88	0.71	1.00	-
Math	94,516	0.82	0.75	0.75	0.64	1.00

At grade 8, students with valid scale scores on the following were included in the analysis: (1) Grade 8 Science, ELA and Math; or (2) Grade 8 Science, ELA and Algebra I; or (3) Grade 8 Science, ELA, and Algebra II; or (4) Grade 8 Science, ELA, and Geometry. The results at grade

8 are displayed in Table 9.4.2. The only difference in calculating the grade 8 intercorrelation matrix in comparison to grade 5 pertained to the math scale scores. Depending on which course a student was enrolled in, there were four different math assessments that grade 8 students could have taken: Math 8, Algebra I, Algebra II, or Geometry. Therefore, instead of one math scale score the grade 8 intercorrelation matrix is based on four distinct math scale scores. It is impossible for students to have scale scores on two different math tests; thus, those cells in the correlation matrix are represented by N/A.

Similar to Grade 5, the correlation between Grade 8 science and ELA writing was .63. The correlations between science and various math test scores ranged from .68 to .79. This is most likely due to the higher- and lower-achieving students taking different assessments, which could decrease the scale score variance for each math test. Thus, the magnitude of the correlations between science and the various math tests appears reasonable when considering that math achievement is more homogeneous within each sub-group than if all students at all ability levels were taking the same assessment.

Table 9.4.2: Grade 8 Intercorrelations by Content Area

Content Area	N	Science	ELA	ELA-R	ELA-W	Math 8	Alg. I	Alg. II	Geo.
Science	99,430	1.00	-	-	-	-	-	-	-
ELA	99,430	0.76	1.00	-	-	-	-	-	-
ELA-Reading	99,430	0.79	0.95	1.00	-	-	-	-	-
ELA-Writing	99,430	0.63	0.92	0.76	1.00	-	-	-	-
Math 8	65,034	0.72	0.65	0.65	0.55	1.00	-	-	-
Algebra I	28,890	0.79	0.65	0.67	0.52	-	1.00	-	-
Algebra II	783	0.68	0.40	0.47	0.24	-	-	1.00	-
Geometry	4,723	0.76	0.55	0.55	0.44	-	-	-	1.00

For Grade 11, students with valid scale scores on both the New Jersey Graduation Proficiency Assessment (NJGPA; the ELA and mathematics components) and the NJSLA–S were included in the analysis. The results for grade 11 are presented in Table 9.4.3. The correlation between science and ELA was .73 while the correlation between science and math was .80. This was consistent with the relationships between science, math and ELA seen in other grades, and is in line with expectations. The relationships between science and ELA reading and ELA writing were .75 and .61, respectively, showing similar relationships seen in grades 5 and 8.

Table 9.4.3: Grade 11 Intercorrelations by Content Area

Content Area	N	Science	ELA	ELA-R	ELA-W	Math
Science	92,295	1.00	-	-	-	-
ELA	92,295	0.73	1.00	-	-	-
ELA Reading	92,295	0.75	0.95	1.00	-	-
ELA Writing	92,295	0.61	0.92	0.78	1.00	-
Math	92,295	0.80	0.72	0.74	0.61	1.00

9.5 Evidence Based on the Consequences of Testing

Standard 1.25 states that “[w]hen unintended consequences result from test use, an attempt should be made to investigate whether such consequences arise from the test’s sensitivity to characteristics other than those it is intended to assess or from the test’s failure to fully represent the intended construct” (p. 30). Lane and Stone (2002, p. 24) list the types of evidence that can be collected to evaluate the consequences of a large-scale statewide accountability assessment program.

- Student, teacher, and administrator motivation and effort
- Curriculum and instructional content and strategies
- Content and format of classroom assessments
- Improved learning for all students
- Professional development support
- Use and nature of test preparation activities
- Student, teacher, administrator, and public awareness and beliefs about the assessment and criteria for judging performance and the use of assessment results

No NJSLA–S validity evidence based on the consequences of testing currently exists. Future NJSLA–S validity studies, including evidence based on consequences, are detailed in Section 9.7.3.

9.6 Other Validity Evidence

Each part within this technical report contributes evidence relevant to validity. The following is a summary of evidence within each part:

Part 1: Introduction—This part describes the purpose of the assessment including:

- intended inferences and uses of test scores
- the relationship between the NJSLS–S and NJSLA–S

Part 2: Test Development—This part describes the processes used to design and develop the NJSLA–S including:

- the steps taken to link test development to the intended inferences and uses of the NJSLA–S
- the training and QC procedures implemented in the item development process
- the use of NJDOE, the NJSAC, and the Sensitivity committee to ensure the work of item writers and content specialists was aligned to the NJSLS–S
- the statistical review of each item after being field tested
- the steps taken to ensure the test construction process matched the NJSLA–S blueprint and statistical constraints

Part 3: Test Administration—This part describes the care that was taken to implement standardized test administration procedures including:

- documents produced to communicate NJSLA–S test administration procedures for all versions of the test
- steps taken to ensure testing materials were handled using safe and secure procedures
- accommodations and accessibility features that were used during the test administration to provide all NJSLA–S test-takers with equal opportunities on the test

Part 4: Scoring—This part describes the procedures that were implemented to verify the accuracy of scoring student responses including:

- confirming all computer-scored answer keys for both MC and TE item types
- development of unique scoring guides for each CR item
- selecting and training the scorers, team leaders, and scoring directors charged with handscoring the CR items
- monitoring handscorers to verify they are implementing the scoring rubric accurately
- verifying that student raw scores and subscores were calculated accurately

Part 5: Standard Setting—This part and the 2019 NJSLA–S Technical Report describe the methods that were undertaken to set the NJSLA–S performance standards including:

- approval of all NJSLA–S Standard-Setting methods by the NJTAC
- development of performance level descriptors
- selection of a representative group of New Jersey educators to serve as standard-setting panelists
- evaluation of the standard-setting meeting by the standard-setting panelists
- external review of the standard-setting meeting by an NJTAC member
- documentation of all results in the NJSLA–S Standard-Setting Report

Part 6: Item and Test Statistics—This part describes the battery of statistics that were used to evaluate the NJSLA–S at both the test and item level including:

- summaries of item performance across grade level, content domain, scientific practice, and item type to verify that the items are appropriate
- measures of test speededness to assess whether students could finish the test in the allotted time
- confirming the test items were not disadvantaging large subgroups of students via DIF statistics
- descriptive statistics of raw and scale scores by test form and subgroups of students to evaluate how appropriate the test is for portions of the population
- evaluating the IRT assumptions of the PCM to ensure it is appropriate for modeling student ability estimates

- evaluating IRT person fit statistics by subgroups of students

Part 7: Equating and Scaling—This part describes the methods used to ensure all students at a given grade level received scale scores that were comparable including:

- documenting the equating and scaling procedures
- descriptions of the special equating(s)

Part 8: Reliability—This part describes the reliability statistics that were calculated to verify the consistency of the NJSLA–S test scores including:

- verifying the reliability at the total score, form, subscore, item type, and subgroup levels
- evaluating graphic displays of IRT reliability such as TIFs and CSEMs
- assessing the consistency of student performance-level classifications
- assessing rater agreement rates for the handscoring of all CR items

9.7 Summary

Messick (1989) defined validity as “an integrated evaluative judgment of the degree to which empirical evidence and theoretical rationales support the adequacy and appropriateness of inferences and actions based on test scores and other modes of assessment” (p. 13). Making an integrated evaluative judgment with such a diverse assortment of evidence is challenging given that the validity process is ongoing and exists throughout the duration of the testing program. Overall, there is ample evidence that the NJSLA–S fosters valid inferences and uses. However, the NJSLA–S validity argument requires continuing attention, and an iterative process of identifying its weakest components, making modifications, and then reevaluating their effectiveness is needed. As Cronbach (1980) said “the job of validation is not to support an interpretation, but to find out what might be wrong with it. A proposition deserves some degree of trust only when it has survived serious attempts to falsify it” (p. 103). The following sections set forth the pros and cons of the NJSLA–S validity evidence by the primary inferences and uses of the test.

9.7.1 Student Performance Level Classifications: Overall Scale Score

The most important inferences made from the NJSLA–S involve the student performance level classifications. Students are classified in Levels 1 through 4; students at or above Level 3 are deemed proficient. All interpretations based on NJSLA–S performance-level classifications should be validated for evaluating student performance as it pertains to the KSAs defined in the NJSLA–S.

Overwhelming validity evidence in support of the proposed performance-level classification interpretations has been presented throughout this document and within the validity section. The NJSLA–S was developed and constructed by well-trained experts with assistance from NJDOE and the NJSAC to specifically measure the wide range of KSAs defined in the NJSLA–S. It was administered under strict standardized processes and procedures. The accuracy of the scoring of all NJSLA–S items was verified. The performance level classifications were determined at standard setting using methodology that was reviewed and approved by the NJTAC. After the test administration, the items were statistically reviewed to ensure they met the assumptions of

the proposed IRT model. Finally, both the overall scale and the performance-level classifications were verified as being internally consistent.

There are some areas in which the validity evidence in support of the performance-level classification inferences could be improved. The validity section on consequences also has no evidence, which is somewhat expected due to the challenge of integrating consequential validity evidence into a coherent validity argument (Cizek, 2016), as well as to the fact that it is hard to identify the long-term consequences of a testing program after its first year of operational use. Also, the Reporting PLDs would be more useful in providing guidance to test score users if they contained both performance level- and grade-specific KSAs. The current versions are generic for each performance level and do not differentiate among grade-level skills.

Overall, the evidence in favor of the valid interpretations of performance-level classification outweighs the areas in which evidence is lacking or non-existent. As a standards-based assessment, the content validity evidence linking the test scores and interpretations to the NJSL-S and the test blueprint are of chief importance (Sireci et al., 2008). Studying the issues noted above would enhance the validity evidence.

9.7.2 Student Performance Level Classifications: Domains and Practices Subscores

Inferences and uses of subscores are of secondary importance to the overall scale score and performance-level classifications. Student subscores are used to classify their performance as Below Expectations, Near/Met Expectations, or Above Expectations. Students do not receive either a raw or a scale score in any of the subscore categories. The validity evidence pertaining to interpretations based on NJSLA-S subscore performance-level classifications is limited, and caution in using the subscores should be emphasized.

Some validity evidence in support of the interpretations of subscores is presented throughout this document. Much of the validity evidence supporting the overall scale score—for instance, the test administration and scoring procedures—also contributes to subscore validity evidence. Aside from that, item development, test construction, and PLD creation were all undertaken with the explicit goal of being able to report student performance in the six subscore categories. The subscore performance-level procedures were approved by the NJTAC, and each subscore raw-to-theta score table was independently calibrated and verified by two MI psychometricians. Psychometrically, the subscores displayed adequate reliability coefficients and CSEMs.

Finally, the connection of the NJSLA-S subscores to the NJSL-S is unclear. The NJSL-S emphasizes the SEPs, DCIs, and CCCs, whereas the NJSLA-S is reporting subscore categories back to students, teachers, and administrators in categories that are clusters of SEPs and DCIs. One of the stated goals of the NJSLA-S is to provide feedback to schools on their overall performance in the six subscore categories, but it is not clear how to use or interpret that information within the framework of the NJSL-S. Constructing links between the NJSL-S and the reporting categories of the NJSLA-S would improve the ability of teachers, schools, and administrators to use and interpret the information in the subscores.

Overall, the intended inferences being made from the NJSLA-S subscores lack enough validity evidence that any interpretations and uses should be made with caution. NJDOE has sagaciously emphasized caution in both their communications with LEAs and in the Score Interpretation

Guide. Future studies of response processes and factor structures, as well as links from the NJSLS–S to the NJSLA–S reporting categories, could provide insights into how to best interpret and use the subscores; as previously noted in Section 2.4, ongoing, iterative improvements to item development and test construction might alleviate the lack of balance between individual scientific practices and the three content domains.

9.7.3 Future NJSLA–S Validity Studies

As was noted earlier, Kane (2006) labeled the process of evaluating validity evidence as validation, and he conceptualized that process as ongoing, ever evolving, and extending through the duration of an assessment program. NJDOE is committed to addressing the limitations within the NJSLA–S validity evidence and iteratively enhancing the validity of the inferences made from its test scores. One future study is planned, and some details are provided in Section 9.7.3.1.

9.7.3.1 Consequences of the NJSLA–S. Two of the goals of the NJSLA–S are to influence adoption of the NJSLS–S curriculum and to inform instruction, which will in turn improve the educational opportunities for New Jersey students. As described in Section 9.5: Evidence Based on the Consequences of Testing, Lane and Stone (2002) list many possible studies of the consequences of testing programs. They generally involve evaluating whether the testing program is having its intended effect and/or whether it is having unintended consequences. Sources of the data come from students, teachers, administrators, and parents. The future study would likely follow recommendations from Lane and Stone to evaluate the consequences of the NJSLA–S as NJDOE is committed to evaluating the effects of the NJSLA–S.

PART 10: REPORTING

Standard 6.10 states that “[w]hen test score information is released, those responsible for testing programs should provide interpretations appropriate to the audience” (p. 119). The NJSLA–S score reports were designed to effectively communicate test scores while avoiding possible misinterpretations or over-interpretation of the figures. This means that the score reports only show scale scores and performance levels rather than raw scores. This section briefly describes the five different reports that were produced for the NJSLA–S. An example of each of the five reports is explained and presented below. More comprehensive descriptions of each component within the reports can be found in the NJSLA–S Score Interpretation Guide (SIG) at the [NJSLA–S website](#) under **NJSLA–Science Guides**. Two versions of the SIG are publicly available. One version is targeted to educators, administrators, and other district personnel who need to understand the score reports. The other version is targeted to parents and focuses on the Individual Student Reports.

10.1 Individual Student Report

The Individual Student Report (ISR) is a two-sided document intended for use by students, parents, teachers, and other school personnel who have to know a student’s strengths and weaknesses in science. It shows the student scale score; the Reporting PLD associated with the student’s performance; data for comparison across the state, district, and school; subscore performance levels; and descriptions of the Near/Met Expectations performance level for each subscore. Figures 10.1.1 and 10.1.2 show examples of the front and back of an ISR. A complete list of Reporting PLDs can be found in Appendix E.



New Jersey Student Learning Assessment - Science (NJSLA-S)
Individual Student Report

This report shows how FIRSTNAME004 performed on the high school science assessment. **This assessment is just one measure of how well your child is performing academically. The results from this assessment should be used in combination with other indicators of achievement in drawing conclusions about your student's performance in science.**

Visit the NJ Parent Portal at nj-results.pearsonaccessnext.com and use this code to access your student's results online.

Rz8Nppysqw7T

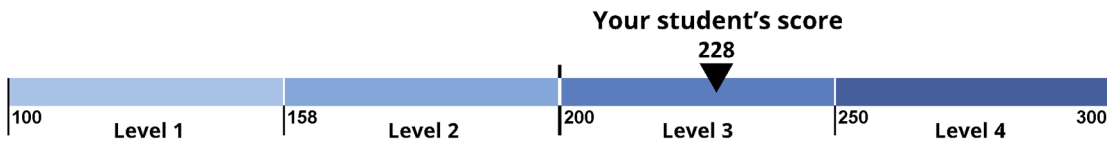
How did FIRSTNAME004 perform on the NJSLA-S?

Your student's score: **228**

Performance: **Level 3**

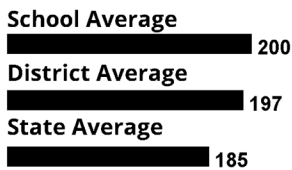
Proficient

- Level 4** (250 – 300) Advanced Proficiency
- Level 3** (200 – 249) Proficient
- Level 2** (158 – 199) Near Proficiency
- Level 1** (100 – 157) Below Proficient

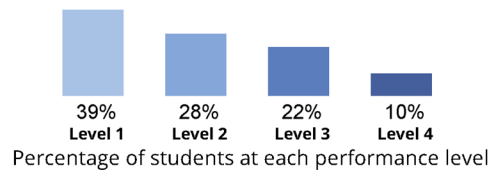


FIRSTNAME004's score on the NJSLA-S indicates that your student is at Level 3.

Students who are at Level 3 demonstrated appropriate grade-level understanding of the New Jersey Student Learning Standards-Science (NJSLA-S) by comprehending information from a variety of sources (e.g., text, charts, graphs, tables) and applying the knowledge gained from scientific investigations to develop accurate explanations and models of observed phenomena. The students often chose and used the appropriate tools to make observations and to gather, classify, and present data. The students used both essential and non-essential information to recognize patterns and relationships between data and designed systems. The students were able to use information to make real-world connections and predictions.



How Students Statewide Performed



See page 2 of this report for specific information on your student's performance using the science domains and practices.

Figure 10.1.1. Sample Individual Student Report–Page 1

How did your student perform using the domains and practices?

The domains are the content components related to specific disciplines of science.

The practices are methods by which scientists investigate and build models and theories about the world.



Earth & Space Science

Your student's performance is **Below Expectations**.

A student designated as Near/Met Expectations demonstrates knowledge of the processes that operate on and within the Earth and also its place in the solar system and galaxy.



Investigating Practices

Your student's performance is **Above Expectations**.

A student designated as Near/Met Expectations asks questions, plans and carries out investigations based on observations of phenomena, and organizes the data effectively.



Life Science

Your student's performance is **Near/Met Expectations**.

A student designated as Near/Met Expectations demonstrates knowledge of patterns, processes, and relationships of living organisms.



Sensemaking Practices

Your student's performance is **Below Expectations**.

A student designated as Near/Met Expectations recognizes patterns and relationships in data to develop explanations or models of the phenomena.



Physical Science

Your student's performance is **Below Expectations**.

A student designated as Near/Met Expectations demonstrates knowledge of the mechanisms of cause and effect in all systems and processes that can be understood through a common set of physical and chemical processes.



Critiquing Practices

Your student's performance is **Near/Met Expectations**.

A student designated as Near/Met Expectations evaluates and creates arguments regarding different explanations and claims to convey a deeper understanding of the natural world.



How will my student's school use the test results?

Results from the test give your student's teacher information about their academic performance. The results also give your school and school district important information to make improvements to the education program.

Learn more about the New Jersey Student Learning Assessment — Science

For more information about the assessment, sample questions, practice tests, and the Score Interpretation Guide (SIG) for this report please visit www.measinc.com/nj/science.

Learn More about the New Jersey Learning Standards

Explore your school website, or ask your principal, for information on your school's annual assessment schedule; the curriculum chosen by your district to give students more hands-on learning experiences that meet state standards; and to learn more about how test results contribute to school improvements. You can also learn more about New Jersey's K-12 standards at <https://www.nj.gov/education/standards/science/Index.shtml>.



New Jersey Student Learning Assessment - Science (NJSLA-S)
Grade 5

Purpose: This report describes student performance in terms of scale score, and using domains and practices, in comparison to school, district and state averages.	TOTAL NUMBER OF STUDENT RECORDS*	NUMBER OF STUDENT WITH VALID SCORES**	AVERAGE SCALE SCORE	Student Performance Using Domains and Practices (Percent)								
				EARTH & SPACE SCIENCE	LIFE SCIENCE	PHYSICAL SCIENCE	INVESTIGATING PRACTICES	SENSEMAKING PRACTICES	CRITIQUIING PRACTICES			
STATE	102,628	101,221	225	36 21 43	24 63 13	33 21 46	36 21 43	24 63 13	33 21 46			
DISTRICT	72	69	201	13 58 29	24 20 56	35 35 30	13 58 29	24 20 56	35 35 30			
SCHOOL	19	15	180	34 42 24	46 37 17	29 60 11	34 42 24	46 37 17	29 60 11			
STUDENT	SID	DOB	SE	ELL	SCALE SCORE	INDIVIDUAL STUDENT PERFORMANCE INDICATOR						
ALASTNAME, FIRSTNAME M.	0123456789	02/02/2009	504	Y	259	4	✓	✓	✓	≈	✓	✓
DLASTNAME, FIRSTNAME M.	0123456789	02/02/2009	IEP	F1	233	3	!	✓	!	!	≈	✓
ELASTNAME, FIRSTNAME M.	0123456789	02/02/2009	B	F2	115	1	✓	!	≈	!	!	!
FLASTNAME, FIRSTNAME M.	0123456789	02/02/2009	N	-	167	2	✓	≈	✓	≈	!	✓
GLASTNAME, FIRSTNAME M.	0123456789	02/02/2009	N	Y	Not Tested - 01							
HLASTNAME, FIRSTNAME M.	0123456789	02/02/2009	N	F4	241	3	✓	✓	≈	≈	✓	!
ILASTNAME, FIRSTNAME M.	0123456789	02/02/2009	IEP	F3	137	1	!	!	✓	!	≈	!
JLASTNAME, FIRSTNAME M.	0123456789	02/02/2009	N	R	172	2	✓	≈	!	✓	!	✓

1 Below Proficient (100-149)	2 Near Proficiency (150-199)	3 Proficient (200-242)	4 Advanced Proficiency (243-300)
------------------------------	------------------------------	------------------------	----------------------------------

! Below Expectations	≈ Near/Met Expectations	✓ Above Expectations
----------------------	-------------------------	----------------------

Districts may assign Not Tested or Void codes for students that did not receive a scale score.

For more information see the Score Interpretation Guide at www.measinc.com/nj/s/science.

* Total Number of Student Records - The number of students registered for the test.

** Number of Students with Valid Scores - The number of students who took the test and completed enough items for the test to be scored.

Figure 10.3.1. Sample Student Roster

10.4 School Summary and District Summary of Schools

The NJSLA–S School Summary and District Summary of Schools reports display aggregate student performance at the state, district, and school levels. The School Summary shows only one school while the District Summary of Schools shows all the schools in a district. Other aggregations include gender, ethnicity/race, disability status, and English learner status. Aggregate student performance is illustrated by the percentages of students with each subscore performance level. Figure 10.4.1 displays an example of the School Summary report. Figure 10.4.2 displays an example of the District Summary of Schools report.

SCHOOL SUMMARY

Grade 5



STATE OF NEW JERSEY
DEPARTMENT OF EDUCATION

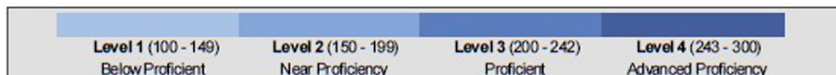
CONFIDENTIAL - DO NOT DISTRIBUTE

SAMPLE DISTRICT NAME
SAMPLE SCHOOL NAME
NEW JERSEY

SPRING 2023

New Jersey Student Learning Assessment - Science (NJSLA-S) Grade 5

Purpose: This report describes group performance in using the domains and practices, in comparison to state and district averages. PERFORMANCE DISTRIBUTION BY %	Number of Students with Valid Scores	Student Performance Using Domains and Practices (Percent)					
		EARTH & SPACE SCIENCE	LIFE SCIENCE	PHYSICAL SCIENCE	INVESTIGATING PRACTICES	SENSEMAKING PRACTICES	CRITIQUING PRACTICES
STATE	99,999						
DISTRICT	5,664						
SAMPLE SCHOOL NAME	204						



For more information see the Score Interpretation Guide at www.measinc.com/nj/s/science.



Figure 10.4.1. Sample School Performance Level Summary Report—Domains and Practices

DISTRICT SUMMARY OF SCHOOLS

Grade 5



STATE OF NEW JERSEY
DEPARTMENT OF EDUCATION

CONFIDENTIAL - DO NOT DISTRIBUTE

SAMPLE DISTRICT NAME

NEW JERSEY

SPRING 2023

New Jersey Student Learning Assessment - Science (NJSLA-S) Grade 5

Purpose: This report describes group performance in using the domains and practices, in comparison to state and district averages.	Number of Students with Valid Scores	Student Performance Using Domains and Practices (Percent)					
		EARTH & SPACE SCIENCE	LIFE SCIENCE	PHYSICAL SCIENCE	INVESTIGATING PRACTICES	SENSEMAKING PRACTICES	CRITIQUING PRACTICES
PERFORMANCE DISTRIBUTION BY %							
STATE	99,999						
		36 21 43	24 63 13	33 21 46	36 21 43	24 63 13	33 21 46
DISTRICT	5,664						
		13 58 29	24 20 56	35 35 30	13 58 29	24 20 56	35 35 30
ABRAHAM LINCOLN MIDDLE SCHOOL	204						
		34 42 24	46 37 17	29 60 11	34 42 24	46 37 17	29 60 11
ADA LOVELACE MIDDLE SCHOOL	198						
		21 79 0	12 57 31	33 40 27	21 79 0	12 57 31	33 40 27
BENJAMIN FRANKLIN MIDDLE SCHOOL	177						
		29 18 53	22 64 14	29 22 49	29 18 53	22 64 14	29 22 49
BOOKER T. WASHINGTON MIDDLE SCHOOL	204						
		11 57 32	28 20 52	35 34 30	11 57 32	28 20 52	35 34 30
CHARLOTTE HAWKINS BROWN MIDDLE SCHOOL	198						
		37 42 21	47 39 14	32 60 8	37 42 21	47 39 14	32 60 8
ELEANOR ROOSEVELT MIDDLE SCHOOL	177						
		29 60 11	12 49 39	35 41 24	29 60 11	12 49 39	35 41 24



For more information see the Score Interpretation Guide at www.measinc.com/nj/s/science.

Figure 10.4.2. Sample District Performance Level Summary Report—Domains and Practices

10.5 School and District Performance Level Summary Reports

The NJSLA–S School and District Performance-Level Summary reports display aggregate student performance for the state and district. The School Performance-Level Summary also shows student performance at the school level. Other aggregations for the district or school include gender, ethnicity/race, disability status, and English learner status. Aggregate student performance is illustrated by the average scale score and the percentages of students in each performance-level classification. Figures 10.5.1 and 10.5.2 display examples of the School and District Performance-Level Summary reports.

SCHOOL PERFORMANCE LEVEL SUMMARY

Grade 5



STATE OF NEW JERSEY
DEPARTMENT OF EDUCATION

CONFIDENTIAL - DO NOT DISTRIBUTE

SAMPLE DISTRICT NAME
SAMPLE SCHOOL NAME
NEW JERSEY
SPRING 2023

New Jersey Student Learning Assessment - Science (NJSLA-S) Grade 5

Purpose: This report describes group achievement in terms of average scale scores and performance levels.	Total Number of Student Records	No Scores Reported	Number of Students with Valid Scores	Average Scale Score	Performance Levels								≥ Level 3	
					Level 1 Below Proficient		Level 2 Near Proficiency		Level 3 Proficient		Level 4 Advanced Proficiency			
					#	%	#	%	#	%	#	%	#	%
State	999,999	999,999	999,999	999	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%
District	999,999	999,999	999,999	999	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%
School	999,999	999,999	999,999	999	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%
Gender														
Female	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Male	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Non-Binary/Undesignated	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Ethnicity/Race														
Hispanic or Latino	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
American Indian or Alaska Native	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Asian	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Black or African-American	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Native Hawaiian or Other Pacific Islander	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
White	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Two or more races	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Not Indicated	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Students with Disabilities														
IEP - Yes	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
504	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
English Language Learner														
Current EL	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Former EL	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Other														
Economically Disadvantaged	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Non-Economically Disadvantaged	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Homeless	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Migrant	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%

For more information see the Score Interpretation Guide at www.measinc.com/nj/science.

Figure 10.5.1. Sample School Performance Level Summary Report

DISTRICT PERFORMANCE LEVEL SUMMARY

Grade 5



STATE OF NEW JERSEY
DEPARTMENT OF EDUCATION

CONFIDENTIAL - DO NOT DISTRIBUTE

SAMPLE DISTRICT NAME

NEW JERSEY

SPRING 2023

New Jersey Student Learning Assessment - Science (NJSLA-S) Grade 5

Purpose: This report describes group achievement in terms of average scale scores and performance levels.	Total Number of Student Records	No Scores Reported	Number of Students with Valid Scores	Average Scale Score	Performance Levels								≥ Level 3	
					Level 1 Below Proficient		Level 2 Near Proficiency		Level 3 Proficient		Level 4 Advanced Proficiency			
					#	%	#	%	#	%	#	%	#	%
State	999,999	999,999	999,999	999	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%
District	999,999	999,999	999,999	999	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%	999,999	999.9%
Gender														
Female	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Male	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Non-Binary/Undesignated	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Ethnicity/Race														
Hispanic or Latino	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
American Indian or Alaska Native	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Asian	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Black or African-American	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Native Hawaiian or Other Pacific Islander	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
White	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Two or more races	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Not Indicated	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Students with Disabilities														
IEP - Yes	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
504	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
English Language Learner														
Current EL	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Former EL	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Other														
Economically Disadvantaged	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Non-Economically Disadvantaged	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Homeless	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%
Migrant	99,999	99,999	99,999	999	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%	99,999	99.9%

For more information see the Score Interpretation Guide at www.measinc.com/nj/science.

Figure 10.5.2. Sample District Performance Level Summary Report

REFERENCES

- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for Educational and Psychological Testing*. Washington, D.C.: American Educational Research Association.
- Angoff, W.H., & Ford, S.F. (1973). Item-race interaction on a test of scholastic aptitude. *Journal of Educational Measurement*, 10(2), 95–106.
- Bentler, P.M. (1990). Comparative fit indexes in structural models. *Psychological Bulletin*, 107, 238-246.
- Beauducel, A., & Herzberg, P.Y. (2006). On the performance of maximum likelihood versus means and variance adjusted weighted least squares estimation in CFA. *Structural Equation Modeling*, 13, 186–203.
- Brennan, R. L. (2004). *Manual for BB-CLASS: A computer program that uses the beta-binomial model for classification consistency and accuracy (version 1)*. CASMA Research Report 9. Iowa City, IA.
- Brown, W. (1910). Some experimental results in the correlation of mental abilities. *British Journal of Psychology*, 3, 296–322.
- Brown, T. A. (2006). *Confirmatory Factor Analysis for Applied Research*. New York, NY: The Guilford Press.
- Brown, M.W., & Cudeck, R. (1993). Alternative ways of assessing model fit. In K.A. Bollen & J.S. Long (Eds.), *Testing structural equation models* (pp. 136–162). Sage.
- Chen, W.-H., & Thissen, D. (1997). Local dependence indexes for item pairs using item response theory. *Journal of Educational and Behavioral Statistics*, 22, 265–289.
- Cheung, G.W., & Rensvold, R.B. (2002). Evaluating goodness-of-fit indexes for testing measurement invariance. *Structural Equation Modeling*, 9, 233–255.
- Cizek, G.J., & Bunch, M.B. (2007). *Standard Setting*. Sage Publications.
- Cizek, G.J. (2016). Validating test score meaning and defending test score use: different aims, different methods. *Assessment in Education: Principles, Policy, & Practice*, 23(2), 212–225.
- Clauser, B. E., & Mazor, K. M. (1998). Using statistical procedures to identify differentially functioning test items. *Educational Measurement: Issues and Practice*, 17, 31–44.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20, 37–46.
- Cohen, J. (1968). Nominal scale agreement with provision for scaled disagreement or partial credit. *Psychological Bulletin*, 70, 213–220.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16, 297–334.
- Cronbach, L. J. (1980). Validity on parole: How can we go straight? *New Directions for Testing and Measurement: Measuring Achievement over a Decade*, 5, 99–108.

- 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). Lawrence Erlbaum.
- Educational Testing Service. (2015). *ETS guidelines for fair tests and communications*. Educational Test Service.
- Engelhard, G. & Wang, J. (2021). *Rasch model for solving measurement problems: Invariant measurement in the social sciences*. Sage Publications, Inc.
- Engelhard, G., & Wind, S. A. (2018). *Invariant measurement with raters and rating scales: Rasch models for rater-mediated assessments*. Routledge.
- Every Student Succeeds Act, 20 U.S.C. § 6301 (2015). <https://www.congress.gov/bill/114th-congress/senate-bill/1177>
- Forte, E. (2013). Re-conceptualizing alignment in the evidence-centered design context [Paper presentation]. The Annual Meeting of the American Educational Research Association, San Francisco, CA, United States.
- Forte, E. (2017). Evaluating alignment in large-scale standards-based assessment systems [White paper]. Council of Chief State School Officers.
- Garrido, L.E., Abad, F.J., & Ponsoda, V. (2016). Are fit indices really fit to estimate the number of factors with categorical variables? Some cautionary findings via monte carlo simulation. *Psychological Methods*, 21(1), 93-111.
- Gorsuch, R. L. (1983). *Factor Analysis*. Lawrence Erlbaum Associates.
- Hambleton, R. K. (1989). Principles and selected applications of item response theory. In R. L. Linn (Ed.), *Educational measurement* (3rd ed.). American Council on Education and Macmillan.
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Kluwer Academic Publishers.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Sage.
- Hambleton, R. K., & van der Linden (1982). Advanced in item response theory and applications: An introduction. *Applied Psychological Measurement*, 6, 373–378.
- Hills, J. R. (1989). Screening for potentially biased items in testing programs. *Educational Measurement: Issues and Practice*, 8, 5–11.
- Hu, L., & Bentler, P.M. (1999). Cutoff criteria for fit indexes in covariance structure analysis: Conventional criteria versus new alternatives. *Structural Equation Modeling*, 6, 1–55.
- Jöreskog, K.G., & Sörbom, D. (1988). *LISREL 7. A guide to the program and applications* (2nd ed.). International Education Services.
- Kane, M. T. (2001). So much remains the same: Conception and status of validation in setting standards. In G. Cizek (Ed.), *Setting performance standards: Concepts, methods and perspectives* (pp. 53–88). Erlbaum

- Kane, M.T. (2006). Validation. In R.L. Brennan (Ed.), *Educational measurement* (pp. 17–64). American Council on Education/Praeger.
- Kane, M.T. (2013). Validating the interpretations and uses of test scores. *Journal of Educational Measurement*, 50, 1–73.
- Kline, R.B. (2011). *Principles and practice of structural equation modeling* (3rd edition). The Guilford Press.
- Koffler, S. (2019). NJSLA–S Cut Score Evaluation. New Jersey Technical Advisory Committee.
- Kolen, M. J., & Brennan, R. L. (2004). *Test equating: Methods and practice*. Springer.
- Kuder, G. F., & Richardson, M. W. (1937). The theory of the estimation of test reliability. *Psychometrika*, 2, 151–160.
- Lane, S., & Stone, C. A. (2002). Strategies for examining the consequences of assessment and accountability programs. *Educational Measurement: Issues and Practice*, 21, 23–30.
- Linacre, J.M. (2002). What do infit and outfit, mean-square and standardized mean? *Rasch Measurement Transactions*, 16(2), 878.
- Linacre, J.M. (2016). *A User's Guide to Winsteps MINISTEP Rasch-Model Computer Programs*.
- Liu, I-M., & Agresti, A. (1996). Mantel-Haenszel-type inference for cumulative odds ratios with a stratified ordinal response. *Biometrics*, 52, 1223–1234.
- 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 mental test scores*. Addison-Wesley.
- Mantel, N., & Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of National Cancer Institute*, 22, 719–748.
- Masters, G. N. (1982). A Rasch model for partial credit scoring. *Psychometrika*, 47, 149–174.
- McNeill, K.L., Katsh-Singer, R., & Pelletier, P. (2015). Assessing science practices: Moving your class along a continuum. *Science Scope*, 39, 21–28
- Messick S. (1989). Validity. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 13–104). American Council on Education.
- Messick, S. (1995). Validity of psychological assessment: Validation of inferences from persons' responses and performances as scientific inquiry into score meaning. *American Psychologist*, 50, 741–749.
- Miller, G., Rotou, O., & Twing, J. (2004). Evaluation of the 0.3 logits screening criterion in common item equating. *Journal of Applied Measurement*, 5, 172–177.
- Muthén, B.O. (1993). Goodness of fit with categorical and other nonnormal variables. In K.A. Bollen & J.S. Long (Eds.), *Testing structural equation models* (pp. 205–243). Sage.
- Muthén, B.O., du Toit, S.H.C, & Spisic, D. (1997). Robust inference using weighted least squares and quadratic estimating equations in latent variable modeling with categorical and continuous outcomes [Unpublished manuscript]. <https://www.statmodel.com/wlscv.shtml>

- National Research Council (2012). A framework for K-12 science education: Practices, crosscutting concepts and core ideas. National Academies Press.
- New Jersey Department of Education (2019). New Jersey Student Learning Assessment–Science (NJSLA–S): Technical report grades 5, 8, and 11 2019. <https://measinc-nj-science.com/sites/default/files/2022-03/2019%20NJSLA-S%20Technical%20Report-Accessible.pdf>
- New Jersey Department of Education (2024). Spring 2024 District Test and District Technology Coordinator Training: New Jersey Student Learning Assessments (NJSLA) and New Jersey Graduation Proficiency Assessment (NJGPA) [Training materials]. https://nj.mypearsonsupport.com/resources/test-administration-resource/2024_NJSLA_NJGPA_DTC_Training_Part1_V1_1.pdf
- NGSS Lead States (2013). Next generation science standards: For states, by states. The National Academies Press.
- Nye, C.D., & Drasgow, F. (2011). Assessing goodness of fit: Simple rules of thumb simply do not work. *Organizational Research Methods*, 14(3), 548-570.
- Ostini, R., & Nering, M. L. (2010). New perspectives and applications. In M. L. Nering & R. Ostini (Ed.), *Handbook of Polytomous Item Response Models* (pp. 3-20). Routledge.
- Padilla, J. L., & Benitez, I. (2014). Validity evidence based on response processes. *Psicothema*, 26, 136–144.
- Penfield, R. D. (2007). An approach for categorizing DIF in polytomous items. *Applied Measurement in Education*, 20, 335–355.
- Penfield, R. D., & Algina, J. (2003). Applying the Liu-Agresti estimator of the cumulative odds ratio to DIF detection in polytomous items. *Journal of Educational Measurement*, 40, 353–370.
- Penfield, R. D., & Algina, J. (2006). A generalized DIF effect variance estimator for measuring unsigned differential test functioning in mixed format tests. *Journal of Educational Measurement*, 43, 295–312.
- Satorra, A., & Bentler, P.M. (1994). Corrections to test statistics and standard errors in covariance structure analysis. In A. von Eye & C.C. Clogg (Eds.), *Latent variables analysis: Applications for developmental research* (pp. 399-419). Sage Publications, Inc.
- Shi, D., & Maydeu-Olivares, A. (2020). The effect of estimation methods on SEM fit indices. *Educational and Psychological Measurement*, 80(3), 421-445.
- Sireci, S. G., Baldwin, P., Martone, A., Zenisky, A. L., Kaira, L., Lam, W., Shea, C. L., Han, K., Deng, N., Delton, J., & Hambleton, R. K. (2008). *Massachusetts Adult Proficiency Tests Technical Manual*. Massachusetts Department of Elementary and Secondary Education.
- Smarter Balanced Assessment Consortium. (2018). Smarter Balanced Assessment Consortium: 2017–18 summative technical report. Retrieved from <https://portal.smarterbalanced.org/library/en/2017-18-summative-assessment-technical-report.pdf>
- Steiger, J.H. (1990). Structural model evaluation and modification: An interval estimation approach. *Multivariate Behavioral Research*, 25, 173-180.

- Swineford, F. (1949). Law School Admission Test - WLS. Educational Testing Service. <https://onlinelibrary.wiley.com/doi/pdf/10.1002/j.2333-8504.1949.tb00915.x>
- Thek, A.D., Hoole, E.R., & Lottridge, S.M. (2006). What are you thinking? Postsecondary student think-alouds of scientific and quantitative reasoning tasks. *The Journal of General Education*, 55(1), 17–39.
- Thek, A., Sundre, D. L., Horst, S. J., & Finney, S. J. (2009). Examining inferences about test-taking motivation: The Student Opinion Scale (SOS) [Paper presentation]. Annual meeting of the Northeastern Educational Research Association, Rocky Hill, CT, United States.
- Thompson, S. J., Johnstone, C. J., & Thurlow, M. L. (2002). Universal Design Applied to Large Scale Assessments (Synthesis Report 44). National Center on Educational Outcomes.
- Traub, R. E., & Rowley, G. L. (1991). NCME instructional module: Understanding reliability. *Educational Measurement: Issues and Practice*, 10(1), 37-45.
- Webb, N. L. (1997). Criteria for alignment of expectations and assessments in mathematics and science education (Research monograph No. 6). Council of Chief State School Officers.
- Webb, N.L. (1999). Alignment study in language arts, mathematics, science, and social studies of state standards and assessments for four states. Council of Chief State School Officers.
- Webb, N. L. (2002). Alignment study in language arts, mathematics, science, and social studies of state standards and assessments for four states. Washington, DC: Council of Chief State School Officers.
- Wright B.D., & Linacre J.M. (1994). Reasonable mean-square fit values. *Rasch Measurement Transactions*, 8(3), 370.
- Wright B.D., & Masters G.N. (1982). *Rating Scale Analysis*. MESA Press.
- Xia, Y., & Yang, Y. (2019). RMSEA, CFI, and TLI in structural equation models with ordered categorical data: The story they tell depends on the estimation methods. *Behavioral Research Methods*, 51, 409-428.
- Yen, W. M. (1984). Effects of local item dependence on the fit and equating performance of the three-parameter logistic model. *Applied Psychological Measurement*, 8, 125–145.
- Zieky, M. (1993). Practical questions in the use of DIF statistics in item development. In P. W. Holland, & H. Wainer (Eds.), *Differential item functioning* (pp. 337–347). Erlbaum.

APPENDIX A: GLOSSARY OF ABBREVIATIONS

Table A.1: Glossary of NJSLA–S Abbreviations

Abbreviation	Definition
ABBI	Assessment Banking for Building Interoperability
AERA	American Educational Research Association
AF&A	Accessibility Features and Accommodations
APA	American Psychological Association
ASL	American Sign Language
CBT	Computer-Based Test
CCC	Crosscutting Concept
CFA	Confirmatory Factor Analysis
CR	Constructed Response
CSEM	Conditional Standard Error of Measurement
CTT	Classical Test Theory
DCI	Disciplinary Core Idea
DIF	Differential Item Functioning
DTC	District Test Coordinator
EconDis	Economically Disadvantaged
EL	English Learner
ESEA	Elementary and Secondary Education Act
ESSA	Every Student Succeeds Act
ICC	Item Characteristic Curve
IIF	Item Information Function
IRT	Item Response Theory
ISR	Individual Student Report
KIS	Key Information Sheet
KSA	Knowledge, Skills, and Abilities
LEA	Local Education Agency
MC	Multiple choice
MH	Mantel-Haenszel
MI	Measurement Incorporated
MSA	Machine-Scorable Assessment
NBP	National Braille Press
NCME	National Council on Measurement in Education
NJASK	New Jersey Assessment of Skills and Knowledge
NJBCT	New Jersey Biology Competency Test
NJBSC	New Jersey Bias and Sensitivity Committee

Abbreviation	Definition
NJDOE	New Jersey Department of Education
NJSAC	New Jersey Science Advisory Committee
NJTAC	New Jersey Technical Advisory Committee
NJSLA–S	New Jersey Student Learning Assessment–Science
NJSLS–S	New Jersey Student Learning Standards for Science
NRC	National Research Council
OIB	Ordered Item Booklet
OPLS	Online Performance-Level Setting
PAN	PearsonAccess ^{next}
PBA	Performance-Based Assessment
PBS	Phenomenon-Based Scenario
PBT	Paper-Based Test
PCA	Principal Components Analysis
PCM	Partial Credit Model
PIA	Preliminary Item Analysis
PLD	Performance-Level Descriptor
<i>rpb</i>	Item-Total Correlation
SEM	Standard Error of Measurement
SEP	Science and Engineering Practice
SIG	Score Interpretation Guide
SRF	Summative Record File
STC	School Test Coordinator
SWD	Students with Disabilities
TA	Test Administrator
TCM	Test Coordinator Manual
TE	Technology-enhanced
TIF	Test Information Function
TLC	Teneo Linguistics Company
TTS	Text-to-Speech

APPENDIX B: NEW JERSEY SCIENCE ADVISORY AND BIAS AND SENSITIVITY COMMITTEES—DISTRICT AND COUNTY REPRESENTATION

Table B.1: Grade 5 NJSAC District and County Representation

Number	District	School County
1	Glen Rock Public School District	Bergen
2	River Edge School District	Bergen
3	Northern Burlington County Regional School District	Burlington
4	Avalon School District	Cape May County
5	Livingston Public Schools	Essex
6	Swedesboro-Woolwich School District	Gloucester
7	Swedesboro-Woolwich School District	Gloucester
8	Jersey City Global CS	Hudson
9	Readington Township School District	Hunterdon
10	Lawrence Township Public School District	Mercer
11	West Windsor-Plainsboro Regional School District	Mercer
12	West Windsor Plainsboro Regional School District	Mercer
13	Metuchen Public School District	Middlesex
14	Rumson Borough School District	Monmouth
15	Washington Township School District	Morris
16	Brick Township Public School District	Ocean
17	Cranford Public School District	Union

Table B.2: Grade 8 NJSAC District and County Representation

Number	District	School County
1	Franklin Lakes School District	Bergen
2	Leonida Public School District	Bergen
3	Lyndhurst School District	Bergen
4	Lyndhurst School District	Bergen
5	Cinnaminson Township Public Schools	Burlington
6	Maria L. Varisco-Rogers Charter School	Essex
7	Clinton Township School District	Hunterdon
8	Melvin H. Kreps Middle School	Mercer
9	West Windsor-Plainsboro Regional School District	Mercer
10	East Brunswick Township School District	Middlesex
11	New Brunswick School District	Middlesex
12	North Brunswick Township School District	Middlesex
13	Matawan Aberdeen Regional School District	Monmouth
14	Mount Olive Township School District	Morris
15	Memorial Middle School	Ocean
16	Passaic City School District	Passaic
17	Berkeley Heights Board of Education	Union
18	Roselle Park Public School District	Union

Table B.3: Grade 11 NJSAC District and County Representation

Number	District	School County
1	Atlantic County Institute of Technology	Atlantic
2	Greater Egg Harbor Regional High School District	Atlantic
3	Moorestown Township Public Schools	Burlington
4	Cherry Hill School District	Camden
5	South Orange/Maplewood	Essex
6	Greater Egg Harbor Regional High School District	Hudson
7	Jersey City Public Schools	Hudson
8	Princeton Public Schools	Mercer
9	West Windsor-Plainsboro Regional School District	Mercer
10	Asbury Park School District	Monmouth
11	Jefferson Township Public Schools	Morris
12	Parsippany-Troy Hills School District	Morris
13	Paramus Public School District	Paramus
14	Passaic Academy for Science and Engineering	Passaic
15	Paterson Public Schools	Passaic
16	Paterson Charter School for Science/Technology	Passaic
17	Pennsville Public School District	Salem
18	Somerset County Vocational & Technical High School	Somerset

Table B.4: NJBSC District and County Representation

Number	District	School County
1	Oakland Public School District, Curriculum Office	Bergen
2	Cherry Hill School District	Camden
3	Millburn Township Public Schools	Essex
4	Jersey City Global Charter	Hudson
5	East Brunswick (Retired)	Middlesex
6	Freehold Township District	Monmouth
7	Morris County School of Technology	Morris
8	Mt. Olive Public School District	Morris
9	Clifton Public School District	Passaic
10	Paterson Public School District	Passaic

APPENDIX C: STATISTICAL REVIEW REFERENCE SHEET

<p><i>p-value</i> (All items: 1-point to 4-point items)</p>	<ul style="list-style-type: none">• A measure of item difficulty based on classical test theory• Proportion correct; the proportion of students who answered a dichotomous (1-point) item correctly• Percentage of maximum score point; item mean divided by the highest attainable score point for a polytomous (2- or 4-point) item• <i>P-values</i> can range from 0 to 1.
<p>*FLAGGED IF:</p>	<p><i>P < .25 (too hard)</i> <i>P > .90 (too easy)</i></p>
<p>RASCH VALUE (All items: 1-point to 4-point items)</p>	<ul style="list-style-type: none">• A measure of item difficulty based on item response theory, with values generally ranging from –3 to +3. Higher values indicate greater difficulty (reverse of <i>p-value</i>)• Items with Rasch values targeted at cut scores for performance categories are especially useful for measurement.
<p>SCORE POINT DISTRIBUTION (2-/4-point items)</p>	<ul style="list-style-type: none">• Percentage of responses at each score point• If any score point has fewer than 10% of responses (2-point item) or 5% of responses (4-point item), the score point is not measuring relevant ability effectively.
<p>*FLAGGED IF:</p>	<p><i>Response percentage < 10% at any score point (2-point items)</i> <i>Response percentage < 5% at any score point (4-point items)</i></p>
<p>ITEM-TOTAL CORRELATION (All items: 1-point to 4-point items)</p>	<ul style="list-style-type: none">• A measure of the degree to which an item discriminates between those students who know the material (using total test score as a proxy for that knowledge) and those who do not• <i>rpb</i>: correlation between an item and the total test score• <i>rpb</i> can range from –1 to +1.
<p>*FLAGGED IF:</p>	<p><i>rpb < .20 (All items, 1-point to 4-point items)</i></p>

DIF CATEGORY
(All items: 1-point
to 4-point items)

- A statistical procedure for detecting potential item bias
- Differential item functioning (DIF) categorization looks at the extent to which an item performs differently across different groups—Male/Female, White/Black, White/Hispanic, and White/Asian—controlling for the groups’ ability (using total test score as a proxy).
- Each item is classified as A, B, or C:
 - A: Item displays negligible DIF; does not need review for bias.
 - B: Item displays moderate DIF; needs review for bias.
 - C: Item displays severe DIF; needs *careful* review for bias.

***FLAGGED IF:**

DIF CATEGORY = B or C

APPENDIX D: 2019 NJSLA–S STANDARD SETTING: EXECUTIVE SUMMARY

Appendix D contains the executive summary from the standard setting report submitted by Measurement Incorporated in 2019. The standard setting study is summarized in greater detail in the 2019 NJSLA–S technical report (NJDOE, 2019).

The New Jersey Student Learning Assessment–Science (NJSLA–S) is the assessment battery New Jersey uses to satisfy reporting requirements for the Every Student Succeeds ACT (ESSA; P.L. 115–94) for science in grades 5, 8, and 11.

The New Jersey Department of Education (NJDOE) conducted standard setting for science tests in grades 5, 8, and 11 during the week of July 23–25, 2019. Educators from throughout the state of New Jersey participated in this three-day meeting. Staff of Measurement Incorporated (MI), the contractor, and Pearson Education, its subcontractor, facilitated the meeting.

The main goals of the meeting were to

1. allow workshop participants (panelists) to gain an understanding of the test contents and performance level descriptors (PLDs),
2. learn a standard-setting procedure known as the Bookmark procedure, and
3. have panelists recommend cut scores for each test that differentiate Level 1 from Level 2, Level 2 from Level 3, and Level 3 from Level 4 performance (i.e., three cut scores to yield four performance levels).

These recommendations are designed to help inform the New Jersey State Board of Education (Board) as it completes its task of establishing performance standards for these assessments.

From July 23 through July 25, 2019, MI/Pearson staff met with representatives of NJDOE and 39 educator-panelists from around the state to recommend performance standards on the three tests.

Process and Procedures

The panelists, nominated by district superintendents, were chosen specifically to represent the demographics and geographic distribution of educators throughout the state. A profile of the 39 panelists is provided in the 2019 NJSLA–S technical report (NJDOE, 2019). Panelists spent the entire first day examining the tests and PLDs under the direction of NJDOE and MI staff. On the second day, following an introduction to the Bookmark standard-setting procedure, the panelists separated into their respective grade-level groups, each led by two facilitators (one psychometrician and one content specialist) from MI/Pearson. Panelists in all groups received a thorough orientation to the standard-setting software and practice exercises to prepare them for their standard-setting task. MI staff provided additional information to panelists as they proceeded through three rounds of recommending cut scores, discussing decisions, and settling on final recommendations.

In accordance with a plan previously approved by NJDOE, MI employed the Bookmark procedure. This procedure is the most widely used standard-setting procedure for statewide assessments and is thoroughly documented in the approved plan and elsewhere (cf. Cizek & Bunch, 2007). In this procedure, panelists review all test items in a specially formatted test

booklet (ordered item booklet, or OIB) that places the easiest item on page one, the most difficult item on the final page, and all items in between ordered by difficulty, based on actual student responses. Using threshold PLDs developed previously by NJDOE (with the assistance of New Jersey educators), panelists place a bookmark at the point in the test booklet where they believe the probability of a student at the threshold of Level 2, Level 3, or Level 4 would begin to have less than a two-thirds chance of answering correctly. These page numbers are then mathematically translated into raw cut scores. The average (median) of the panelists' bookmarked pages becomes the group bookmark, and the associated raw score becomes the cut score for that level for that grade for that round. The procedure is more fully described in Chapter 1 of the report. All reviews were completed within software created by MI and used previously for several other successful standard-setting activities.

Panelists considered each test in three online rounds. During Round 1, each panelist placed three bookmarks, one for Level 2, one for Level 3, and one for Level 4. MI staff analyzed the data for Round 1 and led discussions of the results: difficulties encountered, dispersion of bookmarks for each level, reasons for those dispersions, rationales for individual bookmark placements, and differences in interpretation of the PLDs.

After discussion of Round 1 results, panelists then started Round 2, repeating the process of placing bookmarks as they had in Round 1. After Round 2, MI staff again analyzed the data and presented results to the panelists, along with score distributions showing percentages of students who would be classified at each level on the basis of the Round 2 cut scores (impact data).

After discussion of Round 2 results and impact data, panelists once again placed three bookmarks in Round 3. These bookmarks defined the final cut scores (averaged over all panelists in a given group) to be forwarded to NJDOE. Facilitators then presented Round 3 results to panelists and gave them an opportunity to evaluate the process and outcomes. One panelist in grade 11 had to leave after Round 2.

Results

Final recommended performance standards are reported in Table ES-1. The cut scores include both the raw score associated with the median bookmark and that score expressed in terms of a percentage of the total points possible. The final column in Table ES-1 shows the total number of points possible for each test. There were no cross-grade discussions of cut scores.

Table ES-1: Final Recommendations from Standard-Setting Panelists

Grade	Level	Total Points	Raw Cut Score	Cut Score % Correct
Grade 5	Level 2	60	25	42%
Grade 5	Level 3	60	39	65%
Grade 5	Level 4	60	49	82%
Grade 8	Level 2	70	20	29%
Grade 8	Level 3	70	40	57%
Grade 8	Level 4	70	52	74%
Grade 11	Level 2	78	31	40%
Grade 11	Level 3	78	45	58%
Grade 11	Level 4	78	60	77%

The impact of these cut scores on New Jersey students is summarized in Figure ES-1. Overall, 26.3% of grade 5 students, 17.6% of grade 8 students, and 26.5% of grade 11 students scored at or above Level 3. The numbers of students upon which these percentages are based are not the entire population. By prior agreement between NJDOE and MI, the available data was analyzed during the week prior to standard setting: 64,419 fifth graders, 88,295 eighth graders, and 76,001 eleventh graders. It should be noted that special care was taken to make sure these data were representative of the entire state. Thus, when all of the data are analyzed, it is possible that the percentages in each category could change slightly.

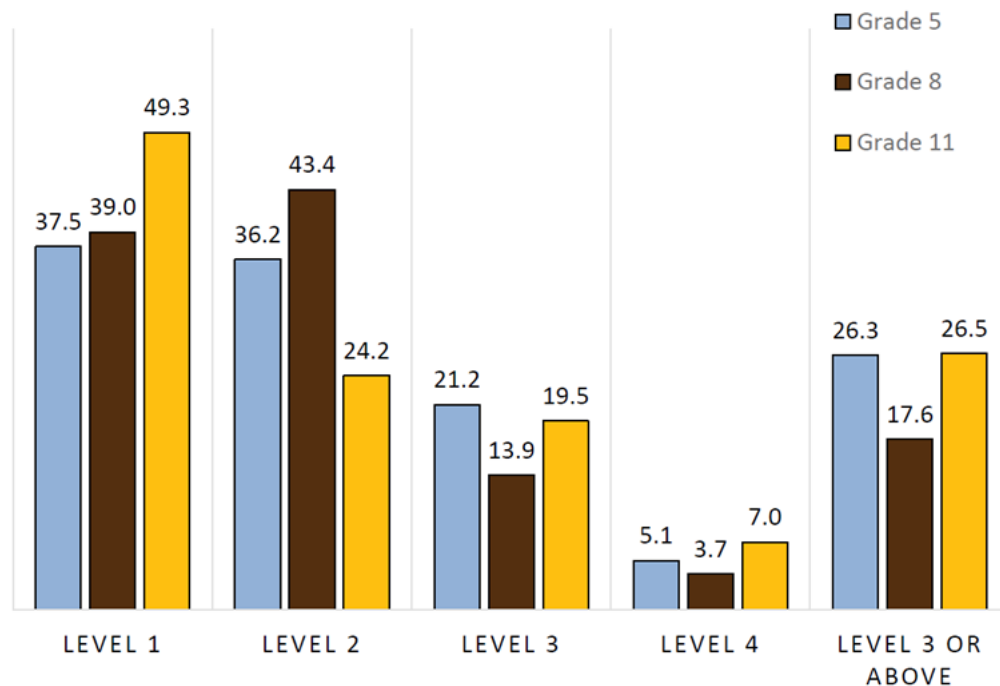


Figure ES-1. Percentages of Students Classified at Each Level after Round 3

Impact of impact data. From Round 2 to Round 3, there was some movement (in both directions) in cut scores. In grade 5; the Level 2 cut score actually went up by 1 raw score point.

At grade 8, the Level 2 cut score went down by 7 raw score points (a difference of two pages in the OIB), but the cut scores for Levels 3 and 4 did not change. One grade 8 panelist commented on the back of the evaluation form that anticipated pressure from local school administrators may have caused some panelists to lower their cut scores for Level 2. Yet, there was no change in the Level 3 or Level 4 cut scores for grade 8. At grade 11, the Level 2 and Level 3 raw cut scores went down by 4 and 2 points, respectively; the Level 4 cut score was unchanged from Round 2 to Round 3.

Evaluation of process and outcomes. The panelists were given an opportunity after presentation of Round 3 results to evaluate the entire process and outcomes. Those conducting the standard setting were especially interested in knowing how reasonable they found the final cut scores to be. Their responses to key statements on the evaluation form are summarized in Table ES-2.

Table ES-2: Responses to Key Evaluation Questions

[Responses: Grade 5–14; Grade 8–12; Grade 11–12]

Statement	% Strongly Disagree			% Disagree			% Uncertain			% Agree			% Strongly Agree		
	5	8	11	5	8	11	5	8	11	5	8	11	5	8	11
The process was fair.	0	0	0	0	0	0	0	0	8	7	17	33	93	83	58
The process was orderly.	0	0	0	0	0	0	0	0	0	7	17	33	93	83	67
My group’s final cut score for Level 2 is reasonable.	0	0	0	0	0	0	0	8	0	14	8	50	86	83	50
My group’s final cut score for Level 3 is reasonable.	0	0	0	0	0	0	0	0	0	14	17	25	86	83	75
My group’s final cut score for Level 4 is reasonable.	0	0	0	0	0	0	0	0	0	21	8	25	79	92	75

These last three statements had a follow-up direction: If you disagree, should it have been higher or lower? Circle one.

Panelists were also encouraged to enter comments on the back of the form, particularly if they disagreed with the reasonableness of any of the cut scores. The open-ended responses to the reasonableness items are summarized in Table ES-3.

Table ES-3: Summary of Reasonableness Ratings and Comments

Statement	Grade 5	Grade 8	Grade 11
My group’s final cut score for Level 2 is reasonable.	No objections; no recommended changes	No objections; one suggestion that impacts data skew Round 3 cuts	No objections; no recommended changes
My group’s final cut score for Level 3 is reasonable.	No objections; no recommended changes	No objections; no recommended changes	No objections; no recommended changes
My group’s final cut score for Level 4 is reasonable.	No objections; one recommendation to raise cut by 1	No objections; no recommended changes	No objections; no recommended changes

Summary and Recommendations

The standard setting for NJSLA–S was conducted in strict accordance with the approved plan. Panelists understood the process well, as indicated by their responses to the Evaluation Form. The standard-setting process for NJSLA–S was sound, both in conception and execution, representative of the highest standards in contemporary educational measurement, and representative of standards operating among state assessment programs nationwide. The cut scores produced after three rounds of test review reflect well the PLDs panelists used to complete the standard-setting task. It is proposed that the cut score recommendations presented here be given strong consideration for approval.

APPENDIX E: NJSLA–S PERFORMANCE-LEVEL DESCRIPTORS

E.1 Policy PLDs

NJSLA–S Policy-Level Performance-Level Descriptors

Level 1	Level 2	Level 3	Level 4
<p>Level 1 students demonstrate minimal understanding of the disciplinary concepts and have difficulty applying the scientific practices. They may have significant difficulty engaging in public discussion on scientific topics and discerning valid and reliable scientific technological information related to their everyday lives even with focused effort achieving minimal success.</p>	<p>Level 2 students demonstrate partial understanding of the disciplinary concepts and performance with the scientific practices. They may have difficulty engaging in public discussion on scientific topics and discerning valid and reliable scientific technological information related to their everyday lives without the focused effort needed to achieve some success.</p>	<p>Level 3 students demonstrate appropriate grade-level understanding of the disciplinary concepts and performance with the scientific practices. They can likely engage in public discussion on scientific topics and discern valid and reliable scientific technological information related to their everyday lives with some success.</p>	<p>Level 4 students demonstrate a deep understanding of the disciplinary concepts and superior performance with the scientific practices. They can likely engage in public discussions on scientific topics and discern valid and reliable scientific and technological information related to their everyday lives with a high degree of success.</p>

E.2 Threshold PLDs

E.2.1 Grade 5 Threshold PLDs

The Threshold Performance-Level Descriptors (PLDs) define the minimum knowledge, skills, and practices that students must display for each Disciplinary Core Idea and Science and Engineering Practice to reach a certain performance level. They expand upon the brief overall PLDs included in the Score Interpretation Guide.

Grade 5 Threshold Performance-Level Descriptors (Physical Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
PS1: Matter and Its Interactions	<ul style="list-style-type: none"> that matter is made of particles that can be identified by their properties and that weight does not change during visible physical changes that the properties of substances may change when combined, but the total weight will stay the same 	<ul style="list-style-type: none"> that matter is made of particles with unique, measurable properties that are conserved when changing state that a change to a substance(s) may or may not result in one or more new substances, but the total weight will remain the same 	<ul style="list-style-type: none"> of distinguishing properties of matter and the relationship between visible and non-visible matter that the outcome of the combination of one or more substances is predictable based on the properties of the substances
PS2: Motion and Stability: Forces and Interactions	<ul style="list-style-type: none"> that objects are acted upon by forces that can cause predictable patterns of motion that the size of a force, the properties of objects, and the position of the objects relative to one another have an effect on their interaction 	<ul style="list-style-type: none"> that an object's motion is a product of the net force acting on the object and can therefore cause predictable motion of how certain relationships among the interactions between objects are interconnected and can explain how the objects ultimately affect each other 	<ul style="list-style-type: none"> of the relationship between net force and motion of an object in predicting future movement that the relationships between the interactions and the properties of objects are dependent upon systems in which the objects exist

Grade 5 Threshold Performance-Level Descriptors (Physical Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
PS3: Energy	<ul style="list-style-type: none"> • that differences in the movement of energy can cause objects to move at different speeds • that energy in various forms can be transferred from place to place • that energy is transferred when objects collide • that energy can be converted into forms for practical use 	<ul style="list-style-type: none"> • that energy can move from place to place in different forms with varying levels of magnitude • that effects of transferred energy are observable • of the relationship between the transfer of energy and the change in motion when objects collide • that there is a relationship between energy and its conversion for practical uses 	<ul style="list-style-type: none"> • that predictions can be made regarding the interactions of objects based on the amount of energy the objects possess • of the transformation from one type of energy to other type(s) of energy • that when objects collide, there are predictable outcomes • that stored energy is converted energy from the Sun

Grade 5 Threshold Performance-Level Descriptors (Physical Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
PS4: Waves and Their Applications in Technologies for Information Transfer	<ul style="list-style-type: none"> • that there are similarities and differences in the patterns of waves • that in order for an object to be seen, light must reflect off the object • that information can be transmitted over long distances using communication methods/devices 	<ul style="list-style-type: none"> • that the characteristics of a wave determine the net motion of the wave • that there exists a relationship among the path of light, light reflection, and the visibility of objects • of how different communication methods/devices operate 	<ul style="list-style-type: none"> • of how changing the amount of energy can change the characteristics of a wave • that a change in the path of light or light reflection will cause a change in the visibility of an object • of the advantages of different communication methods/devices and how those devices transmit digitized information over long distances

Grade 5 Threshold Performance-Level Descriptors (Life Science)
Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>LS1: From Molecules to Organisms: Structures and Processes</p>	<ul style="list-style-type: none"> • of the internal or external structures of plants or animals and their functions • that animals or plants reproduce and have life cycles • that both animals and plants take in materials to survive • that animals have sense receptors that they use to guide their actions 	<ul style="list-style-type: none"> • of internal and external structures of plants and animals and how their functions support survival, growth, behavior, or reproduction • that animals and plants reproduce for continued existence and have life cycles that are unique but have some similarities • of the relationship between plants and animals and the materials they take in for specific various functions • that an animal’s brain processes information received from specialized sense receptors that they use to guide their actions 	<ul style="list-style-type: none"> • of the variation and function of internal and external structures across the plant and animal kingdoms • of the relationships among the components of life cycles • that animals and plants acquire energy from different sources but use the energy for similar functions • that animals respond to environmental changes using sensory information and stored memories

Grade 5 Threshold Performance-Level Descriptors (Life Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>LS2: Ecosystems: Interactions, Energy, and Dynamics</p>	<ul style="list-style-type: none"> • that in a food web, all organisms have a role <p>OR</p> <ul style="list-style-type: none"> • of the requirements of a healthy ecosystem • that materials cycle through an environment • that organisms respond to changes in their environment • that living in groups helps animals 	<ul style="list-style-type: none"> • that organisms have different roles in a food web, with a focus on the cycling of materials • that the health and stability of an ecosystem depend on the overall biodiversity and the availability of resources • of how materials cycle through multiple components of an environment • of organisms responding to changes in their environment • that living in specialized groups helps animals, depending on the situation 	<ul style="list-style-type: none"> • that the materials that animals consume can be traced through multiple levels of the food web back to plants • that the balance of the flow of matter can be disrupted by changes in the ecosystem • of the impact of change on the cycling of matter in a system • of how changes in an environment affect multiple organisms • that the dynamics of a group can change over time
<p>LS3: Heredity: Inheritance and Variation of Traits</p>	<ul style="list-style-type: none"> • that traits and characteristics are based on both inheritance and environmental factors • that organisms have variations in traits 	<ul style="list-style-type: none"> • that while there are similarities in traits between siblings, they each have characteristics that are influenced by the environment • that some traits are inherited in a predictable way while others may be influenced by the environment 	<ul style="list-style-type: none"> • that environmental factors affect traits or functions • that patterns in traits are expressed over multiple generations • that traits, whether inherited or influenced by the environment, have some similarities and some differences

Grade 5 Threshold Performance-Level Descriptors (Life Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>LS4: Biological Evolution: Unity and Diversity</p>	<ul style="list-style-type: none"> • that fossils are evidence of plant and animal life long ago • that variations among organisms help them survive and reproduce • that some organisms can survive in a particular environment while others cannot • that plants and animals are affected by change in their habitat 	<ul style="list-style-type: none"> • that fossils are evidence of varying environments • that certain characteristics are advantageous to the survival of a species • that an environment must meet the needs of an organism for survival • that plants and animals may adapt to changes in their environment 	<ul style="list-style-type: none"> • that fossils are evidence of changing environments over time • that specific variation in a characteristic can influence an organism's survival • that changes in an environment affect an organism's ability to survive • that the effects of habitat change may cause adaptation to occur

Grade 5 Threshold Performance-Level Descriptors (Earth and Space Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
ESS1: Earth's Place in the Universe	<ul style="list-style-type: none"> • that the Sun is an object in the sky and gives off light • that Earth is a rotating body in relative position to the Sun • that Earth's rotation affects day and night • that there are observable patterns in moon phases, shadows, and star patterns • that patterns of rock formations can contain fossils and can change due to Earth forces 	<ul style="list-style-type: none"> • that distance affects relative size • of changes in patterns (daylight hours, shadow length, stars, moon phases) that can be observed during day and night as Earth rotates and orbits around the Sun • that fossil records can help identify rock layer formations because of changes caused by natural processes 	<ul style="list-style-type: none"> • that relative distance affects brightness • that Earth's orbit and rotation at different times of day and year, together with the orbit of the Moon and position of the Sun, create patterns that affect how humans view objects from Earth • that a geological history can be determined by examining rock layers and fossil records
ESS2: Earth's Systems	<ul style="list-style-type: none"> • that Earth's four major systems can interact with each other and that components of the systems can change • that maps can be used to locate Earth's features and processes • that Earth has oceans and areas of freshwater • that weather conditions in different areas change over time • that organisms affect the environment 	<ul style="list-style-type: none"> • of how specific processes change components of Earth's four major systems and, in turn, have an effect on the systems themselves • that maps can be used to determine patterns of Earth's features and processes • of the distribution of water on Earth and its availability and accessibility • that patterns of weather form the basis of climate data • of how organisms affect the environment 	<ul style="list-style-type: none"> • of patterns of processes affecting Earth's four major systems and how changes in those processes will likely affect the components of those systems • that the locations of Earth's features are related to geologic changes • that the water cycle affects the distribution of water on Earth • that climatic patterns can be used to predict future weather conditions of an area • that behavior of organisms in an environment can help predict changes to the physical characteristics of that environment

Grade 5 Threshold Performance-Level Descriptors (Earth and Space Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
ESS3: Earth and Human Activity	<ul style="list-style-type: none"> • that humans use both renewable and nonrenewable resources for fuel and energy and that such use can affect the environment • that humans can identify different types of natural hazards • that humans have different effects on the environment or its resources 	<ul style="list-style-type: none"> • that using fuel from natural sources can be positive and negative in multiple ways • that Earth’s processes create unavoidable hazards and that humans have an important role in designing solutions to reduce negative impact • that individuals and communities can protect and reduce the negative effects that human activities can have on the environment 	<ul style="list-style-type: none"> • that humans have to make informed decisions about which natural resources to use by analyzing their risks and benefits • that there are benefits and risks to human-created solutions designed to lessen the impact of natural hazards • that humans have to make informed decisions based on the positive and negative effects of their activities in an effort to protect Earth

Grade 5 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>Asking Questions (for Science) and Defining Problems (for engineering) (AQDP): A practice of science is to ask and refine questions that lead to descriptions and explanations of how the natural and designed world works and which can be empirically tested.</p>	<ul style="list-style-type: none"> identify or ask relevant questions that are testable and that can show cause-and-effect relationships in the natural or designed world 	<ul style="list-style-type: none"> identify or ask relevant questions that can be investigated describe problems that can be solved predict reasonable outcomes clarify and redesign a solution to a problem 	<ul style="list-style-type: none"> generate questions based on investigations incorporating variables to determine patterns while defining and solving a design problem
<p>Developing and Using Models (DUM): A practice of both science and engineering is to use and construct models as helpful tools for representing ideas and explanations. These tools include diagrams, drawings, physical replicas, mathematical representations, analogies, and computer simulations.</p>	<ul style="list-style-type: none"> describe or use a model to show the relationship among components in a phenomenon 	<ul style="list-style-type: none"> develop or refine a model to minimize limitations, or test cause and effect relationships 	<ul style="list-style-type: none"> evaluate and revise or develop models to show relationships in cause-and-effect systems
<p>Planning and Carrying Out Investigations (PACI): Scientists and engineers plan and carry out investigations in the field or laboratory, working collaboratively as well as individually. Their investigations are systematic and require clarifying what counts as data and identifying variables or parameters.</p>	<ul style="list-style-type: none"> plan an investigation and collect observational data using appropriate methods or tools that help identify outcomes from changing a variable 	<ul style="list-style-type: none"> plan or conduct an investigation by evaluating appropriate methods or tools for collecting data while making predictions about a fair test in which variables are controlled 	<ul style="list-style-type: none"> plan and conduct multiple trials of an investigation to produce data that can be compared to make predictions, to serve as evidence for an explanation of a phenomenon, or to test a design solution

Grade 5 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>Analyzing and Interpreting Data (AID): Scientific investigations produce data that must be analyzed in order to derive meaning. Because data patterns and trends are not always obvious, scientists use a range of tools—including tabulation, graphical interpretation, visualization, and statistical analysis—to identify the significant features and patterns in the data. Scientists identify sources of error in the investigations and calculate the degree of certainty in the results. Modern technology makes the collection of large data sets much easier, providing secondary sources for analysis.</p>	<ul style="list-style-type: none"> organize relevant data to identify similarities or differences and describe how the data can be interpreted to make sense of phenomena 	<ul style="list-style-type: none"> analyze and represent relevant data describing how the data can be interpreted to make sense of phenomena 	<ul style="list-style-type: none"> evaluate and analyze data to refine a problem statement or make sense of phenomena
<p>Using Mathematics and Computational Thinking (UMCT): In both science and engineering, mathematics and computation are fundamental tools for representing physical variables and their relationships. They are used for a range of tasks, such as constructing simulations; statistically analyzing data; and recognizing, expressing, and applying quantitative relationships.</p>	<ul style="list-style-type: none"> identify ways to organize or analyze qualitative or quantitative data 	<ul style="list-style-type: none"> collect and organize data to reveal patterns, determine whether qualitative or quantitative data would be more appropriate 	<ul style="list-style-type: none"> organize complex data sets of qualitative or quantitative data, as determined to be appropriate, for determining relationships and patterns, creating algorithms, or utilizing mathematical representations to support conclusions

Grade 5 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>Constructing Explanations (for Science) and Designing Solutions (for Engineering) (CEDS): The products of science are explanations and the products of engineering are solutions.</p>	<ul style="list-style-type: none"> identify evidence or scientific ideas that support relationships to create solutions to a problem 	<ul style="list-style-type: none"> construct an explanation using evidence which utilizes scientific ideas to solve problems 	<ul style="list-style-type: none"> using evidence, evaluate and refine explanations of relationships among variables in determining the strengths and weaknesses of a design
<p>Engaging in Argument from Evidence (EAE): Argumentation is the process by which explanations and solutions are reached.</p>	<ul style="list-style-type: none"> identify evidence or compare facts in a claim 	<ul style="list-style-type: none"> distinguish among facts to construct, support, or evaluate a claim 	<ul style="list-style-type: none"> make or evaluate a claim using multiple sets of data
<p>Obtaining, Evaluating, and Communicating Information (OEI): Scientists and engineers must be able to communicate clearly and persuasively the ideas and methods they generate. Critiquing and communicating ideas individually and in groups is a critical professional activity.</p>	<ul style="list-style-type: none"> compare and summarize information to communicate basic scientific explanations of a phenomenon 	<ul style="list-style-type: none"> compare and combine information from various sources to communicate scientific explanations in various media 	<ul style="list-style-type: none"> evaluate scientific information to describe evidence and support future investigations

E.2.2 Grade 8 Threshold PLDs

The Threshold Performance-Level Descriptors (PLDs) define the minimum knowledge, skills, and practices that students must display for each Disciplinary Core Idea and Science and Engineering Practice to reach a certain performance level. They expand upon the brief overall PLDs included in the Score Interpretation Guide.

Grade 8 Threshold Performance-Level Descriptors (Physical Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
PS1: Matter and Its Interactions	<ul style="list-style-type: none"> that everything is made from atoms and that the states of matter have some unique characteristics that temperature and/or pressure have an effect on changes of state that chemical reactions create new substances while the mass does not change, and energy is involved 	<ul style="list-style-type: none"> that substances are made from one or more types of atoms and that the particles in the states of matter have unique characteristics that atoms are regrouped and conserved during chemical processes, and energy is either released or stored 	<ul style="list-style-type: none"> that substances can be made from two to thousands of atoms that can be combined in a variety of ways that the same numbers of atoms are regrouped into different molecules to create new substances with different properties, and therefore, the mass does not change
PS2: Motion and Stability: Forces and Interactions	<ul style="list-style-type: none"> that the movement of an object is the sum of its forces that forces among objects are either attractive or repulsive and are dependent upon the distance between the objects 	<ul style="list-style-type: none"> that in every interaction, there is a pair of forces acting on the two interacting objects and that the size of the forces on the first object equals the size of the forces on the second object that the size of the electromagnetic force depends upon the magnitudes of the charges, currents, or magnetic strengths due to the fields created 	<ul style="list-style-type: none"> of the effect of balanced versus unbalanced forces on the motion of objects that there is a relationship among forces, the fields created, and the magnitudes of the charges, currents, or magnetic strengths involved and among the distance between interacting objects and the masses of the interacting objects

Grade 8 Threshold Performance-Level Descriptors (Physical Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
PS3: Energy	<ul style="list-style-type: none"> • to identify kinetic energy, potential energy, temperature, and heat • that if there is a change in motion energy, it is due to energy being transferred in or out of the system • to identify that, during a collision, energy is transferred, and both objects exert a force • to identify reactants needed to make food in plants and the products of cellular respiration 	<ul style="list-style-type: none"> • of the proportional relationships that define kinetic and potential energy and the relationship between temperature and energy • of the relationship between energy and motion and how the amount of energy needed to cause changes is related to the properties of the substance • by describing the interaction between two objects in terms of force and energy transfer • to describe in general the processes of photosynthesis and cellular respiration including their reactants and products 	<ul style="list-style-type: none"> • to explain the relationship among the variables for kinetic and potential energy and explain how temperature is affected by composition, state, and energy of the particles in the system • to explain the flow of energy in a system, the relationship between the properties of a substance, and the energy needed to change the temperature or motion of the particles • to explain why objects exert a force on each other and that energy is transferred during an interaction • to explain the relationship between photosynthesis and cellular respiration and predict effects of a change to the system
PS4: Waves and Their Applications in Technologies for Information Transfer	<ul style="list-style-type: none"> • to identify properties of a simple wave • to identify the effect on a beam of light as it crosses between media and when it interacts with an object • to identify methods and their characteristics for transmitting information 	<ul style="list-style-type: none"> • to describe the properties of a simple wave and how it moves • to describe the effect on light as it crosses between media, the path it follows, and its interaction with objects • by describing how digitized signals are a more reliable way to encode and transmit information than analog signals 	<ul style="list-style-type: none"> • to explain the relationship between the properties of a wave and the requirement of a medium for transmission • by explaining how the properties of an object affect how light interacts with it and that the wave model of light is useful for explaining certain properties of light • to explain why digitized signals are a more reliable way to encode and transmit information than analog signals

Grade 8 Threshold Performance-Level Descriptors (Life Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>LS1: From Molecules to Organisms: Structures and Processes</p>	<ul style="list-style-type: none"> • that cells contain special structures which may be specific to the type of cell in a living unicellular or multicellular organism • of why genetic material is transferred differently in asexual reproduction and sexual reproduction, of how animal behaviors aid in reproduction for both the animal and/or some plants, and discuss genetic factors and local conditions that can affect growth of an organism • that matter and energy cycle through plants, creating sugars which can be broken down or rearranged to release the energy • that sense receptors can send various signals to the brain 	<ul style="list-style-type: none"> • that cells are the smallest unit of life, that living organisms can consist of one or more cells, and that multicellular organisms often contain specialized systems working together, and discuss the functions of special structures within cells • of characteristics, specialized features, and animal behaviors that increase the reproduction chance for both animals and plants, and explain how growth is affected by both genetic and environmental factors • of the process of photosynthesis for the creation of food and of the fact that to use that food, it needs to be broken down through another series of chemical reactions • that nerves transmit sense receptor inputs to be processed in the brain, resulting in memories or responses 	<ul style="list-style-type: none"> • of how parts of a cell function together in a manner similar to how systems interact in multicellular organisms • of characteristics, specialized features, and animal behaviors that increase the reproduction chance for both animals and plants and explain how growth is affected by both genetic and environmental factors • of the relationship between photosynthesis and cellular respiration and of how an organism obtains energy to sustain life • of the different ways a sense receptor reacts to inputs and of the process by which the signal is processed

Grade 8 Threshold Performance-Level Descriptors (Life Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>LS2: Ecosystems: Interactions, Energy, and Dynamics</p>	<ul style="list-style-type: none"> • that organisms are dependent on resources for which they may need to compete • that matter and/or energy are cycled through a food web of an ecosystem • that there are physical and biological components of ecosystems, that changes to those will cause disruption, and that biodiversity is related to species representation and can be used to determine overall health of an ecosystem • that changes in biodiversity have an impact on humans 	<ul style="list-style-type: none"> • of how growth and survival of an organism are dependent on access to limited resources and interactions with other organisms • of how matter and energy transfer between trophic levels • of the dynamic nature of ecosystems and of how biodiversity is used as a measure of an ecosystem’s health • of how changing biodiversity can affect humans and the services humans rely on 	<ul style="list-style-type: none"> • of an organism’s reliance on the environment and of how populations are limited by access to resources, predatory interactions, and competition • of how a food web can model mechanisms for the cycling of matter, including the role of decomposers, which in turn account for the conservation of energy • of the relationship between biodiversity and ecosystem health, and of the predicted outcomes of disturbances to an ecosystem • of why changes in biodiversity affect humans
<p>LS3: Heredity: Inheritance and Variation of Traits</p>	<ul style="list-style-type: none"> • that genes are located on inherited chromosomes and that the gene may be slightly different from the parent’s • that in sexual reproduction, each parent contributes half of the genetic material and that mutations that occur can be beneficial, harmful, or neutral 	<ul style="list-style-type: none"> • that genes control production of proteins and that mutations cause genetic variation • about genetic contributions during sexual reproduction and the general effects that mutations cause 	<ul style="list-style-type: none"> • of how genes control protein production and of what effect mutations could have on this process • of why individuals have two of each chromosome and how mutations may result in structural and functional changes

Grade 8 Threshold Performance-Level Descriptors (Life Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>LS4: Biological Evolution: Unity and Diversity</p>	<ul style="list-style-type: none"> • that fossils can show the evolutionary progression of organisms living today, that organisms may be artificially selected for reproduction based on desired traits, and that while embryos across species may have similarities as they develop, the organisms with more advantageous traits are more likely to survive • that environmental conditions will drive trait commonality in species 	<ul style="list-style-type: none"> • of the uses for the fossil record and of embryological development, including similarities not evident in the fully formed anatomy, where certain traits, whether natural or artificially selected, will provide advantages for survival • of how environmental conditions can change a species over generations and of how distributions of traits reflect adaptation by natural selection 	<ul style="list-style-type: none"> • of evolutionary history based on anatomical similarities and to predict predominance of certain traits in a population • to predict trait distribution in a species based on changing environmental conditions

Grade 8 Threshold Performance-Level Descriptors (Earth and Space Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
ESS1: Earth's Place in the Universe	<ul style="list-style-type: none"> • that the celestial bodies have observable patterns and that we exist in a galaxy called the Milky Way • that gravity acts on objects, that there are eclipses, and that Earth's tilt causes seasons • that fossils are used to date rock layers and that tectonic processes change Earth 	<ul style="list-style-type: none"> • to predict the observed motion of the Sun, Moon, and stars • that gravity is an attractive force, that alignment of the Earth-Moon-Sun causes solar and lunar eclipses, and that changes in seasons are due to intensity of sunlight • that Earth's history can be determined from rock layers and that tectonic processes create and destroy Earth materials 	<ul style="list-style-type: none"> • to explain the predictable observed patterns of the Sun, Moon, and stars • to predict eclipses and seasonal changes based on data or models • that rock layers and fossils only provide relative dates and that the sea floor has different ages
ESS2: Earth's Systems	<ul style="list-style-type: none"> • of where Earth's energy comes from and that Earth processes vary in timeframe and size • that Earth's plates move in different ways • that water cycles in Earth's spheres and affects weather patterns, that ocean water density varies, and that moving water affects landforms • that both living and nonliving factors influence complex weather patterns 	<ul style="list-style-type: none"> • that energy and matter have caused, and continue to cause, changes on Earth • that rocks and fossils help determine how Earth's plates have moved • of the way that water cycles, of the factors that affect the movement of water in Earth's spheres, of the causes of ocean density differences, and of the way that moving water affects landforms • of how weather patterns are influenced by living and nonliving factors that vary with location and of how the ocean is a major driving factor 	<ul style="list-style-type: none"> • of the interaction between Earth's processes driven by differing energy sources to explain Earth's history or predict future geological events • to predict effects of plate movement on Earth's landscape • to predict weather patterns that are the result of the cycling of water and of impacts of density on ocean currents • to predict the effect living and nonliving factors, including the ocean, have on weather and climate

Grade 8 Threshold Performance-Level Descriptors (Earth and Space Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
ESS3: Earth and Human Activity	<ul style="list-style-type: none"> • that resources are not evenly distributed • that natural hazards can be mapped • that human populations may negatively impact resources and that human activity has both positive and negative impacts on different organisms • of climate science and of the fact that human activities have an effect on global temperatures 	<ul style="list-style-type: none"> • that there are renewable and nonrenewable resources • that mapping hazards can help understand geological forces • on how humans have altered the biosphere and that humans are making technological gains to minimize negative impacts • of how human activities affect temperatures and that climate science may help lead to decisions to benefit life on Earth 	<ul style="list-style-type: none"> • of the relationship of past geological processes and the distribution of resources • to predict future hazards based on historical occurrences • to predict whether human activities would be positive or negative and to evaluate solutions based on the rate of resource consumption • to predict when human activities will have significant impacts on the Earth's climate

Grade 8 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>Analyzing and Interpreting Data (AID): Scientific investigations produce data that must be analyzed in order to derive meaning. Because data patterns and trends are not always obvious, scientists use a range of tools—including tabulation, graphical interpretation, visualization, and statistical analysis—to identify the significant features and patterns in the data. Scientists identify sources of error in the investigations and calculate the degree of certainty in the results. Modern technology makes the collection of large data sets much easier, providing secondary sources for analysis.</p>	<ul style="list-style-type: none"> identify and/or interpret data, graphical displays, and/or concepts of statistics and/or their limitations to provide evidence for phenomena 	<ul style="list-style-type: none"> analyze, interpret, and/or use simple data sets and/or concepts of statistics to identify relationships and/or define operational ranges for objects, processes, and/or systems 	<ul style="list-style-type: none"> analyze and interpret complex or multiple data sets and/or construct graphical displays to identify and/or explain relationships, limitations of data, when to use concepts of statistics, and/or to justify operational ranges for objects, processes, and/or systems
<p>Asking Questions (for science) and Defining Problems (for engineering) (AQDP): A practice of science is to ask and refine questions that lead to descriptions and explanations of how the natural and designed world works and which can be empirically tested.</p>	<ul style="list-style-type: none"> identify questions that arise from observations and models in order to clarify information and/or arguments, refine models, and/or determine relationships 	<ul style="list-style-type: none"> ask testable questions that arise from observations of phenomena, models, and/or unexpected results in order to clarify information, evidence, arguments, and/or design problems that can be solved through development of objects/tools, processes, and/or systems 	<ul style="list-style-type: none"> analyze and/or evaluate testable questions that arise from observations of phenomena, models, and/or unexpected results in order to clarify information, evidence, arguments, and/or design problems that can be solved through development of objects/tools, processes, and/or systems

Grade 8 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>Constructing Explanations (for science) and Designing Solutions (for engineering) (CEDS): The products of science are explanations, and the products of engineering are solutions.</p>	<ul style="list-style-type: none"> • identify or revise an explanation and/or design project based on models or representations, or by applying scientific reasoning and/or evidence 	<ul style="list-style-type: none"> • construct, revise, and/or use an explanation based on models or representations, or by applying scientific reasoning and/or evidence, or by undertaking a design project to construct and/or implement a solution 	<ul style="list-style-type: none"> • analyze, construct, and/or elaborate on an explanation based on models or representations by applying scientific reasoning and/or evidence, or by evaluating a design project to construct and/or implement solutions and/or optimize performance
<p>Developing and Using Models (DUM): A practice of both science and engineering is to use and construct models as helpful tools for representing ideas and explanations. These tools include diagrams, drawings, physical replicas, mathematical representations, analogies, and computer simulations.</p>	<ul style="list-style-type: none"> • use a simple model to show relationships, make predictions, or generate data and/or describe its limitations 	<ul style="list-style-type: none"> • develop and/or revise a simple model to show relationships, make predictions, or generate data and/or evaluate its limitations 	<ul style="list-style-type: none"> • develop, revise, and/or evaluate a complex model to show relationships, make predictions, or generate data and/or evaluate its merits and limitations
<p>Engaging in Argument from Evidence (EAE): Argumentation is the process by which explanations and solutions are reached.</p>	<ul style="list-style-type: none"> • identify evidence in arguments to support or refute explanations, • provide critiques of procedures or models, and/or • identify competing design solutions 	<ul style="list-style-type: none"> • identify and/or compare multiple pieces of evidence in arguments, • provide critiques about explanations or questions, and/or • write arguments that support or refute the advertised performance of a device, process, or system 	<ul style="list-style-type: none"> • critique arguments, procedures, or models; • construct and/or use written arguments to support or refute explanations, models, and/or solutions; or • analyze empirical evidence to support written arguments

Grade 8 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>Obtaining, Evaluating, and Communicating Information (OEIC): Scientists and engineers must be able to communicate clearly and persuasively the ideas and methods they generate. Critiquing and communicating ideas individually and in groups is a critical professional activity.</p>	<ul style="list-style-type: none"> • read and use information from multiple simple scientific sources to describe patterns, clarify claims, and/or assess accuracy 	<ul style="list-style-type: none"> • integrate information from multiple, complex, qualitative sources to clarify claims, assess accuracy, and evaluate conclusions 	<ul style="list-style-type: none"> • integrate information from multiple, complex, quantitative sources to describe patterns, clarify claims, assess accuracy, and evaluate conclusions
<p>Planning and Carrying Out Investigations (PACI): Scientists and engineers plan and carry out investigations in the field or laboratory, working collaboratively as well as individually. Their investigations are systematic and require clarifying what counts as data and identifying variables or parameters.</p>	<ul style="list-style-type: none"> • plan and/or conduct an investigation that includes the identification of appropriate tools and methods for collecting data in order to provide evidence or test a design solution 	<ul style="list-style-type: none"> • plan an investigation that includes the identification of variables and/or controls, or indicate how much data is sufficient to serve as evidence necessary to test a design solution, or evaluate an experimental design 	<ul style="list-style-type: none"> • plan and refine an investigation that includes the identification of variables and controls, tools, how data will be collected, and how much data is sufficient to serve as evidence necessary to test a design solution, or revise an experimental design
<p>Using Mathematics and Computational Thinking (UMCT): In both science and engineering, mathematics and computation are fundamental tools for representing physical variables and their relationships. They are used for a range of tasks such as constructing simulations; statistically analyzing data; and recognizing, expressing, and applying quantitative relationships.</p>	<ul style="list-style-type: none"> • identify qualitative and quantitative data and when the use of digital tools is warranted, • select appropriate mathematical representations, and • use algorithms to solve problems and/or address engineering questions 	<ul style="list-style-type: none"> • decide whether to use qualitative or quantitative data, • use digital tools to analyze large data sets, • use mathematical representations, and • explain and/or evaluate algorithms or mathematical concepts for solving problems and/or addressing engineering questions 	<ul style="list-style-type: none"> • explain when to use qualitative or quantitative data, • evaluate digital tools, • explain mathematical representations, and/or • create algorithms to solve problems and/or address engineering questions

E.2.3 Grade 11 Threshold PLDs

The Threshold Performance-Level Descriptors (PLDs) define the minimum knowledge, skills, and practices that students must display for each Disciplinary Core Idea and Science and Engineering Practice to reach a certain performance level. They expand upon the brief overall PLDs included in the Score Interpretation Guide.

Grade 11 Threshold Performance-Level Descriptors (Physical Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
PS1: Matter and Its Interactions	<ul style="list-style-type: none">• of subatomic particles, their interactions, and the involvement of energy in these interactions• of an understanding of how collisions between molecules affect reaction rates• that some reactions are reversible• that atoms are conserved during reactions• that nuclear processes involve energy	<ul style="list-style-type: none">• of atomic properties and patterns through the use of the periodic table, as well as different types of particle interactions and the energy involved• of the factors that affect reaction rates and equilibrium systems• of the energy involved in the rearranging of atoms and molecules• of the different types of reactions and how to make predictions about them• that energy and matter are conserved in nuclear processes	<ul style="list-style-type: none">• of varying atomic structures• of how the periodic table models the patterns of the properties and electron structure of the elements• of how particle interactions affect bulk properties of substances• of how collisions lead to changes in the sum of all the bond energies• of how atom conservation and chemical properties can be used to make predictions on chemical reactions• of multiple nuclear processes

Grade 11 Threshold Performance-Level Descriptors (Physical Science)
Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
PS2: Motion and Stability: Forces and Interactions	<ul style="list-style-type: none"> of quantified acceleration and momentum of types of fields and attractive/repulsive forces of gravitational and/or electric fields that electrical energy can be stored or transmitted 	<ul style="list-style-type: none"> (quantified knowledge) of factors that affect Newton’s second law, single object momentum systems, and conservation of momentum of how interactions happen at a distance due to fields of electrical interactions at the atomic level of the difference between magnetic and electric fields <p>OR</p> <ul style="list-style-type: none"> (quantified knowledge) of Coulomb’s law and Newton’s universal law of gravitation of how electrical energy can be stored in a battery or transmitted by electric currents 	<ul style="list-style-type: none"> (quantified knowledge) of outside interactions that affect the momentum and acceleration of a single- or multiple-object system of how to predict changes in electrical and gravitational forces of how to describe fields as force and energy fields and predict the effect of electrical and/or magnetic fields due to interactions between the two fields
PS3: Energy	<ul style="list-style-type: none"> of how different types of energy can be transferred of systems in which energy is conserved and how the availability of energy restricts what is possible in a closed system of the nature of the relationship between two objects interacting in a field using the energy prospective of how energy can be converted to different forms 	<ul style="list-style-type: none"> of how energy manifests itself at the microscopic and macroscopic scale and how energy transfers in a system (quantified knowledge) of how energy transfers in and out of a system <p>OR</p> <ul style="list-style-type: none"> of possible and impossible events based on energy availability, and defined stable states of how the distance between two objects affects the energy of a field of how energy can be converted to less useful forms of how solar energy can be captured and used for other processes, such as photosynthesis 	<ul style="list-style-type: none"> of the amount of various types of energy in a given situation and how microscopic changes affect macroscopic manifestations of energy of how to evaluate physical changes in a system using the conservation of energy of how to predict changes in energy in a field based on the position and nature of objects of the importance of energy conservation and efficiency

Grade 11 Threshold Performance-Level Descriptors (Physical Science)
Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>PS4: Waves and Their Applications in Technologies for Information Transfer</p>	<ul style="list-style-type: none"> • of how a wave travels through a medium, including understanding of examples of digitized information, and qualitative understanding of superposition principle • of the wave and particle models of electromagnetic radiation, the absorption of electromagnetic radiation, and the relationship between frequency and energy of light • of everyday experiences that involve waves and how wave signals are produced, transmitted, and captured 	<ul style="list-style-type: none"> • (quantified knowledge) of the relationship among frequency, wavelength, and speed in a real-world phenomenon <p>OR</p> <ul style="list-style-type: none"> • of the advantages and disadvantages of digitizing information • of the effect of absorption of electromagnetic waves, features of electromagnetic radiation that can be explained by either the wave or particle model, and the nature of photoelectric materials • of technologies used to produce, transmit, and/or capture signals and technologies used to store and interpret information 	<ul style="list-style-type: none"> • of waves in various media and how combining waves of different frequencies can make a wide variety of patterns and thereby encode and transmit information • of the difference between the wave- and particle-like behavior of electromagnetic radiation and how either the wave or particle model can be used to explain how an electron is emitted and how it can damage living cells • of how technology can be used to store and/or interpret information

Grade 11 Threshold Performance-Level Descriptors (Life Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>LS1: From Molecules to Organisms: Structures and Processes</p>	<ul style="list-style-type: none"> • of how multicellular organisms utilize feedback mechanisms and have specialized cells that are organized and function according to the proteins coded by the DNA • of the role of cellular division (mitosis) in creating genetically identical cells that differentiate into complex multicellular organisms • of photosynthesis and cellular respiration as the chemical processes of life that produce or utilize carbon-based molecules that are recombined into different products of living systems 	<ul style="list-style-type: none"> • of how positive and negative feedback mechanisms are beneficial to multicellular organisms, which have systems of specialized cells that perform essential life functions expressed through proteins coded for by genes • of how mitosis and differentiation produce and maintain complex organisms from a single cell • of the chemistry behind photosynthesis, how cellular respiration uses energy to maintain the organism, and how the products of these processes are used to build larger molecules 	<ul style="list-style-type: none"> • of how changing genes (mutation) can lead to functional changes of a protein and how positive and/or negative feedback helps maintain the equilibrium of an organism • of how genetic material from two variants of each chromosome pair is maintained as a single cell (fertilized egg) grows to a multicellular organism • of the interdependence of photosynthesis and cellular respiration and their role in the growth and maintenance of living systems

Grade 11 Threshold Performance-Level Descriptors (Life Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>LS2: Ecosystems: Interactions, Energy, and Dynamics</p>	<ul style="list-style-type: none"> • of both living and non-living factors that contribute to the carrying capacity of the ecosystem • of how food webs often have photosynthetic producers at the lowest level, how a small amount of matter and energy will transfer upward in the food web reducing the amount of organisms that can exist at higher levels, and how this relates to the carbon cycle • of how ecosystems have interactions that keep the population numbers stable, and ecosystems are resilient to modest changes, but humans can disrupt ecosystems and species survival • of how group behavior has evolved to increase individual and group survival 	<ul style="list-style-type: none"> • of how carrying capacity is affected by challenges and/or availability of resources • of how photosynthesis and cellular respiration are connected and use carbon in maintaining life processes, that the matter and energy of a food web are used and restructured by the organisms in the food web, and that a small amount is used by the next levels of the food web • of complex ecosystem interactions and their effects on population size, including biological and physical disturbances, extreme fluctuations, and the ways human activity can have an effect on an ecosystem • of how group behaviors can increase the chances of survival for individuals and their genetic relatives 	<ul style="list-style-type: none"> • of how carrying capacity affects the population size of a given species within an ecosystem • of how carbon and matter are used in the maintenance of life processes (including photosynthesis and both anaerobic and aerobic respiration) through the food web, including how carbon cycles through Earth's spheres • of how changes to populations and environments caused by human interactions and other physical events within ecosystems can result in changes that affect both the organisms and the environment • of how changes to the group or conditions can affect the survival of individuals and their genetic relatives

Grade 11 Threshold Performance-Level Descriptors (Life Science)
Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
LS3: Heredity: Inheritance and Variation of Traits	<ul style="list-style-type: none"> of how all cells have the same DNA containing genes that are the organisms' characteristics, but not all DNA codes for protein of the processes within meiosis, errors that can occur during DNA replication, and mutations due to environmental factors that can create genetic diversity, which may be passed to future generations 	<ul style="list-style-type: none"> that chromosomes contain genes that code for proteins and regions that do not code for proteins, and that different cells express different genes that while the process of DNA replication is tightly regulated and highly accurate, errors still occur, and combined with mutations due to environmental factors, DNA replication can create genetic diversity that may affect survivability and the transmission of traits to future generations 	<ul style="list-style-type: none"> of the mechanisms of gene regulation and different possible functions of segments of non-protein coding DNA of the mechanisms within meiosis that create genetic diversity, as well as the effects of environmental factors on DNA replication and the impact of the changes to DNA on genetic diversity within populations
LS4: Biological Evolution: Unity and Diversity	<ul style="list-style-type: none"> of the different types of evidence of evolution of how natural selection allows inheritable advantageous traits to become more common if they increase chances of survival within populations that natural selection selects for inheritable traits that provide a survival advantage for a particular environment that changes to the environment may cause the selection of different traits leading to changes in the population known as adaptation that the frequency of traits depends on natural selection forces that can change with a changing environment of how biodiversity increases or decreases and how humans need resources and biodiversity, but are having adverse effects on biodiversity 	<ul style="list-style-type: none"> of how different sources of evidence for evolution can support each other of how gene expression and genetic variation in the individual lead to differences in performance of the individuals in a population, and how positively selected traits are more common in a population because they increase survival that evolution occurs when there is genetic variation, competition, and selective reproduction of organisms with desirable genetic traits that organisms with desirable traits will become more common, but as the environment changes, different traits may provide the selective advantages that some populations may increase while others may go extinct of specific results of human activities that affect the environment and biodiversity and reasons why preservation of biodiversity is desirable 	<ul style="list-style-type: none"> of how DNA sequences, amino acid sequences, and anatomical and embryological evidence support that evolution has occurred of how natural selection leads to different levels of performance of the individual that factors affecting natural selection work together creating changes in the diversity within populations and ecosystems that changing environments cause changes in selection pressures that result in further adaptation or extinction of ways that humans can maintain or increase biodiversity while meeting the needs of humanity and why this is beneficial to life on Earth

Grade 11 Threshold Performance-Level Descriptors (Earth and Space Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
ESS1: Earth's Place in the Universe	<ul style="list-style-type: none"> • of the Big Bang, which allowed for the creation of galaxies and stars, where many elements are created • of identifying properties of orbits, factors that affect the orbit, and how the orbit affects the stellar body • of plate tectonics and erosion, which cause the destruction of early rock records on Earth, and that we have to rely on other objects in the solar system for information on Earth's formation 	<ul style="list-style-type: none"> • that light spectra emitted from a star can give information about its life cycle, composition, and distance • of features of motion of orbital objects, what changes that motion, and the effects of changing the motion of the stellar body • of the fact that while there is a range in the age of the rocks on Earth, the early rock history has been destroyed, and we rely on studying other stellar bodies to explain how the Earth formed 	<ul style="list-style-type: none"> • of the life cycle of stars and explain how the characteristics of a star can support the Big Bang theory • of the laws explaining motions of orbiting objects, their changes, and the changes to the stellar bodies as a result of those changes • of why different areas of the Earth have rocks of different ages and the processes that are erasing the early rock history
ESS2: Earth's Systems	<ul style="list-style-type: none"> • of how Earth has a series of interacting dynamic systems • that Earth's surface is in motion, and that motion can create physical features on the Earth's surface • of the properties of water that are essential to Earth's dynamics • of Earth's atmosphere and how it undergoes temperature changes • that dynamic and delicate feedbacks between the Earth's systems and biosphere exist 	<ul style="list-style-type: none"> • of methods of investigation of Earth's dynamic systems and how the data can be used to describe the effects of these systems • that Earth's surface is in motion due to convection, creating physical features that have changed throughout history • of how the properties of water are essential to Earth's processes • of how Earth's atmosphere undergoes short-term and long-term temperature changes at the global scale due to changes in the biosphere, including human activities • of how dynamic and delicate feedback between the Earth's systems and biosphere causes a continual co-evolution of Earth's surface and the life that exists on it 	<ul style="list-style-type: none"> • of Earth's dynamic systems in explaining the effects of these systems and the development of the currently accepted model of the structure of the planet • of the theory of plate tectonics allowing for the prediction of future plate movements and interpretations of Earth's geologic history • of how the properties of water can be used to explain Earth's processes • of why Earth's atmosphere undergoes short-term and long-term temperature changes at the global scale • of how positive and/or negative feedbacks between the biosphere and other Earth systems cause a continual co-evolution of Earth's surface and the life that exists on it

Grade 11 Threshold Performance-Level Descriptors (Earth and Space Science)

Students should be able to demonstrate knowledge:

DCI	Level 2	Level 3	Level 4
<p>ESS3: Earth and Human Activity</p>	<ul style="list-style-type: none"> • that new technologies have associated costs, risks, and benefits • that natural hazards have shaped human history • that human activities can have both positive and negative impacts on biodiversity • of humans’ abilities to use technology to model, predict, and manage current and future impacts 	<ul style="list-style-type: none"> • that new technologies have associated costs, risks, and benefits at the economic, social, environmental, and/or geopolitical level • of how natural hazards and geological events have shaped human history through changes in the human population including through migration at the local, regional, and/or global scale • that human impacts on biodiversity can be mitigated by the development of new technologies and/or responsible resource management • of technologies that allow modeling, predicting, and managing of current and future impacts on oceans, the atmosphere, and the biosphere 	<ul style="list-style-type: none"> • of new technologies in order to explain their associated costs, risks, and benefits at the economic, social, environmental, and/or geopolitical level • of how natural hazards affect human population and migration at the local, regional, and global scale • of new technologies and responsible resource management to predict their effects on biodiversity • to explain how humans’ abilities to model, predict, and manage current and future impacts have increased alongside the magnitudes of human impacts

Grade 11 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>(Investigating) Asking Questions (for Science) and Defining Problems (for engineering) (AQDP): A practice of science is to ask and refine questions that lead to descriptions and explanations of how the natural and designed worlds work and which can be empirically tested. Engineering questions clarify problems to determine criteria for successful solutions and identify constraints to solve problems about the designed world. Both scientists and engineers also ask questions to clarify ideas.</p> <p>Asking questions and defining problems in 9–12 progresses to formulating, refining, and evaluating empirically testable questions and design problems using models and simulations.</p>	<ul style="list-style-type: none"> ask relevant questions or define problems in different contexts, based on unexpected results, independent and dependent variables, models, theories, etc. 	<ul style="list-style-type: none"> ask relevant and testable questions that arise from careful observation of phenomena, unexpected results, or models or theories for the purpose of determining relationships, providing an explanation, or clarifying and refining a design 	<ul style="list-style-type: none"> analyze, evaluate, and/or revise questions that arise from careful observation of phenomena, unexpected results, or models or theories for the purpose of determining relationships, providing an explanation, or clarifying and refining a design
<p>(Sensemaking) Developing and Using Models (DUM): A practice of both science and engineering is to use and construct models as helpful tools for representing ideas and explanations. These tools include diagrams, drawings, physical replicas, mathematical representations, analogies, and computer simulations. Modeling in 9–12 progresses to using, synthesizing, and developing models to predict and show relationships among variables between systems and their components in the natural and designed worlds.</p>	<ul style="list-style-type: none"> use a model to generate data that test the model’s reliability and/or evaluates its merits and limitations 	<ul style="list-style-type: none"> develop simple models and revise different types of models that test and/or predict relationships among systems/ phenomena based on the models’ merits and limitations 	<ul style="list-style-type: none"> develop or revise complex models that test and/or predict relationships/ phenomena based on the models’ merits and limitations

Grade 11 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>(Investigating) Planning and Carrying Out Investigations (PACI): Scientists and engineers plan and carry out investigations in the field or laboratory, working collaboratively as well as individually. Their investigations are systematic and require clarifying what counts as data and identifying variables or parameters. Planning and carrying out investigations in 9–12 progresses to include investigations that provide evidence for and test conceptual, mathematical, physical, and empirical models.</p>	<ul style="list-style-type: none"> identify ways to conduct an investigation (including making a directional hypothesis) or test a design solution through manipulating variables or acquiring data 	<ul style="list-style-type: none"> plan and/or conduct an investigation (including making a directional hypothesis) or test a design solution through manipulating variables or acquiring data 	<ul style="list-style-type: none"> revise and/or evaluate an investigation in which an independent variable is manipulated or an unsatisfactory performance is found
<p>(Sensemaking) Analyzing and Interpreting Data (AID): Scientific investigations produce data that must be analyzed in order to derive meaning. Because data patterns and trends are not always obvious, scientists use a range of tools—including tabulation, graphical interpretation, visualization, and statistical analysis—to identify the significant features and patterns in the data. Scientists identify sources of error in the investigations and calculate the degree of certainty in the results. Modern technology makes the collection of large data sets much easier, providing secondary sources for analysis. Analyzing data in 9–12 progresses to introducing more detailed statistical analysis, the comparison of data sets for consistency, and the use of models to generate and analyze data.</p>	<ul style="list-style-type: none"> identify the appropriate statistics and/or data, and/or their limitations, when providing evidence for claims, design solutions, or solving problems 	<ul style="list-style-type: none"> apply and/or analyze data and statistics to identify or solve scientific and engineering problems, or to make scientific claims 	<ul style="list-style-type: none"> evaluate the use of data and statistics and/or their limitations to solve problems, make claims, or design solutions

Grade 11 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>(Investigating) Using Mathematics and Computational Thinking (UMCT): In both science and engineering, mathematics and computation are fundamental tools for representing physical variables and their relationships. They are used for a range of tasks such as constructing simulations; statistically analyzing data; and recognizing, expressing, and applying quantitative relationships. Mathematical and computational thinking in 9–12 progresses to using algebraic thinking and analysis, a range of linear and nonlinear functions including trigonometric functions, exponentials and logarithms, and computational tools for statistical analysis to analyze, represent, and model data. Simple computational simulations are created and used based on mathematical models of basic assumptions.</p>	<ul style="list-style-type: none"> • apply/use mathematical concepts to describe conclusions that may require deciding when to use qualitative versus quantitative data 	<ul style="list-style-type: none"> • apply/use mathematical computational representations to see if a model is viable, or decide if qualitative or quantitative data meet criteria for success 	<ul style="list-style-type: none"> • through the use of evaluation of mathematical computations, create a model or justify the choice of qualitative versus quantitative data
<p>(Sensemaking) Constructing Explanations (for science) and Designing Solutions (for engineering) (CEDs): The products of science are explanations and the products of engineering are solutions. Constructing explanations and designing solutions in 9–12 progresses to explanations and designs that are supported by multiple and independent student-generated sources of evidence consistent with scientific ideas, principles, and theories.</p>	<ul style="list-style-type: none"> • identify and describe appropriate data and/or evidence for supporting claims, solving problems, constructing explanations, or designing solutions 	<ul style="list-style-type: none"> • make or revise claims, explanations, or solutions by applying appropriate data and/or evidence 	<ul style="list-style-type: none"> • evaluate, design, or construct claims, explanations, or solutions by applying appropriate data, evidence, and/or scientific theories and laws

Grade 11 SEP Threshold Performance-Level Descriptors
Students should be able to:

SEP	Level 2	Level 3	Level 4
<p>(Critiquing) Engaging in Argument from Evidence (EAE): Argumentation is the process by which explanations and solutions are reached. Engaging in argument from evidence in 9–12 progresses to using appropriate and sufficient evidence and scientific reasoning to defend and critique claims and explanations about the natural and designed worlds. Arguments may also come from current scientific or historical episodes in science.</p>	<ul style="list-style-type: none"> identify and/or describe the main points of an argument or claim that is based on scientific evidence 	<ul style="list-style-type: none"> evaluate and/or defend a claim or argument— or choose between competing arguments— related to currently accepted explanations or solutions 	<ul style="list-style-type: none"> construct and/or critique an argument or claim by using scientific evidence
<p>(Critiquing) Obtaining, Evaluating, and Communicating Information (OEI): Scientists and engineers must be able to communicate clearly and persuasively the ideas and methods they generate. Critiquing and communicating ideas individually and in groups is a critical professional activity. Obtaining, evaluating, and communicating information in 9–12 progresses to evaluating the validity and reliability of the claims, methods, and designs.</p>	<ul style="list-style-type: none"> read and compare sources of information to describe patterns in evidence and/ or evidence for solving problems or answering scientific questions 	<ul style="list-style-type: none"> integrate information from multiple sources to gather valid and reliable evidence for solving problems or answering scientific questions 	<ul style="list-style-type: none"> evaluate information from multiple sources and determine the usefulness of evidence, ensuring it is valid and reliable, for solving problems or answering scientific questions

E.3 Reporting PLDs

E.3.1 Reporting PLDs–Level 1

Students who are at Level 1 demonstrated a minimal understanding of the New Jersey Student Learning Standards for Science (NJSLS–S) by misinterpreting information from a variety of sources (e.g., text, charts, graphs, tables) and inconsistently applying the knowledge gained from scientific investigations to develop incorrect explanations or models of observed phenomena. The students had difficulty choosing and using, even with significant scaffolding, the appropriate tools to make observations and to gather, classify, and present data. The students struggled to use essential information to recognize patterns and relationships between data and designed systems. The students seldom used information to make real-world connections or predictions.

E.3.2 Reporting PLDs–Level 2

Students who are at Level 2 demonstrated a limited grade-level understanding of the New Jersey Student Learning Standards for Science (NJSLS–S) by partially interpreting information from a variety of sources (e.g., text, charts, graphs, tables) and inconsistently applying the knowledge gained from scientific investigations to develop incomplete explanations or models of observed phenomena. The students had some difficulty choosing and using the appropriate tools to make observations and to gather, classify, and present data. The students may be able to use essential information to recognize patterns and relationships between data and designed systems. The students inconsistently used information to make real-world connections and predictions.

E.3.3 Reporting PLDs–Level 3

Students who are at Level 3 demonstrated appropriate grade-level understanding of the New Jersey Student Learning Standards for Science (NJSLS–S) by comprehending information from a variety of sources (e.g., text, charts, graphs, tables) and applying the knowledge gained from scientific investigations to develop accurate explanations and models of observed phenomena. The students often chose and used the appropriate tools to make observations and to gather, classify, and present data. The students used both essential and non-essential information to recognize patterns and relationships between data and designed systems. The students were able to use information to make real-world connections and predictions.

E.3.4 Reporting PLDs–Level 4

Students who are at Level 4 demonstrate advanced understanding of the New Jersey Student Learning Standards for Science (NJSLS–S) by integrating information from a variety of sources (e.g., text, charts, graphs, tables) and analyzing the knowledge gained from scientific investigations to develop sophisticated explanations and models of observed phenomena. The students consistently chose and used the appropriate tools to make observations and to gather, classify, and present relevant data. The students considered both essential and non-essential information to explain patterns and relationships between data and designed systems. The students regularly used information and provided supporting explanations in making real-world connections and predictions.

APPENDIX F: DETAILED TEST MAPS

Table F.1: Grade 5 Test Map—Metadata and Item Statistics

Grade 5											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	<i>rpb</i>	Median Time
1905B007_01	1	TE	AID	LS3.A	SC	Life	Sensemaking	1.017	0.33	0.62	148
1905B007_03	1	TE	UMCT	LS3.A	S, P, and Q	Life	Investigating	0.902	0.34	0.60	97
1905B007_05	1	TE	AID	LS3.A	SC	Life	Sensemaking	0.553	0.41	0.59	80
1905B007_08	4	CR	EAE	LS3.A	PAT	Life	Critiquing	0.393	0.41	0.65	407
1905B007_10	1	TE	CEDS	LS1.C	S & SM	Life	Sensemaking	0.541	0.42	0.66	38
1905B009_01	1	TE	UMCT	LS4.B	S, P, and Q	Life	Investigating	0.350	0.44	0.62	112
1905B009_02	1	TE	AID	LS4.B	S, P, and Q	Life	Sensemaking	-0.809	0.65	0.65	43
1905B009_05	1	TE	AQDP	LS4.B	PAT	Life	Investigating	1.579	0.22	0.26	74
1905M005_01	1	TE	EAE	ESS2.D	S & SM	Earth and Space	Critiquing	0.232	0.46	0.47	118
1905M005_03	1	TE	OECI	ESS2.D	PAT	Earth and Space	Critiquing	0.881	0.37	0.55	85
1905M005_04	1	MC	UMCT	ESS2.D	PAT	Earth and Space	Investigating	-0.282	0.58	0.25	51
1905M008_01	1	MC	EAE	ESS1.A	S, P, and Q	Earth and Space	Critiquing	-0.515	0.65	0.53	134
1905M008_05	1	TE	AID	ESS1.A	PAT	Earth and Space	Sensemaking	0.111	0.53	0.56	78
1905M008_06	1	MC	OECI	ESS1.A	C and E	Earth and Space	Critiquing	-0.794	0.67	0.49	69
1905M040_01	1	TE	PACI	PS2.A	PAT	Physical	Investigating	0.733	0.37	0.47	163
1905M040_03	1	MC	AID	PS2.A	C and E	Physical	Sensemaking	-0.224	0.55	0.44	75
1905M040_05	1	TE	AID	PS2.A	C and E	Physical	Sensemaking	0.775	0.35	0.39	55
1905M044_02	1	MC	OECI	LS2.D	C and E	Life	Critiquing	0.072	0.50	0.55	112
1905M044_03	1	TE	EAE	LS2.D	C and E	Life	Critiquing	1.038	0.32	0.45	63
1905M044_04	1	TE	AID	LS2.D	SF	Life	Sensemaking	0.990	0.32	0.39	89
1905M076_01	1	MC	EAE	PS3.A	E&M	Physical	Critiquing	-0.352	0.58	0.23	75
1905M076_03	1	TE	CEDS	PS2.A	PAT	Physical	Sensemaking	0.472	0.42	0.54	78
1905M076_05	1	TE	CEDS	PS3.B	PAT	Physical	Sensemaking	1.204	0.29	0.51	47
2105M015_04	1	TE	EAE	ESS2.E	S & SM	Earth and Space	Critiquing	1.026	0.31	0.35	69
2105M015_05	1	TE	CEDS	ESS2.E	C and E	Earth and Space	Sensemaking	0.525	0.41	0.09	69

Grade 5											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	<i>rpb</i>	Median Time
2105M015_06	1	MC	AQDP	ESS2.E	C and E	Earth and Space	Investigating	-0.966	0.69	0.36	97
2205B003_01	1	MC	AQDP	PS1.B	C and E	Physical	Investigating	0.508	0.41	0.14	98
2205B003_02	1	TE	PACI	PS1.B	SF	Physical	Investigating	0.809	0.35	0.42	52
2205B003_03	1	TE	EAE	PS1.B	SF	Physical	Critiquing	1.259	0.28	0.25	53
2205B003_04	4	CR	OECI	PS1.B	C and E	Physical	Critiquing	1.656	0.18	0.56	421
2205B003_05	1	TE	CEDS	ESS2.C	SC	Earth and Space	Sensemaking	-0.328	0.58	0.50	35
2205B009_01	1	TE	DUM	ESS2.A	S & SM	Earth and Space	Sensemaking	-0.298	0.57	0.52	50
2205B009_02	1	TE	AQDP	ESS2.A	S & SM	Earth and Space	Investigating	1.339	0.27	0.43	42
2205B009_03	1	TE	PACI	ESS2.A	S & SM	Earth and Space	Investigating	1.430	0.25	0.23	40
2205B009_04	1	TE	DUM	ESS2.A	S & SM	Earth and Space	Sensemaking	0.372	0.44	0.42	35
2205B009_05	4	CR	EAE	ESS2.A	S & SM	Earth and Space	Critiquing	1.091	0.28	0.63	354
2205M004_03	1	TE	CEDS	LS1.D	SF	Life	Sensemaking	0.265	0.46	0.61	142
2205M004_05	1	TE	EAE	LS2.C	C and E	Life	Critiquing	1.162	0.29	0.40	101
2205M004_07	1	TE	DUM	LS1.D	SF	Life	Sensemaking	0.273	0.46	0.32	63
2205M006_01	1	MC	AQDP	LS2.B	S & SM	Life	Investigating	0.030	0.51	0.40	91
2205M006_02	1	TE	EAE	LS2.B	S & SM	Life	Critiquing	0.591	0.40	0.26	79
2205M006_05	1	TE	DUM	LS2.B	S & SM	Life	Sensemaking	0.103	0.49	0.45	56
2205M011_01	1	TE	AQDP	ESS2.B	PAT	Earth and Space	Investigating	1.401	0.26	0.27	43
2205M011_02	1	TE	EAE	ESS2.B	PAT	Earth and Space	Critiquing	0.576	0.40	0.55	55
2205M011_04	1	MC	OECI	ESS2.B	PAT	Earth and Space	Critiquing	0.349	0.44	0.37	87
2205M012_01	1	MC	AQDP	PS3.C	S & SM	Physical	Investigating	-0.534	0.61	0.49	103
2205M012_03	1	MC	OECI	PS3.C	E&M	Physical	Critiquing	0.645	0.39	0.45	78
2205M012_04	1	TE	EAE	PS3.C	C and E	Physical	Critiquing	1.587	0.23	0.06	46
2205M022_01	1	MC	AQDP	PS4.B	S & SM	Physical	Investigating	0.433	0.43	0.44	70
2205M022_03	1	MC	PACI	PS4.B	SF	Physical	Investigating	0.362	0.44	0.47	51
2205M022_05	1	TE	PACI	PS4.B	S, P, and Q	Physical	Investigating	0.447	0.42	0.26	38

Table F.2: Grade 8 Test Map—Metadata and Item Statistics

Grade 8											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	<i>rpb</i>	Median Time
1908B000_03	1	TE	CEDS	PS2.A	C and E	Physical	Sensemaking	1.003	0.17	0.29	91
1908B000_04	1	TE	AID	PS2.B	C and E	Physical	Sensemaking	0.573	0.27	0.40	73
1908B000_08	1	TE	EAE	PS3.A	S, P, and Q	Physical	Critiquing	0.624	0.24	0.40	56
1908B000_11	4	CR	CEDS	PS2.A	C and E	Physical	Sensemaking	-0.027	0.33	0.70	413
1908B000_12	1	TE	AID	PS2.A	SC	Physical	Sensemaking	0.589	0.27	0.38	48
1908M003_02	1	TE	DUM	LS3.A	S, P, and Q	Life	Sensemaking	1.580	0.15	0.24	86
1908M003_07	1	TE	DUM	LS3.B	PAT	Life	Sensemaking	1.413	0.16	0.29	48
1908M003_08	1	MC	CEDS	LS3.B	PAT	Life	Sensemaking	0.397	0.31	0.42	40
1908M005_02	1	TE	UMCT	PS1.A	S & SM	Physical	Investigating	1.080	0.22	0.37	122
1908M005_03	1	TE	UMCT	PS1.B	E&M	Physical	Investigating	1.901	0.11	0.41	86
1908M005_05	1	MC	AQDP	PS1.B	C and E	Physical	Investigating	-0.173	0.39	0.39	53
1908M026_01	1	TE	OECI	ESS3.D	C and E	Earth and Space	Critiquing	1.459	0.15	0.42	163
1908M026_04	1	MC	AID	ESS3.D	PAT	Earth and Space	Sensemaking	-0.096	0.39	0.34	63
1908M026_06	1	MC	EAE	ESS3.D	PAT	Earth and Space	Critiquing	-0.205	0.41	0.30	67
1908M030_01	1	MC	EAE	LS1.A	S & SM	Life	Critiquing	0.269	0.38	0.38	150
1908M030_02	1	TE	AID	LS1.B	C and E	Life	Sensemaking	0.106	0.38	0.40	49
1908M030_05	1	MC	PACI	LS1.D	C and E	Life	Investigating	0.362	0.34	0.39	72
1908M033_02	1	TE	UMCT	ESS2.B	PAT	Earth and Space	Investigating	0.105	0.38	0.33	104
1908M033_03	1	TE	EAE	ESS3.B	PAT	Earth and Space	Critiquing	-0.412	0.47	0.29	93
1908M033_04	1	MC	DUM	ESS2.B	SC	Earth and Space	Sensemaking	-0.857	0.54	0.36	70
2008M000_01	1	TE	DUM	LS1.A	SF	Life	Sensemaking	0.640	0.26	0.19	62
2008M000_02	1	TE	OECI	LS1.A	SF	Life	Critiquing	1.107	0.19	0.46	98
2008M000_04	1	MC	AQDP	LS1.C	E&M	Life	Investigating	0.291	0.32	0.25	42
2008M001_01	1	TE	AID	PS2.B	PAT	Physical	Sensemaking	0.354	0.31	0.59	112
2008M001_05	1	MC	EAE	PS2.B	PAT	Physical	Critiquing	-0.996	0.58	0.40	83

Grade 8											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	<i>rpb</i>	Median Time
2008M001_08	1	MC	PACI	PS2.B	S & SM	Physical	Investigating	0.343	0.31	0.34	44
2008M015_04	1	TE	UMCT	LS3.B	S, P, and Q	Life	Investigating	1.198	0.18	0.31	80
2008M015_05	1	TE	AQDP	LS3.B	C and E	Life	Investigating	0.774	0.24	0.20	48
2008M015_06	1	TE	AID	LS3.B	C and E	Life	Sensemaking	-0.414	0.46	0.49	37
2008M015_09	1	MC	EAE	LS3.B	C and E	Life	Critiquing	0.181	0.34	0.35	48
2108B003_02	1	MC	DUM	PS4.B	SF	Physical	Sensemaking	-1.118	0.60	0.38	65
2108B003_07	1	TE	EAE	PS4.B	SF	Physical	Critiquing	0.667	0.26	0.32	53
2108B003_08	1	TE	PACI	PS4.B	SF	Physical	Investigating	-0.284	0.43	0.26	64
2108B006_01	1	TE	UMCT	ESS1.B	S, P, and Q	Earth and Space	Investigating	0.194	0.34	0.59	94
2108B006_03	1	TE	AQDP	ESS2.C	S & SM	Earth and Space	Investigating	1.232	0.17	0.33	39
2108B006_06	1	TE	CEDS	ESS1.B	PAT	Earth and Space	Sensemaking	1.057	0.20	0.37	55
2108B006_09	1	TE	OECI	ESS2.C	PAT	Earth and Space	Critiquing	0.175	0.34	0.54	23
2108B006_11	3	CR	PACI	ESS2.C	C and E	Earth and Space	Investigating	1.496	0.10	0.53	313
2108B007_01	1	TE	EAE	PS1.B	PAT	Physical	Critiquing	0.257	0.33	0.48	97
2108B007_03	1	MC	AQDP	PS3.A	S, P, and Q	Physical	Investigating	0.980	0.21	0.11	43
2108B007_07	1	TE	CEDS	PS1.B	PAT	Physical	Sensemaking	0.988	0.21	0.25	30
2108B007_08	1	TE	OECI	PS1.B	PAT	Physical	Critiquing	-1.045	0.59	0.40	30
2108M015_02	1	TE	EAE	ESS1.B	PAT	Earth and Space	Critiquing	-0.473	0.47	0.42	92
2108M015_06	1	TE	AID	ESS1.B	PAT	Earth and Space	Sensemaking	-0.273	0.43	0.50	95
2108M015_09	1	MC	AQDP	ESS1.B	PAT	Earth and Space	Investigating	-0.786	0.53	0.49	58
2108M015_10	1	TE	CEDS	ESS1.B	PAT	Earth and Space	Sensemaking	0.314	0.32	0.21	56
2108M027_01	1	TE	AQDP	LS4.A	PAT	Life	Investigating	0.714	0.25	0.46	35
2108M027_03	1	TE	OECI	LS4.A	PAT	Life	Critiquing	-0.495	0.47	0.50	27
2108M027_07	1	MC	AID	LS4.C	C and E	Life	Sensemaking	-0.143	0.41	0.27	64
2208B003_01	1	TE	AID	LS2.C	C and E	Life	Sensemaking	-0.204	0.41	0.48	49
2208B003_05	1	TE	CEDS	LS4.D	E&M	Life	Sensemaking	-0.107	0.40	0.47	40

Grade 8											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	<i>rpb</i>	Median Time
2208B003_07	1	TE	PACI	LS2.C	PAT	Life	Investigating	1.196	0.18	0.47	52
2208B003_09	1	TE	EAE	LS2.C	C and E	Life	Critiquing	0.590	0.27	0.21	38
2208B003_11	3	CR	OECI	LS2.C	SC	Life	Critiquing	-0.220	0.41	0.60	271
2208M021_03	1	TE	OECI	ESS1.C	PAT	Earth and Space	Critiquing	0.581	0.27	0.36	88
2208M021_05	1	TE	AID	ESS1.C	C and E	Earth and Space	Sensemaking	0.104	0.36	0.37	50
2208M021_09	1	MC	AID	ESS1.C	PAT	Earth and Space	Sensemaking	0.443	0.30	0.32	59
2208M021_10	1	TE	EAE	ESS1.C	PAT	Earth and Space	Critiquing	0.065	0.36	0.58	40
2208M028_02	1	MC	AQDP	LS2.A	C and E	Life	Investigating	1.005	0.21	0.11	44
2208M028_06	1	TE	PACI	ESS3.C	C and E	Earth and Space	Investigating	-0.189	0.41	0.25	82
2208M028_07	1	TE	EAE	LS2.C	S, P, and Q	Life	Critiquing	1.242	0.17	0.24	58
2208M028_09	1	MC	EAE	LS2.A	C and E	Life	Critiquing	-0.109	0.40	0.30	42
2208M051_13	1	TE	AQDP	PS4.A	C and E	Physical	Investigating	0.091	0.36	0.37	75
2208M051_16	1	TE	EAE	PS3.B	PAT	Physical	Critiquing	-0.379	0.46	0.44	51
2208M051_17	1	TE	OECI	PS4.A	C and E	Physical	Critiquing	0.344	0.31	0.42	66

Table F.3: Grade 11 Test Map—Metadata and Item Statistics

Grade 11											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	<i>rpb</i>	Median Time
1911B009_01A	1	TE	OECI	LS2.D	S & SM	Life	Critiquing	-0.439	0.55	0.52	83
1911B009_03A	1	TE	AID	LS2.D	PAT	Life	Sensemaking	0.479	0.36	0.62	69
1911B009_05A	1	TE	EAE	LS2.D	S, P, and Q	Life	Critiquing	1.377	0.22	0.45	74
1911B009_07A	4	CR	CEDS	LS2.D	S & SM	Life	Sensemaking	0.683	0.33	0.59	309
1911B009_09A	1	TE	PACI	LS4.C	S & SM	Life	Investigating	1.331	0.20	0.36	87
1911M002_01	1	TE	DUM	ESS2.A	PAT	Earth and Space	Sensemaking	1.627	0.20	0.31	101
1911M002_04	1	MC	EAE	ESS2.A	SC	Earth and Space	Critiquing	0.909	0.28	0.24	87
1911M002_05	1	MC	AQDP	ESS2.B	S, P, and Q	Earth and Space	Investigating	0.166	0.46	0.27	56
1911M023_02	1	MC	AID	LS2.B	PAT	Life	Sensemaking	0.460	0.40	0.51	115
1911M023_05	1	TE	AID	LS2.B	S, P, and Q	Life	Sensemaking	-0.600	0.64	0.52	69
1911M023_06	1	MC	OECI	LS2.B	PAT	Life	Critiquing	0.420	0.39	0.27	73
1911M023_07	1	MC	PACI	LS2.B	PAT	Life	Investigating	0.465	0.38	0.22	58
1911M028_01	1	MC	UMCT	PS3.A	S & SM	Physical	Investigating	-0.617	0.57	0.27	94
1911M028_03	1	MC	UMCT	PS3.A	S & SM	Physical	Investigating	1.373	0.25	0.43	75
1911M028_04	1	MC	UMCT	PS3.B	S & SM	Physical	Investigating	0.103	0.43	0.21	49
1911M028_06	1	MC	UMCT	PS3.B	S & SM	Physical	Investigating	0.576	0.37	0.25	101
1911M079_02	1	MC	AID	ESS1.C	PAT	Earth and Space	Sensemaking	0.249	0.40	0.22	87
1911M079_03	1	TE	DUM	ESS1.C	S, P, and Q	Earth and Space	Sensemaking	-0.025	0.48	0.47	110
1911M079_04	1	TE	OECI	ESS1.C	S, P, and Q	Earth and Space	Critiquing	1.241	0.22	0.45	109
1911M119_02	1	MC	EAE	ESS2.C	PAT	Earth and Space	Critiquing	0.675	0.35	0.27	89
1911M119_05	1	TE	UMCT	ESS2.D	S, P, and Q	Earth and Space	Investigating	0.121	0.45	0.45	105
1911M119_06	1	MC	AID	ESS2.C	S, P, and Q	Earth and Space	Sensemaking	0.312	0.42	0.51	91
1911M124_01	1	TE	AQDP	PS1.A	SF	Physical	Investigating	1.128	0.27	0.20	45
1911M124_02	1	TE	AID	PS1.A	SC	Physical	Sensemaking	0.871	0.31	0.46	99
1911M124_05	1	MC	DUM	PS1.A	SF	Physical	Sensemaking	-0.426	0.57	0.21	54

Grade 11											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	<i>rpb</i>	Median Time
1911M124_10	1	TE	EAE	PS1.A	C and E	Physical	Critiquing	-0.100	0.50	0.47	45
2011M003_01	1	TE	AID	LS2.C	SC	Life	Sensemaking	0.955	0.29	0.49	58
2011M003_03	1	MC	OECI	LS2.C	SC	Life	Critiquing	1.147	0.26	0.29	44
2011M003_04	1	MC	AQDP	LS2.C	SC	Life	Investigating	0.788	0.32	0.39	22
2011M003_05	1	TE	EAE	LS2.C	C and E	Life	Critiquing	0.497	0.38	0.54	58
2011M003_06	1	MC	PACI	LS2.C	C and E	Life	Investigating	0.034	0.47	0.40	58
2011M010_01	1	MC	AQDP	PS4.B	SF	Physical	Investigating	1.055	0.28	0.19	40
2011M010_02	1	TE	EAE	PS4.B	S & SM	Physical	Critiquing	0.543	0.37	0.34	32
2011M010_03	1	MC	DUM	PS4.B	S & SM	Physical	Sensemaking	0.360	0.41	0.27	41
2011M010_05	1	TE	EAE	PS4.B	SF	Physical	Critiquing	1.358	0.23	0.28	45
2011M071_01	1	MC	AID	PS3.D	S, P, and Q	Physical	Sensemaking	-0.598	0.60	0.50	41
2011M071_02	1	TE	EAE	PS3.D	S, P, and Q	Physical	Critiquing	1.334	0.23	0.30	87
2011M071_03	1	TE	EAE	PS3.D	E&M	Physical	Critiquing	0.569	0.36	0.48	53
2011M071_04	1	MC	DUM	ESS3.A	S & SM	Physical	Sensemaking	-0.262	0.53	0.29	70
2011M071_05	1	MC	UMCT	PS3.D	E&M	Physical	Investigating	-0.311	0.54	0.43	53
2111M000_02	1	TE	AID	ESS3.A	C and E	Earth and Space	Sensemaking	-0.791	0.64	0.47	79
2111M000_06	1	MC	PACI	ESS3.A	S & SM	Earth and Space	Investigating	-0.148	0.51	0.33	142
2111M000_08	1	TE	AID	ESS3.C	PAT	Earth and Space	Sensemaking	-0.960	0.67	0.43	61
2111M000_09	1	MC	EAE	ESS3.C	PAT	Earth and Space	Critiquing	-0.075	0.50	0.44	100
2111M004_02	1	TE	AID	LS4.C	S, P, and Q	Life	Sensemaking	2.137	0.13	0.26	31
2111M004_03	1	MC	UMCT	LS2.A	C and E	Life	Investigating	1.189	0.25	0.45	37
2111M004_05	1	MC	CEDS	LS1.A	PAT	Life	Sensemaking	-0.260	0.53	0.28	46
2111M004_06	1	TE	EAE	LS4.C	C and E	Life	Critiquing	1.328	0.23	0.31	49
2211B000_01	1	TE	AQDP	PS2.A	C and E	Physical	Investigating	0.568	0.36	0.42	43
2211B000_03	1	TE	CEDS	PS3.B	S & SM	Physical	Sensemaking	0.536	0.37	0.22	29

Grade 11											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	<i>rpb</i>	Median Time
2211B000_07	1	TE	OECI	PS3.B	S & SM	Physical	Critiquing	1.102	0.27	0.37	31
2211B000_12	3	CR	UMCT	PS3.B	S & SM	Physical	Investigating	1.522	0.19	0.61	194
2211B006_02	1	TE	AQDP	ESS3.B	C and E	Earth and Space	Investigating	0.825	0.31	0.44	64
2211B006_05	1	TE	PACI	ESS3.B	SF	Earth and Space	Investigating	0.971	0.29	0.46	47
2211B006_06	1	TE	PACI	ESS3.B	C and E	Earth and Space	Investigating	-1.220	0.71	0.46	36
2211B006_09	4	CR	EAE	ESS3.B	C and E	Earth and Space	Critiquing	-0.242	0.53	0.61	251
2211B006_12	1	TE	EAE	PS1.A	SF	Physical	Critiquing	0.526	0.37	0.57	43
2211M003_01	1	MC	AQDP	LS1.B	S & SM	Life	Investigating	0.291	0.41	0.36	58
2211M003_02	1	MC	EAE	LS1.B	S & SM	Life	Critiquing	0.568	0.37	0.21	36
2211M003_03	1	TE	EAE	LS1.B	S & SM	Life	Critiquing	0.945	0.29	0.35	40
2211M003_04	1	TE	CEDS	LS1.B	S & SM	Life	Sensemaking	2.320	0.11	0.35	55
2211M003_05	1	TE	PACI	LS3.A	S & SM	Life	Investigating	1.475	0.21	0.45	40
2211M008_01	1	TE	AID	ESS1.B	S & SM	Earth and Space	Sensemaking	1.421	0.22	0.37	97
2211M008_03	1	MC	EAE	ESS1.B	S & SM	Earth and Space	Critiquing	0.242	0.43	0.28	61
2211M008_07	1	TE	PACI	ESS1.B	S & SM	Earth and Space	Investigating	0.896	0.30	0.24	60
HS18060_01	1	TE	OECI	PS1.C	E&M	Physical	Critiquing	1.364	0.23	0.54	146
HS18060_03	1	MC	AID	PS1.C	E&M	Physical	Sensemaking	-0.032	0.52	0.44	41
HS18060_04	1	MC	AID	PS1.C	E&M	Physical	Sensemaking	-0.192	0.56	0.50	34
HS18060_06	1	TE	AID	PS1.C	E&M	Physical	Sensemaking	2.033	0.16	0.35	43

APPENDIX G: SCALE SCORE CUMULATIVE FREQUENCY DISTRIBUTIONS

Table G.1: Grade 5—Scale Score Cumulative Frequency Distribution

Raw Score	Scale Score	All Cum. N	All Cum. %	Female Cum. %	Male Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	White Cum. %
0	100	3	0.00	0.00	0.00	0.01	0.01	0.00	0.00
1	100	45	0.05	0.04	0.06	0.02	0.10	0.06	0.02
2	100	205	0.21	0.15	0.27	0.05	0.41	0.31	0.11
3	100	611	0.63	0.51	0.76	0.12	1.17	1.00	0.30
4	100	1,548	1.61	1.30	1.90	0.24	3.00	2.54	0.73
5	100	3,200	3.32	2.86	3.76	0.69	6.22	5.18	1.51
6	100	5,485	5.69	5.12	6.24	1.24	10.41	8.82	2.68
7	100	8,413	8.73	8.10	9.34	1.93	15.91	13.33	4.27
8	103	11,555	11.99	11.28	12.67	2.65	21.43	18.31	6.03
9	109	14,756	15.31	14.59	16.00	3.50	27.50	23.27	7.73
10	115	18,010	18.68	17.96	19.38	4.50	32.83	28.25	9.73
11	120	21,239	22.03	21.48	22.57	5.42	37.89	33.08	11.92
12	125	24,259	25.17	24.82	25.50	6.46	42.42	37.45	14.14
13	129	27,026	28.04	27.93	28.14	7.53	46.64	41.37	16.17
14	134	29,617	30.73	30.93	30.53	8.67	50.06	45.08	18.19
15	138	32,178	33.38	33.77	33.01	9.83	53.52	48.44	20.40
16	142	34,605	35.90	36.43	35.39	10.92	56.63	51.67	22.54
17	146	37,055	38.44	39.17	37.75	12.15	59.64	54.71	24.85
18	150	39,381	40.86	41.75	40.00	13.38	62.22	57.61	27.13
19	153	41,722	43.28	44.39	42.22	14.78	64.70	60.44	29.49
20	156	43,985	45.63	46.87	44.44	16.28	67.25	62.86	31.93
21	160	46,285	48.02	49.51	46.58	17.82	69.62	65.45	34.40
22	163	48,462	50.28	51.93	48.69	19.45	71.70	67.80	36.85
23	166	50,652	52.55	54.32	50.85	21.08	73.67	70.05	39.45
24	170	52,752	54.73	56.57	52.95	22.71	75.70	72.09	42.00
25	173	54,836	56.89	58.88	54.97	24.45	77.67	74.12	44.47
26	176	56,967	59.10	61.16	57.11	26.57	79.40	76.13	47.10
27	179	59,024	61.23	63.26	59.28	28.48	81.09	78.02	49.69
28	182	60,995	63.28	65.36	61.27	30.40	82.63	79.84	52.13
29	185	63,031	65.39	67.48	63.38	32.46	84.22	81.56	54.82
30	188	65,021	67.45	69.50	65.48	35.04	85.81	83.18	57.30
31	191	67,069	69.58	71.61	67.62	37.53	87.21	84.77	60.01
32	193	69,032	71.62	73.69	69.62	39.95	88.47	86.28	62.66
33	196	70,586	73.23	75.23	71.30	42.03	89.45	87.37	64.81
34	200	72,663	75.38	77.37	73.46	44.81	90.73	89.01	67.55

Raw Score	Scale Score	All Cum. N	All Cum. %	Female Cum. %	Male Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	White Cum. %
35	202	74,587	77.38	79.37	75.46	47.89	91.82	90.23	70.16
36	205	76,400	79.26	81.24	77.35	50.98	92.64	91.35	72.73
37	208	78,164	81.09	82.98	79.27	53.85	93.45	92.46	75.22
38	211	79,868	82.86	84.70	81.08	56.79	94.38	93.38	77.65
39	214	81,491	84.54	86.29	82.86	59.90	95.19	94.21	79.94
40	217	83,054	86.16	87.82	84.57	63.07	96.02	95.00	82.12
41	220	84,521	87.68	89.30	86.13	66.35	96.61	95.69	84.14
42	224	85,995	89.21	90.66	87.82	69.38	97.23	96.45	86.21
43	227	87,337	90.61	91.92	89.34	72.57	97.73	97.02	88.12
44	231	88,601	91.92	93.04	90.83	75.37	98.23	97.55	89.92
45	234	89,802	93.16	94.17	92.19	78.23	98.61	98.03	91.63
46	238	90,904	94.31	95.16	93.48	81.13	98.95	98.42	93.15
47	243	91,912	95.35	96.11	94.62	84.08	99.12	98.71	94.58
48	246	92,794	96.27	96.93	95.62	86.63	99.32	99.00	95.78
49	251	93,579	97.08	97.65	96.53	89.22	99.51	99.27	96.73
50	256	94,284	97.81	98.27	97.37	91.57	99.63	99.48	97.63
51	261	94,887	98.44	98.75	98.14	93.73	99.71	99.65	98.36
52	267	95,357	98.93	99.12	98.74	95.50	99.79	99.78	98.90
53	273	95,723	99.31	99.47	99.15	96.96	99.90	99.87	99.30
54	281	95,981	99.57	99.68	99.47	98.04	99.94	99.93	99.59
55	289	96,157	99.76	99.82	99.70	98.94	99.98	99.96	99.76
56	300	96,290	99.89	99.93	99.86	99.56	99.99	99.98	99.89
57	300	96,356	99.96	99.97	99.95	99.83	100.00	99.99	99.96
58	300	96,381	99.99	99.99	99.99	99.96	100.00	100.00	99.98
59	300	96,390	100.00	100.00	100.00	99.99	100.00	100.00	100.00
60	300	96,392	100.00	100.00	100.00	100.00	100.00	100.00	100.00

Table G.2: Grade 8–Scale Score Cumulative Frequency Distribution

Raw Score	Scale Score	All Cum. N	All Cum. %	Female Cum. %	Male Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	White Cum. %
0	100	5	0.00	0.00	0.01	0.00	0.01	0.01	0.00
1	100	34	0.03	0.02	0.05	0.00	0.05	0.06	0.01
2	100	144	0.14	0.12	0.16	0.02	0.27	0.22	0.06
3	100	396	0.39	0.35	0.43	0.07	0.77	0.52	0.22
4	100	977	0.96	0.83	1.09	0.21	2.01	1.32	0.49
5	100	2,065	2.03	1.78	2.27	0.38	4.27	2.83	1.02
6	102	3,954	3.90	3.40	4.37	0.77	7.91	5.55	1.93
7	108	6,588	6.49	5.71	7.23	1.32	12.69	9.36	3.26
8	114	9,974	9.83	8.80	10.81	2.02	18.31	14.26	5.21
9	119	13,769	13.57	12.34	14.74	2.87	24.50	19.70	7.45
10	124	18,153	17.89	16.67	19.05	4.07	31.18	26.06	10.11
11	128	22,400	22.07	20.85	23.25	5.38	37.99	31.93	12.74
12	132	26,334	25.95	24.94	26.92	6.45	43.86	37.40	15.37
13	136	30,281	29.84	28.92	30.73	7.87	48.91	42.76	18.31
14	140	33,875	33.38	32.71	34.04	9.22	53.61	47.43	21.09
15	143	37,316	36.77	36.40	37.15	10.73	57.79	51.76	23.95
16	146	40,568	39.98	39.79	40.18	12.28	61.41	55.69	26.90
17	150	43,493	42.86	42.93	42.82	13.77	64.66	59.16	29.59
18	152	46,302	45.63	45.70	45.59	15.27	67.43	62.43	32.30
19	155	48,971	48.26	48.50	48.06	16.76	70.23	65.34	34.90
20	158	51,513	50.76	51.19	50.39	18.15	72.71	68.01	37.61
21	161	53,884	53.10	53.57	52.69	19.72	74.86	70.45	40.17
22	163	56,243	55.42	56.03	54.88	21.59	76.90	72.72	42.81
23	166	58,538	57.69	58.43	57.02	23.24	78.80	74.97	45.38
24	168	60,730	59.85	60.70	59.08	25.16	80.50	76.95	47.93
25	171	62,886	61.97	62.94	61.09	27.12	82.18	78.88	50.46
26	173	65,129	64.18	65.26	63.20	29.34	83.60	80.76	53.23
27	176	67,124	66.15	67.34	65.05	31.50	85.10	82.32	55.64
28	178	69,122	68.12	69.34	66.99	33.62	86.58	83.91	58.03
29	180	71,056	70.02	71.29	68.86	35.78	87.72	85.31	60.58
30	183	72,890	71.83	73.20	70.57	38.02	88.95	86.63	62.86
31	185	74,711	73.62	75.07	72.28	40.28	90.17	87.83	65.23
32	187	76,473	75.36	76.81	74.02	42.89	91.13	89.05	67.43
33	189	78,182	77.04	78.55	75.65	45.48	92.01	90.16	69.63
34	192	79,861	78.70	80.22	77.30	47.83	92.91	91.27	71.82
35	194	81,446	80.26	81.81	78.82	50.40	93.73	92.11	74.00
36	196	82,661	81.46	82.98	80.05	52.31	94.27	92.81	75.64

Raw Score	Scale Score	All Cum. N	All Cum. %	Female Cum. %	Male Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	White Cum. %
37	200	84,408	83.18	84.76	81.71	55.46	95.01	93.75	77.97
38	201	85,857	84.61	86.14	83.18	58.29	95.61	94.49	79.88
39	203	87,167	85.90	87.43	84.48	60.79	96.08	95.05	81.75
40	205	88,441	87.15	88.63	85.78	63.31	96.57	95.65	83.47
41	208	89,661	88.36	89.84	86.96	65.71	97.14	96.14	85.15
42	210	90,807	89.48	90.96	88.10	68.10	97.51	96.67	86.70
43	212	91,866	90.53	91.90	89.24	70.52	97.85	97.11	88.12
44	215	92,886	91.53	92.84	90.31	73.04	98.13	97.53	89.45
45	217	93,828	92.46	93.68	91.32	75.34	98.39	97.86	90.72
46	220	94,745	93.37	94.50	92.30	77.72	98.65	98.19	91.92
47	222	95,575	94.18	95.23	93.20	79.98	98.88	98.51	92.95
48	225	96,320	94.92	95.90	93.99	82.02	99.11	98.75	93.91
49	227	97,047	95.63	96.51	94.81	84.04	99.25	98.99	94.83
50	231	97,669	96.25	97.08	95.47	85.89	99.36	99.20	95.59
51	233	98,222	96.79	97.56	96.06	87.61	99.53	99.34	96.26
52	235	98,764	97.33	98.06	96.64	89.45	99.63	99.47	96.92
53	238	99,199	97.75	98.38	97.16	90.80	99.71	99.58	97.48
54	241	99,628	98.18	98.72	97.66	92.24	99.75	99.69	98.03
55	244	99,972	98.52	99.00	98.06	93.41	99.82	99.76	98.45
56	247	100,289	98.83	99.23	98.45	94.56	99.85	99.82	98.82
57	251	100,532	99.07	99.40	98.76	95.55	99.89	99.88	99.07
58	254	100,766	99.30	99.56	99.05	96.59	99.93	99.91	99.32
59	258	100,936	99.47	99.67	99.27	97.30	99.95	99.93	99.51
60	261	101,072	99.60	99.75	99.46	97.97	99.98	99.95	99.63
61	265	101,181	99.71	99.81	99.61	98.56	99.98	99.97	99.72
62	270	101,288	99.81	99.87	99.76	99.05	99.99	99.98	99.84
63	275	101,365	99.89	99.94	99.84	99.40	100.00	99.98	99.91
64	280	101,419	99.94	99.98	99.91	99.70	100.00	99.99	99.94
65	286	101,443	99.97	99.98	99.95	99.80	100.00	99.99	99.97
66	292	101,459	99.98	99.99	99.97	99.91	100.00	100.00	99.98
67	300	101,470	99.99	99.99	99.99	99.95	100.00	100.00	99.99
68	300	101,475	100.00	100.00	100.00	99.98	100.00	100.00	100.00
69	300	101,478	100.00	100.00	100.00	100.00	100.00	100.00	100.00
70	300	101,478	100.00	100.00	100.00	100.00	100.00	100.00	100.00
71	300	101,478	100.00	100.00	100.00	100.00	100.00	100.00	100.00
72	300	101,478	100.00	100.00	100.00	100.00	100.00	100.00	100.00

Table G.3: Grade 11–Scale Score Cumulative Frequency Distribution

Raw Score	Scale Score	All Cum. N	All Cum. %	Female Cum. %	Male Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	White Cum. %
0	100	2	0.00	0.00	0.00	0.00	0.00	0.01	0.00
1	100	6	0.01	0.01	0.01	0.00	0.02	0.01	0.00
2	100	15	0.02	0.02	0.02	0.00	0.04	0.03	0.00
3	100	50	0.05	0.04	0.07	0.02	0.11	0.07	0.03
4	100	125	0.13	0.10	0.16	0.02	0.25	0.20	0.07
5	100	294	0.31	0.24	0.38	0.06	0.63	0.48	0.16
6	100	655	0.70	0.55	0.84	0.08	1.41	1.03	0.38
7	100	1,346	1.43	1.11	1.74	0.18	2.77	2.10	0.84
8	100	2,454	2.61	2.05	3.15	0.45	5.00	3.81	1.53
9	100	3,953	4.20	3.25	5.12	0.85	7.70	6.20	2.53
10	100	5,950	6.33	5.08	7.54	1.30	11.48	9.34	3.85
11	103	8,380	8.91	7.21	10.56	1.87	15.83	13.16	5.54
12	108	11,095	11.80	9.76	13.78	2.63	20.78	17.44	7.36
13	113	13,929	14.81	12.50	17.06	3.54	25.58	21.69	9.46
14	118	16,959	18.04	15.49	20.51	4.58	30.17	26.32	11.80
15	123	19,974	21.24	18.55	23.86	5.75	34.96	30.84	14.11
16	127	22,838	24.29	21.52	26.98	6.79	39.12	35.06	16.53
17	132	25,611	27.24	24.46	29.95	7.90	43.34	38.99	18.84
18	136	28,413	30.22	27.57	32.81	9.00	47.01	43.06	21.33
19	139	31,214	33.20	30.57	35.77	10.25	50.61	47.04	23.87
20	143	33,806	35.96	33.46	38.40	11.66	53.99	50.44	26.33
21	147	36,374	38.69	36.33	41.01	12.91	57.14	53.81	28.88
22	151	38,804	41.27	39.10	43.43	14.34	60.33	56.78	31.31
23	154	41,186	43.80	41.79	45.80	15.83	63.29	59.60	33.76
24	158	43,486	46.25	44.42	48.06	17.34	65.82	62.48	36.08
25	161	45,700	48.61	46.90	50.29	18.94	68.20	65.01	38.52
26	164	47,928	50.97	49.46	52.48	20.70	70.25	67.69	40.91
27	167	50,009	53.19	51.86	54.51	22.28	72.63	70.01	43.18
28	171	51,964	55.27	54.11	56.44	23.82	74.70	72.06	45.39
29	174	53,847	57.27	56.28	58.28	25.26	76.64	73.99	47.58
30	177	55,733	59.28	58.52	60.06	27.01	78.49	75.83	49.82
31	180	57,567	61.23	60.70	61.79	28.88	80.11	77.70	51.93
32	183	59,404	63.18	62.90	63.50	30.76	81.82	79.40	54.18
33	186	61,148	65.04	64.96	65.16	32.53	83.36	80.89	56.38
34	189	62,932	66.93	67.03	66.89	34.61	84.74	82.50	58.60
35	192	64,634	68.74	69.05	68.50	36.69	85.97	83.89	60.82
36	196	65,972	70.17	70.60	69.81	38.19	87.03	85.14	62.44
37	200	67,738	72.04	72.60	71.57	40.36	88.34	86.46	64.80
38	202	69,294	73.70	74.38	73.11	42.59	89.41	87.58	66.89

Raw Score	Scale Score	All Cum. N	All Cum. %	Female Cum. %	Male Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	White Cum. %
39	205	70,783	75.28	76.07	74.57	44.79	90.49	88.66	68.83
40	208	72,289	76.88	77.73	76.12	47.01	91.42	89.81	70.79
41	211	73,733	78.42	79.32	77.61	49.32	92.28	90.71	72.79
42	214	75,093	79.87	80.86	78.96	51.49	92.88	91.60	74.71
43	217	76,393	81.25	82.35	80.23	53.82	93.61	92.47	76.40
44	220	77,747	82.69	83.88	81.59	56.27	94.36	93.28	78.24
45	223	79,010	84.03	85.24	82.92	58.50	95.08	94.01	79.99
46	227	80,103	85.20	86.40	84.08	60.59	95.62	94.58	81.53
47	230	81,291	86.46	87.73	85.28	62.81	96.22	95.27	83.16
48	233	82,346	87.58	88.85	86.40	64.93	96.71	95.82	84.65
49	236	83,399	88.70	90.06	87.43	67.11	97.15	96.32	86.15
50	240	84,425	89.79	91.15	88.52	69.74	97.60	96.78	87.51
51	243	85,342	90.77	92.04	89.57	71.74	97.89	97.20	88.81
52	246	86,249	91.73	92.99	90.55	73.75	98.19	97.58	90.11
53	250	87,167	92.71	93.86	91.62	76.20	98.48	97.93	91.38
54	253	87,965	93.56	94.66	92.52	78.47	98.74	98.26	92.42
55	257	88,708	94.35	95.28	93.47	80.71	98.94	98.50	93.40
56	261	89,422	95.11	95.99	94.28	83.04	99.15	98.70	94.32
57	265	90,133	95.86	96.67	95.10	85.17	99.37	98.91	95.29
58	268	90,746	96.51	97.28	95.80	87.20	99.53	99.12	96.06
59	273	91,294	97.10	97.76	96.47	88.91	99.60	99.29	96.81
60	277	91,803	97.64	98.24	97.07	90.70	99.66	99.42	97.46
61	281	92,228	98.09	98.61	97.60	92.20	99.73	99.53	98.02
62	286	92,648	98.54	98.93	98.16	93.94	99.80	99.64	98.50
63	291	92,986	98.90	99.21	98.59	95.15	99.85	99.74	98.94
64	296	93,261	99.19	99.44	98.96	96.34	99.87	99.81	99.24
65	300	93,487	99.43	99.62	99.25	97.28	99.93	99.89	99.47
66	300	93,665	99.62	99.73	99.51	98.19	99.97	99.92	99.64
67	300	93,788	99.75	99.82	99.69	98.77	99.98	99.95	99.77
68	300	93,891	99.86	99.89	99.83	99.25	99.99	99.98	99.88
69	300	93,945	99.92	99.93	99.90	99.62	100.00	99.98	99.92
70	300	93,987	99.96	99.97	99.96	99.87	100.00	99.98	99.96
71	300	94,008	99.98	99.98	99.98	99.94	100.00	99.99	99.99
72	300	94,018	99.99	99.99	100.00	99.98	100.00	100.00	100.00
73	300	94,023	100.00	100.00	100.00	100.00	100.00	100.00	100.00
74	300	94,023	100.00	100.00	100.00	100.00	100.00	100.00	100.00
75	300	94,023	100.00	100.00	100.00	100.00	100.00	100.00	100.00
76	300	94,023	100.00	100.00	100.00	100.00	100.00	100.00	100.00
77	300	94,023	100.00	100.00	100.00	100.00	100.00	100.00	100.00

APPENDIX H: ITEM PARAMETERS AND MODEL FIT TABLES

Table H.1: Grade 5–IRT Item Parameters and Fit Statistics

NJSLA–S Grade 5							
Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1905B007_01	1.017	0.80	0.71	0.62	1.39	0.00	0.33
1905B007_03	0.902	0.79	0.77	0.60	1.40	0.00	0.34
1905B007_05	0.553	0.82	0.78	0.59	1.42	0.00	0.41
1905B007_08	0.393	1.21	1.16	0.65	0.89	0.00	1.63
1905B007_10	0.541	0.75	0.68	0.66	1.59	0.00	0.42
1905B009_01	0.350	0.78	0.71	0.62	1.55	0.00	0.44
1905B009_02	-0.809	0.71	0.62	0.65	1.58	0.00	0.65
1905B009_05	1.579	1.10	1.27	0.26	0.83	0.03	0.22
1905M005_01	0.232	0.95	0.93	0.47	1.12	0.00	0.46
1905M005_03	0.881	0.89	0.84	0.55	1.23	0.00	0.37
1905M005_04	-0.282	1.20	1.27	0.25	0.51	0.18	0.58
1905M008_01	-0.515	0.82	0.75	0.53	1.42	0.00	0.65
1905M008_05	0.111	0.84	0.80	0.56	1.40	0.00	0.53
1905M008_06	-0.794	0.87	0.80	0.49	1.28	0.00	0.67
1905M040_01	0.733	0.95	0.94	0.47	1.11	0.00	0.37
1905M040_03	-0.224	1.00	0.99	0.44	1.01	0.02	0.55
1905M040_05	0.775	1.02	1.02	0.39	0.96	0.00	0.35
1905M044_02	0.072	0.86	0.82	0.55	1.35	0.00	0.50
1905M044_03	1.038	0.96	0.96	0.45	1.07	0.00	0.32
1905M044_04	0.990	1.02	1.06	0.39	0.94	0.01	0.32
1905M076_01	-0.352	1.23	1.29	0.23	0.44	0.18	0.58
1905M076_03	0.472	0.87	0.83	0.54	1.31	0.00	0.42
1905M076_05	1.204	0.88	0.87	0.51	1.20	0.00	0.29
2105M015_04	1.026	1.06	1.14	0.35	0.86	0.03	0.31
2105M015_05	0.525	1.39	1.57	0.09	0.04	0.21	0.41
2105M015_06	-0.966	1.01	1.07	0.36	0.96	0.00	0.69
2205B003_01	0.508	1.34	1.50	0.14	0.17	0.19	0.41
2205B003_02	0.809	1.00	1.03	0.42	0.99	0.01	0.35
2205B003_03	1.259	1.17	1.28	0.25	0.71	0.05	0.28
2205B003_04	1.656	1.14	1.18	0.56	0.92	0.01	0.70
2205B003_05	-0.328	0.90	0.87	0.50	1.24	0.00	0.58
2205B009_01	-0.298	0.88	0.83	0.52	1.30	0.00	0.57
2205B009_02	1.339	0.96	1.03	0.43	1.04	0.00	0.27
2205B009_03	1.430	1.17	1.36	0.23	0.71	0.06	0.25

NJSLA–S Grade 5							
Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
2205B009_04	0.372	1.02	1.01	0.42	0.97	0.01	0.44
2205B009_05	1.091	1.23	1.15	0.63	0.88	0.00	1.11
2205M004_03	0.265	0.80	0.75	0.61	1.51	0.00	0.46
2205M004_05	1.162	1.00	1.07	0.40	0.97	0.01	0.29
2205M004_07	0.273	1.13	1.20	0.32	0.66	0.05	0.46
2205M006_01	0.030	1.05	1.04	0.40	0.89	0.05	0.51
2205M006_02	0.591	1.20	1.24	0.26	0.56	0.08	0.40
2205M006_05	0.103	0.97	0.96	0.45	1.08	0.00	0.49
2205M011_01	1.401	1.11	1.42	0.27	0.76	0.05	0.26
2205M011_02	0.576	0.87	0.84	0.55	1.30	0.00	0.40
2205M011_04	0.349	1.08	1.13	0.37	0.79	0.08	0.44
2205M012_01	-0.534	0.90	0.87	0.49	1.23	0.00	0.61
2205M012_03	0.645	0.97	0.99	0.45	1.06	0.00	0.39
2205M012_04	1.587	1.32	1.76	0.06	0.48	0.09	0.23
2205M022_01	0.433	1.00	1.00	0.44	1.01	0.01	0.43
2205M022_03	0.362	0.96	0.94	0.47	1.10	0.00	0.44
2205M022_05	0.447	1.20	1.27	0.26	0.51	0.10	0.42

Table H.2: Grade 8–IRT Item Parameters and Fit Statistics

NJSLA–S Grade 8							
Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1908B000_03	1.003	0.93	0.96	0.29	1.07	0.00	0.17
1908B000_04	0.573	0.96	0.94	0.40	1.06	0.00	0.27
1908B000_08	0.624	0.90	0.93	0.40	1.13	0.00	0.24
1908B000_11	-0.027	0.84	0.81	0.70	1.14	0.00	1.34
1908B000_12	0.589	1.00	1.03	0.38	0.99	0.00	0.27
1908M003_02	1.580	1.15	1.27	0.24	0.86	0.02	0.15
1908M003_07	1.413	1.04	1.23	0.29	0.93	0.01	0.16
1908M003_08	0.397	0.98	0.99	0.42	1.02	0.00	0.31
1908M005_02	1.080	1.07	1.15	0.37	0.91	0.02	0.22
1908M005_03	1.901	0.89	0.73	0.41	1.09	0.00	0.11
1908M005_05	-0.173	0.98	0.97	0.39	1.05	0.00	0.39
1908M026_01	1.459	0.95	0.81	0.42	1.07	0.00	0.15
1908M026_04	-0.096	1.04	1.04	0.34	0.90	0.00	0.39
1908M026_06	-0.205	1.09	1.14	0.30	0.74	0.07	0.41
1908M030_01	0.269	1.09	1.13	0.38	0.81	0.06	0.38
1908M030_02	0.106	1.02	1.05	0.40	0.94	0.03	0.38
1908M030_05	0.362	1.05	1.09	0.39	0.90	0.03	0.34
1908M033_02	0.105	1.08	1.10	0.33	0.82	0.04	0.38
1908M033_03	-0.412	1.09	1.11	0.29	0.73	0.05	0.47
1908M033_04	-0.857	1.02	1.03	0.36	0.94	0.01	0.54
2008M000_01	0.640	1.17	1.34	0.19	0.70	0.06	0.26
2008M000_02	1.107	0.89	0.79	0.46	1.14	0.00	0.19
2008M000_04	0.291	1.13	1.18	0.25	0.74	0.05	0.32
2008M001_01	0.354	0.81	0.71	0.59	1.39	0.00	0.31
2008M001_05	-0.996	0.98	0.95	0.40	1.08	0.02	0.58
2008M001_08	0.343	1.04	1.07	0.34	0.91	0.02	0.31
2008M015_04	1.198	1.02	0.96	0.31	0.99	0.00	0.18
2008M015_05	0.774	1.15	1.22	0.20	0.78	0.04	0.24
2008M015_06	-0.414	0.88	0.86	0.49	1.35	0.00	0.46
2008M015_09	0.181	1.03	1.06	0.35	0.92	0.02	0.34
2108B003_02	-1.118	0.97	0.95	0.38	1.08	0.00	0.60
2108B003_07	0.667	1.05	1.13	0.32	0.91	0.02	0.26
2108B003_08	-0.284	1.12	1.15	0.26	0.66	0.06	0.43
2108B006_01	0.194	0.80	0.72	0.59	1.45	0.00	0.34
2108B006_03	1.232	0.99	1.08	0.33	0.99	0.00	0.17
2108B006_06	1.057	0.98	0.94	0.37	1.03	0.00	0.20

NJSLA–S Grade 8							
Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
2108B006_09	0.175	0.85	0.80	0.54	1.33	0.00	0.34
2108B006_11	1.496	0.90	0.78	0.53	1.05	0.00	0.30
2108B007_01	0.257	0.90	0.85	0.48	1.21	0.00	0.33
2108B007_03	0.980	1.22	1.43	0.11	0.70	0.06	0.21
2108B007_07	0.988	1.08	1.27	0.25	0.86	0.03	0.21
2108B007_08	-1.045	0.94	0.96	0.40	1.15	0.00	0.59
2108M015_02	-0.473	0.96	0.94	0.42	1.12	0.00	0.47
2108M015_06	-0.273	0.89	0.86	0.50	1.32	0.00	0.43
2108M015_09	-0.786	0.87	0.85	0.49	1.41	0.00	0.53
2108M015_10	0.314	1.17	1.24	0.21	0.66	0.07	0.32
2108M027_01	0.714	0.91	0.86	0.46	1.14	0.00	0.25
2108M027_03	-0.495	0.87	0.85	0.50	1.39	0.00	0.47
2108M027_07	-0.143	1.11	1.14	0.27	0.70	0.05	0.41
2208B003_01	-0.204	0.90	0.88	0.48	1.26	0.00	0.41
2208B003_05	-0.107	0.92	0.89	0.47	1.22	0.00	0.40
2208B003_07	1.196	0.88	0.78	0.47	1.14	0.00	0.18
2208B003_09	0.590	1.15	1.25	0.21	0.74	0.05	0.27
2208B003_11	-0.220	1.01	1.00	0.60	1.01	0.00	1.22
2208M021_03	0.581	1.01	1.02	0.36	0.98	0.00	0.27
2208M021_05	0.104	1.02	1.04	0.37	0.95	0.01	0.36
2208M021_09	0.443	1.06	1.11	0.32	0.88	0.03	0.30
2208M021_10	0.065	0.81	0.74	0.58	1.45	0.00	0.36
2208M028_02	1.005	1.23	1.44	0.11	0.69	0.06	0.21
2208M028_06	-0.189	1.14	1.19	0.25	0.61	0.08	0.41
2208M028_07	1.242	1.08	1.22	0.24	0.90	0.02	0.17
2208M028_09	-0.109	1.10	1.16	0.30	0.72	0.08	0.40
2208M051_13	0.091	1.02	1.04	0.37	0.95	0.02	0.36
2208M051_16	-0.379	0.95	0.93	0.44	1.17	0.00	0.46
2208M051_17	0.344	0.96	0.95	0.42	1.07	0.00	0.31

Table H.3: Grade 11–IRT Item Parameters and Fit Statistics

NJSLA–S Grade 11							
Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1911B009_01A	-0.439	0.86	0.81	0.52	1.46	0.00	0.55
1911B009_03A	0.479	0.75	0.70	0.62	1.60	0.00	0.36
1911B009_05A	1.377	0.92	0.83	0.45	1.13	0.00	0.22
1911B009_07A	0.683	1.10	1.07	0.59	0.87	0.00	1.32
1911B009_09A	1.331	0.89	0.97	0.36	1.12	0.00	0.20
1911M002_01	1.627	1.07	1.12	0.31	0.91	0.02	0.20
1911M002_04	0.909	1.09	1.20	0.24	0.81	0.04	0.28
1911M002_05	0.166	1.12	1.14	0.27	0.65	0.08	0.46
1911M023_02	0.460	0.90	0.90	0.51	1.23	0.00	0.40
1911M023_05	-0.600	0.80	0.76	0.52	1.58	0.00	0.64
1911M023_06	0.420	1.13	1.15	0.27	0.69	0.06	0.39
1911M023_07	0.465	1.17	1.21	0.22	0.59	0.07	0.38
1911M028_01	-0.617	1.12	1.16	0.27	0.63	0.00	0.57
1911M028_03	1.373	1.00	1.03	0.43	1.00	0.01	0.25
1911M028_04	0.103	1.17	1.20	0.21	0.49	0.11	0.43
1911M028_06	0.576	1.15	1.21	0.25	0.65	0.09	0.37
1911M079_02	0.249	1.15	1.21	0.22	0.56	0.06	0.40
1911M079_03	-0.025	0.91	0.89	0.47	1.29	0.00	0.48
1911M079_04	1.241	0.84	0.79	0.45	1.23	0.00	0.22
1911M119_02	0.675	1.12	1.22	0.27	0.71	0.08	0.35
1911M119_05	0.121	0.93	0.93	0.45	1.19	0.00	0.45
1911M119_06	0.312	0.88	0.86	0.51	1.32	0.00	0.42
1911M124_01	1.128	1.17	1.25	0.20	0.72	0.05	0.27
1911M124_02	0.871	0.92	0.89	0.46	1.15	0.00	0.31
1911M124_05	-0.426	1.16	1.19	0.21	0.50	0.18	0.57
1911M124_10	-0.100	0.91	0.89	0.47	1.29	0.00	0.50
2011M003_01	0.955	0.89	0.87	0.49	1.19	0.00	0.29
2011M003_03	1.147	1.08	1.16	0.29	0.87	0.03	0.26
2011M003_04	0.788	0.99	1.03	0.39	1.00	0.00	0.32
2011M003_05	0.497	0.85	0.82	0.54	1.36	0.00	0.38
2011M003_06	0.034	0.99	1.00	0.40	1.01	0.03	0.47
2011M010_01	1.055	1.18	1.28	0.19	0.68	0.06	0.28
2011M010_02	0.543	1.05	1.07	0.34	0.88	0.02	0.37
2011M010_03	0.360	1.13	1.18	0.27	0.65	0.09	0.41
2011M010_05	1.358	1.06	1.25	0.28	0.87	0.03	0.23
2011M071_01	-0.598	0.86	0.82	0.50	1.42	0.00	0.60

NJSLA–S Grade 11							
Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
2011M071_02	1.334	1.06	1.12	0.30	0.91	0.02	0.23
2011M071_03	0.569	0.91	0.90	0.48	1.21	0.00	0.36
2011M071_04	-0.262	1.10	1.11	0.29	0.69	0.10	0.53
2011M071_05	-0.311	0.94	0.95	0.43	1.18	0.00	0.54
2111M000_02	-0.791	0.87	0.84	0.47	1.34	0.00	0.64
2111M000_06	-0.148	1.04	1.06	0.33	0.86	0.02	0.51
2111M000_08	-0.960	0.90	0.87	0.43	1.24	0.00	0.67
2111M000_09	-0.075	0.95	0.93	0.44	1.17	0.00	0.50
2111M004_02	2.137	1.04	1.16	0.26	0.95	0.01	0.13
2111M004_03	1.189	0.92	0.87	0.45	1.13	0.00	0.25
2111M004_05	-0.260	1.09	1.12	0.28	0.70	0.07	0.53
2111M004_06	1.328	1.04	1.11	0.31	0.93	0.02	0.23
2211B000_01	0.568	0.97	0.94	0.42	1.08	0.00	0.36
2211B000_03	0.536	1.17	1.24	0.22	0.59	0.08	0.37
2211B000_07	1.102	1.00	1.00	0.37	1.00	0.00	0.27
2211B000_12	1.522	0.87	0.79	0.61	1.14	0.00	0.56
2211B006_02	0.825	0.95	0.94	0.44	1.10	0.00	0.31
2211B006_05	0.971	0.91	0.94	0.46	1.13	0.00	0.29
2211B006_06	-1.220	0.87	0.80	0.46	1.27	0.00	0.71
2211B006_09	-0.242	1.12	1.16	0.61	0.93	0.00	2.12
2211B006_12	0.526	0.82	0.78	0.57	1.42	0.00	0.37
2211M003_01	0.291	1.03	1.03	0.36	0.93	0.00	0.41
2211M003_02	0.568	1.18	1.29	0.21	0.55	0.12	0.37
2211M003_03	0.945	1.04	1.02	0.35	0.94	0.00	0.29
2211M003_04	2.320	0.96	0.92	0.35	1.04	0.00	0.11
2211M003_05	1.475	0.91	0.87	0.45	1.12	0.00	0.21
2211M008_01	1.421	1.00	0.97	0.37	1.01	0.00	0.22
2211M008_03	0.242	1.10	1.13	0.28	0.71	0.05	0.43
2211M008_07	0.896	1.14	1.18	0.24	0.74	0.05	0.30
HS18060_01	1.364	0.83	0.71	0.54	1.24	0.00	0.23
HS18060_03	-0.032	0.95	0.94	0.44	1.16	0.00	0.52
HS18060_04	-0.192	0.88	0.85	0.50	1.39	0.00	0.56
HS18060_06	2.033	1.07	1.13	0.35	0.93	0.01	0.16

APPENDIX I: RAW SCORE-TO-SCALE-SCORE-CONVERSION TABLES

Table I.1: Grade 5–Operational

NJSLA–S Grade 5							
Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS CSEM	Lower SS	Upper SS
0	–4.873	1.833	–45.295	100	17	100	117
1	–3.650	1.013	6.638	100	17	100	117
2	–2.930	0.726	37.212	100	17	100	117
3	–2.498	0.600	55.557	100	17	100	117
4	–2.184	0.526	68.890	100	17	100	117
5	–1.934	0.476	79.506	100	17	100	117
6	–1.725	0.440	88.381	100	17	100	117
7	–1.544	0.412	96.067	100	17	100	117
8	–1.383	0.390	102.904	103	17	100	120
9	–1.238	0.372	109.061	109	16	100	125
10	–1.105	0.357	114.709	115	15	100	130
11	–0.982	0.345	119.932	120	15	105	135
12	–0.867	0.334	124.815	125	14	111	139
13	–0.759	0.325	129.402	129	14	115	143
14	–0.656	0.317	133.775	134	13	121	147
15	–0.558	0.309	137.937	138	13	125	151
16	–0.464	0.303	141.928	142	13	129	155
17	–0.374	0.298	145.750	146	13	133	159
18	–0.287	0.293	149.445	150	12	138	162
19	–0.203	0.288	153.012	153	12	141	165
20	–0.121	0.284	156.494	156	12	144	168
21	–0.041	0.280	159.891	160	12	148	172
22	0.036	0.277	163.160	163	12	151	175
23	0.112	0.274	166.388	166	12	154	178
24	0.187	0.272	169.572	170	12	158	182
25	0.260	0.269	172.672	173	11	162	184
26	0.332	0.267	175.730	176	11	165	187
27	0.403	0.266	178.745	179	11	168	190
28	0.473	0.264	181.717	182	11	171	193
29	0.542	0.263	184.647	185	11	174	196
30	0.611	0.262	187.577	188	11	177	199
31	0.680	0.262	190.507	191	11	180	202
32	0.748	0.261	193.395	193	11	182	204
33	0.817	0.262	196.325	196	11	185	207

NJSLA–S Grade 5							
Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS CSEM	Lower SS	Upper SS
34	0.885	0.262	199.212	200	11	189	211
35	0.954	0.263	202.142	202	11	191	213
36	1.024	0.264	205.115	205	11	194	216
37	1.094	0.266	208.087	208	11	197	219
38	1.165	0.268	211.102	211	11	200	222
39	1.237	0.270	214.160	214	11	203	225
40	1.311	0.273	217.302	217	12	205	229
41	1.386	0.276	220.487	220	12	208	232
42	1.463	0.280	223.756	224	12	212	236
43	1.543	0.285	227.154	227	12	215	239
44	1.625	0.290	230.636	231	12	219	243
45	1.711	0.296	234.287	234	13	221	247
46	1.801	0.303	238.109	238	13	225	251
47	1.896	0.312	242.143	243	13	230	256
48	1.996	0.321	246.390	246	14	232	260
49	2.103	0.333	250.933	251	14	237	265
50	2.217	0.346	255.774	256	15	241	271
51	2.342	0.361	261.082	261	15	246	276
52	2.479	0.380	266.900	267	16	251	283
53	2.632	0.403	273.397	273	17	256	290
54	2.806	0.432	280.785	281	18	263	299
55	3.008	0.469	289.363	289	20	269	300
56	3.251	0.520	299.682	300	22	278	300
57	3.559	0.595	312.761	300	22	278	300
58	3.985	0.722	330.850	300	22	278	300
59	4.699	1.010	361.170	300	22	278	300
60	5.918	1.831	412.933	300	22	278	300

Note. Grade 5 theta to scale linear conversion: slope = 42.46393; intercept = 161.6317

Table I.2: Grade 8–Operational

NJSLA–S Grade 8							
Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS SCME	Lower SS	Upper SS
0	-5.294	1.831	-15.712	100	18	100	118
1	-4.076	1.010	30.305	100	18	100	118
2	-3.363	0.721	57.242	100	18	100	118
3	-2.938	0.594	73.298	100	18	100	118
4	-2.630	0.519	84.934	100	18	100	118
5	-2.387	0.469	94.115	100	18	100	118
6	-2.185	0.432	101.747	102	16	100	118
7	-2.012	0.403	108.283	108	15	100	123
8	-1.858	0.380	114.101	114	14	100	128
9	-1.721	0.362	119.277	119	14	105	133
10	-1.596	0.346	123.999	124	13	111	137
11	-1.480	0.333	128.382	128	13	115	141
12	-1.374	0.321	132.386	132	12	120	144
13	-1.274	0.311	136.164	136	12	124	148
14	-1.179	0.303	139.753	140	11	129	151
15	-1.090	0.295	143.116	143	11	132	154
16	-1.005	0.288	146.327	146	11	135	157
17	-0.923	0.282	149.425	150	11	139	161
18	-0.845	0.277	152.372	152	10	142	162
19	-0.770	0.272	155.205	155	10	145	165
20	-0.697	0.268	157.963	158	10	148	168
21	-0.626	0.264	160.646	161	10	151	171
22	-0.557	0.261	163.253	163	10	153	173
23	-0.489	0.258	165.822	166	10	156	176
24	-0.423	0.256	168.315	168	10	158	178
25	-0.359	0.254	170.733	171	10	161	181
26	-0.295	0.252	173.151	173	10	163	183
27	-0.232	0.250	175.531	176	9	167	185
28	-0.170	0.249	177.873	178	9	169	187
29	-0.108	0.248	180.216	180	9	171	189
30	-0.047	0.247	182.520	183	9	174	192
31	0.014	0.246	184.825	185	9	176	194
32	0.074	0.246	187.092	187	9	178	196
33	0.134	0.245	189.359	189	9	180	198
34	0.195	0.245	191.663	192	9	183	201
35	0.255	0.245	193.930	194	9	185	203

NJSLA–S Grade 8							
Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS SCME	Lower SS	Upper SS
36	0.315	0.245	196.197	196	9	187	205
37	0.375	0.246	198.464	200	9	191	209
38	0.436	0.246	200.768	201	9	192	210
39	0.497	0.247	203.073	203	9	194	212
40	0.558	0.248	205.377	205	9	196	214
41	0.620	0.249	207.720	208	9	199	217
42	0.682	0.250	210.062	210	9	201	219
43	0.744	0.251	212.404	212	9	203	221
44	0.808	0.253	214.822	215	10	205	225
45	0.872	0.254	217.240	217	10	207	227
46	0.937	0.256	219.696	220	10	210	230
47	1.003	0.258	222.189	222	10	212	232
48	1.070	0.260	224.721	225	10	215	235
49	1.139	0.263	227.327	227	10	217	237
50	1.208	0.265	229.934	231	10	221	241
51	1.280	0.269	232.654	233	10	223	243
52	1.353	0.272	235.412	235	10	225	245
53	1.428	0.276	238.246	238	10	228	248
54	1.505	0.280	241.155	241	11	230	252
55	1.585	0.285	244.177	244	11	233	255
56	1.668	0.291	247.313	247	11	236	258
57	1.755	0.297	250.600	251	11	240	262
58	1.845	0.305	254.000	254	12	242	266
59	1.940	0.313	257.589	258	12	246	270
60	2.041	0.322	261.405	261	12	249	273
61	2.149	0.333	265.485	265	13	252	278
62	2.264	0.346	269.830	270	13	257	283
63	2.389	0.362	274.553	275	14	261	289
64	2.527	0.380	279.766	280	14	266	294
65	2.680	0.403	285.547	286	15	271	300
66	2.853	0.431	292.082	292	16	276	300
67	3.055	0.468	299.714	300	18	282	300
68	3.298	0.519	308.895	300	18	282	300
69	3.605	0.594	320.493	300	18	282	300
70	4.030	0.721	336.550	300	18	282	300
71	4.742	1.010	363.449	300	18	282	300
72	5.960	1.831	409.465	300	18	282	300

Note. Grade 8 theta to scale linear conversion: slope = 37.78004; intercept = 184.2960

Table I.3: Grade 11–Operational

NJSLA–S Grade 11							
Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS CSEM	Lower SS	Upper SS
0	-5.173	1.831	-98.329	100	18	100	118
1	-3.955	1.010	-33.995	100	18	100	118
2	-3.241	0.721	3.717	100	18	100	118
3	-2.816	0.594	26.165	100	18	100	118
4	-2.509	0.519	42.381	100	18	100	118
5	-2.266	0.468	55.216	100	18	100	118
6	-2.065	0.431	65.832	100	18	100	118
7	-1.892	0.402	74.970	100	18	100	118
8	-1.740	0.379	82.999	100	18	100	118
9	-1.604	0.360	90.182	100	18	100	118
10	-1.480	0.344	96.732	100	18	100	118
11	-1.367	0.330	102.700	103	17	100	120
12	-1.261	0.319	108.299	108	17	100	125
13	-1.163	0.308	113.475	113	16	100	129
14	-1.071	0.300	118.335	118	16	102	134
15	-0.983	0.292	122.983	123	15	108	138
16	-0.900	0.285	127.367	127	15	112	142
17	-0.821	0.279	131.539	132	15	117	147
18	-0.744	0.274	135.606	136	14	122	150
19	-0.671	0.269	139.462	139	14	125	153
20	-0.599	0.265	143.265	143	14	129	157
21	-0.530	0.262	146.910	147	14	133	161
22	-0.462	0.258	150.501	151	14	137	165
23	-0.396	0.256	153.987	154	14	140	168
24	-0.332	0.253	157.368	158	13	145	171
25	-0.268	0.251	160.748	161	13	148	174
26	-0.205	0.249	164.076	164	13	151	177
27	-0.144	0.248	167.298	167	13	154	180
28	-0.083	0.246	170.520	171	13	158	184
29	-0.022	0.245	173.742	174	13	161	187
30	0.038	0.244	176.911	177	13	164	190
31	0.097	0.243	180.027	180	13	167	193
32	0.156	0.243	183.143	183	13	170	196
33	0.215	0.242	186.260	186	13	173	199
34	0.273	0.242	189.323	189	13	176	202
35	0.332	0.241	192.439	192	13	179	205
36	0.390	0.241	195.503	196	13	183	209
37	0.448	0.241	198.566	200	13	187	213
38	0.506	0.241	201.630	202	13	189	215

NJSLA–S Grade 11							
Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS CSEM	Lower SS	Upper SS
39	0.564	0.241	204.693	205	13	192	218
40	0.622	0.241	207.757	208	13	195	221
41	0.681	0.242	210.873	211	13	198	224
42	0.739	0.242	213.937	214	13	201	227
43	0.798	0.243	217.053	217	13	204	230
44	0.857	0.244	220.169	220	13	207	233
45	0.917	0.245	223.339	223	13	210	236
46	0.977	0.246	226.508	227	13	214	240
47	1.038	0.247	229.730	230	13	217	243
48	1.099	0.248	232.952	233	13	220	246
49	1.161	0.250	236.226	236	13	223	249
50	1.224	0.252	239.554	240	13	227	253
51	1.288	0.253	242.934	243	13	230	256
52	1.352	0.256	246.315	246	14	232	260
53	1.418	0.258	249.801	250	14	236	264
54	1.486	0.261	253.393	253	14	239	267
55	1.555	0.264	257.037	257	14	243	271
56	1.625	0.267	260.734	261	14	247	275
57	1.698	0.271	264.590	265	14	251	279
58	1.772	0.275	268.499	268	15	253	283
59	1.849	0.280	272.566	273	15	258	288
60	1.929	0.285	276.791	277	15	262	292
61	2.012	0.291	281.175	281	15	266	296
62	2.099	0.298	285.771	286	16	270	300
63	2.190	0.305	290.577	291	16	275	300
64	2.286	0.314	295.648	296	17	279	300
65	2.387	0.324	300.982	300	17	283	300
66	2.495	0.335	306.687	300	17	283	300
67	2.612	0.348	312.867	300	17	283	300
68	2.738	0.363	319.522	300	17	283	300
69	2.877	0.382	326.864	300	17	283	300
70	3.031	0.404	334.998	300	17	283	300
71	3.206	0.433	344.241	300	17	283	300
72	3.409	0.470	354.963	300	17	283	300
73	3.653	0.520	367.851	300	17	283	300
74	3.961	0.595	384.119	300	17	283	300
75	4.387	0.722	406.620	300	17	283	300
76	5.101	1.010	444.333	300	17	283	300
77	6.320	1.831	508.719	300	17	283	300

Note. Grade 11 theta to scale linear conversion: slope = 52.81995; intercept = 174.9036

APPENDIX J: RAW SCORE-TO-THETA SUBSCORE TABLES

Table J.1: Grade 5 Earth and Space Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.987	1.849	-6.761	-1.214	Below
1	-2.723	1.043	-4.287	-1.159	Below
2	-1.942	0.768	-3.093	-0.790	Below
3	-1.446	0.652	-2.423	-0.469	Below
4	-1.066	0.585	-1.944	-0.188	Below
5	-0.750	0.542	-1.562	0.063	Below
6	-0.474	0.510	-1.239	0.291	Below
7	-0.227	0.485	-0.955	0.501	Below
8	-0.001	0.466	-0.700	0.697	Below
9	0.208	0.450	-0.467	0.884	Below
10	0.406	0.440	-0.253	1.065	Near/Met
11	0.596	0.434	-0.055	1.247	Near/Met
12	0.784	0.435	0.133	1.436	Near/Met
13	0.976	0.442	0.313	1.638	Near/Met
14	1.177	0.456	0.492	1.861	Near/Met
15	1.395	0.480	0.675	2.115	Near/Met
16	1.641	0.514	0.871	2.411	Above
17	1.928	0.561	1.086	2.770	Above
18	2.282	0.632	1.333	3.231	Above
19	2.755	0.753	1.625	3.884	Above
20	3.515	1.033	1.966	5.065	Above
21	4.766	1.844	2.000	7.531	Above

Table J.2: Grade 5 Life Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.704	1.841	-6.466	-0.943	Below
1	-2.461	1.027	-4.002	-0.920	Below
2	-1.712	0.746	-2.831	-0.593	Below
3	-1.249	0.627	-2.189	-0.309	Below
4	-0.899	0.560	-1.740	-0.059	Below
5	-0.610	0.519	-1.388	0.169	Below
6	-0.355	0.492	-1.093	0.383	Below
7	-0.123	0.474	-0.834	0.588	Below
8	0.096	0.463	-0.598	0.790	Below
9	0.307	0.456	-0.377	0.991	Near/Met
10	0.514	0.454	-0.167	1.194	Near/Met
11	0.720	0.455	0.037	1.402	Near/Met
12	0.929	0.461	0.238	1.620	Near/Met
13	1.146	0.472	0.438	1.854	Near/Met
14	1.377	0.490	0.642	2.112	Near/Met
15	1.630	0.518	0.853	2.407	Above
16	1.919	0.560	1.079	2.759	Above
17	2.269	0.628	1.327	3.210	Above
18	2.734	0.747	1.613	3.855	Above
19	3.485	1.028	1.942	5.027	Above
20	4.729	1.841	1.967	7.491	Above

Table J.3: Grade 5 Physical Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.533	1.850	-6.307	-0.759	Below
1	-2.267	1.044	-3.833	-0.702	Below
2	-1.483	0.770	-2.638	-0.329	Below
3	-0.984	0.655	-1.966	-0.002	Below
4	-0.600	0.590	-1.485	0.285	Below
5	-0.277	0.549	-1.100	0.546	Below
6	0.008	0.520	-0.773	0.789	Below
7	0.268	0.500	-0.482	1.017	Near/Met
8	0.509	0.484	-0.217	1.235	Near/Met
9	0.737	0.472	0.029	1.446	Near/Met
10	0.957	0.465	0.260	1.653	Near/Met
11	1.171	0.462	0.478	1.864	Near/Met
12	1.385	0.466	0.686	2.085	Near/Met
13	1.608	0.480	0.889	2.327	Above
14	1.849	0.505	1.092	2.607	Above
15	2.125	0.548	1.303	2.948	Above
16	2.463	0.619	1.534	3.392	Above
17	2.920	0.744	1.804	4.036	Above
18	3.669	1.029	2.126	5.213	Above
19	4.916	1.843	2.151	7.680	Above

Table J.4: Grade 5 Sensemaking Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.828	1.847	-6.598	-1.057	Below
1	-2.568	1.040	-4.127	-1.008	Below
2	-1.792	0.765	-2.938	-0.645	Below
3	-1.299	0.650	-2.274	-0.325	Below
4	-0.921	0.586	-1.800	-0.041	Below
5	-0.601	0.547	-1.422	0.220	Below
6	-0.316	0.523	-1.100	0.469	Below
7	-0.050	0.508	-0.813	0.712	Below
8	0.204	0.501	-0.548	0.956	Near/Met
9	0.455	0.501	-0.297	1.206	Near/Met
10	0.708	0.508	-0.053	1.470	Near/Met
11	0.973	0.522	0.190	1.755	Near/Met
12	1.256	0.546	0.438	2.075	Near/Met
13	1.574	0.584	0.698	2.450	Near/Met
14	1.950	0.647	0.979	2.921	Above
15	2.439	0.762	1.295	3.582	Above
16	3.210	1.038	1.654	4.767	Above
17	4.467	1.846	1.698	7.236	Above

Table J.5: Grade 5 Critiquing Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.766	1.842	-6.529	-1.003	Below
1	-2.521	1.028	-4.064	-0.979	Below
2	-1.771	0.746	-2.890	-0.653	Below
3	-1.310	0.624	-2.246	-0.374	Below
4	-0.966	0.554	-1.796	-0.136	Below
5	-0.686	0.507	-1.447	0.074	Below
6	-0.447	0.473	-1.157	0.264	Below
7	-0.235	0.448	-0.906	0.437	Below
8	-0.044	0.427	-0.684	0.597	Below
9	0.131	0.410	-0.484	0.746	Below
10	0.294	0.397	-0.301	0.889	Below
11	0.447	0.386	-0.132	1.026	Near/Met
12	0.592	0.378	0.025	1.160	Near/Met
13	0.733	0.373	0.174	1.293	Near/Met
14	0.871	0.371	0.316	1.427	Near/Met
15	1.009	0.371	0.452	1.565	Near/Met
16	1.147	0.374	0.587	1.707	Near/Met
17	1.289	0.379	0.720	1.858	Near/Met
18	1.436	0.388	0.853	2.018	Near/Met
19	1.591	0.401	0.989	2.193	Above
20	1.759	0.420	1.130	2.389	Above
21	1.946	0.445	1.278	2.613	Above
22	2.158	0.480	1.439	2.878	Above
23	2.412	0.530	1.617	3.206	Above
24	2.731	0.605	1.824	3.637	Above
25	3.169	0.731	2.072	4.267	Above
26	3.899	1.019	2.371	5.428	Above
27	5.131	1.837	2.376	7.886	Above

Table J.6: Grade 5 Investigating Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.675	1.854	-6.456	-0.893	Below
1	-2.397	1.052	-3.974	-0.819	Below
2	-1.597	0.780	-2.767	-0.426	Below
3	-1.081	0.667	-2.082	-0.080	Below
4	-0.679	0.606	-1.587	0.230	Below
5	-0.336	0.568	-1.188	0.516	Below
6	-0.027	0.545	-0.845	0.791	Below
7	0.263	0.533	-0.536	1.062	Near/Met
8	0.543	0.528	-0.249	1.335	Near/Met
9	0.823	0.531	0.026	1.620	Near/Met
10	1.110	0.542	0.297	1.924	Near/Met
11	1.415	0.564	0.569	2.261	Near/Met
12	1.752	0.600	0.852	2.652	Above
13	2.146	0.661	1.155	3.137	Above
14	2.651	0.773	1.492	3.810	Above
15	3.439	1.045	1.872	5.006	Above
16	4.706	1.850	1.931	7.481	Above

Table J.7: Grade 8 Earth and Space Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.179	1.844	-6.945	-1.413	Below
1	-2.928	1.034	-4.478	-1.377	Below
2	-2.164	0.756	-3.298	-1.030	Below
3	-1.686	0.639	-2.644	-0.727	Below
4	-1.321	0.573	-2.181	-0.462	Below
5	-1.018	0.531	-1.815	-0.221	Below
6	-0.752	0.503	-1.506	0.003	Below
7	-0.509	0.484	-1.234	0.217	Below
8	-0.282	0.470	-0.987	0.424	Below
9	-0.065	0.461	-0.757	0.627	Near/Met
10	0.145	0.456	-0.539	0.829	Near/Met
11	0.352	0.453	-0.328	1.031	Near/Met
12	0.557	0.453	-0.123	1.236	Near/Met
13	0.762	0.455	0.080	1.445	Near/Met
14	0.971	0.460	0.282	1.661	Near/Met
15	1.186	0.469	0.483	1.889	Above
16	1.412	0.483	0.687	2.137	Above
17	1.657	0.508	0.895	2.418	Above
18	1.934	0.547	1.113	2.754	Above
19	2.267	0.613	1.347	3.187	Above
20	2.714	0.734	1.612	3.815	Above
21	3.446	1.020	1.916	4.975	Above
22	4.678	1.837	1.922	7.435	Above

Table J.8: Grade 8 Life Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.113	1.837	-6.869	-1.358	Below
1	-2.881	1.021	-4.411	-1.350	Below
2	-2.145	0.737	-3.250	-1.040	Below
3	-1.696	0.615	-2.618	-0.774	Below
4	-1.362	0.545	-2.180	-0.545	Below
5	-1.091	0.500	-1.840	-0.341	Below
6	-0.857	0.469	-1.560	-0.154	Below
7	-0.649	0.447	-1.318	0.021	Below
8	-0.456	0.431	-1.103	0.190	Below
9	-0.276	0.420	-0.906	0.354	Below
10	-0.103	0.412	-0.721	0.515	Near/Met
11	0.065	0.408	-0.547	0.677	Near/Met
12	0.230	0.406	-0.378	0.839	Near/Met
13	0.395	0.406	-0.214	1.004	Near/Met
14	0.561	0.408	-0.052	1.173	Near/Met
15	0.729	0.413	0.110	1.348	Near/Met
16	0.902	0.420	0.273	1.532	Near/Met
17	1.082	0.429	0.438	1.726	Above
18	1.272	0.442	0.609	1.935	Above
19	1.475	0.459	0.786	2.163	Above
20	1.696	0.482	0.972	2.419	Above
21	1.943	0.514	1.172	2.713	Above
22	2.228	0.558	1.391	3.066	Above
23	2.577	0.627	1.636	3.518	Above
24	3.041	0.747	1.921	4.162	Above
25	3.792	1.028	2.250	5.334	Above
26	5.035	1.841	2.274	7.797	Above

Table J.9: Grade 8 Physical Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.254	1.844	-7.021	-1.488	Below
1	-3.003	1.033	-4.552	-1.455	Below
2	-2.246	0.750	-3.371	-1.120	Below
3	-1.779	0.627	-2.719	-0.838	Below
4	-1.432	0.555	-2.264	-0.600	Below
5	-1.152	0.507	-1.912	-0.392	Below
6	-0.913	0.473	-1.622	-0.204	Below
7	-0.701	0.449	-1.375	-0.028	Below
8	-0.507	0.433	-1.156	0.142	Below
9	-0.325	0.423	-0.959	0.310	Below
10	-0.148	0.417	-0.774	0.478	Near/Met
11	0.025	0.416	-0.599	0.649	Near/Met
12	0.199	0.418	-0.428	0.825	Near/Met
13	0.375	0.423	-0.259	1.009	Near/Met
14	0.556	0.430	-0.088	1.201	Near/Met
15	0.745	0.440	0.085	1.405	Near/Met
16	0.944	0.453	0.265	1.624	Near/Met
17	1.157	0.471	0.452	1.863	Above
18	1.389	0.493	0.649	2.129	Above
19	1.647	0.525	0.860	2.434	Above
20	1.945	0.569	1.091	2.799	Above
21	2.306	0.637	1.350	3.262	Above
22	2.784	0.756	1.649	3.919	Above
23	3.549	1.035	1.996	5.102	Above
24	4.804	1.845	2.035	7.572	Above

Table J.10: Grade 8 Sensemaking Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.277	1.841	-7.038	-1.516	Below
1	-3.036	1.026	-4.575	-1.497	Below
2	-2.291	0.742	-3.404	-1.178	Below
3	-1.836	0.618	-2.763	-0.908	Below
4	-1.500	0.546	-2.318	-0.681	Below
5	-1.229	0.497	-1.975	-0.483	Below
6	-0.999	0.463	-1.694	-0.304	Below
7	-0.796	0.439	-1.454	-0.138	Below
8	-0.612	0.421	-1.244	0.021	Below
9	-0.439	0.409	-1.053	0.175	Below
10	-0.275	0.402	-0.878	0.328	Below
11	-0.115	0.398	-0.712	0.481	Near/Met
12	0.042	0.397	-0.553	0.637	Near/Met
13	0.200	0.399	-0.398	0.798	Near/Met
14	0.360	0.403	-0.243	0.964	Near/Met
15	0.525	0.409	-0.088	1.138	Near/Met
16	0.695	0.417	0.070	1.321	Near/Met
17	0.873	0.428	0.231	1.515	Near/Met
18	1.062	0.442	0.399	1.725	Above
19	1.265	0.460	0.575	1.955	Above
20	1.487	0.484	0.762	2.212	Above
21	1.736	0.515	0.963	2.509	Above
22	2.024	0.561	1.183	2.865	Above
23	2.375	0.630	1.431	3.319	Above
24	2.843	0.750	1.719	3.967	Above
25	3.597	1.030	2.052	5.142	Above
26	4.844	1.842	2.080	7.607	Above

Table J.11: Grade 8 Critiquing Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.370	1.841	-7.131	-1.609	Below
1	-3.128	1.027	-4.668	-1.587	Below
2	-2.379	0.745	-3.497	-1.261	Below
3	-1.918	0.624	-2.854	-0.982	Below
4	-1.573	0.555	-2.405	-0.741	Below
5	-1.291	0.510	-2.056	-0.527	Below
6	-1.048	0.479	-1.766	-0.330	Below
7	-0.830	0.456	-1.515	-0.145	Below
8	-0.629	0.441	-1.290	0.032	Below
9	-0.440	0.430	-1.085	0.206	Below
10	-0.258	0.424	-0.893	0.378	Below
11	-0.080	0.420	-0.710	0.551	Near/Met
12	0.097	0.420	-0.534	0.727	Near/Met
13	0.274	0.423	-0.360	0.909	Near/Met
14	0.455	0.429	-0.188	1.098	Near/Met
15	0.642	0.437	-0.013	1.298	Near/Met
16	0.839	0.450	0.164	1.513	Near/Met
17	1.048	0.466	0.349	1.747	Near/Met
18	1.276	0.489	0.542	2.009	Above
19	1.529	0.520	0.749	2.309	Above
20	1.822	0.565	0.975	2.669	Above
21	2.178	0.633	1.228	3.127	Above
22	2.650	0.752	1.521	3.778	Above
23	3.408	1.032	1.861	4.956	Above
24	4.658	1.843	1.893	7.422	Above

Table J.12: Grade 8 Investigating Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.825	1.846	-6.593	-1.056	Below
1	-2.569	1.037	-4.124	-1.014	Below
2	-1.801	0.759	-2.939	-0.662	Below
3	-1.318	0.642	-2.280	-0.355	Below
4	-0.950	0.576	-1.813	-0.087	Below
5	-0.644	0.533	-1.443	0.155	Below
6	-0.377	0.504	-1.132	0.379	Below
7	-0.134	0.483	-0.858	0.590	Near/Met
8	0.091	0.467	-0.610	0.792	Near/Met
9	0.304	0.456	-0.380	0.988	Near/Met
10	0.508	0.449	-0.164	1.181	Near/Met
11	0.707	0.444	0.042	1.373	Near/Met
12	0.903	0.442	0.240	1.565	Near/Met
13	1.098	0.443	0.434	1.763	Above
14	1.297	0.449	0.624	1.970	Above
15	1.503	0.460	0.813	2.192	Above
16	1.722	0.478	1.005	2.438	Above
17	1.962	0.506	1.204	2.721	Above
18	2.239	0.550	1.415	3.063	Above
19	2.578	0.619	1.649	3.506	Above
20	3.033	0.741	1.921	4.145	Above
21	3.776	1.025	2.238	5.314	Above
22	5.016	1.840	2.256	7.777	Above

Table J.13: Grade 11 Earth and Space Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.153	1.845	-6.920	-1.387	Below
1	-2.903	1.032	-4.450	-1.355	Below
2	-2.149	0.746	-3.268	-1.030	Below
3	-1.690	0.620	-2.619	-0.761	Below
4	-1.354	0.545	-2.171	-0.537	Below
5	-1.085	0.496	-1.829	-0.342	Below
6	-0.857	0.463	-1.551	-0.162	Below
7	-0.653	0.442	-1.316	0.010	Below
8	-0.463	0.430	-1.108	0.181	Below
9	-0.281	0.424	-0.917	0.355	Below
10	-0.102	0.423	-0.737	0.533	Near/Met
11	0.078	0.426	-0.561	0.717	Near/Met
12	0.262	0.432	-0.386	0.910	Near/Met
13	0.452	0.440	-0.208	1.112	Near/Met
14	0.650	0.451	-0.025	1.326	Near/Met
15	0.859	0.463	0.164	1.554	Near/Met
16	1.081	0.480	0.361	1.801	Near/Met
17	1.322	0.502	0.569	2.075	Above
18	1.588	0.532	0.790	2.386	Above
19	1.893	0.575	1.031	2.755	Above
20	2.260	0.642	1.298	3.222	Above
21	2.743	0.759	1.604	3.881	Above
22	3.511	1.036	1.957	5.066	Above
23	4.767	1.846	1.998	7.535	Above

Table J.14: Grade 11 Life Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.857	1.840	-6.617	-1.096	Below
1	-2.615	1.026	-4.155	-1.076	Below
2	-1.869	0.744	-2.985	-0.752	Below
3	-1.408	0.623	-2.344	-0.473	Below
4	-1.065	0.554	-1.896	-0.233	Below
5	-0.783	0.509	-1.547	-0.020	Below
6	-0.541	0.478	-1.258	0.177	Below
7	-0.323	0.456	-1.007	0.361	Below
8	-0.123	0.440	-0.782	0.537	Near/Met
9	0.065	0.428	-0.577	0.707	Near/Met
10	0.244	0.420	-0.385	0.874	Near/Met
11	0.418	0.414	-0.203	1.039	Near/Met
12	0.588	0.411	-0.028	1.204	Near/Met
13	0.756	0.409	0.142	1.370	Near/Met
14	0.923	0.410	0.308	1.538	Near/Met
15	1.092	0.413	0.473	1.712	Near/Met
16	1.265	0.418	0.638	1.892	Above
17	1.443	0.427	0.803	2.083	Above
18	1.630	0.439	0.972	2.288	Above
19	1.829	0.455	1.147	2.512	Above
20	2.046	0.478	1.329	2.764	Above
21	2.290	0.510	1.525	3.055	Above
22	2.573	0.556	1.739	3.407	Above
23	2.919	0.626	1.980	3.859	Above
24	3.384	0.748	2.262	4.506	Above
25	4.137	1.030	2.592	5.681	Above
26	5.383	1.843	2.619	8.147	Above

Table J.15: Grade 11 Physical Science Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.150	1.841	-6.911	-1.389	Below
1	-2.907	1.028	-4.448	-1.365	Below
2	-2.157	0.747	-3.277	-1.037	Below
3	-1.693	0.627	-2.633	-0.753	Below
4	-1.345	0.558	-2.182	-0.509	Below
5	-1.060	0.513	-1.829	-0.291	Below
6	-0.814	0.481	-1.536	-0.092	Below
7	-0.594	0.458	-1.281	0.093	Below
8	-0.393	0.441	-1.053	0.268	Below
9	-0.205	0.427	-0.846	0.436	Below
10	-0.026	0.417	-0.653	0.600	Near/Met
11	0.145	0.410	-0.471	0.760	Near/Met
12	0.310	0.405	-0.297	0.918	Near/Met
13	0.473	0.402	-0.129	1.075	Near/Met
14	0.633	0.400	0.033	1.233	Near/Met
15	0.794	0.401	0.193	1.394	Near/Met
16	0.955	0.403	0.351	1.559	Near/Met
17	1.118	0.407	0.508	1.729	Above
18	1.287	0.414	0.666	1.907	Above
19	1.461	0.423	0.827	2.095	Above
20	1.645	0.435	0.992	2.298	Above
21	1.842	0.453	1.163	2.520	Above
22	2.057	0.476	1.343	2.770	Above
23	2.297	0.507	1.536	3.058	Above
24	2.577	0.553	1.748	3.406	Above
25	2.920	0.623	1.986	3.853	Above
26	3.379	0.744	2.263	4.494	Above
27	4.124	1.026	2.586	5.663	Above
28	5.365	1.840	2.605	8.124	Above

Table J.16: Grade 11 Sensemaking Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.311	1.842	-7.073	-1.548	Below
1	-3.066	1.029	-4.609	-1.523	Below
2	-2.314	0.748	-3.436	-1.192	Below
3	-1.848	0.628	-2.790	-0.906	Below
4	-1.498	0.559	-2.337	-0.660	Below
5	-1.212	0.514	-1.983	-0.441	Below
6	-0.964	0.483	-1.689	-0.240	Below
7	-0.742	0.460	-1.433	-0.051	Below
8	-0.538	0.444	-1.204	0.128	Below
9	-0.346	0.432	-0.994	0.302	Below
10	-0.164	0.424	-0.799	0.472	Below
11	0.013	0.418	-0.614	0.640	Near/Met
12	0.187	0.415	-0.436	0.809	Near/Met
13	0.358	0.414	-0.263	0.978	Near/Met
14	0.529	0.414	-0.092	1.151	Near/Met
15	0.702	0.417	0.076	1.327	Near/Met
16	0.877	0.421	0.245	1.509	Near/Met
17	1.057	0.428	0.416	1.699	Near/Met
18	1.244	0.437	0.589	1.899	Above
19	1.440	0.449	0.766	2.113	Above
20	1.648	0.466	0.950	2.347	Above
21	1.875	0.489	1.143	2.608	Above
22	2.129	0.521	1.348	2.910	Above
23	2.423	0.566	1.574	3.272	Above
24	2.782	0.636	1.827	3.736	Above
25	3.259	0.757	2.124	4.395	Above
26	4.027	1.037	2.471	5.582	Above
27	5.284	1.847	2.514	8.054	Above

Table J.17: Grade 11 Critiquing Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-3.742	1.833	-6.491	-0.993	Below
1	-2.523	1.010	-4.037	-1.008	Below
2	-1.812	0.718	-2.889	-0.735	Below
3	-1.392	0.590	-2.277	-0.506	Below
4	-1.087	0.519	-1.865	-0.309	Below
5	-0.842	0.474	-1.554	-0.131	Below
6	-0.631	0.447	-1.301	0.039	Below
7	-0.439	0.430	-1.084	0.206	Below
8	-0.259	0.420	-0.889	0.371	Below
9	-0.085	0.415	-0.708	0.538	Near/Met
10	0.087	0.413	-0.533	0.707	Near/Met
11	0.257	0.413	-0.363	0.878	Near/Met
12	0.429	0.415	-0.194	1.051	Near/Met
13	0.602	0.418	-0.025	1.229	Near/Met
14	0.778	0.422	0.145	1.412	Near/Met
15	0.959	0.428	0.317	1.602	Near/Met
16	1.146	0.437	0.491	1.801	Above
17	1.342	0.449	0.669	2.015	Above
18	1.550	0.464	0.853	2.246	Above
19	1.775	0.486	1.046	2.504	Above
20	2.026	0.516	1.251	2.800	Above
21	2.314	0.560	1.474	3.154	Above
22	2.664	0.628	1.722	3.607	Above
23	3.130	0.748	2.008	4.251	Above
24	3.881	1.028	2.339	5.423	Above
25	5.124	1.841	2.363	7.885	Above

Table J.18: Grade 11 Investigating Score Table

Raw Score	Theta	CSEM	Lower	Upper	Level
0	-4.058	1.847	-6.828	-1.289	Below
1	-2.801	1.037	-4.357	-1.246	Below
2	-2.033	0.758	-3.170	-0.896	Below
3	-1.554	0.639	-2.511	-0.596	Below
4	-1.191	0.570	-2.046	-0.337	Below
5	-0.893	0.525	-1.680	-0.106	Below
6	-0.635	0.493	-1.375	0.105	Below
7	-0.403	0.470	-1.108	0.302	Below
8	-0.190	0.453	-0.870	0.490	Below
9	0.009	0.441	-0.652	0.670	Near/Met
10	0.199	0.432	-0.448	0.847	Near/Met
11	0.383	0.425	-0.255	1.021	Near/Met
12	0.562	0.422	-0.071	1.194	Near/Met
13	0.739	0.421	0.108	1.370	Near/Met
14	0.916	0.422	0.283	1.549	Near/Met
15	1.096	0.426	0.457	1.734	Near/Met
16	1.280	0.433	0.630	1.929	Above
17	1.471	0.444	0.806	2.137	Above
18	1.675	0.459	0.986	2.363	Above
19	1.895	0.480	1.174	2.615	Above
20	2.139	0.511	1.373	2.906	Above
21	2.422	0.555	1.589	3.255	Above
22	2.767	0.624	1.831	3.702	Above
23	3.227	0.744	2.110	4.343	Above
24	3.973	1.026	2.434	5.511	Above
25	5.213	1.840	2.453	7.972	Above

APPENDIX K: SUBSCORE PROFICIENCY CLASSIFICATIONS

Table K.1: Grade 5 Content Disaggregated Subscore Proficiency Classifications

Group	N	Earth and Space Science			Life Science			Physical Science		
		%Below	%Near/Met	%Above	%Below	%Near/Met	%Above	%Below	%Near/Met	%Above
All Students	96,392	56.35	31.84	11.82	54.72	29.50	15.78	52.86	36.79	10.35
Male	49,082	53.44	33.24	13.32	54.16	28.89	16.95	51.72	36.69	11.59
Female	47,299	59.36	30.38	10.26	55.31	30.12	14.58	54.05	36.88	9.06
Am. Indian	169	54.44	30.18	15.38	53.85	30.18	15.98	49.70	36.69	13.61
Asian	10,765	25.68	43.51	30.81	23.53	38.38	38.09	24.29	46.94	28.77
Black	14,028	76.08	20.19	3.73	74.98	19.97	5.05	71.85	24.92	3.23
Hispanic	31,700	73.11	22.51	4.38	71.69	21.82	6.49	68.04	27.95	4.01
Pacific Islander	179	54.19	31.84	13.97	48.60	34.64	16.76	53.07	36.87	10.06
White	36,375	44.18	40.57	15.25	42.38	36.82	20.80	41.63	45.67	12.69
EL–Yes	9,158	92.29	7.16	0.55	90.87	7.93	1.20	86.93	12.42	0.66
EL–No	87,234	52.57	34.43	13.00	50.93	31.76	17.31	49.29	39.35	11.37
EconDis–Yes	36,109	76.53	20.26	3.21	75.05	20.10	4.85	71.48	25.38	3.14
EconDis–No	60,283	44.26	38.77	16.97	42.55	35.12	22.33	41.71	43.62	14.67
SWD–Yes	20,003	75.05	19.82	5.13	74.89	17.99	7.11	72.81	22.91	4.28
SWD–No	76,389	51.45	34.98	13.57	49.44	32.51	18.05	47.64	40.42	11.94
CBT	75,574	50.95	35.35	13.70	49.02	32.66	18.31	47.50	40.51	11.99
PBT	143	87.41	9.79	2.80	86.01	9.09	4.90	86.71	10.49	2.80
TTS	17,960	73.21	21.09	5.70	73.24	19.38	7.38	70.24	24.78	4.98
SP	1,473	94.43	5.36	0.20	90.16	8.21	1.63	86.01	13.24	0.75
SP TTS	983	94.10	5.39	0.51	88.50	10.38	1.12	85.15	14.75	0.10
Human Reader	201	86.07	12.94	1.00	88.56	9.95	1.49	84.08	13.93	1.99

Table K.2: Grade 5 Practice Disaggregated Subscore Proficiency Classifications

Group	N	Investigating			Sensemaking			Critiquing		
		%Below	%Near/Met	%Above	%Below	%Near/Met	%Above	%Below	%Near/Met	%Above
All Students	96,392	53.37	36.22	10.42	50.80	36.20	12.99	58.54	31.54	9.92
Male	49,082	52.73	36.19	11.08	48.58	36.67	14.74	57.09	31.72	11.19
Female	47,299	54.02	36.24	9.74	53.11	35.71	11.18	60.04	31.35	8.61
Am. Indian	169	46.15	40.83	13.02	49.70	34.32	15.98	56.80	31.95	11.24
Asian	10,765	24.85	47.21	27.94	20.62	46.60	32.78	26.88	45.12	28.00
Black	14,028	71.86	24.78	3.36	71.51	24.02	4.48	78.66	18.79	2.54
Hispanic	31,700	68.83	27.33	3.85	68.03	26.93	5.04	74.88	21.64	3.48
Pacific Islander	179	54.19	30.17	15.64	46.93	37.43	15.64	54.75	34.08	11.17
White	36,375	42.18	44.61	13.21	37.81	45.37	16.82	46.91	40.53	12.56
EL–Yes	9,158	87.37	12.12	0.51	89.86	9.40	0.74	92.30	7.28	0.41
EL–No	87,234	49.80	38.75	11.46	46.70	39.02	14.28	54.99	34.09	10.92
EconDis–Yes	36,109	72.12	24.94	2.94	71.60	24.50	3.90	78.46	19.11	2.43
EconDis–No	60,283	42.13	42.97	14.90	38.34	43.22	18.44	46.60	38.99	14.41
SWD–Yes	20,003	72.69	22.83	4.48	70.94	23.08	5.97	77.84	18.00	4.16
SWD–No	76,389	48.30	39.72	11.97	45.53	39.64	14.83	53.48	35.09	11.43
CBT	75,574	48.02	39.85	12.13	44.93	40.05	15.02	53.15	35.30	11.56
PBT	143	84.62	9.79	5.59	84.62	11.19	4.20	88.81	6.99	4.20
TTS	17,960	70.48	24.74	4.77	69.29	24.31	6.40	76.08	19.39	4.53
SP	1,473	88.19	11.61	0.20	91.45	8.01	0.54	90.90	8.69	0.41
SP TTS	983	87.79	11.90	0.31	90.03	9.56	0.41	91.25	8.55	0.20
Human Reader	201	83.58	15.42	1.00	85.07	12.94	1.99	90.55	7.96	1.49

Table K.3: Grade 8 Content Disaggregated Subscore Proficiency Classifications

Group	N	Earth and Space Science			Life Science			Physical Science		
		%Below	%Near/Met	%Above	%Below	%Near/Met	%Above	%Below	%Near/Met	%Above
All Students	101,478	65.35	27.42	7.22	65.29	27.69	7.02	66.00	27.13	6.87
Male	52,212	63.66	28.05	8.29	65.64	26.85	7.51	64.48	27.16	8.36
Female	49,201	67.18	26.73	6.09	64.96	28.55	6.49	67.65	27.07	5.27
Am. Indian	154	66.23	27.92	5.84	66.23	26.62	7.14	68.83	27.27	3.90
Asian	10,718	33.33	44.58	22.09	32.27	44.75	22.98	32.94	45.88	21.19
Black	14,998	83.22	15.13	1.65	82.79	15.39	1.82	85.44	13.22	1.34
Hispanic	32,921	80.74	16.91	2.35	79.80	17.99	2.21	81.94	16.12	1.94
Pacific Islander	206	57.28	32.52	10.19	55.34	37.38	7.28	56.80	34.95	8.25
White	39,768	55.12	35.80	9.09	56.26	35.38	8.37	55.03	36.08	8.88
EL–Yes	7,151	95.69	4.18	0.13	94.66	5.15	0.20	95.65	4.15	0.20
EL–No	94,327	63.05	29.19	7.76	63.06	29.40	7.54	63.76	28.87	7.38
EconDis–Yes	35,709	82.80	15.36	1.85	82.13	16.07	1.80	84.31	14.15	1.54
EconDis–No	65,769	55.88	33.97	10.14	56.15	34.00	9.85	56.07	34.17	9.76
SWD–Yes	20,520	82.56	14.47	2.97	82.33	14.68	2.99	82.60	14.53	2.87
SWD–No	80,958	60.99	30.71	8.30	60.97	30.99	8.04	61.80	30.32	7.88
CBT	84,298	61.57	30.22	8.22	61.61	30.37	8.02	62.10	30.04	7.86
PBT	72	93.06	5.56	1.39	84.72	12.50	2.78	88.89	11.11	0.00
TTS	14,433	81.66	15.56	2.78	81.48	16.04	2.48	83.23	14.42	2.35
SP	1,876	95.90	4.00	0.11	93.92	6.02	0.05	95.52	4.37	0.11
SP TTS	729	96.16	3.84	0.00	92.46	7.41	0.14	95.34	4.53	0.14
Human Reader	47	93.62	6.38	0.00	95.74	4.26	0.00	97.87	2.13	0.00

Table K.4: Grade 8 Practice Disaggregated Subscore Proficiency Classifications

Group	N	Investigating			Sensemaking			Critiquing		
		%Below	%Near/Met	%Above	%Below	%Near/Met	%Above	%Below	%Near/Met	%Above
All Students	101,478	62.99	30.26	6.75	66.34	26.33	7.33	66.54	26.73	6.73
Male	52,212	62.79	29.81	7.41	64.08	26.96	8.96	66.43	26.22	7.35
Female	49,201	63.24	30.71	6.05	68.77	25.63	5.60	66.69	27.24	6.07
Am. Indian	154	67.53	25.97	6.49	70.13	24.03	5.84	67.53	24.68	7.79
Asian	10,718	31.27	46.36	22.36	33.07	45.15	21.79	33.50	44.74	21.77
Black	14,998	80.09	18.24	1.67	85.60	12.83	1.57	84.55	13.86	1.59
Hispanic	32,921	77.20	20.76	2.04	82.24	15.69	2.07	81.92	16.00	2.09
Pacific Islander	206	53.88	34.95	11.17	59.22	32.04	8.74	61.17	31.07	7.77
White	39,768	53.91	38.05	8.04	55.52	34.83	9.66	56.66	35.16	8.18
EL–Yes	7,151	94.06	5.72	0.22	95.86	3.93	0.21	95.89	3.97	0.14
EL–No	94,327	60.64	32.12	7.24	64.10	28.02	7.87	64.32	28.45	7.23
EconDis–Yes	35,709	79.15	19.15	1.70	84.62	13.69	1.69	84.05	14.34	1.62
EconDis–No	65,769	54.22	36.29	9.49	56.41	33.19	10.40	57.04	33.45	9.51
SWD–Yes	20,520	80.90	16.41	2.69	82.59	14.19	3.23	83.30	13.97	2.73
SWD–No	80,958	58.46	33.77	7.78	62.22	29.40	8.38	62.29	29.96	7.75
CBT	84,298	59.17	33.10	7.73	62.48	29.15	8.37	62.84	29.49	7.67
PBT	72	80.56	18.06	1.39	87.50	12.50	0.00	90.28	8.33	1.39
TTS	14,433	79.55	18.21	2.24	83.33	14.02	2.65	82.82	14.65	2.52
SP	1,876	94.03	5.81	0.16	95.68	4.21	0.11	94.62	5.28	0.11
SP TTS	729	93.55	6.31	0.14	96.16	3.70	0.14	94.65	5.21	0.14
Human Reader	47	93.62	6.38	0.00	95.74	4.26	0.00	95.74	4.26	0.00

Table K.5: Grade 11 Content Disaggregated Subscore Proficiency Classifications

Group	N	Earth and Space Science			Life Science			Physical Science		
		%Below	%Near/Met	%Above	%Below	%Near/Met	%Above	%Below	%Near/Met	%Above
All Students	94,023	49.40	38.79	11.81	49.53	37.17	13.30	52.27	32.77	14.96
Male	47,959	51.25	36.40	12.35	51.73	34.68	13.60	52.44	30.48	17.08
Female	45,924	47.53	41.26	11.21	47.31	39.75	12.94	52.15	35.13	12.72
Am. Indian	141	60.99	30.50	8.51	53.90	36.88	9.22	60.28	25.53	14.18
Asian	10,003	21.18	47.41	31.41	21.17	43.98	34.85	22.81	37.99	39.20
Black	12,731	68.13	28.20	3.68	68.13	27.86	4.01	70.23	25.03	4.74
Hispanic	28,687	64.34	31.27	4.39	65.48	29.43	5.09	67.92	26.27	5.81
Pacific Islander	313	42.49	46.65	10.86	38.98	48.56	12.46	46.96	36.74	16.29
White	40,005	40.13	45.27	14.61	39.79	43.65	16.56	43.16	38.41	18.43
EL–Yes	5,290	87.92	11.78	0.30	90.57	9.22	0.21	88.39	10.85	0.76
EL–No	88,733	47.10	40.40	12.49	47.09	38.84	14.08	50.12	34.08	15.81
EconDis–Yes	28,095	65.68	30.35	3.97	66.37	28.78	4.85	69.00	25.46	5.54
EconDis–No	65,928	42.46	42.39	15.14	42.36	40.75	16.90	45.14	35.89	18.97
SWD–Yes	18,600	68.77	25.76	5.47	68.26	25.04	6.70	70.67	21.96	7.37
SWD–No	75,423	44.62	42.01	13.37	44.91	40.16	14.92	47.73	35.44	16.83
CBT	84,251	47.17	40.22	12.61	47.22	38.60	14.18	50.18	33.87	15.94
PBT	242	76.86	21.49	1.65	78.51	18.18	3.31	76.86	19.83	3.31
TTS	7,732	63.76	30.20	6.04	64.30	28.62	7.07	66.01	26.01	7.98
SP	1,474	88.20	11.60	0.20	91.04	8.89	0.07	87.86	11.60	0.54
SP TTS	273	89.74	9.89	0.37	88.64	11.36	0.00	85.71	14.29	0.00
Human Reader	26	88.46	11.54	0.00	92.31	7.69	0.00	96.15	0.00	3.85

Table K.6: Grade 11 Practice Disaggregated Subscore Proficiency Classifications

Group	N	Investigating			Sensemaking			Critiquing		
		%Below	%Near/Met	%Above	%Below	%Near/Met	%Above	%Below	%Near/Met	%Above
All Students	94,023	53.00	35.88	11.13	52.03	33.32	14.65	52.05	33.32	14.63
Male	47,959	53.88	33.66	12.46	52.91	31.79	15.30	54.01	30.41	15.58
Female	45,924	52.13	38.16	9.71	51.18	34.89	13.94	50.08	36.33	13.59
Am. Indian	141	55.32	34.04	10.64	59.57	28.37	12.06	56.74	34.04	9.22
Asian	10,003	23.71	44.51	31.78	22.37	40.29	37.34	23.29	39.18	37.53
Black	12,731	71.31	25.58	3.10	71.71	23.75	4.54	69.48	25.50	5.02
Hispanic	28,687	67.50	28.61	3.89	69.05	25.56	5.39	66.82	27.28	5.89
Pacific Islander	313	49.52	40.26	10.22	42.81	42.81	14.38	46.65	39.62	13.74
White	40,005	44.47	41.97	13.56	41.47	39.92	18.61	43.50	38.53	17.97
EL–Yes	5,290	85.65	13.93	0.42	92.72	6.90	0.38	89.64	9.92	0.43
EL–No	88,733	51.05	37.18	11.76	49.61	34.89	15.50	49.81	34.72	15.47
EconDis–Yes	28,095	68.52	27.89	3.59	70.30	24.66	5.04	67.89	26.44	5.67
EconDis–No	65,928	46.38	39.28	14.34	44.25	37.01	18.74	45.30	36.25	18.44
SWD–Yes	18,600	70.36	24.35	5.28	71.40	21.42	7.18	70.40	22.31	7.30
SWD–No	75,423	48.72	38.72	12.57	47.26	36.25	16.49	47.53	36.04	16.43
CBT	84,251	51.07	37.04	11.89	49.64	34.73	15.63	49.92	34.50	15.58
PBT	242	76.03	20.66	3.31	81.82	16.12	2.07	77.27	19.42	3.31
TTS	7,732	65.87	28.53	5.60	67.37	24.94	7.70	66.00	26.16	7.84
SP	1,474	84.67	14.99	0.34	94.50	5.16	0.34	89.01	10.45	0.54
SP TTS	273	83.88	16.12	0.00	93.41	6.59	0.00	88.64	10.99	0.37
Human Reader	26	92.31	7.69	0.00	92.31	7.69	0.00	84.62	15.38	0.00

APPENDIX L: EXECUTIVE SUMMARY OF THE NJSLA–S ALIGNMENT EVALUATION STUDY

Appendix L contains the executive summary from the alignment evaluation study report submitted by edCount, LLC in September 2022.

Introduction

The New Jersey Student Learning Assessment–Science (NJSLA–S) assesses students in grades 5, 8, and 11 on their understanding and explanations of scientific phenomena and scenarios. In spring 2019, the NJSLA–S was administered for the first time. Due to the coronavirus pandemic, statewide assessments were cancelled for the 2019–2020 and 2020–2021 school years; thus, the 2022 year marked the second administration of this assessment.

The NJSLA–S is composed of two parts: a performance-based assessment (PBA) and a machine scorable assessment (MSA). Each item within the NJSLA–S represents an interaction of disciplinary core ideas (DCIs), science and engineering practices (SEPs), and crosscutting concepts (CCCs).

The New Jersey Department of Education (NJDOE) commissioned edCount, LLC (edCount), to conduct an independent evaluation of the alignment quality of the NJSLA–S for grades 5, 8, and 11 with the New Jersey Student Learning Standards for Science (NJSLS–S) in 2022. This report documents the methodology for and results of this independent alignment evaluation. The NJDOE intends to use the information gained via this evaluation to inform decisions about future item and assessment development and for federal peer review purposes.

Evaluation Methodology

Evidence of alignment quality is critical to validity evaluation for standards-based assessments (Forte, 2017; Webb, 1997, 1999). Such evidence must draw upon an examination of how a test has been designed and developed, as well as instances of the test itself (Forte, 2013). As is the case for all validity evidence, evidence of alignment quality is necessary to support the interpretation and use of test scores. A well-aligned test is one that elicits a sample of student performance that is adequate to support inferences about student achievement in relation to the standards-based domains on which the test is based.

None of the traditional alignment methods are suited to meet the challenge of evaluating the multidimensional science standards within the NJSLS–S. These methods, such as Webb (1999), involve panelists’ ratings of content and cognitive complexity and the analysis of those ratings in relation to overall criteria for Domain Concurrence, Balance of Representation, Range of Knowledge, and Depth of Knowledge (DOK), which will not alone address the needs of the NJSLS–S.

To address the unique aspects of the three-dimensional nature of the NJSLS–S and the NJSLA–S, edCount addresses the following alignment questions.

1. To what extent do the blueprints support the consistent creation of test forms that reflect the standards and the score scale?

2. To what extent do the Performance-Level Descriptors (PLDs) reflect meaningful and appropriate score interpretations across the full range of the score scale?
3. To what extent does the set of phenomena, tasks, and items reflect the blueprints and provide performance opportunities across the full range of the score scale?

This approach evaluates the quality of alignment of the assessment to the multidimensional standards and provides evidence of the extent to which the assessment supports inferences about student achievement in relation to the standards-based domains.

Evaluation Findings and Recommendations

Evaluation Question 1 addresses the blueprints (and not test items). This question focuses on the extent to which the blueprints support the consistent creation of test forms that reflect the NJSLS–S and the score scale. edCount evaluators found that the blueprint development of the NJSLS–S is well-documented across all grades, including a clear description of the review and revision process by stakeholders. Each blueprint also meets the alignment criteria of strong evidence of alignment for Domain Concurrence, Balance of Representation, and Phenomena Design.

edCount commends the NJDOE for the use of the emerging best practice of PLDs as a cognitive complexity framework for forms development and for NJDOE’s plans to include range PLD expectations within the test blueprints, as well as on the close monitoring of test content, including longitudinal representation of content and item types by form. Another commendable practice the NJDOE used is the thoroughness of the phenomena design guidance provided within the test development documentation.

To further supplement these practices, edCount recommends that the NJDOE consider including guidance on the balance of score points within the test blueprint, in addition to guidance around the number of items by domain. edCount also recommends that the NJDOE consider including guidance on the longitudinal sampling of content, such as the number of forms to be developed before all assessed DCIs have been represented on a form at least once.

Evaluation Question 2 addresses the PLDs. Evaluators found that the PLDs were developed to represent clear and appropriate expectations for performance on the assessments. Evaluators also found documentation that indicated that the item review process included item alignment to PLDs as part of the new item development process. The PLDs for all three grade levels exhibit strong evidence of alignment with the NJSLS–S in terms of PLD-Domain Concurrence. Panelists noted that every reporting category/domain from the standards is fully represented by the PLDs. The PLDs describe increasingly sophisticated and reasonable levels of performance for the concepts defined in the standards. However, the grade 11 panelists noted some uneven progressions between Levels 3 and 4 of the Earth and Space Science portion of the PLDs, due primarily to vagueness in the PLD wording. Overall, the PLDs for all three evaluated forms meet the criteria for “adequately differentiated.”

edCount commends the NJDOE on the development of extremely detailed range PLDs and on leveraging these PLDs during item development. edCount also commends the NJDOE on the

inclusion of New Jersey educators and experts in the assessment field in the PLD development process.

As a result of these findings for Evaluation Question 2, edCount recommends the NJDOE review the grade 11 Earth and Space Science PLDs to ensure that sufficient progression is clarified in the language.

Evaluation Question 3 focuses on the relationship among the dimensions of the standards, phenomena, and items that contribute to students' scores. For the first part of this question, which examines the development process for the phenomena, tasks, clusters, and items on the test form, edCount evaluators found that the test forms were developed to ensure that they consisted of item clusters tied to phenomena and the overall test forms reflect each respective blueprint. edCount commends the NJDOE on the test form design and the inclusion of clusters of items designed around the same phenomena.

In terms of the extent to which phenomena represent the intended concepts and problems to be solved, panelists found the phenomena from all three test forms are engaging and display strong evidence of alignment. edCount commends the NJDOE on the use of state-specific phenomena and relevant everyday phenomena, which contribute to student engagement. Panelists evaluating the grade 5 form judged the phenomena to be "highly accessible," though panelists evaluating the grades 8 and 11 forms judged these phenomena to be "somewhat accessible," citing some distracting or confusing elements, inclusion of content more appropriately suited to a different grade level, or the requirement of skills other than science knowledge, specifically reading or mathematics skills.

Across the test forms, panelists aligned each item to a DCI, SEP, and CCC, with the option to indicate "no alignment" for each of these dimensions. edCount evaluators used this information to determine the alignment with intended targets. All three forms meet the criteria for strong evidence of alignment with the intended targets, indicating that panelists judged more than 75 percent of the items on the form to align to the intended DCI; panelists also identified 100 percent of items on all three forms as aligning to additional dimensions of the standards (SEP and CCC). edCount commends the NJDOE for the strong representation of multidimensionality within the test forms, reflecting the nature of the NJSLS-S.

All test forms meet expectations for Domain Concurrence, Range of Knowledge, and Balance of Representation. Further, all three test forms display strong evidence of alignment in terms of all three criteria above, with the exception of the grade 5 test form, which meets the criteria for moderate evidence of alignment in terms of Range of Knowledge, given that 27 percent of the Earth and Space Science DCIs are represented on this form. edCount commends the NJDOE for their strong plan for monitoring sampling of all DCIs in these grade levels across forms.

For all three forms, panelists evaluated the items on the form as being cognitively challenging, though panelists noted in all three forms that, while a range of cognitive challenge levels is present within the form, items tend to skew toward the higher levels of cognitive challenge, with less representation at the lower levels.

Mapping items to PLDs is not required in the federal peer review elements but provides critical insight into how well the set of items on which students' scores are based reflects the

descriptions of performance and skills within the PLDs. edCount commends the NJDOE on the inclusion of PLD levels within their item and test development processes. For this component of alignment, edCount evaluators examined panelist alignments of items to PLD levels to determine the extent to which the distribution is adequate to support score interpretations for all performance levels. The grade 11 test form shows moderate evidence of alignment in terms of PLD range, with the distribution of items across the PLDs unevenly supporting adequate score interpretations for all performance levels. The grades 5 and 8 test forms show limited evidence of alignment in terms of PLD range; the distribution of items across the PLDs is inadequate to support score interpretations for the lower performance levels for both of these forms. These findings are consistent with panelist comments on the cognitive challenge level of the NJSLA–S forms.

Given the findings for Evaluation Question 3, edCount recommends the NJDOE 1) consider including the intended representation of score point values by reporting category, as well as item numbers, within test development documentation; 2) consider reviewing the performance levels of items on the assessment, to ensure that the forms support score interpretations across all four performance levels; and 3) consider reviewing phenomena that panelists identified as not meeting the highest expectations for accessibility.

For the three NJSLA–S forms reviewed, all forms meet expectations across most evaluation criteria addressing both test development and alignment outcomes. Test development activities follow industry practices, and edCount commends the NJDOE for the inclusion of key stakeholders throughout the process. While the alignment outcomes for the test forms reviewed are overwhelmingly positive, edCount encourages the NJDOE to consider the findings and recommendations in their ongoing improvement efforts.

These findings are notable given the depth and breadth of the methodology used, which exceeds the requirements laid out for state assessments through federal peer review. A test form is the product of a complex, multi-faceted development process, and the levels of alignment for the NJSLA–S forms are the outcome of a clear and standardized test development process. We commend the NJDOE on the development of test forms that meet the majority of the rigorous expectations of this alignment evaluation.

APPENDIX M: EXECUTIVE SUMMARY OF EVALUATION OF THE COGNITIVE PROCESS STUDY

Appendix M contains the executive summary from the cognitive process study report submitted by edCount, LLC in February 2023.

Introduction

The New Jersey Department of Education (NJDOE) commissioned edCount, LLC, (edCount) to conduct an independent evaluation of the degree to which the items on the New Jersey Student Learning Assessment–Science (NJSLA–S) in grades 5, 8, and 11 elicit the intended response processes as represented in the New Jersey Student Learning Standards for Science (NJSLS–S) in 2023. Critical Element 3.2 of the state assessment peer review guidance requires states to provide evidence that their assessments tap the intended cognitive processes appropriate for each grade level as represented in the state’s academic content standards (U.S. Department of Education, 2018). The NJDOE intends to use the results of this evaluation to inform decisions about future item and assessment development and for federal peer review purposes.

To serve this purpose, edCount conducted a cognitive lab study of a strategic sample of grade 5, 8, and 11 NJSLA–S items with a sample of New Jersey students. The goal of the cognitive lab study was to investigate the degree to which:

1. the items elicit the intended construct-relevant response processes appropriate for the grade level;
2. the items include any construct-irrelevant attributes that interfere with students’ demonstration of their knowledge and skills; and
3. the items require complex demonstrations or applications of knowledge and skills.

edCount managed all materials and logistics for the cognitive lab study, provided evaluators to conduct the study, conducted all data analyses, and produced a final written report documenting the study and its findings. edCount worked with the NJDOE to confirm all data collection protocols and cognitive lab methodology, and finalize the sample of items to be evaluated, as well as the characteristics of the sample of students. Additionally, edCount presented the proposal for the study to New Jersey’s Technical Advisory Committee (TAC). Finally, edCount worked with the selected NJDOE-approved districts on the selection of students for inclusion in the study, as well as the timing of the data collection event.

edCount recommends that the NJDOE also use the results of this evaluation to inform decisions about future item and assessment development.

Evaluation Methodology

Evidence of intended and elicited cognitive processes is critical to a validity evaluation for standards-based assessments. Such evidence must draw upon an examination of how the test was designed and developed and how the items elicit the cognitive processes they intend to measure. As in the case for all validity evidence, evidence of cognitive processes is necessary to

support the interpretation and use of test scores, to support inferences of student performance as an accurate reflection of the way the content is assessed.

edCount’s approach to evaluating cognitive processes between the aforementioned NJSLA–S and the NJSL–S encompasses the collection and evaluation of a comprehensive body of evidence. This evidence aligns with the demands of both the federal peer review criteria and the Standards for Educational and Psychological Testing (*The Standards*; AERA, APA, & NCME, 2014).

The methodology for this evaluation was carefully developed based on studies that have utilized cognitive labs to capture a verbal report of a problem solver’s account of his or her own mental processing (Baxter & Glaser, 1998; Ericsson & Simon, 1993) using concurrent and retrospective accounts (Cohen, 1987; Leighton, 2004). Evaluators asked students to describe their thinking concurrently while working through the sample items and had them complete retrospective cognitive interview accounts to clarify the concurrent account and allow additional time for probing questions such as students’ level of familiarity with the topics/phenomena included within the item set.

To address the unique aspects of the multidimensional nature of the NJSLA–S and the NJSL–S, edCount addressed the following questions:

- 1. To what extent do the items on the NJSLA–S elicit the intended cognitive processes as represented in the NJSL–S?**
 - a. To what degree do the items on the NJSLA–S elicit the intended construct-relevant response processes appropriate for the grade level?
 - b. To what degree do the items include any construct-irrelevant attributes that interfere with students’ demonstration of their knowledge and skills?
- 2. How do students interact with the item types within the NJSLA–S?**
 - a. To what degree do students interact with the stimuli and assessment activities as intended?
 - b. Do any aspects of the item interfere with students’ ability to respond, and if so, how?

Evaluation Findings and Recommendations

Evaluation Question 1: Findings and Recommendations

Overall, students engaged with items from the NJSLA–S in the manner intended. Students made valid attempts at the items, fully understood key information most of the time, and grappled in problem-solving in an appropriate manner. The data results in Chapter 3 of the report prepared by edCount detail how students engaged with the items and the associated SEP reporting categories of Critiquing, Investigating, and Sensemaking. When students had difficulty in engaging with items, it was largely during the problem-solving process, which is an appropriate construct-relevant challenge for students to be facing.

While very few items were found to be problematic for students, edCount recommends the NJDOE review the items some students found to be confusing or unclear to determine if there is a more direct manner in which to pose the question.

For items containing multiple stimuli, edCount recommends NJDOE consider scaffolding the stimuli so as not to overwhelm students with information that may not be needed to answer some questions. As the students work through a set of items associated with a phenomenon, additional information can be introduced.

edCount commends the overall construction of the NJSLA–S items, as students were challenged, yet able to apply appropriate reasoning and problem-solving. The vast majority of students made valid attempts on the items (97 percent) and demonstrated full understanding of key information in 89 percent of item attempts. For problem-solving, students executed the problem-solving process without issues in 76 percent of the attempts.

edCount commends the format and diversity of items. Throughout the study, students engaged with the test as intended. In post cognitive lab interviews, no item was universally liked or disliked by all students. In fact, the questions that students most often mentioned tended to be polarizing, both favorite and least favorite, and reflected a matter of personal preference rather than any issues with item construction.

Evaluation Question 2: Findings and Recommendations

Students found the technological interface of the test to be accessible and faced no trouble in moving through all test items. Many students deftly moved between items based on their comfort level with the assigned items. Few cited any difficulties with the test format or the interface of the various item types. In instances where students discussed these item types in the post cognitive lab interviews, it was clear that their comments were a matter of preference rather than confusion.

The analysis on item types executed in Chapter 3 showed that there were some differences between different item types. Students performed more poorly on multi-select selected response items as well as drag-and-drop/order items. While not always the case, some of the challenges noted in these items were related to gaps in content knowledge or academic vocabulary.

Because some item types, such as drag-and-drop and multi-select selected response, resulted in lower response accuracy, edCount recommends a review of these item types across the item bank to determine if this trend is consistent or an artifact of the items selected for this study.

The test administration platform did not interfere with students' ability to respond on the NJSLA–S. Across all students and items, no events occurred, and there were no comments made about how the technological interface of the test was a challenge to navigate.

edCount commends the thoroughness with which the NJSLA–S items were developed. Fewer than 2 percent of item attempts were marred by any difficulty in understanding the item, and this trend was not concentrated in any area (DCI or SEP). Students reported being appropriately challenged by the items, and comments were generally positive.

REFERENCES

- American Educational Research Association, American Psychological Association, National Council on Measurement in Education, & Joint Committee on Standards for Educational and Psychological Testing. (2014). Standards for educational and psychological testing. Washington, DC: American Educational Research Association.
- Baxter, G. P., & Glaser, R. (1998) Investigating the cognitive complexity of science assessments. *Educational Measurement: Issues and Practice*, 17(3), 37–45.
- Cohen, A. (1987). Using verbal reports in research on language learning. In: C. Færch & G. Kasper (Eds.), *Introspection in second language research* (pp. 82–95). Clevedon: Multilingual Matters.
- Ericsson, K. A., & Simon H. A. (1993). *Protocol analysis: verbal reports as data*. Cambridge, MA: MIT Press.
- Leighton, J.P. (2021). Rethinking think-alouds: The often-problematic collection of response process data. *Applied Measurement in Education*, 31(1), 61-74.

APPENDIX N: OBSERVED *P*-VALUES FOR THE FIT AND UNDERFIT SUBGROUPS OF STUDENTS

Table N.1: Grade 5 Fit and Underfit *p*-values

Parent UIN	Domain	Practice	Item Type	ACC		EcoDisad		EL		SWD	
				P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
2205B009_05	Earth and Space	Critiquing	CR	0.179	0.074	0.169	0.073	0.077	0.036	0.170	0.065
1905M008_01	Earth and Space	Critiquing	MC	0.554	0.167	0.571	0.170	0.427	0.134	0.553	0.168
1905M008_06	Earth and Space	Critiquing	MC	0.581	0.168	0.600	0.178	0.489	0.177	0.576	0.154
2205M011_04	Earth and Space	Critiquing	MC	0.383	0.248	0.361	0.257	0.324	0.226	0.388	0.276
2105M015_04	Earth and Space	Critiquing	TE	0.247	0.214	0.237	0.201	0.181	0.183	0.259	0.212
1905M005_01	Earth and Space	Critiquing	TE	0.358	0.134	0.366	0.126	0.249	0.096	0.361	0.130
1905M005_03	Earth and Space	Critiquing	TE	0.255	0.082	0.233	0.073	0.103	0.049	0.283	0.089
2205M011_02	Earth and Space	Critiquing	TE	0.271	0.118	0.279	0.120	0.164	0.099	0.285	0.124
2105M015_06	Earth and Space	Investigating	MC	0.646	0.268	0.644	0.267	0.562	0.254	0.641	0.249
1905M005_04	Earth and Space	Investigating	MC	0.542	0.325	0.552	0.338	0.506	0.311	0.532	0.311
2205M011_01	Earth and Space	Investigating	TE	0.184	0.322	0.175	0.331	0.149	0.325	0.190	0.326
2205B009_02	Earth and Space	Investigating	TE	0.191	0.197	0.178	0.188	0.116	0.187	0.186	0.185
2205B009_03	Earth and Space	Investigating	TE	0.195	0.273	0.189	0.264	0.164	0.290	0.201	0.244
1905M008_05	Earth and Space	Sensemaking	TE	0.408	0.079	0.425	0.090	0.265	0.050	0.413	0.085
2105M015_05	Earth and Space	Sensemaking	TE	0.391	0.340	0.393	0.353	0.384	0.315	0.391	0.345
2205B003_05	Earth and Space	Sensemaking	TE	0.480	0.174	0.461	0.167	0.280	0.116	0.525	0.185
2205B009_01	Earth and Space	Sensemaking	TE	0.464	0.171	0.455	0.169	0.310	0.110	0.490	0.174
2205B009_04	Earth and Space	Sensemaking	TE	0.342	0.158	0.371	0.170	0.247	0.138	0.341	0.147
1905B007_08	Life	Critiquing	CR	0.296	0.060	0.303	0.068	0.195	0.037	0.303	0.057
1905M044_02	Life	Critiquing	MC	0.414	0.137	0.401	0.131	0.306	0.131	0.387	0.125
2205M006_02	Life	Critiquing	TE	0.347	0.189	0.347	0.208	0.296	0.173	0.355	0.182
2205M004_05	Life	Critiquing	TE	0.211	0.187	0.201	0.187	0.151	0.173	0.222	0.196
1905M044_03	Life	Critiquing	TE	0.230	0.138	0.227	0.145	0.161	0.121	0.223	0.135
2205M006_01	Life	Investigating	MC	0.433	0.188	0.431	0.207	0.328	0.168	0.442	0.193
1905B007_03	Life	Investigating	TE	0.210	0.108	0.202	0.101	0.120	0.083	0.222	0.114

Parent UIN	Domain	Practice	Item Type	ACC		EcoDisad		EL		SWD	
				P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
1905B009_01	Life	Investigating	TE	0.300	0.045	0.313	0.046	0.159	0.026	0.299	0.041
1905B009_05	Life	Investigating	TE	0.164	0.250	0.162	0.246	0.119	0.246	0.153	0.237
1905B007_01	Life	Sensemaking	TE	0.210	0.050	0.208	0.050	0.112	0.040	0.211	0.045
1905B007_05	Life	Sensemaking	TE	0.288	0.091	0.273	0.087	0.155	0.060	0.303	0.093
1905B007_10	Life	Sensemaking	TE	0.296	0.062	0.278	0.051	0.167	0.037	0.316	0.063
2205M006_05	Life	Sensemaking	TE	0.396	0.121	0.409	0.129	0.284	0.098	0.398	0.110
2205M004_03	Life	Sensemaking	TE	0.329	0.091	0.332	0.094	0.197	0.070	0.322	0.088
2205M004_07	Life	Sensemaking	TE	0.403	0.202	0.402	0.214	0.348	0.228	0.395	0.196
1905B009_02	Life	Sensemaking	TE	0.502	0.080	0.537	0.100	0.344	0.059	0.500	0.079
1905M044_04	Life	Sensemaking	TE	0.247	0.167	0.250	0.156	0.169	0.164	0.250	0.162
2205B003_04	Physical	Critiquing	CR	0.111	0.074	0.103	0.089	0.053	0.052	0.103	0.071
2205M012_03	Physical	Critiquing	MC	0.293	0.142	0.293	0.140	0.218	0.124	0.299	0.145
1905M076_01	Physical	Critiquing	MC	0.549	0.323	0.560	0.339	0.500	0.318	0.541	0.318
2205M012_04	Physical	Critiquing	TE	0.185	0.368	0.186	0.373	0.167	0.352	0.185	0.358
2205B003_03	Physical	Critiquing	TE	0.215	0.198	0.236	0.216	0.170	0.204	0.216	0.196
2205M012_01	Physical	Investigating	MC	0.520	0.170	0.530	0.177	0.415	0.160	0.513	0.162
2205B003_01	Physical	Investigating	MC	0.372	0.345	0.385	0.344	0.346	0.309	0.381	0.346
2205M022_01	Physical	Investigating	MC	0.369	0.180	0.354	0.174	0.273	0.167	0.344	0.161
2205M022_03	Physical	Investigating	MC	0.349	0.151	0.351	0.150	0.251	0.144	0.344	0.131
1905M040_01	Physical	Investigating	TE	0.261	0.134	0.253	0.132	0.178	0.123	0.270	0.127
2205B003_02	Physical	Investigating	TE	0.271	0.187	0.270	0.181	0.205	0.167	0.271	0.175
2205M022_05	Physical	Investigating	TE	0.376	0.252	0.392	0.235	0.340	0.237	0.364	0.233
1905M040_03	Physical	Sensemaking	MC	0.469	0.229	0.469	0.222	0.389	0.225	0.461	0.222
1905M040_05	Physical	Sensemaking	TE	0.271	0.121	0.264	0.135	0.192	0.116	0.269	0.122
1905M076_03	Physical	Sensemaking	TE	0.310	0.092	0.301	0.091	0.177	0.073	0.316	0.107
1905M076_05	Physical	Sensemaking	TE	0.192	0.112	0.173	0.116	0.106	0.107	0.198	0.118

Note. **P_FT** = *p-values* for the group of students not showing underfit; **P_UF** = *p-values* for the group of students showing underfit.

Table N.2: Grade 8 Fit and Underfit *p*-values

Parent UIN	Domain	Practice	Item Type	ACC		EcoDisad		EL		SWD	
				P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
1908M026_06	Earth and Space	Critiquing	MC	0.353	0.211	0.358	0.198	0.290	0.181	0.360	0.231
1908M033_03	Earth and Space	Critiquing	TE	0.391	0.156	0.417	0.187	0.323	0.139	0.380	0.155
1908M026_01	Earth and Space	Critiquing	TE	0.087	0.067	0.089	0.075	0.028	0.051	0.105	0.075
2108M015_02	Earth and Space	Critiquing	TE	0.373	0.126	0.383	0.145	0.309	0.138	0.383	0.123
2208M021_03	Earth and Space	Critiquing	TE	0.197	0.144	0.208	0.134	0.156	0.122	0.201	0.134
2208M021_10	Earth and Space	Critiquing	TE	0.217	0.047	0.247	0.049	0.111	0.043	0.232	0.047
2108B006_09	Earth and Space	Critiquing	TE	0.217	0.064	0.207	0.063	0.075	0.039	0.249	0.072
2108B006_11	Earth and Space	Investigating	CR	0.049	0.032	0.050	0.039	0.016	0.016	0.053	0.036
2108M015_09	Earth and Space	Investigating	MC	0.440	0.132	0.433	0.119	0.356	0.115	0.424	0.111
1908M033_02	Earth and Space	Investigating	TE	0.302	0.162	0.313	0.173	0.245	0.180	0.294	0.155
2208M028_06	Earth and Space	Investigating	TE	0.360	0.225	0.362	0.208	0.317	0.193	0.374	0.228
2108B006_01	Earth and Space	Investigating	TE	0.195	0.041	0.224	0.047	0.096	0.036	0.193	0.034
2108B006_03	Earth and Space	Investigating	TE	0.109	0.187	0.115	0.177	0.070	0.156	0.120	0.193
1908M033_04	Earth and Space	Sensemaking	MC	0.441	0.210	0.475	0.228	0.332	0.175	0.453	0.223
1908M026_04	Earth and Space	Sensemaking	MC	0.305	0.134	0.329	0.134	0.253	0.142	0.308	0.134
2208M021_09	Earth and Space	Sensemaking	MC	0.244	0.166	0.234	0.167	0.219	0.168	0.250	0.162
2108M015_06	Earth and Space	Sensemaking	TE	0.336	0.077	0.323	0.099	0.201	0.084	0.359	0.086
2108M015_10	Earth and Space	Sensemaking	TE	0.235	0.166	0.276	0.174	0.198	0.133	0.273	0.165
2208M021_05	Earth and Space	Sensemaking	TE	0.290	0.147	0.285	0.144	0.230	0.141	0.306	0.151
2108B006_06	Earth and Space	Sensemaking	TE	0.142	0.119	0.126	0.110	0.088	0.119	0.141	0.112
2208B003_11	Life	Critiquing	CR	0.313	0.081	0.316	0.082	0.241	0.086	0.303	0.068
1908M030_01	Life	Critiquing	MC	0.331	0.196	0.323	0.179	0.296	0.176	0.307	0.183
2208M028_09	Life	Critiquing	MC	0.359	0.227	0.354	0.243	0.338	0.216	0.379	0.241
2008M015_09	Life	Critiquing	MC	0.274	0.168	0.275	0.145	0.229	0.154	0.285	0.156

Parent UIN	Domain	Practice	Item Type	ACC		EcoDisad		EL		SWD	
				P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
2208M028_07	Life	Critiquing	TE	0.125	0.220	0.123	0.217	0.097	0.235	0.127	0.216
2008M000_02	Life	Critiquing	TE	0.098	0.067	0.113	0.075	0.035	0.049	0.123	0.068
2208B003_09	Life	Critiquing	TE	0.236	0.192	0.228	0.200	0.212	0.224	0.226	0.175
2108M027_03	Life	Critiquing	TE	0.370	0.095	0.384	0.109	0.269	0.095	0.347	0.087
1908M030_05	Life	Investigating	MC	0.267	0.155	0.272	0.150	0.192	0.166	0.283	0.153
2208M028_02	Life	Investigating	MC	0.180	0.221	0.181	0.229	0.130	0.235	0.204	0.222
2008M000_04	Life	Investigating	MC	0.277	0.191	0.281	0.185	0.271	0.227	0.264	0.178
2208B003_07	Life	Investigating	TE	0.101	0.080	0.104	0.090	0.051	0.060	0.111	0.087
2008M015_04	Life	Investigating	TE	0.116	0.086	0.133	0.102	0.076	0.089	0.111	0.076
2008M015_05	Life	Investigating	TE	0.178	0.177	0.210	0.181	0.142	0.168	0.180	0.158
2108M027_01	Life	Investigating	TE	0.156	0.094	0.171	0.098	0.089	0.094	0.167	0.086
1908M003_08	Life	Sensemaking	MC	0.230	0.138	0.226	0.129	0.180	0.128	0.255	0.141
2108M027_07	Life	Sensemaking	MC	0.349	0.188	0.357	0.203	0.312	0.217	0.344	0.190
1908M030_02	Life	Sensemaking	TE	0.289	0.173	0.308	0.180	0.219	0.187	0.309	0.176
2008M000_01	Life	Sensemaking	TE	0.224	0.245	0.238	0.242	0.218	0.223	0.230	0.237
2208B003_01	Life	Sensemaking	TE	0.300	0.124	0.311	0.119	0.196	0.104	0.323	0.123
2208B003_05	Life	Sensemaking	TE	0.296	0.108	0.305	0.111	0.205	0.107	0.309	0.099
1908M003_02	Life	Sensemaking	TE	0.086	0.168	0.098	0.198	0.051	0.166	0.100	0.178
1908M003_07	Life	Sensemaking	TE	0.098	0.281	0.103	0.264	0.072	0.262	0.108	0.279
2008M015_06	Life	Sensemaking	TE	0.324	0.104	0.370	0.119	0.250	0.097	0.315	0.094
2008M001_05	Physical	Critiquing	MC	0.507	0.241	0.519	0.231	0.447	0.232	0.498	0.244
2108B003_07	Physical	Critiquing	TE	0.206	0.193	0.192	0.179	0.161	0.179	0.223	0.195
1908B000_08	Physical	Critiquing	TE	0.167	0.132	0.165	0.142	0.134	0.130	0.176	0.135
2208M051_16	Physical	Critiquing	TE	0.340	0.103	0.337	0.109	0.209	0.074	0.381	0.106
2208M051_17	Physical	Critiquing	TE	0.216	0.115	0.232	0.115	0.153	0.111	0.228	0.106

Parent UIN	Domain	Practice	Item Type	ACC		EcoDisad		EL		SWD	
				P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
2108B007_01	Physical	Critiquing	TE	0.231	0.067	0.242	0.082	0.164	0.064	0.224	0.069
2108B007_08	Physical	Critiquing	TE	0.493	0.172	0.516	0.180	0.359	0.116	0.485	0.154
1908M005_05	Physical	Investigating	MC	0.301	0.109	0.315	0.129	0.212	0.107	0.300	0.113
2008M001_08	Physical	Investigating	MC	0.260	0.151	0.253	0.150	0.214	0.140	0.275	0.154
2108B007_03	Physical	Investigating	MC	0.183	0.243	0.188	0.260	0.183	0.257	0.170	0.237
2108B003_08	Physical	Investigating	TE	0.357	0.189	0.380	0.198	0.259	0.153	0.387	0.200
1908M005_02	Physical	Investigating	TE	0.147	0.185	0.153	0.186	0.104	0.192	0.144	0.178
1908M005_03	Physical	Investigating	TE	0.049	0.085	0.055	0.103	0.024	0.076	0.064	0.075
2208M051_13	Physical	Investigating	TE	0.276	0.153	0.293	0.151	0.218	0.140	0.290	0.139
1908B000_11	Physical	Sensemaking	CR	0.207	0.029	0.214	0.031	0.121	0.017	0.229	0.029
2108B003_02	Physical	Sensemaking	MC	0.514	0.230	0.540	0.227	0.414	0.223	0.534	0.215
1908B000_03	Physical	Sensemaking	TE	0.126	0.145	0.121	0.137	0.090	0.140	0.134	0.150
1908B000_04	Physical	Sensemaking	TE	0.181	0.098	0.184	0.105	0.124	0.086	0.200	0.110
1908B000_12	Physical	Sensemaking	TE	0.198	0.145	0.193	0.159	0.154	0.163	0.227	0.164
2008M001_01	Physical	Sensemaking	TE	0.158	0.029	0.171	0.034	0.072	0.024	0.180	0.032
2108B007_07	Physical	Sensemaking	TE	0.162	0.242	0.159	0.223	0.124	0.241	0.171	0.241

Note. **P_FT** = *p-values* for the group of students not showing underfit; **P_UF** = *p-values* for the group of students showing underfit.

Table N.3: Grade 11 Fit and Underfit *p*-values

Parent UIN	Domain	Practice	Item Type	ACC		EcoDisad		EL		SWD	
				P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
2211B006_09	Earth and Space	Critiquing	CR	0.452	0.103	0.474	0.098	0.321	0.088	0.425	0.089
2111M000_09	Earth and Space	Critiquing	MC	0.416	0.193	0.404	0.181	0.355	0.224	0.419	0.189
1911M002_04	Earth and Space	Critiquing	MC	0.228	0.288	0.229	0.253	0.182	0.264	0.250	0.271
1911M119_02	Earth and Space	Critiquing	MC	0.294	0.271	0.291	0.278	0.237	0.283	0.306	0.277
2211M008_03	Earth and Space	Critiquing	MC	0.407	0.216	0.399	0.207	0.394	0.254	0.374	0.186
1911M079_04	Earth and Space	Critiquing	TE	0.135	0.100	0.146	0.111	0.063	0.077	0.150	0.122
2111M000_06	Earth and Space	Investigating	MC	0.434	0.203	0.441	0.203	0.364	0.198	0.436	0.199
1911M002_05	Earth and Space	Investigating	MC	0.416	0.238	0.413	0.243	0.378	0.255	0.398	0.243
1911M119_05	Earth and Space	Investigating	TE	0.352	0.139	0.366	0.147	0.261	0.114	0.361	0.147
2211M008_07	Earth and Space	Investigating	TE	0.273	0.211	0.263	0.182	0.223	0.211	0.259	0.174
2211B006_02	Earth and Space	Investigating	TE	0.240	0.129	0.228	0.127	0.159	0.133	0.257	0.119
2211B006_05	Earth and Space	Investigating	TE	0.208	0.159	0.204	0.156	0.122	0.158	0.220	0.173
2211B006_06	Earth and Space	Investigating	TE	0.633	0.192	0.653	0.215	0.514	0.185	0.632	0.227
1911M119_06	Earth and Space	Sensemaking	MC	0.305	0.163	0.301	0.163	0.182	0.124	0.323	0.189
1911M079_02	Earth and Space	Sensemaking	MC	0.337	0.260	0.340	0.248	0.256	0.203	0.345	0.259
2111M000_02	Earth and Space	Sensemaking	TE	0.496	0.170	0.540	0.192	0.288	0.122	0.515	0.172
2111M000_08	Earth and Space	Sensemaking	TE	0.543	0.189	0.587	0.200	0.370	0.124	0.566	0.192
1911M002_01	Earth and Space	Sensemaking	TE	0.131	0.171	0.140	0.169	0.069	0.163	0.147	0.177
1911M079_03	Earth and Space	Sensemaking	TE	0.371	0.126	0.401	0.114	0.216	0.071	0.368	0.109
2211M008_01	Earth and Space	Sensemaking	TE	0.141	0.145	0.146	0.131	0.086	0.094	0.145	0.141
1911M023_06	Life	Critiquing	MC	0.343	0.216	0.357	0.216	0.296	0.226	0.338	0.190
2211M003_02	Life	Critiquing	MC	0.326	0.335	0.330	0.343	0.274	0.305	0.351	0.355
2011M003_03	Life	Critiquing	MC	0.213	0.225	0.210	0.202	0.160	0.183	0.221	0.216
1911B009_01A	Life	Critiquing	TE	0.449	0.156	0.451	0.154	0.283	0.115	0.447	0.179

Parent UIN	Domain	Practice	Item Type	ACC		EcoDisad		EL		SWD	
				P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
1911B009_05A	Life	Critiquing	TE	0.145	0.096	0.150	0.103	0.063	0.086	0.151	0.099
2211M003_03	Life	Critiquing	TE	0.231	0.149	0.239	0.134	0.133	0.112	0.236	0.149
2011M003_05	Life	Critiquing	TE	0.277	0.125	0.277	0.114	0.169	0.097	0.280	0.128
2111M004_06	Life	Critiquing	TE	0.180	0.171	0.193	0.208	0.119	0.153	0.178	0.189
1911M023_07	Life	Investigating	MC	0.361	0.248	0.361	0.218	0.342	0.231	0.347	0.241
2211M003_01	Life	Investigating	MC	0.350	0.192	0.372	0.173	0.301	0.185	0.342	0.169
2011M003_06	Life	Investigating	MC	0.391	0.288	0.398	0.252	0.247	0.242	0.409	0.285
2011M003_04	Life	Investigating	MC	0.258	0.226	0.248	0.211	0.162	0.181	0.259	0.227
2111M004_03	Life	Investigating	MC	0.186	0.118	0.192	0.109	0.110	0.097	0.182	0.125
1911B009_09A	Life	Investigating	TE	0.138	0.181	0.128	0.164	0.072	0.119	0.161	0.200
2211M003_05	Life	Investigating	TE	0.131	0.123	0.136	0.125	0.063	0.089	0.149	0.118
1911B009_07A	Life	Sensemaking	CR	0.232	0.058	0.249	0.065	0.103	0.019	0.250	0.076
1911M023_02	Life	Sensemaking	MC	0.276	0.181	0.282	0.197	0.145	0.142	0.291	0.212
2111M004_05	Life	Sensemaking	MC	0.499	0.293	0.501	0.290	0.464	0.298	0.485	0.265
1911B009_03A	Life	Sensemaking	TE	0.231	0.064	0.235	0.068	0.100	0.049	0.237	0.084
1911M023_05	Life	Sensemaking	TE	0.491	0.148	0.523	0.146	0.269	0.086	0.516	0.138
2211M003_04	Life	Sensemaking	TE	0.065	0.190	0.060	0.184	0.018	0.158	0.066	0.168
2011M003_01	Life	Sensemaking	TE	0.202	0.122	0.204	0.118	0.101	0.097	0.215	0.123
2111M004_02	Life	Sensemaking	TE	0.083	0.269	0.085	0.253	0.044	0.257	0.088	0.242
HS18060_01	Physical	Critiquing	TE	0.142	0.076	0.138	0.060	0.062	0.046	0.148	0.062
1911M124_10	Physical	Critiquing	TE	0.389	0.192	0.390	0.185	0.249	0.152	0.415	0.190
2211B000_07	Physical	Critiquing	TE	0.194	0.161	0.194	0.153	0.108	0.157	0.207	0.179
2011M071_02	Physical	Critiquing	TE	0.187	0.165	0.184	0.178	0.130	0.165	0.180	0.174
2011M071_03	Physical	Critiquing	TE	0.272	0.153	0.278	0.141	0.151	0.112	0.269	0.154
2211B006_12	Physical	Critiquing	TE	0.260	0.107	0.270	0.089	0.144	0.079	0.259	0.110

Parent UIN	Domain	Practice	Item Type	ACC		EcoDisad		EL		SWD	
				P_FT	P_UF	P_FT	P_UF	P_FT	P_UF	P_FT	P_UF
2011M010_02	Physical	Critiquing	TE	0.284	0.186	0.299	0.198	0.201	0.137	0.309	0.190
2011M010_05	Physical	Critiquing	TE	0.176	0.274	0.178	0.278	0.146	0.270	0.195	0.293
2211B000_12	Physical	Investigating	CR	0.121	0.039	0.120	0.034	0.059	0.023	0.117	0.036
1911M028_01	Physical	Investigating	MC	0.530	0.266	0.522	0.230	0.471	0.236	0.516	0.244
1911M028_03	Physical	Investigating	MC	0.205	0.186	0.165	0.168	0.178	0.206	0.171	0.178
1911M028_04	Physical	Investigating	MC	0.391	0.254	0.394	0.250	0.352	0.244	0.404	0.238
1911M028_06	Physical	Investigating	MC	0.333	0.258	0.320	0.257	0.323	0.264	0.313	0.250
2011M071_05	Physical	Investigating	MC	0.483	0.231	0.485	0.187	0.378	0.171	0.450	0.207
2011M010_01	Physical	Investigating	MC	0.258	0.228	0.264	0.230	0.244	0.237	0.246	0.236
1911M124_01	Physical	Investigating	TE	0.220	0.219	0.232	0.201	0.160	0.168	0.231	0.207
2211B000_01	Physical	Investigating	TE	0.283	0.105	0.278	0.104	0.165	0.081	0.273	0.097
HS18060_03	Physical	Sensemaking	MC	0.420	0.203	0.438	0.232	0.304	0.185	0.425	0.240
HS18060_04	Physical	Sensemaking	MC	0.462	0.235	0.471	0.228	0.350	0.204	0.455	0.232
1911M124_05	Physical	Sensemaking	MC	0.547	0.323	0.548	0.347	0.513	0.334	0.534	0.336
2011M071_04	Physical	Sensemaking	MC	0.510	0.292	0.512	0.290	0.487	0.301	0.488	0.274
2011M071_01	Physical	Sensemaking	MC	0.507	0.201	0.519	0.189	0.398	0.157	0.501	0.197
2011M010_03	Physical	Sensemaking	MC	0.351	0.292	0.352	0.278	0.323	0.285	0.363	0.286
HS18060_06	Physical	Sensemaking	TE	0.101	0.174	0.101	0.188	0.050	0.158	0.101	0.186
1911M124_02	Physical	Sensemaking	TE	0.208	0.096	0.220	0.111	0.093	0.068	0.227	0.123
2211B000_03	Physical	Sensemaking	TE	0.335	0.253	0.320	0.221	0.284	0.217	0.343	0.229

Note. **P_FT** = *p-values* for the group of students not showing underfit; **P_UF** = *p-values* for the group of students showing underfit

APPENDIX O: PARAMETER ESTIMATES FROM THE CONFIRMATORY FACTOR ANALYSES FOR THE 2023 NJSLA–S TESTS

Table O.1: Grade 5 Domain Model Parameter Estimates

Earth Space Science	Estimate (standardized)	R-Squared	Standard Error	z
1905M008_01	0.724	0.525	0.003	280.960
1905M008_05	0.738	0.545	0.003	290.327
1905M008_06	0.684	0.468	0.003	240.680
2105M015_06	0.502	0.252	0.004	136.144
2105M015_05	0.116	0.013	0.004	26.8990
2105M015_04	0.472	0.223	0.004	121.544
1905M005_01	0.619	0.383	0.003	198.879
1905M005_03	0.726	0.527	0.003	266.165
1905M005_04	0.333	0.111	0.004	84.6760
2205B003_05	0.659	0.434	0.003	225.515
2205M011_04	0.485	0.235	0.004	134.361
2205M011_02	0.713	0.508	0.003	252.868
2205M011_01	0.364	0.133	0.005	80.3304
2205B009_01	0.667	0.445	0.003	235.220
2205B009_02	0.578	0.335	0.004	155.195
2205B009_03	0.313	0.098	0.004	70.6840
2205B009_04	0.544	0.296	0.003	160.471
2205B009_05	0.733	0.538	0.002	336.396
Life Science	Estimate (standardized)	R-Squared	Standard Error	z
1905B007_01	0.810	0.655	0.002	349.585
1905B007_03	0.781	0.610	0.003	301.870
1905B007_05	0.759	0.576	0.003	302.071
1905B007_08	0.717	0.514	0.002	360.551
1905B007_10	0.837	0.700	0.002	408.575
2205M006_01	0.518	0.268	0.003	153.024
2205M006_02	0.340	0.116	0.004	86.502
2205M006_05	0.584	0.341	0.003	181.298
2205M004_03	0.778	0.605	0.002	331.180
2205M004_05	0.535	0.286	0.004	143.259
2205M004_07	0.418	0.175	0.004	110.983
1905B009_01	0.799	0.639	0.002	357.005
1905B009_02	0.872	0.760	0.002	475.842
1905B009_05	0.367	0.135	0.004	83.307
1905M044_02	0.711	0.506	0.003	268.190
1905M044_03	0.602	0.363	0.003	178.558
1905M044_04	0.525	0.276	0.004	144.707

Physical Science	Estimate (standardized)	R-Squared	Standard Error	z
1905M040_01	0.616	0.380	0.003	190.444
1905M040_03	0.577	0.333	0.003	179.148
1905M040_05	0.523	0.274	0.004	147.132
2205M012_01	0.661	0.436	0.003	224.461
2205M012_03	0.598	0.358	0.003	181.922
2205M012_04	0.076	0.006	0.005	16.033
2205B003_01	0.179	0.032	0.004	42.191
2205B003_02	0.548	0.300	0.004	155.338
2205B003_03	0.350	0.122	0.004	84.022
2205B003_04	0.691	0.478	0.003	272.448
2205M022_01	0.569	0.324	0.003	171.595
2205M022_03	0.610	0.373	0.003	194.492
2205M022_05	0.338	0.114	0.004	85.539
1905M076_01	0.298	0.089	0.004	75.711
1905M076_03	0.705	0.497	0.003	252.510
1905M076_05	0.691	0.478	0.003	219.227

Note. All parameter estimates were significant at $p < .001$; R-squared is the squared standardized estimate and is interpreted as the proportion of variance in the item explained by the latent subscore group (Kline, 2011).

Table O.2: Grade 8 Domain Model Parameter Estimates

Earth Space Science	Estimate (standardized)	R-Squared	Standard Error	z
1908M033_02	0.435	0.189	0.004	116.435
1908M033_03	0.382	0.146	0.004	100.991
1908M033_04	0.477	0.228	0.003	138.596
1908M026_01	0.627	0.393	0.004	161.802
1908M026_04	0.443	0.196	0.004	120.141
1908M026_06	0.393	0.154	0.004	103.104
2208M028_06	0.329	0.108	0.004	83.593
2108M015_06	0.650	0.423	0.003	223.841
2108M015_09	0.652	0.426	0.003	222.697
2108M015_10	0.287	0.082	0.004	68.763
2108M015_02	0.550	0.303	0.003	171.560
2208M021_03	0.491	0.242	0.004	127.687
2208M021_05	0.491	0.241	0.004	134.589
2208M021_09	0.442	0.196	0.004	112.478
2208M021_10	0.743	0.552	0.003	285.193
2108B006_01	0.764	0.584	0.003	296.508
2108B006_03	0.491	0.241	0.004	109.766
2108B006_06	0.536	0.287	0.004	133.312
2108B006_09	0.701	0.491	0.003	243.530
2108B006_11	0.706	0.498	0.003	240.758
Life Science	Estimate (standardized)	R-Squared	Standard Error	z
1908M030_01	0.499	0.249	0.004	138.545
1908M030_02	0.530	0.281	0.003	151.711
1908M030_05	0.520	0.271	0.004	143.633
2208M028_02	0.163	0.027	0.005	34.705
2208M028_07	0.368	0.135	0.005	76.997
2208M028_09	0.398	0.158	0.004	103.650
2008M000_02	0.658	0.433	0.004	183.757
2008M000_04	0.340	0.115	0.004	82.986
2008M000_01	0.270	0.073	0.005	59.943
2208B003_01	0.635	0.403	0.003	211.358
2208B003_05	0.615	0.378	0.003	195.354
2208B003_07	0.678	0.460	0.004	186.370
2208B003_09	0.297	0.088	0.004	68.478
2208B003_11	0.673	0.452	0.002	304.369
1908M003_02	0.372	0.139	0.005	76.257
1908M003_07	0.445	0.198	0.005	91.829
1908M003_08	0.561	0.314	0.004	159.331

2008M015_04	0.463	0.214	0.004	109.834
2008M015_05	0.287	0.082	0.004	65.560
2008M015_06	0.638	0.407	0.003	208.105
2008M015_09	0.473	0.223	0.004	126.428
2108M027_07	0.364	0.132	0.004	93.636
2108M027_01	0.632	0.399	0.003	183.158
2108M027_03	0.649	0.422	0.003	218.901

Physical Science	Estimate (standardized)	R-Squared	Standard Error	z
2108B003_02	0.524	0.275	0.003	162.369
2108B003_08	0.348	0.121	0.004	89.924
2108B003_07	0.452	0.204	0.004	110.297
1908B000_03	0.437	0.191	0.005	95.233
1908B000_04	0.543	0.295	0.004	146.188
1908B000_08	0.559	0.312	0.004	143.319
1908B000_11	0.771	0.595	0.002	431.725
1908B000_12	0.515	0.266	0.004	134.023
1908M005_02	0.532	0.283	0.004	133.031
1908M005_03	0.676	0.457	0.004	155.138
1908M005_05	0.507	0.257	0.004	143.554
2008M001_05	0.535	0.286	0.003	168.818
2008M001_01	0.766	0.587	0.003	295.965
2008M001_08	0.459	0.211	0.004	119.642
2208M051_13	0.490	0.240	0.004	135.586
2208M051_16	0.573	0.329	0.003	180.145
2208M051_17	0.563	0.317	0.003	161.488
2108B007_01	0.640	0.410	0.003	204.429
2108B007_08	0.537	0.288	0.003	157.762
2108B007_03	0.169	0.029	0.005	35.636
2108B007_07	0.372	0.139	0.005	80.845

Note. All parameter estimates were significant at $p < .001$; R-squared is the squared standardized estimate and is interpreted as the proportion of variance in the item explained by the latent subscore group (Kline, 2011).

Table O.3: Grade 11 Domain Model Parameter Estimates

Earth Space Science	Estimate (standardized)	R-Squared	Standard Error	z
2111M000_06	0.440	0.193	0.004	118.810
2111M000_02	0.642	0.412	0.003	214.218
2111M000_09	0.577	0.333	0.003	179.757
2111M000_08	0.596	0.355	0.003	185.437
1911M002_01	0.458	0.210	0.004	103.654
1911M002_04	0.335	0.112	0.004	74.815
1911M002_05	0.356	0.126	0.004	90.458
1911M119_06	0.671	0.450	0.003	223.175
1911M119_02	0.361	0.130	0.004	86.318
1911M119_05	0.586	0.343	0.003	175.468
1911M079_02	0.296	0.088	0.004	71.518
1911M079_03	0.613	0.376	0.003	192.913
1911M079_04	0.640	0.410	0.004	168.825
2211M008_01	0.524	0.274	0.004	129.244
2211M008_03	0.364	0.132	0.004	91.943
2211M008_07	0.328	0.108	0.004	77.210
2211B006_02	0.577	0.333	0.004	159.979
2211B006_05	0.622	0.387	0.004	172.225
2211B006_06	0.646	0.418	0.003	215.614
2211B006_09	0.678	0.460	0.002	309.365
Life Science	Estimate (standardized)	R-Squared	Standard Error	z
1911B009_01A	0.687	0.472	0.003	253.523
1911B009_03A	0.806	0.650	0.002	329.176
1911B009_05A	0.633	0.400	0.004	173.082
1911B009_07A	0.638	0.407	0.002	275.169
1911B009_09A	0.513	0.263	0.005	112.466
1911M023_02	0.666	0.443	0.003	217.491
1911M023_05	0.694	0.482	0.003	248.230
1911M023_06	0.352	0.124	0.004	87.764
1911M023_07	0.284	0.081	0.004	68.173
2211M003_01	0.465	0.216	0.004	124.552
2211M003_03	0.468	0.219	0.004	119.172
2211M003_04	0.572	0.327	0.005	115.560
2211M003_02	0.291	0.085	0.004	68.385
2211M003_05	0.644	0.415	0.004	169.881
2011M003_06	0.523	0.274	0.003	151.310
2011M003_01	0.657	0.432	0.003	193.732
2011M003_03	0.403	0.162	0.004	92.935

2011M003_04	0.527	0.278	0.004	138.844
2011M003_05	0.699	0.488	0.003	237.381
2111M004_05	0.368	0.135	0.004	95.651
2111M004_02	0.417	0.174	0.005	79.571
2111M004_03	0.618	0.382	0.004	169.605
2111M004_06	0.434	0.189	0.004	98.416
Physical Science	Estimate (standardized)	R-Squared	Standard Error	z
1911M028_01	0.349	0.122	0.004	91.325
1911M028_03	0.601	0.361	0.004	158.105
1911M028_04	0.279	0.078	0.004	68.087
1911M028_06	0.332	0.110	0.004	80.447
HS18060_01	0.756	0.572	0.003	245.410
HS18060_03	0.587	0.344	0.003	181.617
HS18060_04	0.659	0.434	0.003	234.292
HS18060_06	0.529	0.280	0.005	114.533
1911M124_02	0.617	0.381	0.003	179.744
1911M124_05	0.284	0.080	0.004	72.038
1911M124_10	0.621	0.385	0.003	201.757
1911M124_01	0.284	0.081	0.004	64.201
2211B000_01	0.542	0.294	0.004	151.947
2211B000_07	0.515	0.265	0.004	131.015
2211B000_03	0.298	0.089	0.004	71.049
2211B000_12	0.727	0.528	0.002	302.092
2011M071_04	0.378	0.143	0.004	99.830
2011M071_01	0.666	0.444	0.003	234.901
2011M071_02	0.425	0.180	0.004	97.043
2011M071_03	0.635	0.403	0.003	193.968
2011M071_05	0.561	0.314	0.003	164.288
2211B006_12	0.732	0.536	0.003	257.220
2011M010_01	0.259	0.067	0.004	57.557
2011M010_02	0.447	0.200	0.004	115.284
2011M010_03	0.357	0.128	0.004	88.614
2011M010_05	0.402	0.162	0.005	86.636

Note. All parameter estimates were significant at $p < .001$; R-squared is the squared standardized estimate and is interpreted as the proportion of variance in the item explained by the latent subscore group (Kline, 2011).