

New Jersey

Student Learning Assessment—Science

(NJSLA–S)

**TECHNICAL REPORT**  
**Grades 5, 8, and 11**  
**2022**

**March 2023**  
PTM XXXX.XX



**State of New Jersey**  
**Department of Education**

Copyright © 2023 by New Jersey Department of Education  
All rights reserved

## State Board of Education

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

Dr. Angelica Allen-McMillan, 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 2022 NJSLA–S Test Construction.....	20
2.4.1 Grade 5 Test Construction .....	22
2.4.2 Grade 8 Test Construction .....	24
2.4.3 Grade 11 Test Construction .....	26
2.5 2022 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 .....	34
3.4.2 Accommodations.....	35
PART 4: SCORING.....	38
4.1 Machine-Scored Items.....	38
4.1.1 Adjudication .....	38

## Table of Contents

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 #1 .....	44
4.3.2 QC #2 .....	45
4.3.3 QC #3 .....	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.....	75
6.2.4 Item Characteristic Curves – CR Items .....	76
6.3 Student Test Performance .....	82
6.3.1 Scale Score Distribution by Form.....	82
6.3.2 Scale Score Distributions by Demographic group.....	83
6.3.3 Subscore Proficiency Classification .....	83
PART 7: EQUATING AND SCALING.....	84
7.1 Summary of Equating and Scaling Procedures .....	84
7.1.1 Rounding Rules.....	86
7.2 Accommodative Form Equivalence .....	87
7.2.1 Special Equatings .....	87
7.3 Subscore Performance Levels.....	88
PART 8: RELIABILITY.....	89
8.1 Classical Test Theory Reliability Estimates.....	89
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. ....	89
8.1.1 Reliability and Measurement Error .....	89
8.1.2 Raw Score Internal Consistency.....	90
8.2 Item Response Theory Reliability .....	94
8.2.1 Test Information Functions .....	94

## Table of Contents

8.2.2 Conditional Standard Error of Measurement .....	97
8.2.3 Item Maps.....	100
8.3 Reliability of Performance Classifications.....	104
8.3.1 Conditional Standard Error of Measurement at Each Cut-Score .....	104
8.3.2 Classification Consistency Indices.....	105
8.4 Reliability of Subscore Performance Classifications .....	105
8.5 Rater Reliability .....	107
PART 9: VALIDITY .....	108
9.1 Evidence Based on Test Content.....	108
9.1.1 Alignment Study .....	109
9.2 Evidence Based on Response Processes.....	110
9.2.1 Cognitive Lab Study .....	111
9.3 Evidence Based on Internal Structure .....	111
9.3.1 Intercorrelations .....	112
9.3.2 Other Internal Structure Evidence.....	112
9.4 Evidence Based on Relationships to Other Variables .....	112
9.5 Evidence Based on the Consequences of Testing .....	114
9.6 Other Validity Evidence .....	115
9.7 Summary .....	116
9.7.1 Student Performance Level Classifications: Overall Scale Score.....	117
9.7.2 Student Performance Level Classifications: Domains and Practices Subscores ....	117
9.7.3 Future NJSLA–S Validity Studies .....	118
PART 10: REPORTING.....	120
10.1 Individual Student Report.....	120
10.2 Student Label .....	123
10.3 Student Roster .....	123
10.4 School Summary and District Summary of Schools.....	125
10.5 School and District Performance Level Summary Reports .....	128
References.....	131
APPENDIX A: Glossary of Abbreviations .....	136
APPENDIX B: New Jersey Science Advisory and Bias and Sensitivity Committees–District and County Representation.....	138
APPENDIX C: Statistical Review Reference Sheet .....	141
APPENDIX D: 2019 NJSLA–S Standard Setting: Executive Summary.....	143
APPENDIX E: NJSLA–S Performance Level Descriptors .....	148
E.1 Policy PLDs.....	148
E.2 Threshold PLDs .....	149
E.2.1 Grade 5 Threshold PLDs.....	149
E.2.2 Grade 8 Threshold PLDs.....	161
E.2.3 Grade 11 Threshold PLDs.....	171

## Table of Contents

E.3 Reporting PLDs .....	189
E.3.1 Reporting PLDs – Level 1.....	189
E.3.2 Reporting PLDs – Level 2.....	189
E.3.3 Reporting PLDs – Level 3.....	189
E.3.4 Reporting PLDs – Level 4.....	189
APPENDIX F: Detailed Test Maps .....	190
APPENDIX G: Scale Score Cumulative Frequency Distributions.....	199
APPENDIX H: Item Parameters and Model Fit Tables .....	206
APPENDIX I: Raw Score-to-Scale Score Conversion Tables .....	212
APPENDIX J: Raw Score-to-Theta Subscore Tables .....	219
APPENDIX K: Subscore Proficiency Classifications.....	237
APPENDIX L: Executive Summary of the NJSLA–S Alignment Evaluation Study .....	243

## 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 DCI .....	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
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: 2022 NJSLA–S Grade 5 Item and Point Totals by Reporting Category.....	22
Table 2.4.3: 2022 NJSLA–S Grade 5 DCIs .....	22
Table 2.4.4: 2022 NJSLA–S Grade 5 SEPs .....	23
Table 2.4.5: 2022 NJSLA–S Grade 5 CCCs .....	23
Table 2.4.6: 2022 NJSLA–S Grade 5 Test Construction Statistics .....	23
Table 2.4.7: 2022 NJSLA–S Grade 5 Test Construction DIF Classifications .....	24
Table 2.4.8: 2022 NJSLA–S Grade 8 Item and Point Totals by Reporting Category.....	24
Table 2.4.9: 2022 NJSLA–S Grade 8 DCIs .....	25
Table 2.4.10: 2022 NJSLA–S Grade 8 SEPs .....	25
Table 2.4.11: 2022 NJSLA–S Grade 8 CCCs .....	25
Table 2.4.12: 2022 NJSLA–S Grade 8 Test Construction Statistics .....	26
Table 2.4.13: 2022 NJSLA–S Grade 8 Test Construction DIF Classifications .....	26
Table 2.4.14: 2022 NJSLA–S Grade 11 Item and Point Totals by Reporting Category.....	26
Table 2.4.15: 2022 NJSLA–S Grade 11 DCIs .....	27
Table 2.4.16: 2022 NJSLA–S Grade 11 SEPs .....	27
Table 2.4.17: 2022 NJSLA–S Grade 11 CCCs .....	27
Table 2.4.18: 2022 NJSLA–S Grade 11 Test Construction Statistics .....	28
Table 2.4.19: 2022 NJSLA–S Grade 11 Test Construction DIF Classifications .....	28
Table 3.1.1: NJSLA–S 2022 Grades 5, 8, and 11 Science Testing Window.....	29
Table 3.2.1: CBT School Test Coordinator Checklist.....	31
Table 3.2.2: PBT School Test Coordinator Checklist.....	31

## Tables and Figures

Table 4.2.1: Scoring Personnel by Grade .....	42
Table 4.2.2: Automatic Rescore Results.....	44
Table 6.1.2: Grade 5 Item Discrimination Distribution and Summary Statistics .....	50
Table 6.1.3: Grade 8 Item Difficulty 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 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: Percent 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: Summary of Yen’s Q3 Statistics .....	76
Table 6.2.12: Constructed-Response Point Distribution Percentages .....	77
Figure 6.2.7. ICC Plot for Grade 5 Constructed-Response Item 1 .....	77
Figure 6.2.8. ICC Plot for Grade 5 Constructed-Response Item 2 .....	78
Figure 6.2.9. ICC Plot for Grade 5 Constructed-Response Item 3 .....	78
Figure 6.2.10. ICC Plot for Grade 8 Constructed-Response Item 1 .....	79

## Tables and Figures

Figure 6.2.11. ICC Plot for Grade 8 Constructed-Response Item 2 .....	79
Figure 6.2.12. ICC Plot for Grade 8 Constructed-Response Item 3 .....	80
Figure 6.2.13. ICC Plot for Grade 11 Constructed-Response Item 1 .....	80
Figure 6.2.14. ICC Plot for Grade 11 Constructed-Response Item 2 .....	81
Figure 6.2.15. ICC Plot for Grade 11 Constructed-Response Item 3 .....	81
Table 6.3.1: Descriptive Statistics of Students' Test Performance by Form .....	82
Table 7.1.1: Scale Score Ranges for Proficiency Levels by Grade.....	85
Table 7.1.2: Slope and Intercept of Theta-to-Scale Score Transformations and Performance Level Cut Scores by Grade .....	86
Table 8.1.1: Coefficient Alpha and SEM by Form.....	91
Table 8.1.2: Coefficient Alpha and SEM by Reporting Category .....	92
Table 8.1.3: Coefficient Alpha and SEM by Demographic Group.....	93
Table 8.1.4: Coefficient Alpha and SEM by Item Type .....	94
Figure 8.2.1. Grade 5 Test Information Function .....	96
Figure 8.2.2. Grade 8 Test Information Function .....	96
Figure 8.2.3. Grade 11 Test Information Function .....	97
Figure 8.2.4. Grade 5 Conditional Standard Error of Measurement.....	98
Figure 8.2.5. Grade 8 Conditional Standard Error of Measurement.....	98
Figure 8.2.6. Grade 11 Conditional Standard Error of Measurement.....	99
Figure 8.2.7 Grade 5 Item Difficulty and Student Ability Distributions .....	101
Figure 8.2.8. Grade 8 Item Difficulty and Student Ability Distributions .....	102
Figure 8.2.9. Grade 11 Item Difficulty and Student Ability Distributions .....	103
Table 8.3.1: Cut Scores with Conditional Standard Error of Measurement.....	104
Table 8.3.2: Performance Level Classification Consistency.....	105
Table 8.4.1: Subscore Performance Classification Consistency and Conditional Standard Error of Measurement .....	106
Table 8.5.1: Inter-rater Agreement Rate of Constructed-Response Items .....	107
Table 9.2.1: Range PLD Alignment by DCI, SEP, and Grade Level.....	111
Table 9.4.1: Grade 5 Intercorrelations by Content Area .....	113
Table 9.4.2: Grade 8 Intercorrelations by Content Area .....	114
Figure 10.2.1. Sample Student Label .....	123
Figure 10.3.1. Sample Student Roster .....	124
Figure 10.4.1. Sample School Performance Level Summary Report – Domains and Practices ..	126
Figure 10.4.2. Sample District Performance Level Summary Report – Domains and Practices..	127
Figure 10.5.1. Sample School Performance Level Summary Report.....	129

## Tables and Figures

Figure 10.5.2. Sample District Performance Level Summary Report.....	130
Table A.1: Glossary of NJSLA–S Abbreviations .....	136
Table B.1: Grade 5 NJSAC District and County Representation .....	138
Table B.2: Grade 8 NJSAC District and County Representation .....	139
Table B.3: Grade 11 NJSAC District and County Representation .....	140
Table B.4: NJBSC District and County Representation .....	140
Table ES-1 Final Recommendations from Standard-Setting Panelists .....	145
Figure ES-1. Percentages of students classified at each level after Round 3 .....	145
Table ES-2 Responses to Key Evaluation Questions.....	146
Table ES-3 Summary of Reasonableness Ratings and Comments .....	147
Table F.1: Grade 5 Test Map – Metadata and Item Statistics.....	190
Table F.2: Grade 8 Test Map – Metadata and Item Statistics.....	193
Table F.3: Grade 11 Test Map – Metadata and Item Statistics.....	196
Table G.1: Grade 5 – Scale Score Cumulative Frequency Distribution.....	199
Table G.2: Grade 8 – Scale Score Cumulative Frequency Distribution.....	201
Table G.3: Grade 11 – Scale Score Cumulative Frequency Distribution.....	203
Table H.1: Grade 5 – IRT Item Parameters and Fit Statistics.....	206
Table H.1: Grade 5 – IRT Item Parameters and Fit Statistics.....	207
Table H.2: Grade 8 – IRT Item Parameters and Fit Statistics.....	208
Table H.2: Grade 8 – IRT Item Parameters and Fit Statistics.....	209
Table H.3: Grade 11 – IRT Item Parameters and Fit Statistics.....	210
Table H.3: Grade 11 – IRT Item Parameters and Fit Statistics.....	211
Table I.1: Grade 5 – Operational.....	212
Table I.2: Grade 8 – Operational.....	214
Table I.3: Grade 11 – Operational.....	216
Table J.1: Grade 5 Earth and Space Science Score Table .....	219
Table J.2: Grade 5 Life Science Score Table.....	220
Table J.3: Grade 5 Physical Science Score Table.....	221
Table J.4: Grade 5 Sensemaking Score Table .....	222
Table J.5: Grade 5 Critiquing Score Table.....	223
Table J.6: Grade 5 Investigating Score Table .....	224
Table J.7: Grade 8 Earth and Space Science Score Table .....	225
Table J.8: Grade 8 Life Science Score Table .....	226
Table J.9: Grade 8 Physical Science Score Table .....	227

## Tables and Figures

Table J.10: Grade 8 Sensemaking Score Table .....	228
Table J.11: Grade 8 Critiquing Score Table.....	229
Table J.12: Grade 8 Investigating Score Table .....	230
Table J.13: Grade 11 Earth and Space Science Score Table .....	231
Table J.14: Grade 11 Life Science Score Table.....	232
Table J.15: Grade 11 Physical Science Score Table.....	233
Table J.16: Grade 11 Sensemaking Score Table .....	234
Table J.17: Grade 11 Critiquing Score Table.....	235
Table J.18: Grade 11 Investigating Score Table .....	236
Table K.1: Grade 5 Content Disaggregated Subscore Proficiency Classifications.....	237
Table K.2: Grade 5 Practice Disaggregated Subscore Proficiency Classifications.....	238
Table K.3: Grade 8 Content Disaggregated Subscore Proficiency Classifications.....	239
Table K.4: Grade 8 Practice Disaggregated Subscore Proficiency Classifications.....	240
Table K.5: Grade 11 Content Disaggregated Subscore Proficiency Classifications.....	241
Table K.6: Grade 11 Practice Disaggregated Subscore Proficiency Classifications.....	242

## PART 1: INTRODUCTION

The purpose of this Technical Report is to provide information about the technical characteristics of the 2022 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 2022 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 April 25, 2022, through June 3, 2022.

- Analyses in Part 5 of this report, Standard Setting, are based on test results from a priority sample due to the short time between the test administration and the 2019 NJSLA–S Standard Setting meeting. The priority sample was representative of the entire state student population in terms of various demographic information including gender, ethnicity, English learner status, disability status, etc.
- 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) 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 *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), 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 and from students’ interactions with teachers on a daily basis, as well as from 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 2022 was the second 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 as well as the topics that are 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 three PS DCIs along with the associated detailed topics for each.

**Table 1.2.3: Physical Science DCI**

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, Katch-Singer, and Pelletier (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 (OEI)	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 (that is, the total number of points a student earned on the test as a whole), 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 that students should have at each of the performance levels, 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 to be below the state minimum level of proficiency. They demonstrate minimal or partial understanding of the DCIs and SEPs. Students at this performance level may need additional instructional support, which could be in the form of 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 2022. 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 either one or two points; each constructed-response (CR) item is worth 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	<b>Investigating</b> <i>AQDP, PACI, UMCT</i>	1–2	3–5	1–2	4–7	1–2	4–8
PS	<b>Sensemaking</b> <i>DUM, AID, CEDS</i>	1–2	3–5	1–2	4–7	1–2	4–8
PS	<b>Critiquing</b> <i>EAE, OECI</i>	1–2	3–5	1–2	4–7	1–2	4–8
PS	<b>Total Items</b>	3–5	11–13	3–5	14–18	3–5	15–21
LS	<b>Investigating</b> <i>AQDP, PACI, UMCT</i>	1–2	3–5	1–2	4–7	1–2	4–8
LS	<b>Sensemaking</b> <i>DUM, AID, CEDS</i>	1–2	3–5	1–2	4–7	1–2	4–8
LS	<b>Critiquing</b> <i>EAE, OECI</i>	1–2	3–5	1–2	4–7	1–2	4–8
LS	<b>Total Items</b>	3–5	11–13	3–5	14–18	3–5	15–21
ESS	<b>Investigating</b> <i>AQDP, PACI, UMCT</i>	1–2	3–5	1–2	4–7	1–2	4–8
ESS	<b>Sensemaking</b> <i>DUM, AID, CEDS</i>	1–2	3–5	1–2	4–7	1–2	4–8
ESS	<b>Critiquing</b> <i>EAE, OECI</i>	1–2	3–5	1–2	4–7	1–2	4–8
ESS	<b>Total Items</b>	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.

<b>MSA: 4 stimuli, 3 items each</b>				<b>PBA: 1 stimulus, 4 items</b>
<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"> <li>• 2 one-point TEs</li> <li>• 1 two-point TEs</li> <li>• 1 four-point CR</li> </ul>
Total # items, MSA: 12 Total points, MSA: 12				Total # items, PBA: 4 Total points, PBA: 8
Total # items, Unit: 16 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 that 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 six 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 2022, NJDOE required that the PBA section contain 7 to 8 total points, with four of those points coming from the CR item.

### 2.1.3 Item Types

Three 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. Some TE items are scored dichotomously; others are rubric-dependent and can be worth up to two points.
- 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 describes each NJSLA–S item type.

**Table 2.1.2: NJSLA–S Item Types**

<b>Item Type</b>	<b>Description</b>
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 (PBS) 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 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, and had to review and approve each item that appears on the NJSLA–S 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, Johnstone, & Thurlow, 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 PBS 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 10 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 2019 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, as well as the counties they come from, are presented in Appendix B.

The NJSAC are the final authority on universal design principle #2: precisely defined constructs. They 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 NJSLA–S test. The NJSAC item reviews predominantly took place in-person at locations approved by NJDOE. Occasionally, it was necessary to conduct the meetings via 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 onto 10 to 18 different 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 gets a copy of 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 and Pearson’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 that is 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 the sections below.

**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**

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

### 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 (see detailed discussions in Part 6: Item and Test Statistics), the items are replaced. The content and psychometric review then resumes until the test matches NJSLA–S’ content and statistical constraints.

### 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 2022 NJSLA–S Test Construction

Overall, the test construction process achieved forms that matched the balance required by the test blueprint. The science content domains were well-balanced at each grade level. Moreover, all grade levels had three eight-point PBA sections representing each of the three content domains, and all of these met the requirement that no more than 50% of the MSA items be MC. However, some constraints were more difficult to achieve. At all three grade levels, it was challenging to identify enough Investigating items—that were also acceptable from a content and statistical perspective—to balance out the three scientific practice reporting categories.

A final test construction content constraint that was not met was the balance between the three content domains across the three scientific practices reporting categories, as shown in the test blueprint in Section 2.1.1. The items associated with each scientific practice were meant to be balanced across all content domains. Table 2.4.1 shows this lack of balance. At each grade level, one content domain was over-represented for each scientific practice. For instance, of the 19 Investigating points available on the grade 11 test, 12 were aligned to the Physical Science content domain, whereas only 2 of 19 were aligned to Earth and Space Science.

Based on the recommendation of content specialists, psychometricians, NJDOE, and the NJSAC, rule adjustments were made to accommodate the differences among the myriad content standards. Balancing the scientific practices across all content domains has been made a primary objective of current NJSLA–S item development. All items currently under development, as described in this paragraph, are scheduled to be field tested in the spring of 2023 and incorporated into the next round of test construction. To iteratively improve the ability of the NJSLA–S to foster valid interpretations and uses of test scores, the test construction issues noted above were addressed by adjusting item development procedures and revising the PBA point total constraints. First, NJSLA–S item development has focused on increasing the proportion of Investigating and, to a lesser extent, Critiquing items. Next, the rules requiring that each PBA have items totaling eight points have been relaxed so that a PBA can total as few as six points. Furthermore, the rules requiring that each CR be worth four points were updated to allow the value of a CR to range from two to four points.

**Table 2.4.1: Points Available by Domain and Practice**

<b>Grade</b>	<b>Practice</b>	<b>Earth</b>	<b>Life</b>	<b>Physical</b>
5	Investigating	3	5	6
5	Sensemaking	10	11	7
5	Critiquing	6	6	6
8	Investigating	5	11	5
8	Sensemaking	13	8	11
8	Critiquing	4	8	7
11	Investigating	2	5	12
11	Sensemaking	12	14	9
11	Critiquing	10	7	7

### 2.4.1 Grade 5 Test Construction

For grade 5, out of 60 total score points, the three content domains ranged from 19 to 22 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 14 out of 60 points being allocated to the Investigating reporting category. Despite being less than ideal, the 14 points were still enough to produce reliable measures of student Investigating abilities. Other content considerations that were met included: MC items only made-up 15 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: 2022 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	9	1	15	19
Life	4	14	1	19	22
Physical	6	9	1	16	19
Total – Domains	15	32	3	50	60
Investigating	2	11	0	13	14
Sensemaking	8	16	1	25	28
Critiquing	5	5	2	12	18
Total – Practices	15	32	3	50	60

**Table 2.4.3: 2022 NJSLA–S Grade 5 DCIs**

DCI	Items	Points
ESS1	6	6
ESS2	7	11
ESS3	2	2
LS1	4	4
LS2	8	8
LS3	4	7
LS4	3	3
PS1	3	3
PS2	3	3
PS3	7	10
PS4	3	3

**Table 2.4.4: 2022 NJSLA–S Grade 5 SEPs**

SEP	Items	Points
AQDP	4	4
PACI	3	3
UMCT	6	7
DUM	7	7
AID	10	13
CEDS	8	8
EAE	9	15
OECI	3	3

**Table 2.4.5: 2022 NJSLA–S Grade 5 CCCs**

CCC	Items	Points
C & E	9	9
E & M	6	9
Patterns	16	22
S & SM	8	9
S, P, & Q	5	5
SC	3	3
SF	3	3

The statistical constraints for the 2022 Grade 5 NJSLA–S operational test form were met. All items had item-total correlations above the 0.2 threshold, and each of the model fit statistics averaged close to their ideal values. The median test time of 104.14 minutes was close to the 105-minute threshold, and out of 200 DIF classifications, there were zero "C" values and only five "B" values. 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: 2022 NJSLA–S Grade 5 Test Construction Statistics**

Statistic	Average	Target	Flags
Item Difficulty	0.20	N/A	N/A
IT Correlation	0.44	> 0.35	0
Infit	0.98	1.00	0
Outfit	0.99	1.00	4
PCM Discrim.	1.03	1.00	4
Lower Asymptote	0.03	0.00	3
Median Time	104.14	< 105	N/A

**Table 2.4.7: 2022 NJSLA–S Grade 5 Test Construction DIF Classifications**

<b>Groups</b>	<b>A</b>	<b>B</b>	<b>C</b>
Male/Female	47	3	0
White/Black	48	2	0
White/Hispanic	50	0	0
White/Asian	50	0	0

**2.4.2 Grade 8 Test Construction**

At grade 8, the least balanced content domain was Earth and Space Science, and it still made-up 22 points of the 72 total score points. Each content domain had one PBA section devoted to it. The scientific practices were less balanced with only 19 out of 72 points allocated to the Critiquing reporting category. Despite being less than ideal, 19 points provided enough information to produce reliable measures of student Critiquing abilities. Other content considerations that were met included: MC items only made-up 22 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: 2022 NJSLA–S Grade 8 Item and Point Totals by Reporting Category**

<b>Domains/Practices</b>	<b>MC Items</b>	<b>TE Items</b>	<b>CR Items</b>	<b>Items</b>	<b>Points</b>
Earth and Space	5	11	1	17	22
Life	12	10	1	23	27
Physical	5	14	1	20	23
Total – Domains	22	35	3	60	72
Investigating	10	9	0	19	21
Sensemaking	8	15	2	25	32
Critiquing	4	11	1	16	19
Total – Practices	22	35	3	60	72

**Table 2.4.9: 2022 NJSLA–S Grade 8 DCIs**

DCI	Items	Points
ESS1	6	10
ESS2	5	5
ESS3	6	7
LS1	8	8
LS2	6	10
LS3	5	5
LS4	4	4
PS1	6	6
PS2	4	7
PS3	6	6
PS4	4	4

**Table 2.4.10: 2022 NJSLA–S Grade 8 SEPs**

SEP	Items	Points
AQDP	5	5
PACI	5	6
UMCT	9	10
DUM	8	9
AID	7	7
CEDS	10	16
EAE	10	13
OECI	6	6

**Table 2.4.11: 2022 NJSLA–S Grade 8 CCCs**

CCC	Items	Points
C & E	15	23
E & M	7	7
Patterns	11	11
S & SM	9	12
S, P, & Q	7	7
SC	5	6
SF	6	6

The statistical constraints for the 2022 Grade 8 NJSLA–S operational test form were met. Five grade 8 items were flagged for having item-total correlations below the .2 threshold, including one with the lowest value of 0.12, while all the other four with a value above 0.17. The infit, outfit, and PCM discrimination model fit statistics each averaged close to their ideal values of 1.00. The median test time of 92.43 minutes was below the 105-minute threshold, and out of 240 DIF classifications, there were zero "C" values and only four "B" values. 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: 2022 NJSLA–S Grade 8 Test Construction Statistics**

Statistic	Average	Target	Flags
Item Difficulty	0.21	N/A	N/A
IT Correlation	0.38	> 0.35	5
Infit	0.99	1.00	0
Outfit	1.01	1.00	4
PCM Discrim.	1.00	1.00	4
Lower Asymptote	0.03	0.00	4
Median Time	92.43	< 105	N/A

**Table 2.4.13: 2022 NJSLA–S Grade 8 Test Construction DIF Classifications**

Groups	A	B	C
Male/Female	59	1	0
White/Black	58	2	0
White/Hispanic	60	0	0
White/Asian	59	1	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 less balanced with only 19 out of 78 points being allocated to the Investigating reporting category. Despite being less than ideal, the 19 points provided sufficient information to produce reliable measures of student Critiquing abilities. Other content considerations that were met included: MC items only made-up 28 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. The seven CCCs were well-balanced. 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 at grade 11.

**Table 2.4.14: 2022 NJSLA–S Grade 11 Item and Point Totals by Reporting Category**

Domains/Practices	MC Items	TE Items	CR Items	Items	Points
Earth and Space	6	14	1	21	24
Life	10	12	1	23	26
Physical	12	12	1	25	28
Total – Domains	28	38	3	69	78
Investigating	12	7	0	19	19
Sensemaking	10	21	1	32	35
Critiquing	6	10	2	18	24
Total – Practices	28	38	3	69	78

**Table 2.4.15: 2022 NJSLA–S Grade 11 DCIs**

DCI	Items	Points
ESS1	5	5
ESS2	7	7
ESS3	9	12
LS1	4	4
LS2	8	11
LS3	7	7
LS4	4	4
PS1	7	7
PS2	3	3
PS3	10	10
PS4	5	8

**Table 2.4.16: 2022 NJSLA–S Grade 11 SEPs**

SEP	Items	Points
AQDP	6	6
PACI	5	5
UMCT	8	8
DUM	7	7
AID	16	16
CEDS	9	12
EAE	13	19
OECI	5	5

**Table 2.4.17: 2022 NJSLA–S Grade 11 CCCs**

CCC	Items	Points
C & E	12	15
E & M	9	9
Patterns	9	12
S & SM	13	16
S, P, & Q	9	9
SC	10	10
SF	7	7

The 2022 Grade 11 NJSLA–S operational test form construction saw seven items flagged for outfit, five items flagged for PCM discrimination, and seven items flagged for lower asymptote. Nevertheless, the flagged items had an outfit, PCM discrimination or lower asymptote value near the thresholds. The items flagged for outfit would not distort or degrade the measures (Engelhard & Wang, 2021; Linacre, 2012) and the average value of infit or outfit, was close to their ideal values of 1.00. Two grade 11 items were flagged for having item-total correlations below the 0.2 threshold, with values hovering about 0.18. The median test time was only 59.51 minutes, which was about 90 minutes below the 150-minute constraint, indicating the

high school students’ lack of motivation to perform their best on the assessment. Of 272 DIF classifications, there were zero “C” values and only six “B” values. 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: 2022 NJSLA–S Grade 11 Test Construction Statistics**

Statistic	Average	Target	Flags
Item Difficulty	0.35	N/A	N/A
IT Correlation	0.43	> 0.35	2
Infit	1.00	1.00	0
Outfit	1.03	1.00	7
PCM Discrim.	0.99	1.00	5
Lower Asymptote	0.03	0.00	7
Median Time	59.51	< 105	N/A

**Table 2.4.19: 2022 NJSLA–S Grade 11 Test Construction DIF Classifications**

Groups	A	B	C
Male/Female	67	2	0
White/Black	68	1	0
White/Hispanic	68	1	0
White/Asian	68	2	0

## 2.5 2022 NJSLA–S State of the Item Bank

Upon the completion of the 2022 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 2023 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 on 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” (NJDOE, 2022, slide 2). Information about the administration of NJSLA–S is available 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. The NJSLA TCM can be read in full at the [NJSLA–S website](#) under **Documents and Downloads**.

- 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

Table 3.1.1 shows the NJSLA–S 2022 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 2022 Grades 5, 8, and 11 Science Testing Window**

Grade	CBT	PBT	Testing Time
5	4/25/22–6/3/22	4/25/22–5/6/22	45 minutes/per unit
8	4/25/22–6/3/22	4/25/22–5/6/22	45 minutes/per unit
11	4/25/22–6/3/22	4/25/22–5/6/22	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<sup>next</sup> to manage each test session; they can monitor the progress of each of their students and lock and unlock units. PearsonAccess<sup>next</sup> 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**

<b>Tasks</b>	<b>TCM Section(s)</b>
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 <sup>next</sup> .	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 <sup>next</sup> .	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**

<b>Tasks</b>	<b>TCM Section(s)</b>
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<sup>next</sup> 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., posters, models, displays, teaching aids) 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. Previously in New Jersey, federal requirements were met by testing grades 4 and 8 students with the NJASK test; grade 11 students were tested via the NJBCT. 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**. 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 are 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 NJ 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 2022 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 of 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. Paper-Based Test.** The conversion of the NJSLA–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 don't 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 2022, 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 to ensure 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. MI’s psychometric department 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 percentage of students selecting each response, and the item-total correlation associated with each response option that was selected more than 100 times. Item means and item-total correlations are also calculated at the item level, and items are flagged for aberrant behavior. 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 specialist 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 NJ 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, NJ 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, NJ 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 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 2022 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**

<b>Grade</b>	<b>Scorers</b>	<b>Team Leaders</b>	<b>Scoring Director</b>
5	183	3	1
8	188	3	1
11	144	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, we conduct a blank check 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	Percent Changed (of those Eligible)
5	2,014	320	15.9
8	1,388	227	16.3
11	1,476	255	17.2

### 4.3 Quality Control

NJDOE conducted three Quality Control (QC) checks to confirm that the processing of student tests and test scores was done correctly. MI selected a sample of several hundred student tests for NJDOE to manually review and independently score. NJDOE staff then compared those scores to the data files 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. 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 sections detail the timing and processes used in the three QCs to ensure student test scores are accurate.

#### 4.3.1 QC #1

After testing and scoring were complete, NJDOE conducted the first QC check to verify preliminary data in a data file generated by Pearson. This activity ran July 19–August 5, 2022. To begin the QC process, MI provided NJDOE with 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, NJDOE verified the student information on the KIS against the student information in the test registration system. Next, NJDOE reviewed the student’s responses to each selected-response item, scored them against the key, and recorded the score on the KIS. This task was somewhat complex for technology-enhanced selected-response items. 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.

#### **4.3.2 QC #2**

QC #2 was conducted from August 15–19, 2022. 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 from QC #1 to determine the 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 QC #3**

QC #3 was conducted from August 22–26, 2022. 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 only show aggregated data. These values were verified in QC #3 by calculating the same figures from the raw data in the data files. This step is not necessarily limited to the schools in the QC sample and the work can be done simultaneously with QC #2, although any problems discovered with the data may require the work to be repeated.

## 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. 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) and item response theory (IRT) statistics, along with the content constraints. The following sections detail how well the 2022 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 were 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 an item was performed in an unintended manner and could negatively impact the reliability or the validity of test score interpretations, it was recommended to NJDOE that the item be removed from the assessment. Descriptions of each type of statistic appear in the following sections.

#### 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 to 1.00; an item with a high  $p$ -value is easier to answer correctly, whereas one 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 0–2 point TE and 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), then flagging the item if its converted p-value falls outside of the .25 to .90 range. It should be noted that the flagging criteria 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, because if items do not discriminate between students with high levels of ability in comparison to students with low levels of ability, then both reliability and the validity of test score interpretations can suffer. CTT 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 negative are considered for removal from the test because they could be harming 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 earth and space science tended to be more slightly challenging than in life science and physical science. At grade 8, only four MC items and three TE items had p-values above .50, indicating most of the items at the harder end of the scale. At grade 11, zero MC or TE items had p-values above .75, indicating no items on the easier end of the scale. At grades 8 and 11, the CR item in life science tended to be more challenging than in the other content domains.

The average item-total correlations (*rpb*) were .62, .64, and .66 for grades 5, 8, and 11 CR items, respectively, indicating positive discrimination powers. 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 had an item-total correlation below .20, while seven (47%) MC and 25 (78%) TE items had item-total correlations above .40. At grade 8, three MC and three TE items had item-total correlations below .20, while 3 (14%) MC items and 16 (43%) TE items had discriminations above .40. At grade 11, two MC and one TE items had item-total correlations below .20, while 13 (46%) MC and 26 (68%) TE items had discriminations above .40.

**Table 6.1.1: Grade 5 Item Difficulty 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	<b>NJSLA-S</b>	<b>15</b>	<b>0</b>	<b>5</b>	<b>10</b>	<b>0</b>	<b>0</b>	<b>.52</b>	<b>.10</b>	<b>.53</b>
	Earth and Space	5	0	2	3	0	0	.49	.15	.55
	Life	4	0	2	2	0	0	.53	.08	.51
	Physical	6	0	1	5	0	0	.54	.05	.53
	Critiquing	5	0	1	4	0	0	.53	.09	.52
	Investigating	2	0	1	1	0	0	.52	.03	.52
	Sensemaking	8	0	3	5	0	0	.51	.12	.53
TE	<b>NJSLA-S</b>	<b>32</b>	<b>1</b>	<b>23</b>	<b>7</b>	<b>1</b>	<b>0</b>	<b>.44</b>	<b>.14</b>	<b>.41</b>
	Earth and Space	9	0	6	2	1	0	.45	.15	.47
	Life	14	1	9	4	0	0	.46	.15	.42
	Physical	9	0	8	1	0	0	.39	.11	.36
	Critiquing	5	0	3	2	0	0	.46	.16	.43
	Investigating	11	1	9	1	0	0	.36	.12	.33
	Sensemaking	16	0	11	4	1	0	.48	.14	.47
CR	<b>NJSLA-S</b>	<b>3</b>	<b>1</b>	<b>2</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>.32</b>	<b>.08</b>	<b>.34</b>
	Earth and Space	1	1	0	0	0	0	.23	N/A	.23
	Life	1	0	1	0	0	0	.38	N/A	.38
	Physical	1	0	1	0	0	0	.34	N/A	.34
	Critiquing	2	0	2	0	0	0	.36	.03	.36
	Investigating	0	0	0	0	0	0	N/A	N/A	N/A
	Sensemaking	1	1	0	0	0	0	.23	N/A	.23

**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	<b>NJSLA-S</b>	<b>15</b>	<b>1</b>	<b>2</b>	<b>5</b>	<b>5</b>	<b>2</b>	<b>.40</b>	<b>.10</b>	<b>.39</b>
	Earth and Space	5	1	1	1	0	2	.37	.16	.32
	Life	4	0	1	3	0	0	.35	.05	.36
	Physical	6	0	0	1	5	0	.44	.04	.45
	Critiquing	5	0	0	1	2	2	.47	.09	.47
	Investigating	2	1	0	1	0	0	.29	.13	.29
	Sensemaking	8	0	2	3	3	0	.38	.07	.39
TE	<b>NJSLA-S</b>	<b>32</b>	<b>0</b>	<b>3</b>	<b>4</b>	<b>11</b>	<b>14</b>	<b>.48</b>	<b>.11</b>	<b>.48</b>
	Earth and Space	9	0	0	1	2	6	.51	.07	.54
	Life	14	0	2	1	4	7	.50	.14	.50
	Physical	9	0	1	2	5	1	.42	.08	.41
	Critiquing	5	0	0	1	1	3	.49	.07	.50
	Investigating	11	0	2	3	2	4	.45	.13	.46
	Sensemaking	16	0	1	0	8	7	.50	.11	.48
CR	<b>NJSLA-S</b>	<b>3</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>3</b>	<b>.62</b>	<b>.04</b>	<b>.63</b>
	Earth and Space	1	0	0	0	0	1	.57	N/A	.57
	Life	1	0	0	0	0	1	.63	N/A	.63
	Physical	1	0	0	0	0	1	.66	N/A	.66
	Critiquing	2	0	0	0	0	2	.64	.02	.64
	Investigating	0	0	0	0	0	0	N/A	N/A	N/A
	Sensemaking	1	0	0	0	0	1	.57	N/A	.57

**Table 6.1.3: Grade 8 Item Difficulty 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	<b>NJSLA-S</b>	<b>22</b>	<b>2</b>	<b>16</b>	<b>4</b>	<b>0</b>	<b>0</b>	<b>.39</b>	<b>.09</b>	<b>.39</b>
	Earth and Space	5	0	4	1	0	0	.45	.06	.44
	Life	12	2	8	2	0	0	.37	.10	.35
	Physical	5	0	4	1	0	0	.39	.08	.38
	Critiquing	4	0	4	0	0	0	.36	.04	.36
	Investigating	10	2	6	2	0	0	.37	.11	.35
	Sensemaking	8	0	6	2	0	0	.43	.08	.44
TE	<b>NJSLA-S</b>	<b>35</b>	<b>14</b>	<b>18</b>	<b>2</b>	<b>1</b>	<b>0</b>	<b>.33</b>	<b>.21</b>	<b>.26</b>
	Earth and Space	11	4	6	1	0	0	.40	.29	.34
	Life	10	4	6	0	0	0	.33	.17	.33
	Physical	14	6	6	1	1	0	.28	.18	.26
	Critiquing	11	4	7	0	0	0	.27	.11	.26
	Investigating	9	3	5	1	0	0	.43	.33	.35
	Sensemaking	15	7	6	1	1	0	.32	.18	.26
CR	<b>NJSLA-S</b>	<b>3</b>	<b>1</b>	<b>2</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>.27</b>	<b>.08</b>	<b>.28</b>
	Earth and Space	1	0	1	0	0	0	.28	N/A	.28
	Life	1	1	0	0	0	0	.19	N/A	.19
	Physical	1	0	1	0	0	0	.35	N/A	.35
	Critiquing	1	1	0	0	0	0	.19	N/A	.19
	Investigating	0	0	0	0	0	0	N/A	N/A	N/A
	Sensemaking	2	0	2	0	0	0	.31	.05	.31

**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	<b>NJSLA-S</b>	<b>22</b>	<b>3</b>	<b>4</b>	<b>12</b>	<b>2</b>	<b>1</b>	<b>.32</b>	<b>.11</b>	<b>.32</b>
	Earth and Space	5	0	2	3	0	0	.31	.03	.31
	Life	12	2	1	6	2	1	.32	.14	.35
	Physical	5	1	1	3	0	0	.30	.09	.34
	Critiquing	4	1	1	2	0	0	.26	.06	.27
	Investigating	10	2	1	4	2	1	.33	.16	.38
	Sensemaking	8	0	2	6	0	0	.33	.04	.32
TE	<b>NJSLA-S</b>	<b>35</b>	<b>3</b>	<b>6</b>	<b>10</b>	<b>13</b>	<b>3</b>	<b>.36</b>	<b>.10</b>	<b>.39</b>
	Earth and Space	11	1	3	3	3	1	.33	.11	.32
	Life	10	1	2	1	4	2	.39	.12	.41
	Physical	14	1	1	6	6	0	.37	.08	.39
	Critiquing	11	0	2	2	6	1	.40	.07	.41
	Investigating	9	1	0	3	3	2	.40	.12	.45
	Sensemaking	15	2	4	5	4	0	.32	.10	.33
CR	<b>NJSLA-S</b>	<b>3</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>3</b>	<b>.64</b>	<b>.06</b>	<b>.62</b>
	Earth and Space	1	0	0	0	0	1	.58	N/A	.58
	Life	1	0	0	0	0	1	.62	N/A	.62
	Physical	1	0	0	0	0	1	.71	N/A	.71
	Critiquing	1	0	0	0	0	1	.62	N/A	.62
	Investigating	0	0	0	0	0	0	N/A	N/A	N/A
	Sensemaking	2	0	0	0	0	2	.65	.09	.65

**Table 6.1.5: Grade 11 Item Difficulty 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	<b>NJSLA-S</b>	<b>28</b>	<b>1</b>	<b>20</b>	<b>7</b>	<b>0</b>	<b>0</b>	<b>.44</b>	<b>.12</b>	<b>.41</b>
	Earth and Space	6	0	5	1	0	0	.41	.15	.38
	Life	10	0	7	3	0	0	.47	.13	.48
	Physical	12	1	8	3	0	0	.42	.11	.41
	Critiquing	6	0	4	2	0	0	.42	.13	.42
	Investigating	12	1	8	3	0	0	.43	.12	.41
	Sensemaking	10	0	8	2	0	0	.45	.14	.43
TE	<b>NJSLA-S</b>	<b>38</b>	<b>12</b>	<b>13</b>	<b>13</b>	<b>0</b>	<b>0</b>	<b>.38</b>	<b>.16</b>	<b>.35</b>
	Earth and Space	14	4	7	3	0	0	.37	.14	.33
	Life	12	3	3	6	0	0	.43	.17	.47
	Physical	12	5	3	4	0	0	.35	.19	.32
	Critiquing	10	4	3	3	0	0	.38	.18	.49
	Investigating	7	3	2	2	0	0	.34	.19	.31
	Sensemaking	21	5	8	8	0	0	.40	.16	.42
CR	<b>NJSLA-S</b>	<b>3</b>	<b>0</b>	<b>3</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>.34</b>	<b>.05</b>	<b>.32</b>
	Earth and Space	1	0	1	0	0	0	.39	N/A	.39
	Life	1	0	1	0	0	0	.30	N/A	.30
	Physical	1	0	1	0	0	0	.32	N/A	.32
	Critiquing	2	0	2	0	0	0	.36	.05	.36
	Investigating	0	0	0	0	0	0	N/A	N/A	N/A
	Sensemaking	1	0	1	0	0	0	.30	N/A	.30

**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	<b>NJSLA–S</b>	<b>28</b>	<b>2</b>	<b>6</b>	<b>7</b>	<b>12</b>	<b>1</b>	<b>.36</b>	<b>.10</b>	<b>.38</b>
	Earth and Space	6	1	1	3	1	0	.32	.11	.34
	Life	10	0	1	2	6	1	.41	.09	.42
	Physical	12	1	4	2	5	0	.34	.10	.32
	Critiquing	6	1	0	2	3	0	.36	.11	.37
	Investigating	12	0	4	3	5	0	.35	.09	.36
	Sensemaking	10	1	2	2	4	1	.37	.11	.38
TE	<b>NJSLA–S</b>	<b>38</b>	<b>1</b>	<b>3</b>	<b>8</b>	<b>11</b>	<b>15</b>	<b>.45</b>	<b>.12</b>	<b>.46</b>
	Earth and Space	14	0	2	3	5	4	.43	.11	.44
	Life	12	0	0	2	4	6	.50	.09	.51
	Physical	12	1	1	3	2	5	.41	.15	.44
	Critiquing	10	0	0	0	5	5	.50	.06	.49
	Investigating	7	1	1	1	1	3	.40	.19	.44
	Sensemaking	21	0	2	7	5	7	.44	.11	.46
CR	<b>NJSLA–S</b>	<b>3</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>3</b>	<b>.66</b>	<b>.05</b>	<b>.67</b>
	Earth and Space	1	0	0	0	0	1	.70	N/A	.70
	Life	1	0	0	0	0	1	.61	N/A	.61
	Physical	1	0	0	0	0	1	.67	N/A	.67
	Critiquing	2	0	0	0	0	2	.69	.02	.69
	Investigating	0	0	0	0	0	0	N/A	N/A	N/A
	Sensemaking	1	0	0	0	0	1	.61	N/A	.61

### 6.1.2 Speededness

The consequence of time limits on examinees' scores is called speededness. A measure of the speededness of a test is the number of items that are not attempted by students and hence their test scores are impacted. In each separately timed subsection of a test, if a student does not attempt the last item, it can be assumed that the student may have run out of time. 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; if a 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 and 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	7 MC, 9 TE, 1 CR	45
5	2	5 MC, 11 TE, 1 CR	45
5	3	3 MC, 12 TE, 1 CR	45
8	1	7 MC, 12 TE, 1 CR	45
8	2	7 MC, 12 TE, 1 CR	45
8	3	8 MC, 11 TE, 1 CR	45
11	1	7 MC, 15 TE, 1 CR	60
11	2	12 MC, 10 TE, 1 CR	60
11	3	9 MC, 13 TE, 1 CR	60

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

Grade	Unit	Location	% Student
5	1	17	3.3
5	2	16	0.4
5	3	15	0.2
8	1	19	2.5
8	2	20	1.3
8	3	19	0.7
11	1	23	1.0
11	2	23	0.5
11	3	23	0.7

### 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 to 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 with the exception of a very small number of items classified as “C” including one TE item for White/Black at grade 5, one TE item for CBT/TTS at grade 8, and one MC item for EL-No/EL-Yes at grade 11. 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 and 11. The “C” DIF items will be re-investigated when more test data become available.

**Table 6.1.10: Grade 5 DIF Classification by Item Type**

<b>Grade</b>	<b>Group</b>	<b>Item Type</b>	<b>A</b>	<b>B</b>	<b>C</b>
5	Male/Female	MC	15	0	0
5	Male/Female	TE	30	2	0
5	Male/Female	CR	3	0	0
5	Male/Female	Total	48	2	0
5	White/Black	MC	15	0	0
5	White/Black	TE	30	1	1
5	White/Black	CR	3	0	0
5	White/Black	Total	48	1	1
5	White/Hispanic	MC	15	0	0
5	White/Hispanic	TE	32	0	0
5	White/Hispanic	CR	3	0	0
5	White/Hispanic	Total	50	0	0
5	White/Asian	MC	15	0	0
5	White/Asian	TE	32	0	0
5	White/Asian	CR	3	0	0
5	White/Asian	Total	50	0	0
5	EL-No/EL-Yes	MC	15	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	50	0	0
5	SWD-No/SWD-Yes	MC	15	0	0
5	SWD-No/SWD-Yes	TE	32	0	0
5	SWD-No/SWD-Yes	CR	3	0	0
5	SWD-No/SWD-Yes	Total	50	0	0
5	EconDis-No/EconDis-Yes	MC	15	0	0
5	EconDis-No/EconDis-Yes	TE	32	0	0
5	EconDis-No/EconDis-Yes	CR	3	0	0
5	EconDis-No/EconDis-Yes	Total	50	0	0
5	CBT/TTS	MC	15	0	0
5	CBT/TTS	TE	32	0	0
5	CBT/TTS	CR	3	0	0
5	CBT/TTS	Total	50	0	0

**Table 6.1.11: Grade 8 DIF Classification by Item Type**

<b>Grade</b>	<b>Group</b>	<b>Item Type</b>	<b>A</b>	<b>B</b>	<b>C</b>
8	Male/Female	MC	22	0	0
8	Male/Female	TE	34	1	0
8	Male/Female	CR	3	0	0
8	Male/Female	Total	59	1	0
8	White/Black	MC	22	0	0
8	White/Black	TE	34	1	0
8	White/Black	CR	3	0	0
8	White/Black	Total	59	1	0
8	White/Hispanic	MC	22	0	0
8	White/Hispanic	TE	34	1	0
8	White/Hispanic	CR	3	0	0
8	White/Hispanic	Total	59	1	0
8	White/Asian	MC	22	0	0
8	White/Asian	TE	35	0	0
8	White/Asian	CR	3	0	0
8	White/Asian	Total	60	0	0
8	EL-No/EL-Yes	MC	22	0	0
8	EL-No/EL-Yes	TE	33	2	0
8	EL-No/EL-Yes	CR	3	0	0
8	EL-No/EL-Yes	Total	58	2	0
8	SWD-No/SWD-Yes	MC	22	0	0
8	SWD-No/SWD-Yes	TE	35	0	0
8	SWD-No/SWD-Yes	CR	3	0	0
8	SWD-No/SWD-Yes	Total	60	0	0
8	EconDis-No/EconDis-Yes	MC	22	0	0
8	EconDis-No/EconDis-Yes	TE	35	0	0
8	EconDis-No/EconDis-Yes	CR	3	0	0
8	EconDis-No/EconDis-Yes	Total	60	0	0
8	CBT/TTS	MC	22	0	0
8	CBT/TTS	TE	34	0	1
8	CBT/TTS	CR	3	0	0
8	CBT/TTS	Total	59	0	1

**Table 6.1.12: Grade 11 DIF Classification by Item Type**

<b>Grade</b>	<b>Group</b>	<b>Item Type</b>	<b>A</b>	<b>B</b>	<b>C</b>
11	Male/Female	MC	27	1	0
11	Male/Female	TE	36	2	0
11	Male/Female	CR	2	1	0
11	Male/Female	Total	65	4	0
11	White/Black	MC	28	0	0
11	White/Black	TE	38	0	0
11	White/Black	CR	3	0	0
11	White/Black	Total	69	0	0
11	White/Hispanic	MC	28	0	0
11	White/Hispanic	TE	38	0	0
11	White/Hispanic	CR	3	0	0
11	White/Hispanic	Total	69	0	0
11	White/Asian	MC	27	1	0
11	White/Asian	TE	38	0	0
11	White/Asian	CR	3	0	0
11	White/Asian	Total	68	1	0
11	EL-No/EL-Yes	MC	27	0	1
11	EL-No/EL-Yes	TE	35	3	0
11	EL-No/EL-Yes	CR	1	2	0
11	EL-No/EL-Yes	Total	63	5	1
11	SWD-No/SWD-Yes	MC	28	0	0
11	SWD-No/SWD-Yes	TE	38	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	28	0	0
11	EconDis-No/EconDis-Yes	TE	38	0	0
11	EconDis-No/EconDis-Yes	CR	3	0	0
11	EconDis-No/EconDis-Yes	Total	69	0	0
11	CBT/TTS	MC	28	0	0
11	CBT/TTS	TE	38	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. 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 a mathematical equation for the relationship of the probability of a student response to a test item and student latent ability level 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 that we are assuming the test items estimate is student understanding of the New Jersey science curriculum: the NJSLA–S, and the directly observable behaviors 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 tests. Under IRT, item parameters (e.g., item difficulty) are independent of the students who took the test; 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 in order 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, the data collected can then be used for psychometric analyses such as equating.

Comprising dichotomous and polytomous items with varying score points (e.g., 0–2 point TE or 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; Rasch models all 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 exist 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, 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  $\theta_j$  is the student proficiency score,  $\beta_i$  denotes the difficulty or location parameter for item  $i$  and  $\tau_{in}$  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  $\theta_j$  to obtain a correct response to item  $i$  and  $\beta_i$  denotes the item difficulty parameter for item  $i$ .

If the PCM's assumptions are met, it is likely a good IRT model to use with the NJSLA–S. Then, at each grade level, 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. Thus, it is imperative that assumptions be checked.

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 0–4 point CR items. Overall, the results of the 2022 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 a strong linear relationship 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 .75. 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	.93	1	-	-
	Life	.95	.82	1	-
	Physical	.93	.80	.82	1
8	Earth and Space	.90	1	-	-
	Life	.93	.75	1	-
	Physical	.93	.76	.80	1
11	Earth and Space	.94	1	-	-
	Life	.95	.85	1	-
	Physical	.94	.82	.84	1

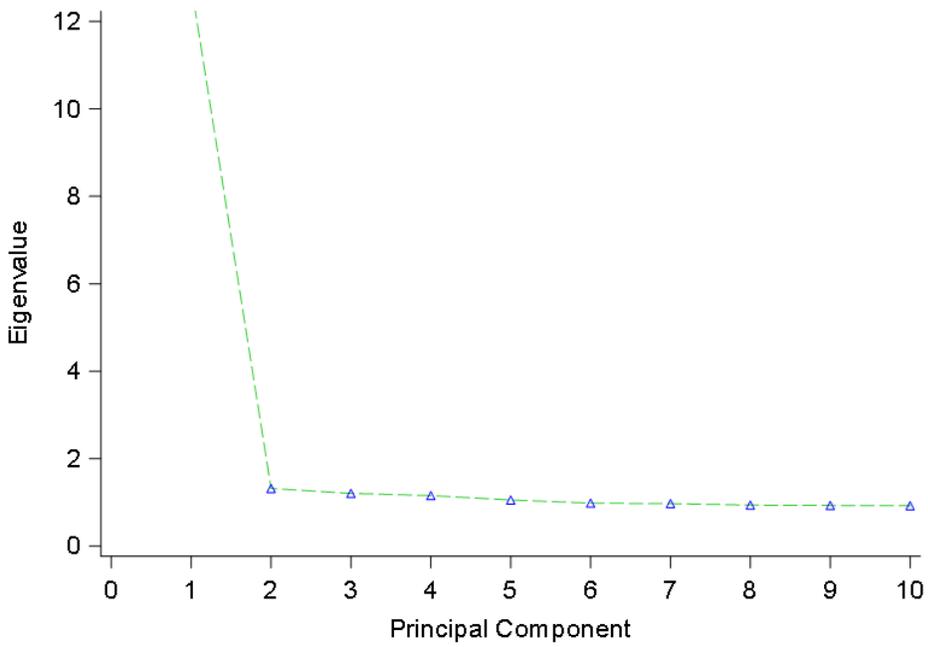
**Table 6.2.2: Correlation Matrix for Practices**

Grade	Practice	NJSLA–S	Investigation	Sensemaking	Critiquing
5	Investigating	.90	1	-	-
	Sensemaking	.97	.82	1	-
	Critiquing	.94	.78	.86	1
8	Investigating	.91	1	-	-
	Sensemaking	.95	.79	1	-
	Critiquing	.90	.75	.78	1
11	Investigating	.90	1	-	-
	Sensemaking	.97	.82	1	-
	Critiquing	.95	.81	.88	1

**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 (Brown, 2006) by converting them into uncorrelated principal components. 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.

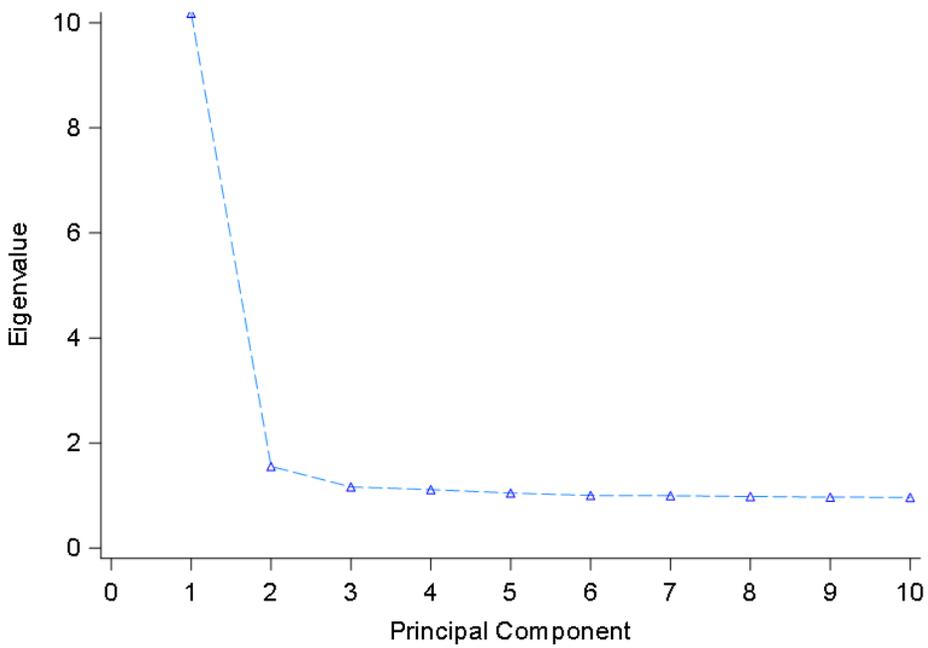
Figures 6.2.1 through 6.2.3 show the scree plots for grades 5, 8, and 11, respectively. As exhibited in these plots, the second most prominent eigenvalue for each grade level is below 2, whereas the most prominent eigenvalues range from approximately 10–14. Each grade’s scree plot shows that only one major dimension is practically 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*

### Grade 11 Scree Plot

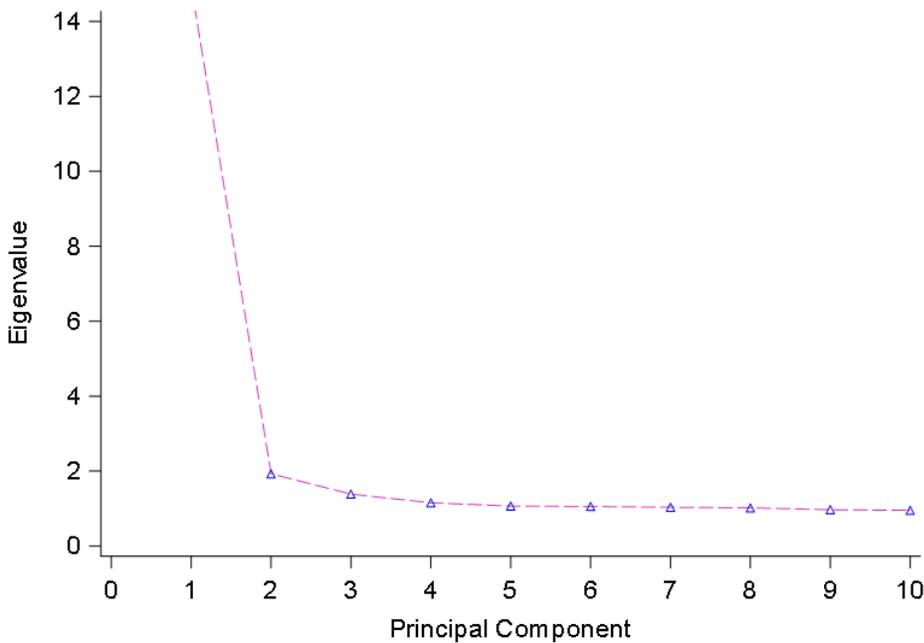


Figure 6.2.3. Grade 11 Scree Plot

#### 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. 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 2022 NJSLA–S test data via Winsteps 3.92.1 (Linacre, 2016). Overall, there is ample evidence that all the grade items fit the assumptions of the PCM. The grade 8 item performance was remarkable with a lower percentage of items flagged than either of the two other grades for each of the four model fit categories.

**6.2.2.1 Item infit and outfit.** Rasch infit and outfit statistics range from 0 to infinity with 1 representing ideal model fit. 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 are measuring near their ability level (Wright and Masters, 1982). Outfit statistics are heavily influenced by unexpected student responses to items that are either relatively easy or relatively hard.

Table 6.2.4 provides a summary of item infit and outfit statistics at each grade level. Only a few items across all grades were flagged for problematic infit statistics. The NJSLA–S outfit statistics were less positive with approximately 18%, 5%, and 12% of all items being flagged for grades 5, 8, and 11, respectively. The problematic outfit statistics, however, are less of a threat to the validity of test score interpretations than are 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	0.69	1.33	2 out of 50	4.0
5	Outfit	1.01	0.59	1.54	9 out of 50	18.0
8	Infit	1.00	0.81	1.35	1 out of 60	1.7
8	Outfit	1.02	0.71	1.54	3 out of 60	5.0
11	Infit	1.00	0.76	1.34	1 out of 69	1.4
11	Outfit	1.03	0.65	1.66	8 out of 69	11.6

**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 assumption of equal discrimination can be tested with the Rasch discrimination statistic. 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. Only four (8%), three (5%), and five (7.2%) items were flagged for having a value outside the 0.5 to 1.5 threshold for grades 5, 8, and 11, respectively. While several of these items were also flagged for item outfit, 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	1.00	0.22	1.62	4 out of 50	8.0
8	Discrimination	0.99	0.40	1.65	3 out of 60	5.0
11	Discrimination	1.00	0.33	1.65	5 out of 69	7.2

**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 guessing assumption is met if items remain within a reasonable threshold. The assumption of guessing can be tested with the Rasch lower asymptote statistic. 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 .1.

Table 6.2.6 provides a summary of the lower asymptote statistics at each grade level. Each grade level saw only four items flagged for having a lower asymptote value outside of the .1 threshold. These items were also flagged for either 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 50	8.0
8	Lower Asymptote	.02	.00	.16	4 out of 60	6.8
11	Lower Asymptote	.03	.00	.13	4 out of 69	5.8

**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 measures of 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. 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, 8.77% of grade 5 students, 4.17% of grade 8 students, and 6.18% of grade 11 students were flagged for person infit statistics. In particular, 12.82% (25 out of 195) grade 5 Pacific Islander students and 9.92% (12 out of 121) grade 11 American Indian students were flagged for person infit. As shown in Table 6.2.9, the grade 5 PBT forms flagged 12.77% of students for person infit. However, there were only 47 PBT test takers.

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 slightly higher percentage of students that were flagged for person outfit statistics relative to other demographic groups at their grade level except for grade 5 Pacific Islander. Within those groups, the students that were flagged also tended to be lower performing. They also were more likely to have taken accommodated forms, which themselves

had much higher percentages of students flagged for person outfit than did the CBT forms. As stated earlier, aberrant person infit statistics are more of a threat to the validity of the inferences than are aberrant person outfit 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. That being said, a deeper investigation into the person outfit statistics for English learners, students with disabilities, and the accommodated forms is warranted and described in greater detail in Section 9.7.3: Future Studies.

**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	<b>NJSLA–S</b>	<b>96,288</b>	<b>162.94</b>	<b>1.02</b>	<b>0.53</b>	<b>2.95</b>	<b>8,441</b>	<b>8.77</b>	<b>166.80</b>
5	Male	49,121	164.78	1.02	0.54	2.95	4,317	8.79	168.85
5	Female	47,156	161.02	1.02	0.53	2.67	4,124	8.75	164.65
5	Am. Indian	150	160.18	1.05	0.73	2.33	14	9.33	146.64
5	Asian	10,437	201.35	1.02	0.54	2.72	979	9.38	203.82
5	Black	13,848	138.38	1.02	0.53	2.67	1,136	8.20	141.14
5	Hispanic	31,155	144.16	1.02	0.55	2.95	2,690	8.63	147.08
5	Pacific Islander	195	171.83	1.04	0.61	1.91	25	12.82	166.12
5	White	37,680	176.03	1.02	0.54	2.95	3,352	8.90	179.97
5	EL – Yes	8,533	121.73	1.03	0.55	2.60	692	8.11	123.51
5	EL – No	87,752	166.95	1.02	0.53	2.95	7,749	8.83	170.66
5	EconDis – Yes	32,449	139.79	1.02	0.53	2.95	2,781	8.57	142.63
5	EconDis – No	63,836	174.71	1.02	0.54	2.95	5,660	8.87	178.67
5	SWD – Yes	19,887	138.76	1.02	0.54	2.58	1,662	8.36	141.95
5	SWD – No	76,399	169.24	1.02	0.53	2.95	6,779	8.87	172.89
8	<b>NJSLA–S</b>	<b>101,373</b>	<b>160.84</b>	<b>1.00</b>	<b>0.58</b>	<b>2.69</b>	<b>4,231</b>	<b>4.17</b>	<b>162.31</b>
8	Male	51,908	160.91	1.00	0.58	2.64	2,134	4.11	162.97
8	Female	49,405	160.75	1.00	0.60	2.69	2,094	4.24	161.57
8	Am. Indian	149	162.18	1.01	0.68	1.60	4	2.68	160.75
8	Asian	10,766	190.28	1.00	0.62	2.64	457	4.24	191.11
8	Black	15,042	143.25	1.00	0.60	2.36	600	3.99	144.05
8	Hispanic	31,831	147.46	1.00	0.58	2.69	1,304	4.10	149.71
8	Pacific Islander	191	166.60	0.98	0.71	1.59	3	1.57	139.00
8	White	40,913	169.49	1.00	0.61	2.23	1,733	4.24	170.41
8	EL – Yes	5,908	131.35	1.00	0.62	2.23	243	4.11	132.72
8	EL – No	95,454	162.67	1.00	0.58	2.69	3,988	4.18	164.11

Grade	Group	N	Mean Scale Score	Mean Person Infit	Person Infit Min	Person Infit Max	N Flagged	% Flagged	Flagged Mean Scale Score
8	EconDis – Yes	31,966	145.19	1.00	0.60	2.69	1,331	4.16	145.88
8	EconDis – No	69,392	168.05	1.00	0.53	2.64	2,900	4.18	169.85
8	SWD – Yes	20,268	144.40	1.00	0.60	2.36	818	4.04	145.83
8	SWD – No	81,097	164.95	1.00	0.58	2.69	3,413	4.21	166.26
11	<b>NJSLA–S</b>	<b>92,925</b>	<b>166.18</b>	<b>1.02</b>	<b>0.58</b>	<b>3.10</b>	<b>5,745</b>	<b>6.18</b>	<b>168.39</b>
11	Male	47,026	167.19	1.03	0.58	3.10	2,923	6.22	170.20
11	Female	45,829	165.11	1.02	0.59	2.68	2,818	6.15	166.48
11	Am. Indian	121	156.12	1.04	0.65	1.67	12	9.92	169.00
11	Asian	10,025	208.35	1.03	0.61	2.47	667	6.65	206.61
11	Black	12,436	139.32	1.03	0.60	2.68	751	6.04	141.69
11	Hispanic	27,349	145.66	1.03	0.60	2.63	1,662	6.08	148.65
11	Pacific Islander	202	174.30	1.02	0.72	1.94	8	3.96	167.00
11	White	40,898	177.26	1.02	0.58	3.10	2,520	6.16	179.03
11	EL – Yes	4,830	116.83	1.03	0.65	2.53	281	5.82	117.02
11	EL – No	88,095	168.88	1.02	0.58	3.10	5,464	6.20	171.03
11	EconDis – Yes	24,929	142.99	1.03	0.60	2.63	1,472	5.90	146.83
11	EconDis – No	67,996	174.68	1.02	0.58	3.10	4,273	6.28	175.81
11	SWD – Yes	18,225	140.88	1.03	0.61	2.59	1,095	6.01	145.56
11	SWD – No	74,697	172.35	1.02	0.58	3.10	4,650	6.23	173.76

**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	<b>NJSLA–S</b>	<b>96,288</b>	<b>162.94</b>	<b>1.01</b>	<b>0.15</b>	<b>4.25</b>	<b>10,922</b>	<b>11.34</b>	<b>163.62</b>
5	Male	49,121	164.78	1.01	0.16	4.25	5,717	11.64	164.92
5	Female	47,156	161.02	1.01	0.15	4.25	5,204	11.04	162.18
5	Am. Indian	150	160.18	1.03	0.26	1.67	15	10.00	151.40
5	Asian	10,437	201.35	1.01	0.16	3.07	1,224	11.73	206.94
5	Black	13,848	138.38	1.02	0.23	4.25	1,608	11.61	136.99
5	Hispanic	31,155	144.16	1.02	0.19	3.52	3,579	11.49	142.01
5	Pacific Islander	195	171.83	1.01	0.51	2.29	26	13.33	172.58
5	White	37,680	176.03	1.01	0.15	4.25	4,158	11.04	178.45
5	EL – Yes	8,533	121.73	1.03	0.23	2.9	1,056	12.38	119.95
5	EL – No	87,752	166.95	1.01	0.15	4.25	9,866	11.24	168.29
5	EconDis – Yes	32,449	139.79	1.02	0.19	4.25	3,767	11.61	137.56
5	EconDis – No	63,836	174.71	1.01	0.15	4.25	7,155	11.21	177.34
5	SWD – Yes	19,887	138.76	1.02	0.20	3.27	2,339	11.76	135.94
5	SWD – No	76,399	169.24	1.01	0.15	4.25	8,583	11.23	171.16
8	<b>NJSLA–S</b>	<b>101,373</b>	<b>160.84</b>	<b>1.01</b>	<b>0.18</b>	<b>3.56</b>	<b>7,970</b>	<b>7.86</b>	<b>158.86</b>
8	Male	51,908	160.91	1.02	0.18	3.52	4,154	8.00	159.32
8	Female	49,405	160.75	1.01	0.35	3.56	3,809	7.71	158.33
8	Am. Indian	149	162.18	1.01	0.64	1.71	8	5.37	158.50
8	Asian	10,766	190.28	1.01	0.18	2.34	777	7.22	189.85
8	Black	15,042	143.25	1.02	0.18	2.4	1,250	8.31	142.34
8	Hispanic	31,831	147.46	1.02	0.36	3.52	2,590	8.14	145.85
8	Pacific Islander	191	166.60	1.00	0.72	1.55	12	6.28	163.08
8	White	40,913	169.49	1.01	0.35	3.56	3,126	7.64	167.98
8	EL – Yes	5,908	131.35	1.02	0.48	3.39	512	8.67	130.00
8	EL – No	95,454	162.67	1.01	0.18	3.56	7,457	7.81	160.85
8	EconDis – Yes	31,966	145.19	1.02	0.40	3.52	2,590	8.10	143.50
8	EconDis – No	69,392	168.05	1.01	0.18	3.56	5,378	7.75	166.28

Grade	Group	N	Mean Scale Score	Mean Person Outfit	Person Outfit Min	Person Outfit Max	N Flagged	% Flagged	Flagged Mean Scale Score
8	SWD – Yes	20,268	144.40	1.02	0.18	3.12	1,712	8.45	143.58
8	SWD – No	81,097	164.95	1.01	0.18	3.56	6,257	7.72	163.05
11	<b>NJSLA–S</b>	<b>92,925</b>	<b>166.18</b>	<b>1.03</b>	<b>0.22</b>	<b>3.29</b>	<b>8,706</b>	<b>9.37</b>	<b>162.52</b>
11	Male	47,026	167.19	1.03	0.22	3.13	4,535	9.64	163.98
11	Female	45,829	165.11	1.02	0.30	3.29	4,167	9.09	160.85
11	Am. Indian	121	156.12	1.03	0.73	1.82	11	9.09	174.55
11	Asian	10,025	208.35	1.02	0.44	2.88	900	8.98	208.98
11	Black	12,436	139.32	1.03	0.24	3.02	1,178	9.47	133.42
11	Hispanic	27,349	145.66	1.03	0.22	2.92	2,630	9.62	140.86
11	Pacific Islander	202	174.30	1.01	0.70	1.97	12	5.94	173.17
11	White	40,898	177.26	1.02	0.30	3.29	3,776	9.23	175.18
11	EL – Yes	4,830	116.83	1.04	0.52	2.72	515	10.66	114.40
11	EL – No	88,095	168.88	1.03	0.22	3.29	8,191	9.30	165.55
11	EconDis – Yes	24,929	142.99	1.03	0.24	2.92	2,395	9.61	137.92
11	EconDis – No	67,996	174.68	1.02	0.22	3.29	6,311	9.28	171.86
11	SWD – Yes	18,225	140.88	1.03	0.18	3.29	1,824	10.01	136.21
11	SWD – No	74,697	172.35	1.02	0.18	3.11	6,882	9.21	169.50

**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	77,23	168.70	1.02	0.53	2.95	6,913	8.95	172.13
5	PBT	47	146.72	1.03	0.65	1.73	6	12.77	153.83
5	TTS	16,87	141.59	1.02	0.54	2.95	1,376	8.15	144.42
5	SP	1,204	122.30	1.02	0.55	2.31	86	7.14	123.77
5	SP TTS	561	121.58	1.02	0.66	1.88	30	5.35	124.73
5	HR	308	129.60	1.02	0.62	2.17	27	8.77	128.52
8	CBT	84,15	164.08	1.00	0.58	2.69	3,499	4.16	165.59
8	PBT	56	143.55	0.96	0.71	1.44	2	3.57	142.00
8	TTS	15,04	147.04	1.00	0.60	2.00	634	4.21	148.64
8	SP	1,494	130.38	1.00	0.63	1.90	59	3.95	134.54
8	SP TTS	416	130.79	1.01	0.62	1.71	23	5.53	137.30
8	HR	173	129.89	1.00	0.67	1.60	13	7.51	124.23
11	CBT	85,01	168.73	1.02	0.58	3.10	5,274	6.20	170.64
11	PBT	62	143.39	1.05	0.79	1.68	5	8.06	117.80
11	TTS	6,525	142.76	1.03	0.63	2.36	392	6.01	147.90
11	SP	1,149	118.70	1.02	0.68	2.53	66	5.74	119.18
11	SP TTS	108	121.11	1.00	0.67	1.39	4	3.70	121.25
11	HR	56	109.13	1.02	0.72	1.62	3	5.36	100.00

\*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	77,23	168.70	1.01	0.15	4.25	8,731	11.30	169.98
5	PBT	47	146.72	0.99	0.61	1.42	6	12.77	141.50
5	TTS	16,87	141.59	1.02	0.27	3.27	1,923	11.39	140.63
5	SP	1,204	122.30	1.03	0.28	2.28	153	12.71	120.35
5	SP TTS	561	121.58	1.02	0.23	1.93	61	10.87	114.87
5	HR	308	129.60	1.02	0.52	2.15	45	14.61	131.62
8	CBT	84,15	164.08	1.01	0.18	3.56	6,560	7.80	162.05
8	PBT	56	143.55	1.01	0.63	1.49	8	14.29	158.50
8	TTS	15,04	147.04	1.02	0.18	3.12	1225	8.14	146.02
8	SP	1,494	130.38	1.01	0.55	2.09	121	8.10	129.16
8	SP TTS	416	130.79	1.03	0.65	1.99	41	9.86	129.54
8	HR	173	129.89	1.03	0.52	1.67	13	7.51	133.15
11	CBT	85,01	168.73	1.03	0.22	3.29	7,929	9.33	165.38
11	PBT	62	143.39	1.04	0.64	1.61	6	9.68	158.33
11	TTS	6,525	142.76	1.03	0.39	2.82	645	9.89	136.29
11	SP	1,149	118.70	1.03	0.59	2.15	109	9.49	118.75
11	SP TTS	108	121.11	1.03	0.68	2.32	10	9.26	112.60
11	HR	56	109.13	1.04	0.74	1.49	5	8.93	100.20

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

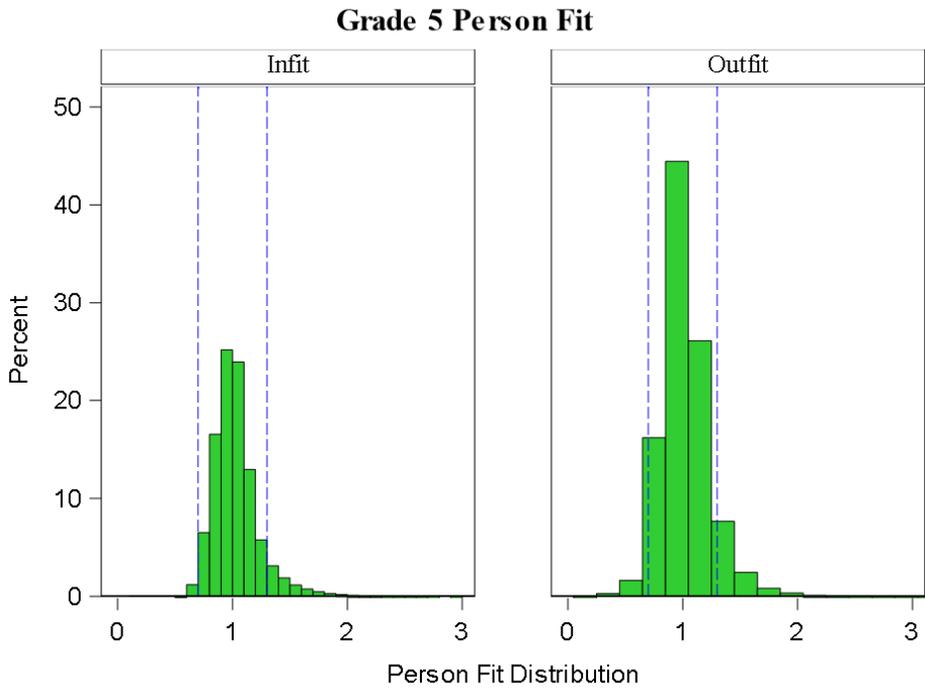


Figure 6.2.4. Grade 5 Person Infit and Outfit Distributions

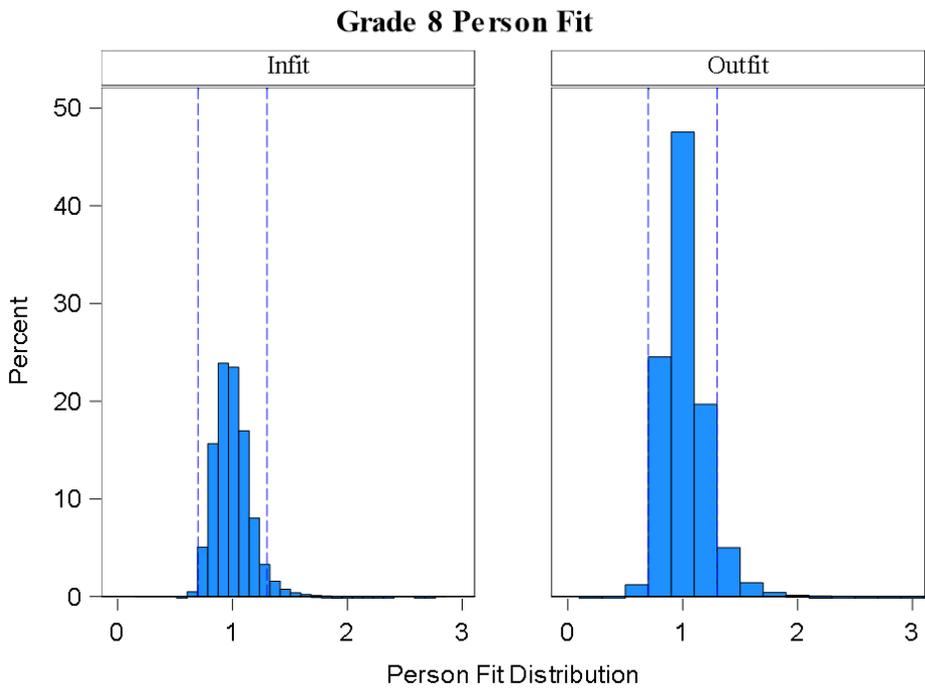


Figure 6.2.5. Grade 8 Person Infit and Outfit Distributions

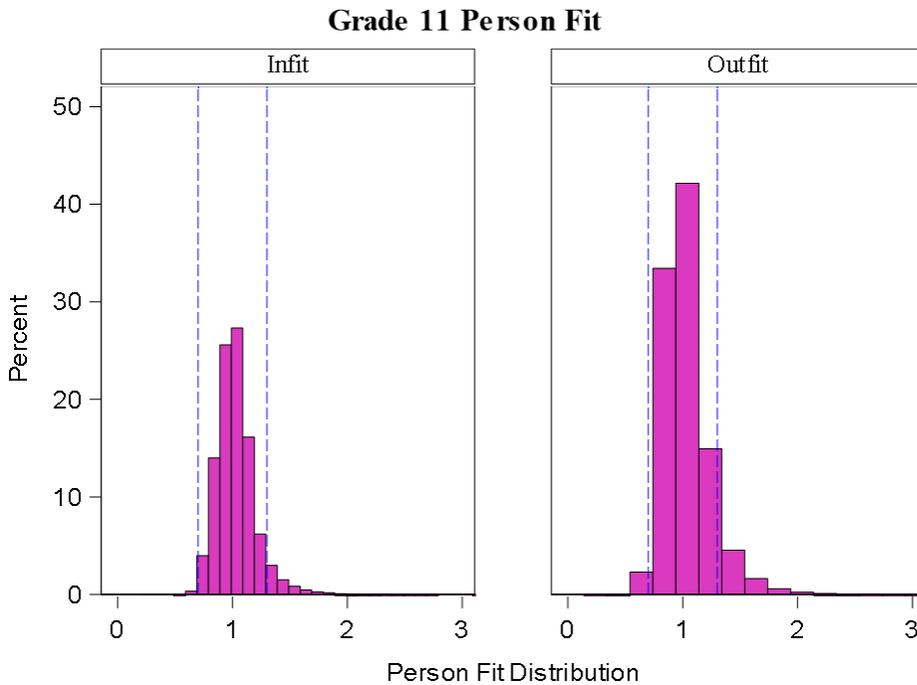


Figure 6.2.6. Grade 11 Person Infit and Outfit Distributions

### 6.2.3 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 that could pose a threat to the validity of inferences made from test scores, 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  $\hat{\theta}_j$  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 6.2.11 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 6.2.11: Summary of Yen’s Q3 Statistics**

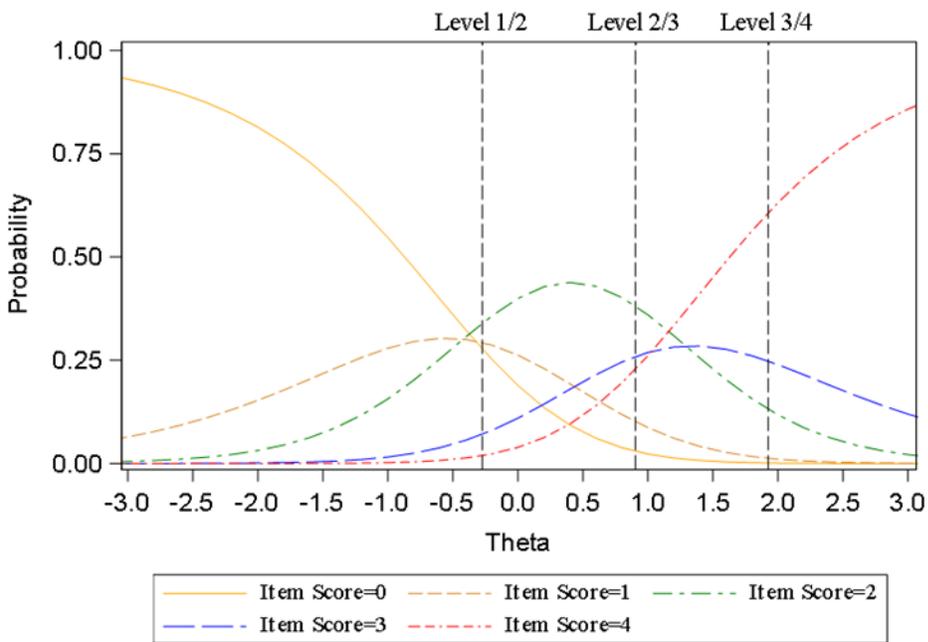
<b>Grade</b>	<b>Mean</b>	<b>Min</b>	<b>Max</b>	<b>Outside <math>-.2</math> to <math>.2</math></b>	<b>% Flagged</b>
5	-.02	-.12	.35	3 out of 1,225	.24
8	-.02	-.13	.25	1 out of 1,770	.06
11	-.01	.11	.24	2 out of 2,346	.09

#### **6.2.4 Item Characteristic Curves – CR Items**

Under IRT, the item characteristic curves (ICC; i.e., 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 0–4 point, hand-scored, constructed-response items are presented in Figures 6.2.7 through 6.2.15 below. The vertical dashed lines represent from left to right the Level 2–4 cut scores on the theta scale. Also, Table 6.2.12 shows the percentages of students receiving each score point for all nine CR items.

**Table 6.2.12: Constructed-Response Point Distribution Percentages**

Grade	Item	%0	%1	%2	%3	%4
5	CR Item 1	33.77	18.06	23.15	12.27	12.76
5	CR Item 2	26.24	25.32	36.84	7.79	3.81
5	CR Item 3	40.53	38.24	14.95	3.10	3.17
8	CR Item 1	37.95	19.03	17.27	15.84	9.90
8	CR Item 2	57.33	18.08	16.34	6.95	1.30
8	CR Item 3	31.94	44.06	12.86	4.24	6.89
11	CR Item 1	42.72	9.04	18.77	7.56	21.92
11	CR Item 2	38.81	22.94	23.28	8.47	6.50
11	CR Item 3	41.16	14.62	24.80	13.90	5.51



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

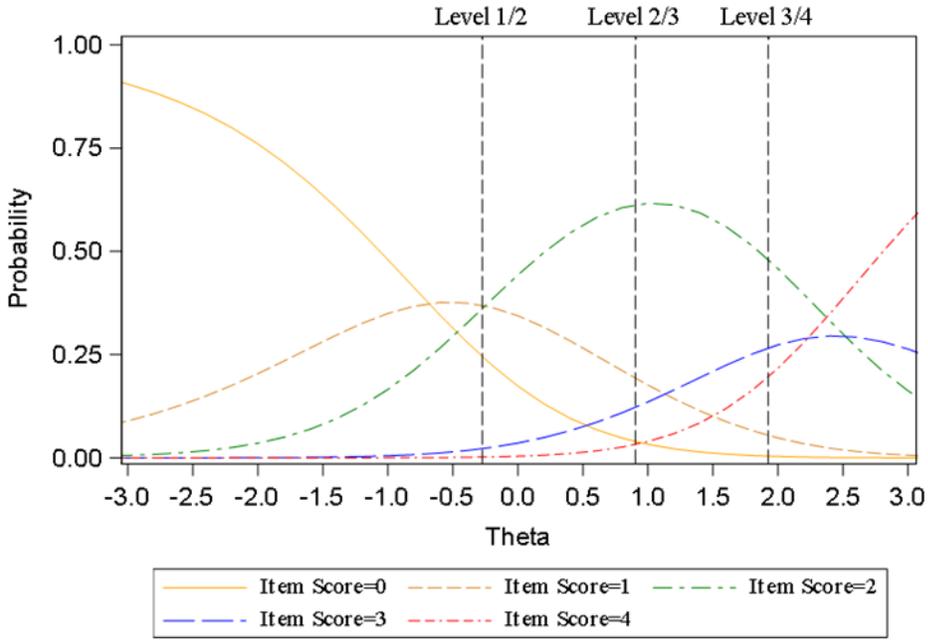


Figure 6.2.8. ICC Plot for Grade 5 Constructed-Response Item 2

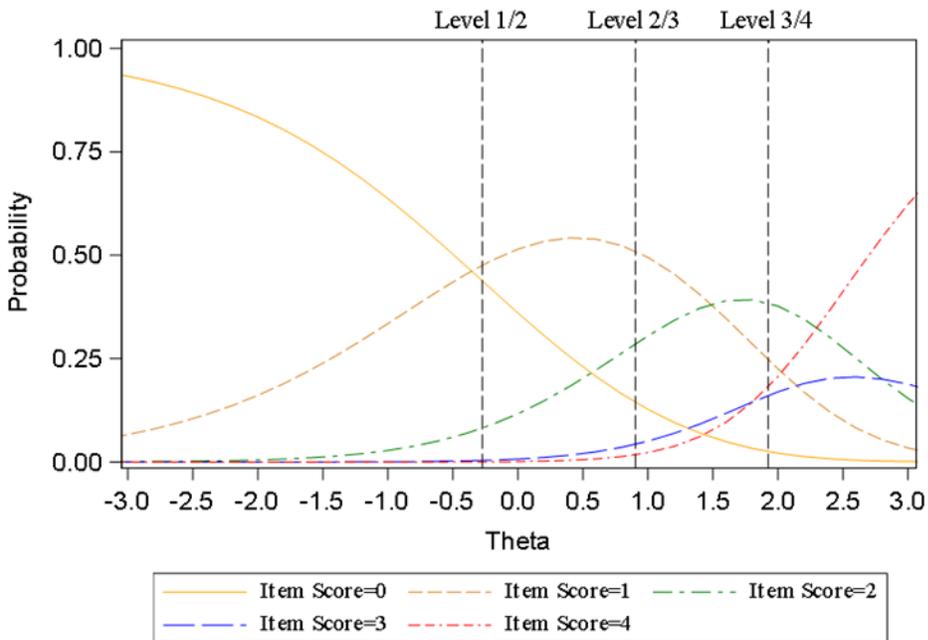


Figure 6.2.9. ICC Plot for Grade 5 Constructed-Response Item 3

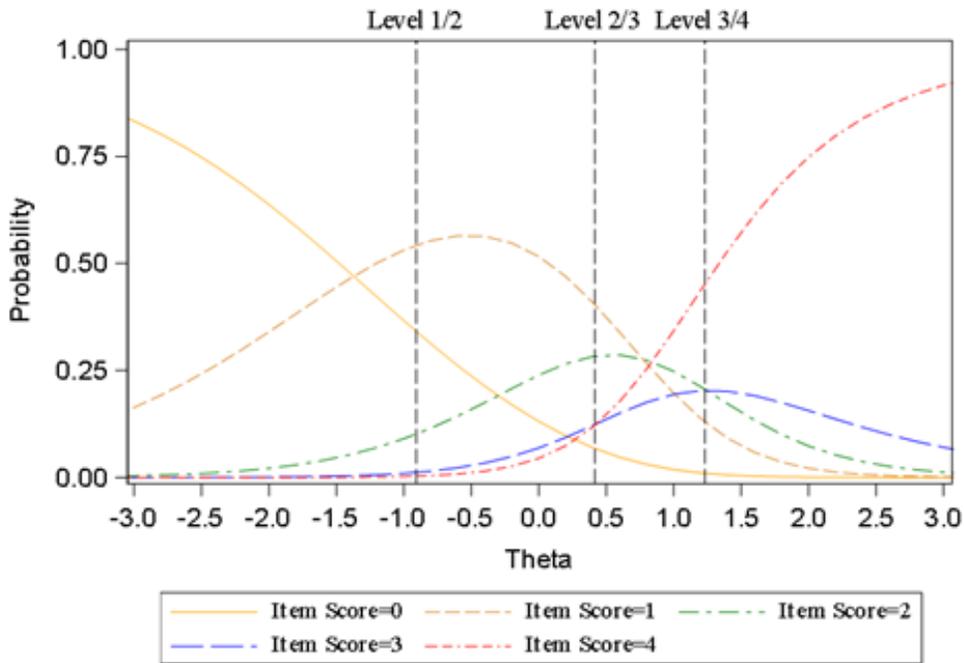


Figure 6.2.10. ICC Plot for Grade 8 Constructed-Response Item 1

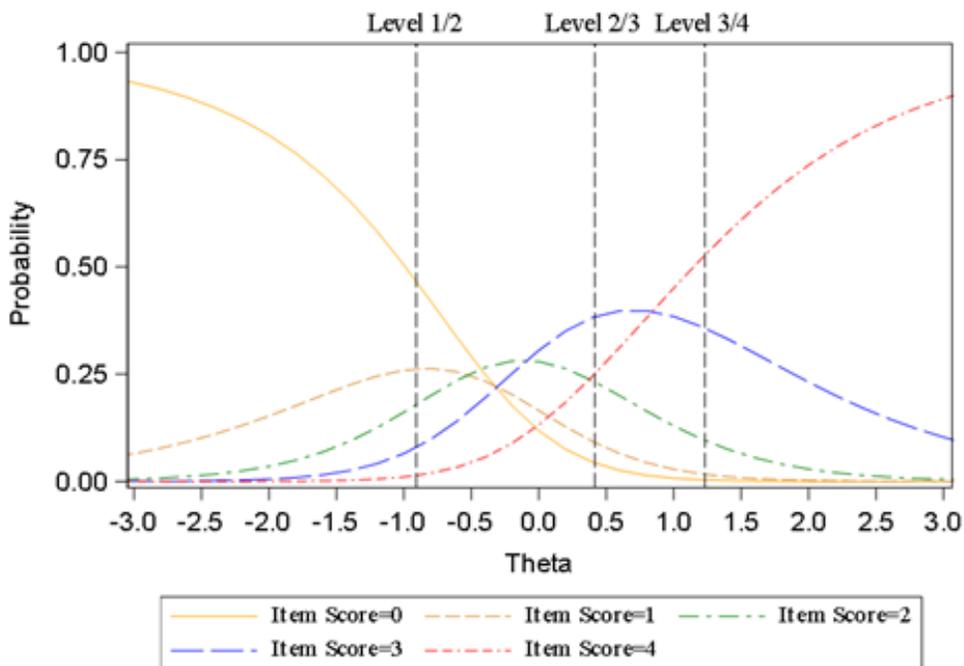


Figure 6.2.11. ICC Plot for Grade 8 Constructed-Response Item 2

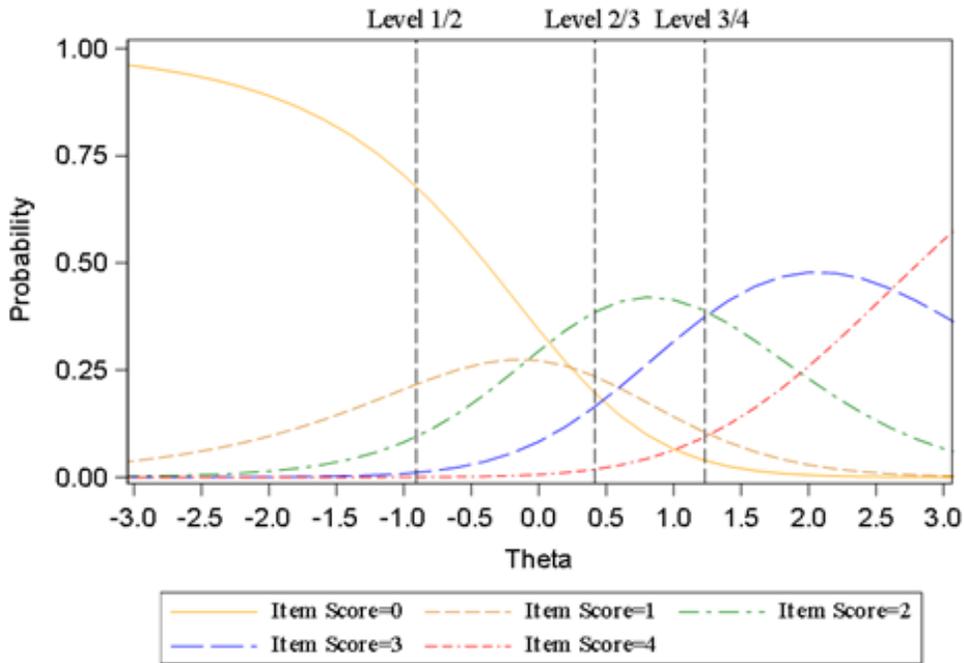


Figure 6.2.12. ICC Plot for Grade 8 Constructed-Response Item 3

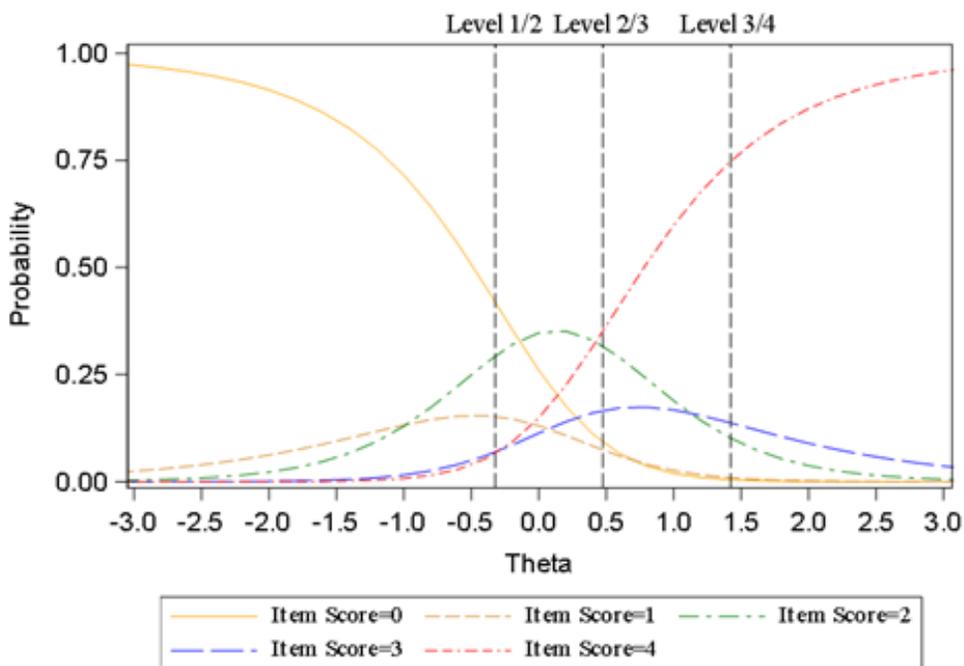


Figure 6.2.13. ICC Plot for Grade 11 Constructed-Response Item 1

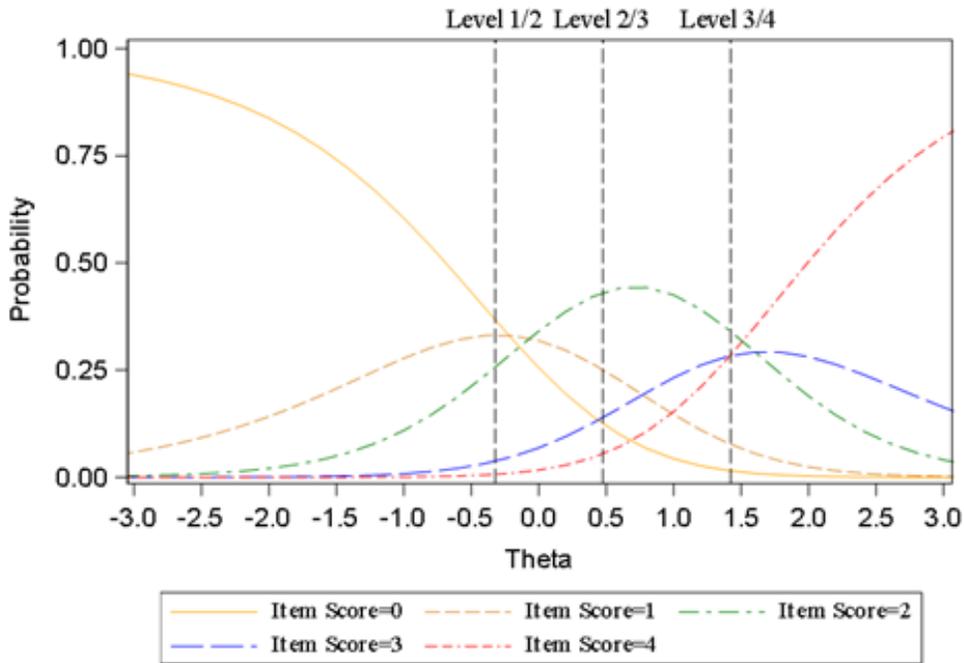


Figure 6.2.14. ICC Plot for Grade 11 Constructed-Response Item 2

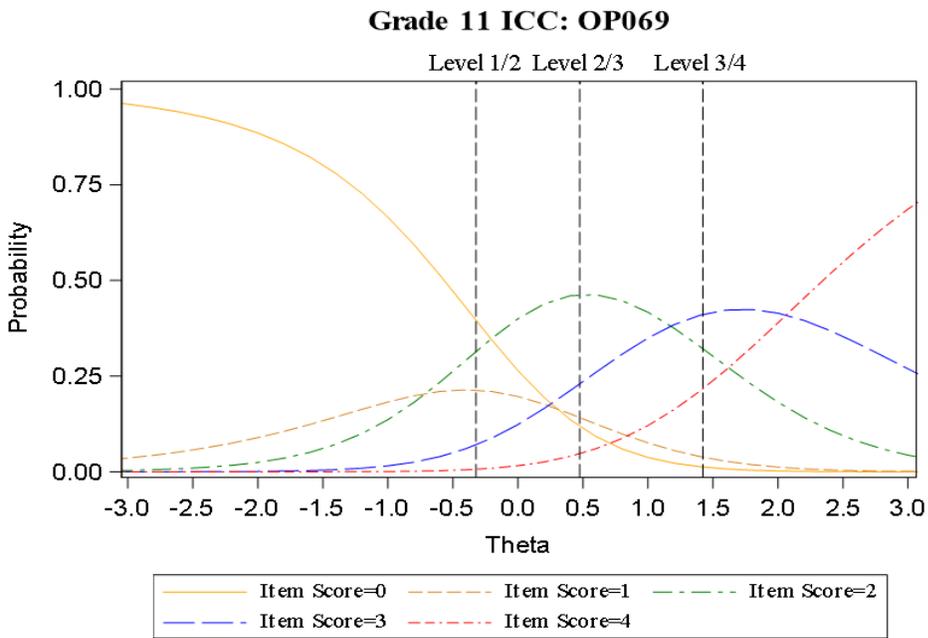


Figure 6.2.15. ICC Plot for Grade 11 Constructed-Response Item 3

### 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 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 2 and 4 cut score ranges are more complex and 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	77,237	168.70	47.19	100	300	36.03	35.10	20.44	8.44
5	PBT	47	146.72	45.71	100	272	59.57	23.40	14.89	2.13
5	TTS	16,877	141.59	43.31	100	300	62.10	24.90	9.72	3.28
5	SP	1,204	122.30	28.24	100	259	81.98	15.53	2.24	0.25
5	SP TTS	561	121.58	25.61	100	226	81.82	17.11	1.07	0.00
5	HR	308	129.60	31.69	100	259	75.65	20.78	2.92	0.65
8	CBT	84,156	164.08	34.08	100	300	36.59	45.91	13.39	4.11
8	PBT	56	143.55	26.54	100	215	67.86	25.00	7.14	0.00
8	TTS	15,049	147.04	29.75	100	286	59.09	33.80	5.83	1.28
8	SP	1,494	130.38	19.86	100	220	82.40	17.27	0.33	0.00
8	SP TTS	416	130.79	19.93	100	210	83.17	16.35	0.48	0.00
8	HR	173	129.89	19.12	100	210	87.28	11.56	1.16	0.00
11	CBT	85,010	168.73	53.37	100	300	44.03	25.42	21.54	9.01
11	PBT	62	143.39	49.79	100	297	64.52	22.58	4.84	8.06
11	TTS	6,525	142.76	45.19	100	300	65.18	20.49	11.52	2.80
11	SP	1,149	118.70	24.76	100	291	89.56	9.40	0.96	0.09
11	SP TTS	108	121.11	28.77	100	254	87.96	9.26	1.85	0.93
11	HR	56	109.13	17.20	100	165	94.64	5.36	0.00	0.00

\* 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.

### 6.3.3 Subscore Proficiency Classification

There are no scale scores for the various subscores. As is explained in Part 7, 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, Life Science saw the highest percentage (59.59%) of students classified as Below Expectations; among the three practices, Sensemaking saw the highest percentage (60.02%) of students classified as Below Expectations. Grade 8 had the same pattern as grade 5. At grade 8, the highest percentages of students classified as Below Expectations were observed for Life Science (70.06%) and Sensemaking (71.23%). At grade 11, the Below Expectations percentages varied from 53.58% for Earth and Space Science to 56.33% for Physical Science, while the Below Expectations percentages were relatively consistent across the practices with 54.79% for Sensemaking to 56.00% for Critiquing. Overall, for each content domain and practice at all grade levels, there was more variance in the Below Expectations percentages of students within each demographic group.

## 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 and 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 2022 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 2022 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, Rotou, & Twing, 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 the evaluation methods or the number of dropped items reaches 20% of the original set of anchor items. For the 2022 NJSLA–S equating, at grades 5 and 8, no anchor items were dropped from the anchor item sets. At grade 11, the anchor evaluation processes were terminated at round 5 and yielded four (14%) 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 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 as follows.

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 state-wide assessments; in the past, with only three performance levels, scale scores of 200 and 250 represented proficient and advanced proficient performance, respectively (NJDOE, 2017). 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 2 and Level 3 cut scores be anchored during the linear transformation and at grade 11, the Level 3 and Level 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 3). 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 ( $\theta$ ) 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  $\theta_1$  and  $\theta_2$  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  $\theta_1$  and  $\theta_2$ , 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	-0.27392	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 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 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 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 2 and 3 for grade 11 as well as Levels 2, 3, and 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 (Level 3 cut score 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 (2017) 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. If this assumption is violated for any accommodative form(s), special equatings performed for those forms are described in the following section. For the 2022 NJSLA–S administration, no special equating was needed and conducted.

### 7.2.1 Special Equatings

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 report student performance in three content domains/disciplinary core ideas (DCI) including Earth and Space, Life, and Physical, and 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 ( $\theta$ ), 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  $\theta^*$ ). Based on the subscore table, the CSEM associated with  $\theta^*$  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  $\theta$  is at or above ( $\theta^* + 1.5 CSEM^*$ ) units.
- A raw subscore is classified as “Below” if its associated  $\theta$  is below ( $\theta^* - 1.5 CSEM^*$ ) units.
- A raw subscore is classified as “Near/Met” if its associated  $\theta$  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 ( $\sigma_X^2$ ) equals the sum of true score variance ( $\sigma_T^2$ ) and error variance ( $\sigma_e^2$ ) 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 ( $\rho_{X, X_s}$ ) is defined as the correlation between observed scores ( $X_1, X_2$ ) on parallel forms, which is equal to true score variance ( $\sigma_T^2$ ) divided by observed score variance ( $\sigma_X^2$ ):

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

With just a few algebraic steps, the CTT definition of the standard error of measurement (SEM,  $\sigma_e$ ) 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,  $\sigma_i^2$  is the item score variance of item  $i$ , and  $\sigma_X^2$  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. 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 .83 to .93. 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 towards the low end of the ability spectrum, reliability ranged from .69 to .91. The relatively low-reliability measures for the Spanish, Spanish TTS, and Human Reader forms are

due to those populations doing very poorly on the test, which limits the amounts of variance in test scores. The grade 11 alpha coefficients ranged from .69 for the Human Reader form to .94 for the CBT form. However, there were only 56 Human Reader form test takers who did poorly on the test.

**Table 8.1.1: Coefficient Alpha and SEM by Form**

Grade	Form*	N	Mean Raw Score	SD	Alpha	SEM
5	CBT	77,237	27.59	13.08	.93	3.47
5	PBT	47	21.43	12.88	.93	3.41
5	TTS	16,877	19.86	12.33	.93	3.29
5	SP	1,204	14.40	8.34	.86	3.08
5	SP TTS	561	14.30	7.50	.83	3.11
5	HR	308	16.64	9.09	.87	3.22
8	CBT	84,156	24.60	12.58	.91	3.82
8	PBT	56	17.32	9.04	.85	3.50
8	TTS	15,049	18.59	10.34	.88	3.56
8	SP	1,494	13.09	5.76	.69	3.19
8	SP TTS	416	13.19	5.85	.70	3.21
8	HR	173	12.87	5.80	.72	3.09
11	CBT	85,010	31.58	16.38	.94	4.08
11	PBT	62	23.82	15.08	.93	3.97
11	TTS	6,525	23.59	13.99	.92	3.85
11	SP	1,149	16.18	7.70	.80	3.40
11	SP TTS	108	16.86	9.05	.86	3.42
11	HR	56	12.82	5.71	.69	3.19

\*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 .74 was for Critiquing at grade 8. For content domains, the lowest subscore reliability of .71 was for Earth and Space at grade 8, which had the least number of items among the content domains for grade 8.

**Table 8.1.2: Coefficient Alpha and SEM by Reporting Category**

Grade	Reporting Category	Total # Items	#MC Items	#TE1 Items	#TE2 Item	#CR Items	Max Points	Alpha	SEM
5	<b>Total</b>	50	15	31	1	3	60	.93	3.44
5	Earth and Space	15	5	8	1	1	19	.82	1.83
5	Life	19	4	14	0	1	22	.84	2.18
5	Physical	16	6	9	0	1	19	.81	1.92
5	Sensemaking	25	8	16	0	1	28	.88	2.26
5	Critiquing	12	5	5	0	2	18	.79	2.03
5	Investigating	13	2	10	1	0	14	.76	1.60
8	<b>Total</b>	60	22	32	3	3	72	.91	3.78
8	Earth and Space	17	5	9	2	1	22	.71	2.13
8	Life	23	12	9	1	1	27	.80	2.25
8	Physical	20	5	14	0	1	23	.77	2.17
8	Sensemaking	25	8	14	1	2	32	.79	2.63
8	Critiquing	16	4	11	0	1	19	.74	1.85
8	Investigating	19	10	7	2	0	21	.76	2.00
11	<b>Total</b>	69	28	38	0	3	78	.94	4.06
11	Earth and Space	21	6	14	0	1	24	.81	2.40
11	Life	23	10	12	0	1	26	.87	2.25
11	Physical	25	12	12	0	1	28	.82	2.38
11	Sensemaking	32	10	21	0	1	35	.88	2.61
11	Critiquing	18	6	10	0	2	24	.82	2.47
11	Investigating	19	12	7	0	0	19	.77	1.88

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 .86, 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 ( $\alpha_{EL-Yes} = .72$ ). The pattern for grade 11 was the same as for grade 5. The coefficient alpha values for all groups at grade 11 were above .91 except for the English learners ( $\alpha_{EL-Yes} = .82$ ).

**Table 8.1.3: Coefficient Alpha and SEM by Demographic Group**

Grade	Group	N	Mean	SD	Alpha	SEM
5	NJSLA-S	96,288	25.95	13.31	.93	3.44
5	Male	49,121	26.44	13.63	.94	3.40
5	Female	47,156	25.45	12.94	.93	3.47
5	Am. Indian	150	25.29	13.62	.94	3.46
5	Asian	10,437	36.45	12.21	.92	3.44
5	Black	13,848	19.01	11.26	.92	3.27
5	Hispanic	31,155	20.71	11.55	.92	3.35
5	Pacific Islander	195	28.50	12.81	.93	3.49
5	White	37,680	29.70	12.40	.92	3.48
5	EL - Yes	8,533	14.20	8.35	.86	3.09
5	EL - No	87,752	27.10	13.14	.93	3.46
5	EconDis - Yes	63,836	29.26	13.11	.93	3.47
5	EconDis - No	32,449	19.46	11.11	.91	3.31
5	SWD - Yes	19,887	19.02	12.20	.93	3.23
5	SWD - No	76,399	27.76	12.99	.93	3.47
8	NJSLA-S	101,371	23.46	12.44	.91	3.78
8	Male	51,908	23.58	13.03	.92	3.78
8	Female	49,405	23.33	11.80	.90	3.78
8	Am. Indian	149	24.06	12.99	.91	3.86
8	Asian	10,766	34.47	13.49	.91	3.97
8	Black	15,042	17.20	8.99	.85	3.47
8	Hispanic	31,831	18.58	9.51	.86	3.55
8	Pacific Islander	191	25.48	12.31	.90	3.86
8	White	40,913	26.49	12.20	.90	3.87
8	EL-Yes	5,908	13.36	6.05	.72	3.20
8	EL-No	95,454	24.09	12.47	.91	3.80
8	EconDis-Yes	69,392	26.08	12.92	.91	3.86
8	EconDis-No	31,966	17.79	9.02	.85	3.50
8	SWD-Yes	20,268	17.70	10.10	.88	3.50
8	SWD-No	81,097	24.90	12.56	.91	3.83
11	NJSLA-S	92,925	30.79	16.36	.94	4.06
11	Male	47,026	31.07	17.10	.94	4.03
11	Female	45,829	30.50	15.55	.93	4.09
11	Am. Indian	121	27.63	15.72	.94	3.97
11	Asian	10,025	43.67	16.46	.94	4.08
11	Black	12,436	22.50	12.92	.91	3.83
11	Hispanic	27,349	24.48	13.64	.92	3.91
11	Pacific Islander	202	33.28	16.05	.93	4.10

Grade	Group	N	Mean	SD	Alpha	SEM
11	White	40,898	34.23	15.82	.93	4.11
11	EL–Yes	4,830	15.49	7.99	.82	3.39
11	EL–No	88,095	31.63	16.28	.94	4.08
11	EconDis–Yes	67,996	33.41	16.55	.94	4.10
11	EconDis–No	24,929	23.66	13.43	.92	3.89
11	SWD–Yes	18,225	22.96	14.42	.93	3.81
11	SWD–No	74,697	32.70	16.23	.94	4.10

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 increases, 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	15	15	7.80	3.49	.75	1.73
5	TE	32	33	14.35	8.04	.91	2.38
5	CR	3	12	3.80	2.78	.72	1.48
8	MC	22	22	8.58	4.11	.74	2.12
8	TE	35	36	11.61	6.53	.85	2.52
8	CR	3	12	3.28	2.88	.74	1.48
11	MC	28	28	12.20	5.56	.82	2.36
11	TE	38	38	14.54	8.49	.91	2.50
11	CR	3	12	4.06	3.40	.76	1.67

## 8.2 Item Response Theory Reliability

The reliability of the scale scores ascertained from the partial credit model (PCM; Masters, 1982) discussed in Section 6.1 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 2022 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 specifically at the most important decision points (i.e., the Level 2–4 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  $\theta$ ,  $I_i(\theta)$  is the amount of information for item  $i$  at an ability level of  $\theta$ , 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, ...,  $m_i$ , 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  $\theta$  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 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 2 and 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 3 cut score, and whether there was a precipitous drop in information at the Level 2 and 4 cut scores. At grade 5 the TIF peaked in between the Level 2 and 3 cut scores. There was a large drop in information at the Level 4 cut. At grade 8 the TIF peaked almost directly on the Level 3 cut score; there was a large drop in information at the Level 2 cut. Similarly, the grade 11 TIF peaked almost directly at the Level 3 cut score. However, there was a large drop in information at the Level 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 11 need more information around the Level 4 cut score on future tests.

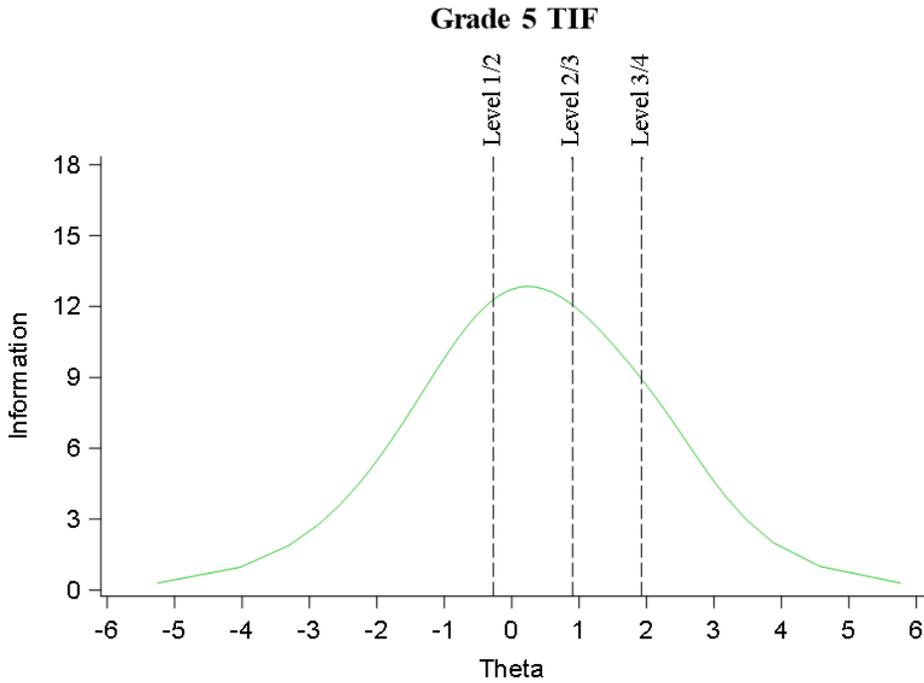


Figure 8.2.1. Grade 5 Test Information Function

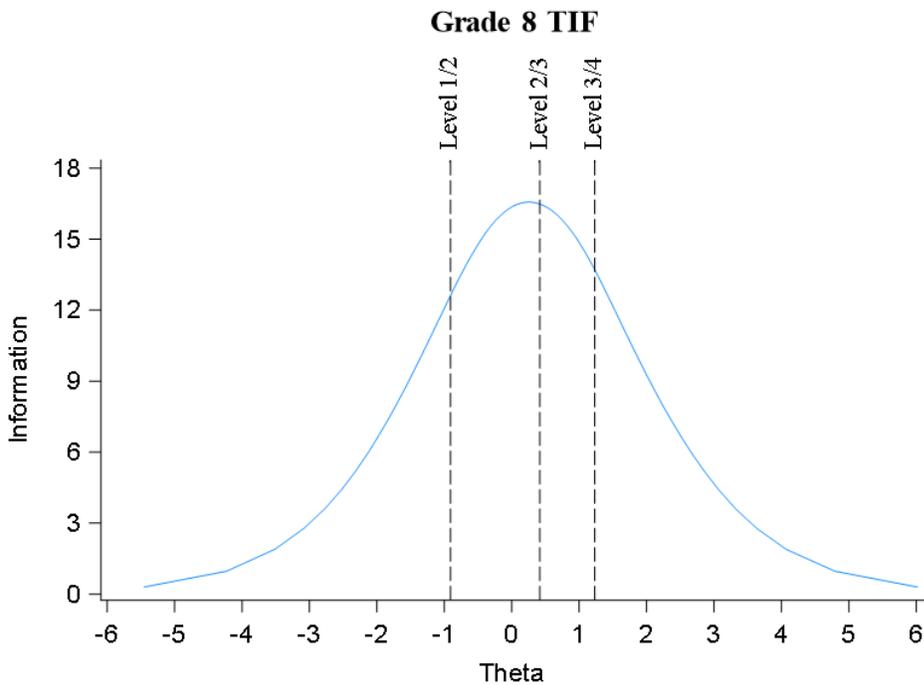


Figure 8.2.2. Grade 8 Test Information Function

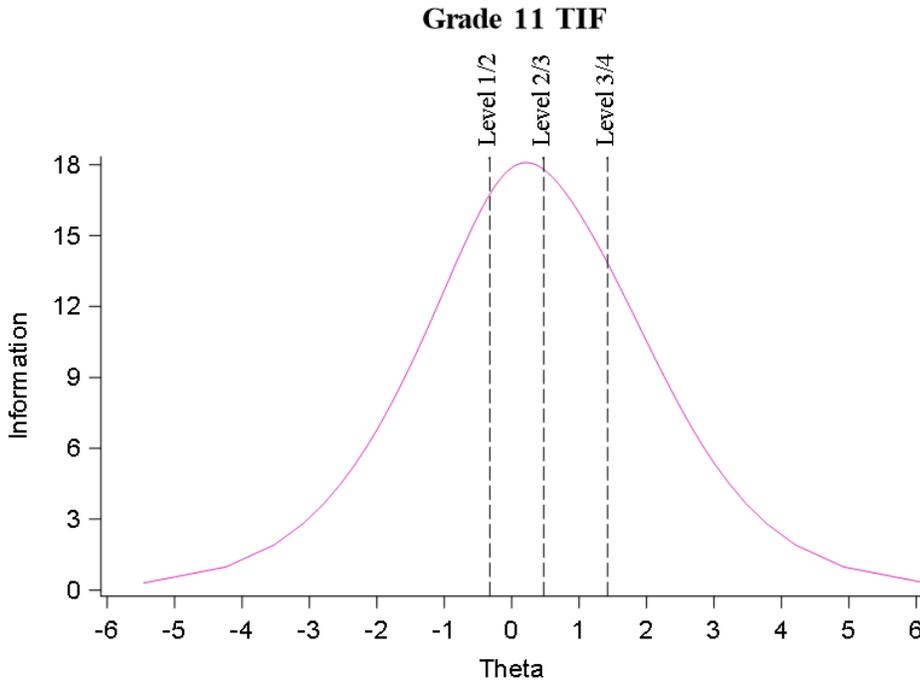


Figure 8.2.3. Grade 11 Test Information Function

### 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  $\theta$ . 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  $\theta$  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 8.2.5 through 8.2.7 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 Level 2, 3, and 4 cut scores are about .25 for each grade level, indicating ability scores around the three cut scores are estimated with precision.

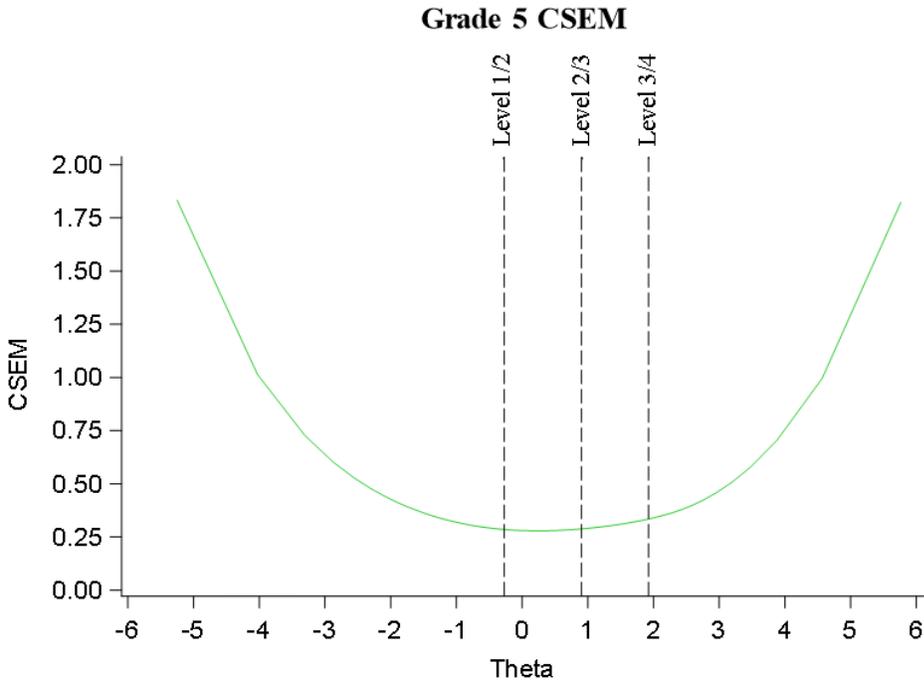


Figure 8.2.4. Grade 5 Conditional Standard Error of Measurement

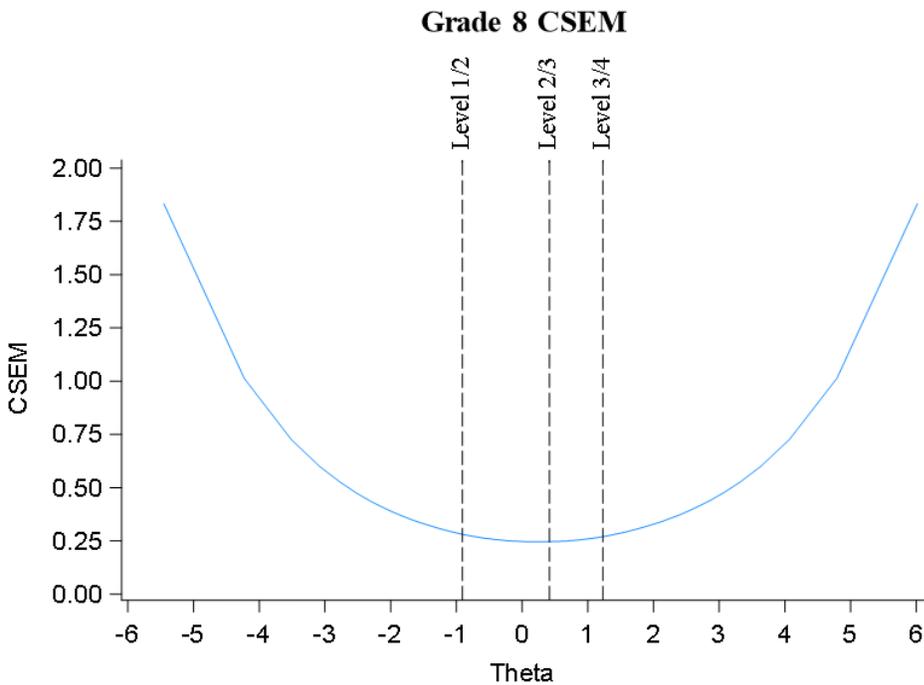


Figure 8.2.5. Grade 8 Conditional Standard Error of Measurement

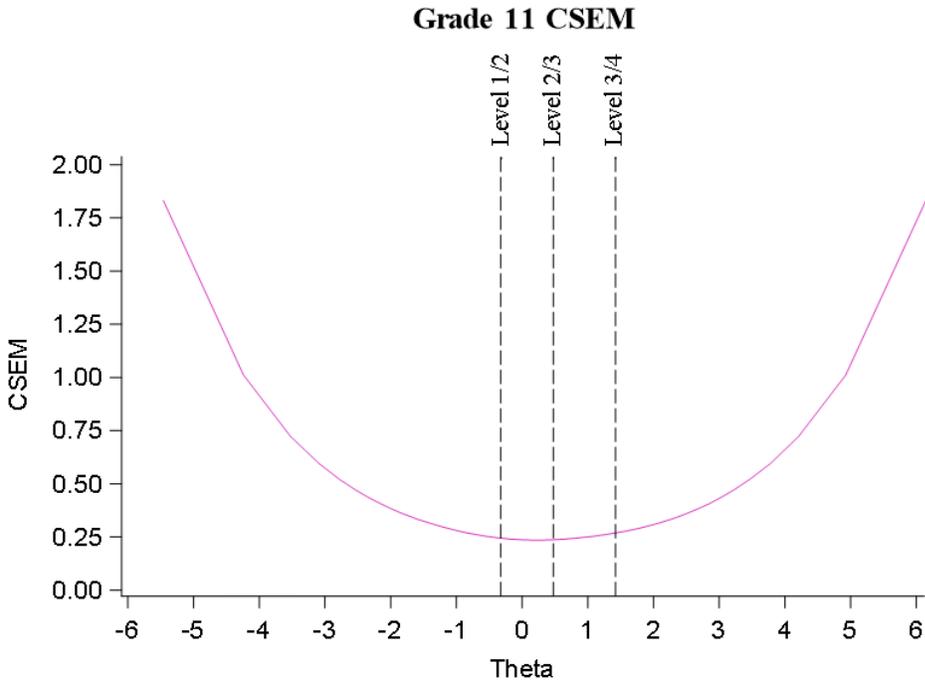


Figure 8.2.6. Grade 11 Conditional Standard Error of Measurement

### 8.2.3 Item Maps

An item map exhibits the distribution of person ability measures 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 2022 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 at or near the Level 3 cut score; the theta distributions peaked around the Level 2 cut score. At grade 5, there were few students above the Level 3 cut score and very few 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. The grade 11 item difficulty and student ability distributions aligned at the decision points on the scale very well.

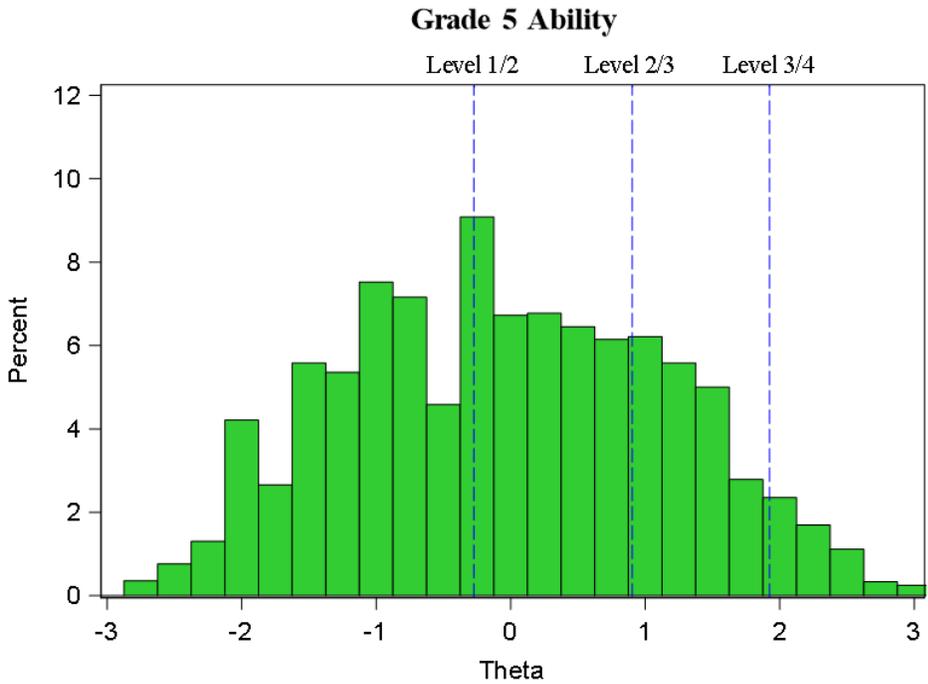
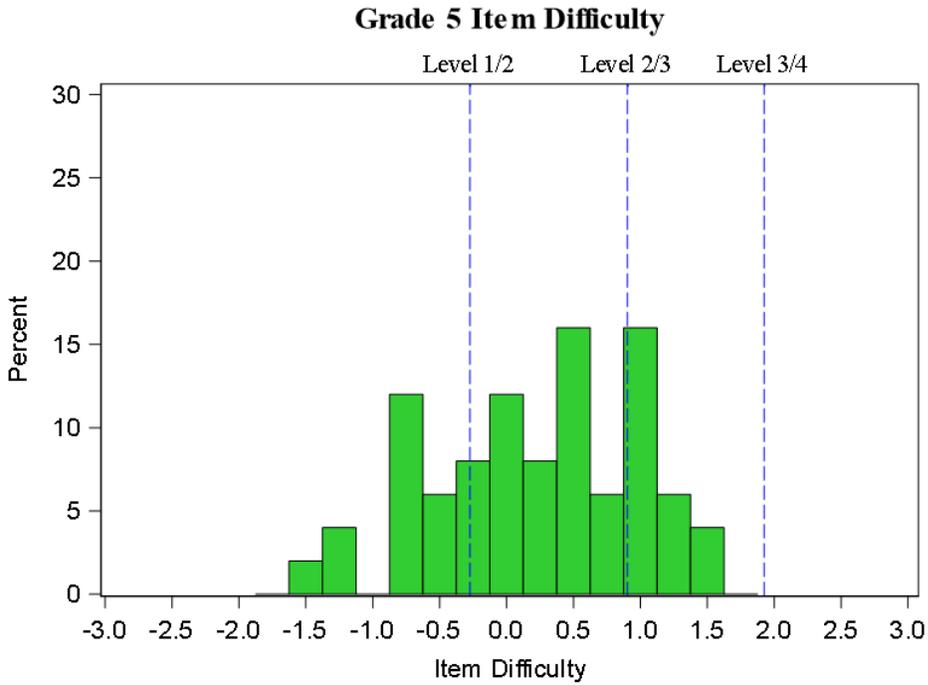


Figure 8.2.7 Grade 5 Item Difficulty and Student Ability Distributions

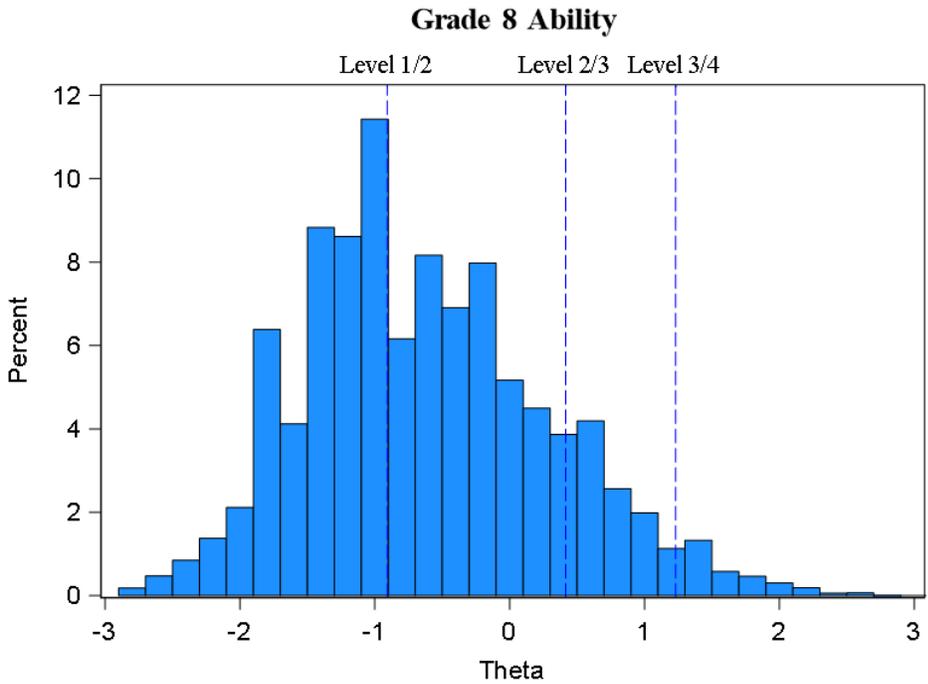
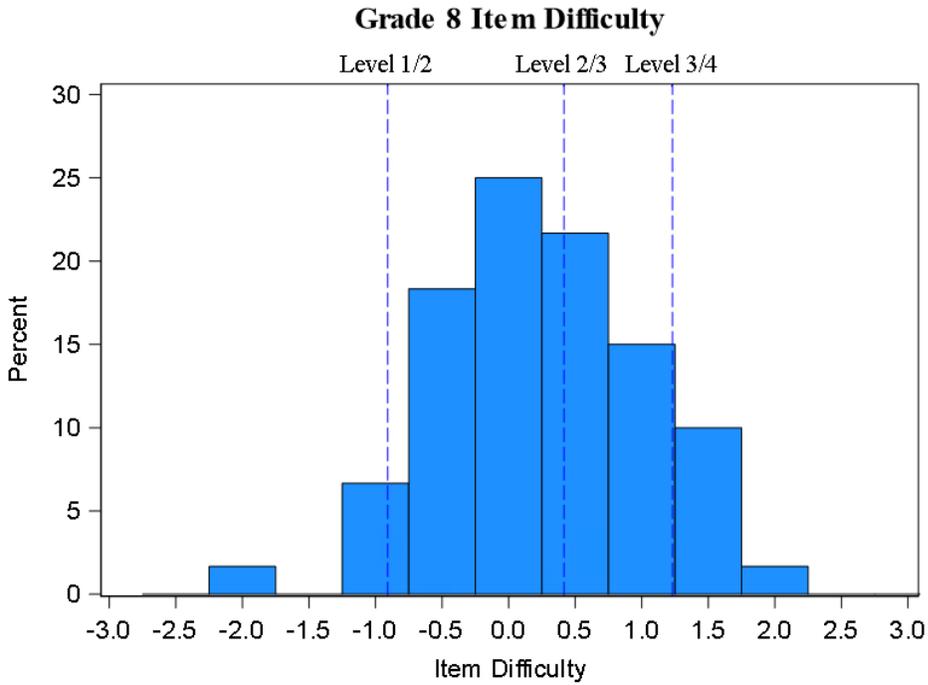


Figure 8.2.8. Grade 8 Item Difficulty and Student Ability Distributions

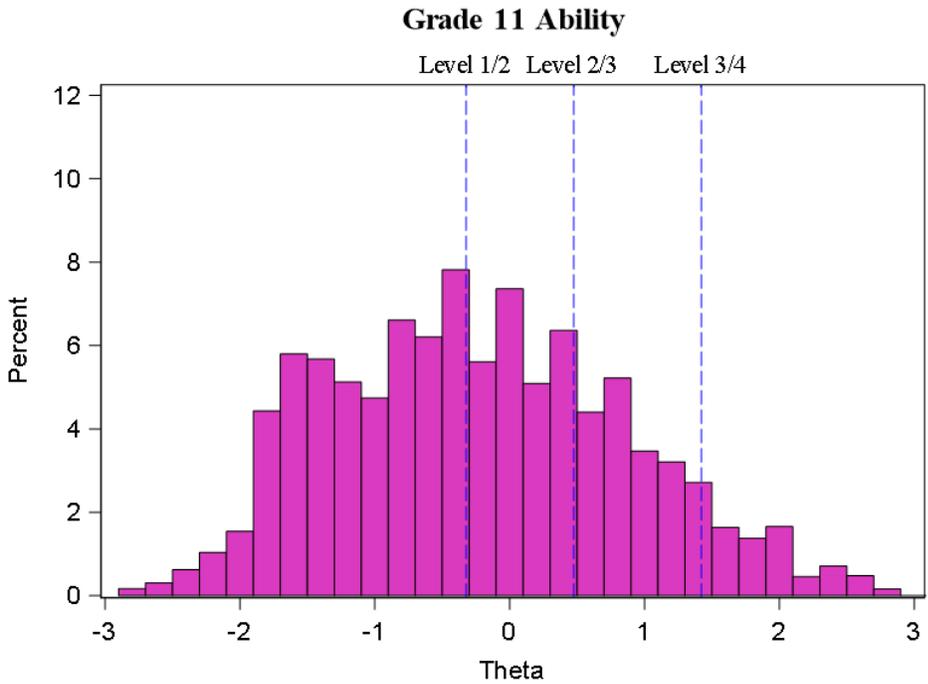
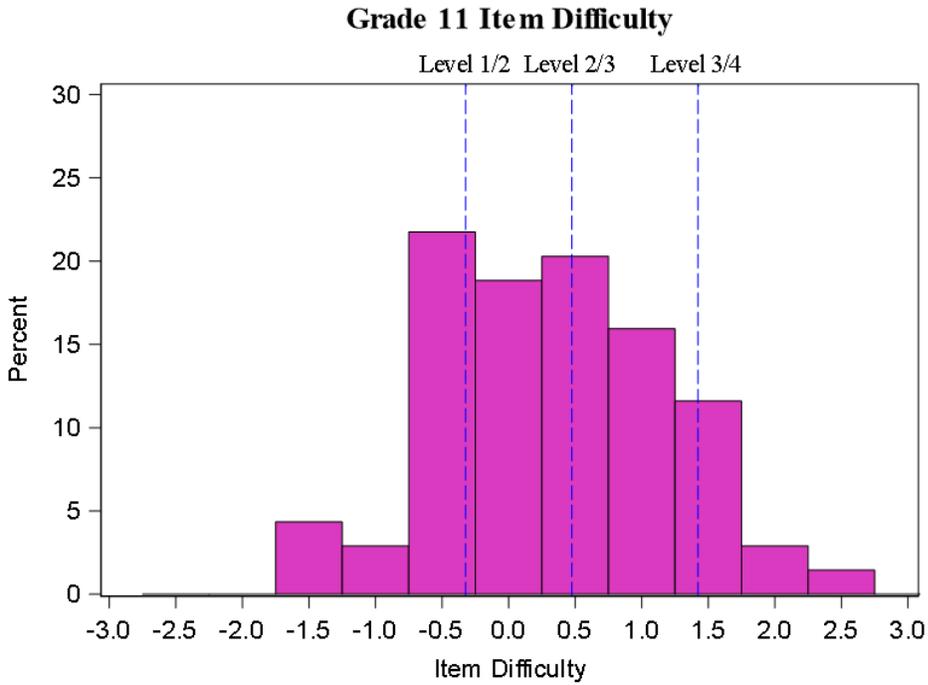


Figure 8.2.9. Grade 11 Item Difficulty and Student Ability Distributions

### 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 2022 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 ( $\theta$ ) scale. Also, the CSEM can be converted and placed on reported scales as needed and appropriate.

The 2022 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 grade 5 the Level 3 cut score’s CSEM was slightly higher than at Level 2, meaning that there was slightly less error in the scale score at 150 than at 200. At grades 8 and 11 the cut score with the least amount of error was the Level 3 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. Any overlap between the upper or lower bounds and one of the other cut scores could indicate reliability problems among the performance-level classifications. For all three cut scores at each grade level, there was no overlap between either their upper or lower boundaries and another cut score.

**Table 8.3.1: Cut Scores with Conditional Standard Error of Measurement**

Grade	Level	Cut Scale Score	CSEM	Lower Bound	Upper Bound
5	Level 2	150	12.1	125.7	174.3
5	Level 3	200	12.2	175.5	224.5
5	Level 4	243	13.9	215.1	270.9
8	Level 2	150	10.7	128.7	171.3
8	Level 3	200	9.3	181.4	218.6
8	Level 4	231	10.2	210.6	251.4
11	Level 2	158	13.0	132.0	184.0
11	Level 3	200	12.5	175.0	225.0
11	Level 4	250	14.2	221.6	278.4

### 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:

$$K = \frac{P_o - P_c}{1 - P_c}, \quad \text{Equation 8.11}$$

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 .75 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 3 cut score because it demarcates the point along the scale where students are deemed proficient or not. The decision consistency at the Level 3 cut score or above was remarkable at .84 to .87, indicating an 84% to 87% probability of being correctly classified as Level 3 or above. The overall NJSLA-S performance classification should be interpreted as being consistent.

**Table 8.3.2: Performance Level Classification Consistency**

Grade	Level 2 Cut	Level 3 Cut	Level 4 Cut	Kappa	$P_o$	$P_o$ for Level 3 or above
5	150	200	243	.63	.75	.84
8	150	200	231	.63	.77	.87
11	158	200	250	.63	.75	.85

### 8.4 Reliability of Subscore Performance Classifications

The methodology used to create the subscore performance level classifications was dependent on the CSEMs in each subscore's 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 <sub>0</sub>	Level	Raw Subscore Cut	Theta	CSEM
5	Earth and Space	.52	.74	Near/Met	9	0.379	0.519
				Above	14	1.787	0.551
5	Life	.56	.76	Near/Met	12	0.393	0.445
				Above	18	1.772	0.557
5	Physical	.53	.74	Near/Met	9	0.237	0.496
				Above	15	1.930	0.598
5	Investigating	.50	.73	Near/Met	6	0.317	0.572
				Above	11	2.080	0.678
5	Sensemaking	.60	.79	Near/Met	15	0.392	0.416
				Above	22	1.761	0.484
5	Critiquing	.48	.70	Near/Met	9	0.280	0.502
				Above	14	1.711	0.598
8	Earth and Space	.48	.74	Near/Met	9	-0.241	0.464
				Above	16	1.195	0.478
8	Life	.59	.81	Near/Met	11	-0.089	0.405
				Above	18	1.074	0.430
8	Physical	.55	.79	Near/Met	9	-0.175	0.435
				Above	16	1.213	0.481
8	Investigating	.54	.77	Near/Met	9	-0.272	0.465
				Above	16	1.336	0.534
8	Sensemaking	.56	.81	Near/Met	14	-0.039	0.370
				Above	22	1.071	0.391
8	Critiquing	.51	.75	Near/Met	6	-0.257	0.499
				Above	12	1.125	0.497
11	Earth and Space	.51	.72	Near/Met	10	0.025	0.405
				Above	17	1.251	0.462
11	Life	.61	.77	Near/Met	12	0.002	0.417
				Above	19	1.289	0.460
11	Physical	.57	.75	Near/Met	11	-0.071	0.408
				Above	18	1.109	0.428
11	Investigating	.52	.73	Near/Met	8	-0.061	0.500
				Above	14	1.499	0.557
11	Sensemaking	.62	.78	Near/Met	15	-0.015	0.363
				Above	23	1.048	0.377
11	Critiquing	.56	.74	Near/Met	10	0.015	0.399
				Above	17	1.246	0.469

## 8.5 Rater Reliability

For constructed-response (CR) items, raters used item-specific scoring rubrics with a score range of 0 to 4. 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 a second read was non-adjacent, then the scores for the response were erased and the paper was re-scored. Thus, all scores in the 10% read-behinds 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. The latter is a variation of Cohen's Kappa that is specifically designed for ordinal variables. As shown in Table 8.5.1, the exact agreement rates ranged from 66.7% to 71.8% for grade 5, from 69.1% to 85.7% for grade 8, and from 72.5% to 78.9% for grade 11. While there was only one grade 5 CR item that had weighted Kappa above .90, all the grades 8 and 11 CR items had weighted Kappa above .90. Overall, rater agreement on the NJSLA-S 0–4 point 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	66.7	33.3	.91
5	CR 2	71.8	28.2	.87
5	CR 3	67.6	32.4	.81
8	CR 1	78.4	21.6	.91
8	CR 2	69.1	30.9	.92
8	CR 3	85.7	14.3	.94
11	CR 1	72.5	27.5	.91
11	CR 2	78.9	21.1	.96
11	CR 3	73.7	26.3	.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 NJSLA–S, and the tests must fulfill the requirements under ESSA (2015).

Adequate representation of the content domains defined in the NJSLA–S is assured through the use of a test blueprint and a responsible test construction process as was described in Part 2. The NJSLA–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 three points of being perfectly balanced. The scientific practices (i.e., Investigating, Sensemaking, and Critiquing) were less balanced. At each grade level, the Sensemaking scientific practice was over-represented and Critiquing and Investigating were under-represented. 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 percent 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 Part 5.2.1: Performance Level Descriptors, 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 in order 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 a given 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, Hoole & Lottridge, 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 we are interested in measuring, which introduces construct irrelevant variance (Thelk, Hoole, Lottridge & Finney 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 NJ 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. Data from this study are currently being analyzed and results will be outlined in a technical report for the NJDOE.

### 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, unexpected patterns of correlations provide 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 .75 at grade 8 between Life science and Earth and Space science. Among the scientific practices, the lowest correlation was .75 at grade 8 between the Investigating and Critiquing scientific practice categories.

### **9.3.2 Other Internal Structure Evidence**

Evidence of the internal structure of the NJSLA–S was also presented via a principal component analysis (PCA). Its 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 a measure of internal consistency that provides a sign as to whether the internal structure of the NJSLA–S is unidimensional. The grade level reliability coefficients presented in Section 8.1 were strong, ranging from .91 to .94. At the subscore level the reliability coefficients were relatively impressive, with the lowest of .71 for grade 8 Earth and Space.

## **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 study only included grades 5 and 8, because at grade 11 the NJSLA–S is the only large-scale statewide assessment delivered to all students. 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 below in Table 9.4.1. The intercorrelation matrix was calculated by correlating students' valid scale scores in ELA, math, and science. 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 are consistent with expectations and showed a correlation of approximately .80; the correlation between science and ELA writing was .64.

**Table 9.4.1: Grade 5 Intercorrelations by Content Area**

Content Area	N	Science	ELA	ELA-R	ELA-W	Math
Science	94,634	1.00	-	-	-	-
ELA	94,634	.80	1.00	-	-	-
ELA Reading	94,634	.82	.95	1.00	-	-
ELA Writing	94,634	.64	.88	.72	1.00	-
Math	94,634	.84	.76	.75	.65	1.00

The results at grade 8 are displayed below 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. Thus, 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 between Grade 8 science and ELA writing was .64. The correlations between science and various math test scores ranged from .62 to .80. This is most likely due to the higher- and lower-achieving students taking different assessments, which could have the effect of decreasing the scale score variance for each math test. Thus, the magnitude of the correlations between science and the various math tests all appear reasonable when considering that math achievement is more homogeneous within each subgroup 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	99,360	1.00	-	-	-	-	-	-	-
ELA	99,360	.75	1.00	-	-	-	-	-	-
ELA-Reading	99,360	.77	.94	1.00	-	-	-	-	-
ELA-Writing	99,360	.64	.92	.77	1.00	-	-	-	-
Math 8	64,116	.70	.63	.63	.54	1.00	-	-	-
Algebra I	29,851	.80	.67	.68	.57	N/A	1.00	-	-
Algebra II	654	.62	.30	.34	.19	N/A	N/A	1.00	-
Geometry	4,739	.73	.51	.54	.39	N/A	N/A	N/A	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 following 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 exists at the moment. Future NJSLA–S validity studies, including evidence based on consequences are detailed below in Section 9.7.3.

## 9.6 Other Validity Evidence

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

Part 1: Introduction—This section 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 section describes the processes used to design and develop the NJSLA–S including:

- the steps taken to link test development to the NJSLA–S’ intended inferences and uses
- 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 section 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 section 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 section and the 2019 NJSLA–S Technical Report describes 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 NJSLSA–S Standard-Setting Report

Part 6: Item and Test Statistics—This section describes the battery of statistics that were used to evaluate the NJSLSA–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 section 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 equatings

Part 8: Reliability—This section describes the reliability statistics that were calculated to verify the consistency of the NJSLSA–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 NJSLSA–S fosters valid inferences and uses. However, the NJSLSA–S validity argument requires continuing attention, and an iterative process of identifying its weakest components, making modifications, and then re-evaluating 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 NJSLA–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 NJSL–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. Numerous future validity studies are planned; they are detailed in the sections below.

**9.7.3.1 Confirmatory Factor Analysis.** The validity evidence based on internal structure is comprehensive and decisive as it pertains to the unidimensionality of the NJSLA–S. However, confirmation of the existence of the theoretical internal structure of the subscores requires a confirmatory factor analysis (CFA), as opposed to the principal components analysis used for monitoring the unidimensionality. CFA is a powerful tool for providing validity evidence that the internal structure of a construct fits the theoretical model (Brown, 2006). The theory behind the NJSLA–S asserts that the DCIs and SEPs can be grouped into, respectively, three content domains and three categories of scientific practices. A CFA can provide insight into whether the DCI and SEP groupings are justified. As noted in Section 9.7.2., another use for a CFA includes evaluating the interdependency among the three content domains and scientific practices.

**9.7.3.2 Accommodated Test Form Equivalence.** Aside from the CFA suggested above, more evidence pertaining to the equivalence of the accommodated test forms could be acquired by a more detailed review of the test results. It is known that some forms had a disproportionately large number of students flagged for person infit and outfit statistics (see Section 6.2.2: Partial Credit Model Fit Statistics). What is not currently known is which items were causing the misfit for those groups. By increasing the depth of the measurement invariance analysis for the item difficulty parameters, the items causing the misfit could be identified and assessed for patterns to determine whether certain characteristics of the items were more likely to lead to certain forms having more students flagged for person infit and outfit. To the extent that information is gleaned from the deeper analysis, then in the spirit of iteratively improving the NJSLA–S, that information could be incorporated into the item development processes to ensure the validity of the NJSLA–S test score inferences.

**9.7.3.3 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. 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 need to know a student’s strengths and needs 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 the ISRs. 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 FIRSTNAME001 performed on the elementary 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](http://nj-results.pearsonaccessnext.com) and use this code to access your student's results online.

4wdmR5FPW4h6

5

How did FIRSTNAME001 perform on the NJSLA-S?

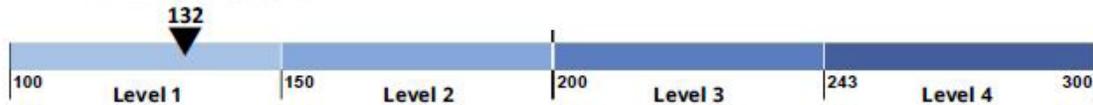
Your student's score: 132

Performance: Level 1

Below Proficient

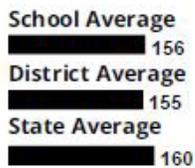
- Level 4** (243 – 300) Advanced Proficiency
- Level 3** (200 – 242) Proficient
- Level 2** (150 – 199) Near Proficiency
- Level 1** (100 – 149) Below Proficient

Your student's score

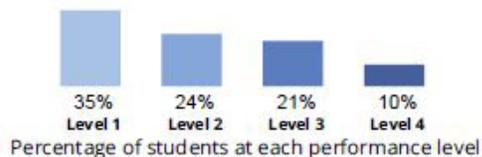


FIRSTNAME001's score on the NJSLA-S indicates that your student is at Level 1.

Students who are at Level 1 demonstrated a minimal understanding of the New Jersey Student Learning Standards-Science (NJSLA-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.



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 **Near/Met 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 **Near/Met 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 **Above 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 **Above 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 **Above 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 **Below 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](http://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/aps/cccs/science/>.

Figure 10.1.2. Sample Individual Student Report – Page 2





STATE OF NEW JERSEY  
DEPARTMENT OF EDUCATION

# STUDENT ROSTER

Grade 5

SAMPLE DISTRICT NAME  
SAMPLE SCHOOL NAME  
NEW JERSEY  
SPRING 2022

## 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.

STATE	DISTRICT	SCHOOL	STUDENT	SID	DOB	SE	ELL	TOTAL NUMBER OF STUDENT RECORDS*	NUMBER OF STUDENT WITH VALID SCORES**	AVERAGE SCALE SCORE	Student Performance Using Domains and Practices (Percent)					CRITIQUING PRACTICES				
											EARTH & SPACE SCIENCE	LIFE SCIENCE	PHYSICAL SCIENCE	INVESTIGATING PRACTICES	SENSEMAKING PRACTICES					
								102,628	101,221	225	36	21	43	33	21	46	33	21	46	
								72	69	201	13	58	29	24	20	56	35	35	30	
								19	15	180	34	42	24	46	37	17	29	60	11	
INDIVIDUAL STUDENT PERFORMANCE INDICATOR																				
										SCALE SCORE										
										259	4	✓	✓	✓	✓	✓	✓	✓	✓	✓
										233	3	!	✓	!	!	!	!	!	!	!
										115	1	✓	!	!	!	!	!	!	!	!
										167	2	✓	!	!	!	!	!	!	!	!
										Not Tested -1										
										241	3	✓	✓	✓	✓	✓	✓	✓	✓	✓
										137	1	!	!	!	!	!	!	!	!	!
										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	~	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/science](http://www.measinc.com/nj/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



CONFIDENTIAL - DO NOT DISTRIBUTE

SAMPLE DISTRICT NAME  
 SAMPLE SCHOOL NAME  
 NEW JERSEY  
 SPRING 2022

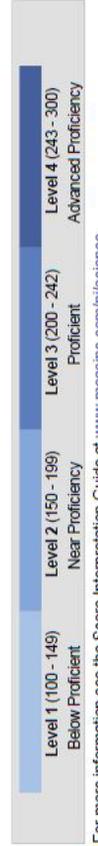
## New Jersey Student Learning Assessment - Science (NJSLA-S) Grade 5

Purpose: This report describes group performance in using the domains and practices, in comparison to state and district averages.

PERFORMANCE DISTRIBUTION BY %

STATE	29	26	28	17
DISTRICT	16	26	37	21
SAMPLE SCHOOL NAME	32	28	18	22

Number of Students with Valid Scores	Student Performance Using Domains and Practices (Percent)					
	EARTH & SPACE SCIENCE	LIFE SCIENCE	PHYSICAL SCIENCE	INVESTIGATING PRACTICES	SENSEMAKING PRACTICES	CRITIQUIING PRACTICES
99,999	36   21   43	24   63   13	33   21   46	36   21   43	24   63   13	33   21   46
5,664	13   58   29	24   20   56	35   35   30	13   58   29	24   20   56	35   35   30
204	34   42   24	46   37   17	29   60   11	34   42   24	46   37   17	29   60   11



For more information see the Score Interpretation Guide at [www.measuringinc.com/nj/science](http://www.measuringinc.com/nj/science)

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

# DISTRICT SUMMARY OF SCHOOLS

Grade 5

CONFIDENTIAL - DO NOT DISTRIBUTE



SAMPLE DISTRICT NAME

NEW JERSEY

SPRING 2022

## 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 %	Student Performance Using Domains and Practices (Percent)							
	Number of Students with Valid Scores	EARTH & SPACE SCIENCE	LIFE SCIENCE	PHYSICAL SCIENCE	INVESTIGATING PRACTICES	SENSEMAKING PRACTICES	CRITIQUIING PRACTICES	
<b>STATE</b>	<b>99,999</b>							
29 26 28 17		36 21 43	24 63 13	33 21 46	36 21 43	24 63 13	33 21 46	
<b>DISTRICT</b>	<b>5,664</b>							
16 26 37 21		13 58 29	24 20 56	35 35 30	13 58 29	24 20 56	35 35 30	
<b>ABRAHAM LINCOLN MIDDLE SCHOOL</b>	<b>204</b>							
32 28 18 22		34 42 24	46 37 17	29 60 11	34 42 24	46 37 17	29 60 11	
<b>ADA LOVELACE MIDDLE SCHOOL</b>	<b>198</b>							
23 42 35 0		21 79 0	12 57 31	33 40 27	21 79 0	12 57 31	33 40 27	
<b>BENJAMIN FRANKLIN MIDDLE SCHOOL</b>	<b>177</b>							
35 33 21 11		29 18 53	22 64 14	29 22 49	29 18 53	22 64 14	29 22 49	
<b>BOOKER T. WASHINGTON MIDDLE SCHOOL</b>	<b>204</b>							
2 39 27 32		11 57 32	28 20 52	35 34 30	11 57 32	28 20 52	35 34 30	
<b>CHARLOTTE HAWKINS BROWN MIDDLE SCHOOL</b>	<b>198</b>							
47 17 25 11		37 42 21	47 39 14	32 60 8	37 42 21	47 39 14	32 60 8	
<b>ELEANOR ROOSEVELT MIDDLE SCHOOL</b>	<b>177</b>							
23 25 37 15		29 60 11	12 49 39	35 41 24	29 60 11	12 49 39	35 41 24	

Level 1 (100 - 149)  
Below Proficient

Level 2 (150 - 199)  
Near Proficiency

Level 3 (200 - 242)  
Proficient

Level 4 (243 - 300)  
Advanced Proficiency

Below Expectations

Near/Met Expectations

Above Expectations

For more information see the Score Interpretation Guide at [www.measinc.com/nj/science](http://www.measinc.com/nj/science). Page 1 of 1

nmddywy-Batch-1234-5678-1234567

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 2022

## 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 1		Level 2		Level 3		Level 4		≥ Level 3			
					#	%	#	%	#	%	#	%	#	%		
<b>State</b>	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%	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%	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%	999,999	999.9%
<b>Gender</b>																
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%	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%	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%	99,999	99.9%
<b>Ethnicity/Race</b>																
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%	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%	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%	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%	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%	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%	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%	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%	99,999	99.9%
<b>Students with Disabilities</b>																
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%	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%	99,999	99.9%
<b>English Language Learner</b>																
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%	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%	99,999	99.9%
<b>Other</b>																
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%	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%	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%	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%	99,999	99.9%

For more information see the Score Interpretation Guide at [www.measinc.com/nj/science](http://www.measinc.com/nj/science).

Figure 10.5.1. Sample School 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.
- 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.
- 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.
- Cizek, G.J., & Bunch, M.B. (2007). *Standard Setting*. Thousand Oaks, CA: 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, 3 1–44.
- 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.
- Cronbach, L.J., Schönemann, P., & McKie, D. (1965). Alpha coefficient for stratified-parallel tests. *Educational and Psychological Measurement*, 25, 291–312.
- Dorans, N. J., & Holland, P. W. (1993). DIF detection and description: Mantel-Haenszel and standardization. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 35–66). Hillsdale, NJ: Lawrence Erlbaum.
- edCount. (2020). A proposal to conduct an independent evaluation of the alignment quality of the New Jersey Student Learning Assessment in Science. Alexandria, VA: edCount.

- Educational Testing Service. (2015). *ETS guidelines for fair tests and communications*. Princeton, NJ: Author.
- Egan, K. L. (2001). *Validity and defensibility of cutscores established by the Bookmark standard setting method*. Paper presented at the 2001 Council of Chief State School Officers Conference on Large-Scale Assessment, Houston.
- Egan, K. L., Schneider, M. C., & Ferrara, S. (2012). Performance level descriptors: History, practice, and a proposed framework. In G. J. Cizek (Eds.), *Setting Performance Standards: Foundations, Methods, and Innovations* (pp. 79-106). New York, NY: 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 presented at the Annual Meeting of the American Educational Research Association, San Francisco, CA.
- Forte, E. (2017). Evaluating alignment in large-scale standards-based assessment systems. Washington, DC: Technical Issues in Large Scale Assessment SCASS of CCSSO.
- Gorsuch, R. L. (1983). *Factor Analysis*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Hambleton, R. K. (1989). Principles and selected applications of item response theory. In R. L. Linn (Ed.), *Educational measurement* (3rd ed.). New York: American Council on Education and Macmillan.
- Hambleton, R. K., Pitoniak, M. J., & Copella, J. M. (2012). Essential steps in setting performance standards on educational tests and strategies for assessing the reliability of results. In G. J. Cizek (Ed.), *Setting Performance Standards: Foundations, Methods, and Innovations* (pp. 79–106). New York, NY: Routledge.
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Norwell, MA: Kluwer Academic Publishers.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Newbury Park, CA: 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.
- 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). Mahwah, N.J.: Erlbaum
- Kane, M.T. (2006). Validation. In R.L. Brennan (Ed.), *Educational measurement* (pp. 17–64). Washington, DC: American Council on Education/Praeger.

- Kane, M.T. (2013). Validating the interpretations and uses of test scores. *Journal of Educational Measurement*, 50, 1–73.
- 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*. NY: 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.
- Lewis, D.M., Mitzel, H.C., & Green, D. R. (1996, June). *Standard setting: A Bookmark approach*. Symposium presented at the Council of the Chief State Officers National Conference on Large-Scale Assessment, Phoenix, AZ.
- Lewis, D.M., Mitzel, H.C., Mercado, R.L., & Shultz, E.M. (2012). The Bookmark Standard Setting Procedure. In G. J. Cizek (Eds.), *Setting Performance Standards: Foundations, Methods, and Innovations* (pp. 225–254). New York, NY: Routledge.
- Lin, J. (2004). *The Bookmark Standard Setting Procedure: Strengths and Weaknesses*. The Centre for Research in Applied Measurement and Evaluation. The University of Alberta. Alberta, Canada.
- Linacre, J. M. (2012). *A User's Guide to Winsteps MINISTEP Rasch-Model Computer Programs*. Chicago, IL.
- 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*. Reading MA: 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* (3<sup>rd</sup> ed., pp. 13–104). New York, NY: 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.
- National Research Council (2012). *A framework for K-12 science education: Practices, crosscutting concepts and core ideas*. Washington, DC: National Academies Press.
- New Jersey Department of Education (2014). NJASK Technical Report. Trenton, NJ.
- New Jersey Department of Education (2017). NJASK Technical Report. Trenton, NJ.
- New Jersey Department of Education (2019). *2019 statewide assessment district test coordinator and district technology coordinator training* [PowerPoint slides]. Trenton, NJ.
- 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). New York, NY: Routledge.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Copenhagen: Danish Institute for Educational Research.
- 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.
- Schneider M. C., & Egan K. L. (2014). *A Handbook for Creating Range and Target Performance Level Descriptors*. Retrieved from: [https://www.nciea.org/sites/default/files/publications/Handbook\\_091914.pdf](https://www.nciea.org/sites/default/files/publications/Handbook_091914.pdf)
- 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*. Boston, MA: Massachusetts Department of Elementary and Secondary Education.
- Sireci, S.G., Hauger, J., Lewis, C., Wells, C, Zenisky, A.L., & Delton, J. (2009). Evaluation of the Standard Setting on the 2005 Grade 12 National Assessment of Educational Progress Mathematics Test. In Buckendahl, C.W., et. al., *Evaluation of the National Assessment of Educational Progress*. Washington, D.C.: U.S. Department of 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>

- 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 presented at the annual meeting of the Northeastern Educational Research Association, Rocky Hill, CT
- Thompson, S. J., Johnstone, C. J., & Thurlow, M. L. (2002). *Universal Design Applied to Large Scale Assessments*. Synthesis Report.
- Traub, R. E., & Rowley, G. L. (1991). Understanding reliability. *Instructional topics in educational measurement*, 37–45.
- Webb, N. L. (1997). *Criteria for alignment of expectations and assessments in mathematics and science education* (Council of Chief State School Officers and National Institute for Science Education Research Monograph No. 6). Madison, WI: University of Wisconsin, Wisconsin Center for Education Research.
- Webb, N.L. (1999). 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.
- 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*. Chicago: MESA Press.
- 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). Hillsdale, NJ: Erlbaum.

## APPENDIX A: Glossary of Abbreviations

**Table A.1: Glossary of NJSLA–S Abbreviations**

<b>Abbreviation</b>	<b>Definition</b>
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 Inc.
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

<b>Abbreviation</b>	<b>Definition</b>
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
NJSL–S	New Jersey Student Learning Standards – Science
NRC	National Research Council
OIB	Ordered Item Booklet
OPLS	Online Performance Level Setting
PAN	PearsonAccess <sup>next</sup>
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
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**

<b>Number</b>	<b>District</b>	<b>School County</b>
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**

<b>Number</b>	<b>District</b>	<b>School County</b>
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**

<b>Number</b>	<b>District</b>	<b>School County</b>
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

<b>P-VALUE</b> <b>(All items: 1-point to 4-point items)</b>	<ul style="list-style-type: none"><li>• A measure of item difficulty based on classical test theory</li><li>• Proportion correct; the proportion of students who answered a dichotomous (1-point) item correctly</li><li>• Percentage of maximum score point; item mean divided by the highest attainable score point for a polytomous (2- or 4-point) item</li><li>• P-values can range from 0 to 1.</li></ul>
<b>*FLAGGED IF:</b>	<b><math>P &lt; .25</math> (too hard)</b> <b><math>P &gt; .90</math> (too easy)</b>
<b>RASCH VALUE</b> <b>(All items: 1-point to 4-point items)</b>	<ul style="list-style-type: none"><li>• A measure of item difficulty based on item response theory, with values generally ranging from <math>-3</math> to <math>+3</math>. Higher values indicate greater difficulty (reverse of P-value)</li><li>• Items with Rasch values targeted at cut scores for performance categories are especially useful for measurement.</li></ul>
<b>SCORE POINT DISTRIBUTION</b> <b>(2-/4-point items)</b>	<ul style="list-style-type: none"><li>• Percentage of responses at each score point</li><li>• 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.</li></ul>
<b>*FLAGGED IF:</b>	<b>Response percentage &lt; 10% at any score point (2-point items)</b> <b>Response percentage &lt; 5% at any score point (4-point items)</b>
<b>ITEM-TOTAL CORRELATION</b> (All items: 1-point to 4-point items)	<ul style="list-style-type: none"><li>• 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</li><li>• RPB: correlation between an item and the total test score</li><li>• RPB can range from <math>-1</math> to <math>+1</math>.</li></ul>
<b>*FLAGGED IF:</b>	<b><math>RPB &lt; .20</math> (All items, 1-point to 4-point items)</b>

**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

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 report (Table 1.2). 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.

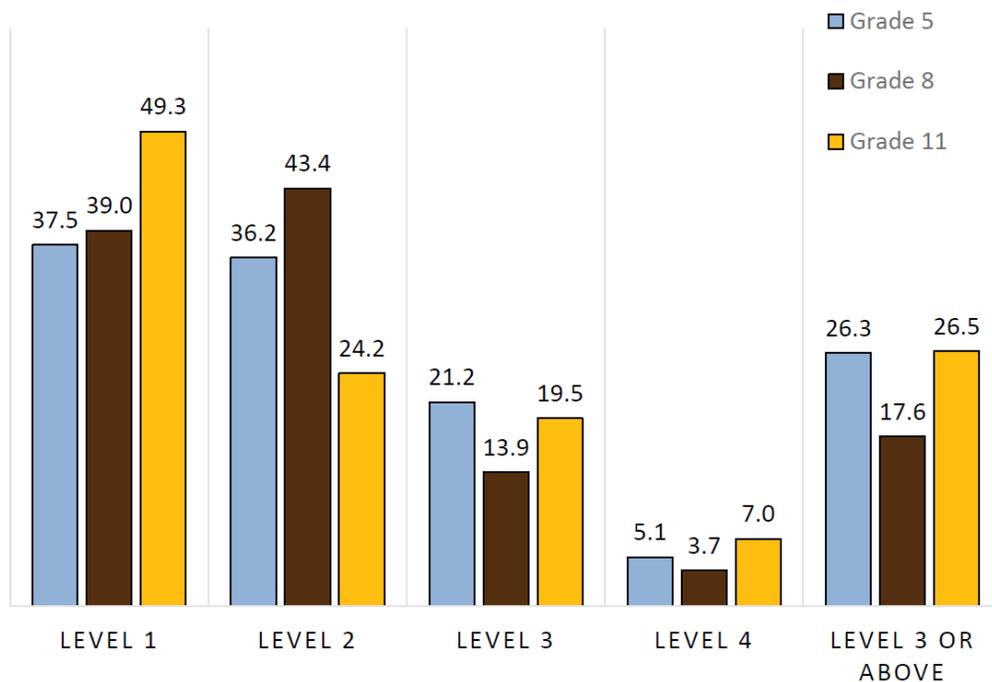
## **Results**

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, we analyzed data available as of 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. In particular, we wished to know 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**

<b>Statement</b>	<b>Grade 5</b>	<b>Grade 8</b>	<b>Grade 11</b>
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. We therefore recommend that the cut score recommendations presented here be given strong consideration for approval.

## APPENDIX E: NJSLA–S Performance Level Descriptors

### E.1 Policy PLDs

**DRAFT**

#### 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
<b>PS1: Matter and Its Interactions</b>	<ul style="list-style-type: none"> <li>that matter is made of particles that can be identified by their properties and that weight doesn't change during visible physical changes</li> <li>that the properties of substances may change when combined, but the total weight will stay the same</li> </ul>	<ul style="list-style-type: none"> <li>that matter is made of particles with unique, measurable properties that are conserved when changing state</li> <li>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</li> </ul>	<ul style="list-style-type: none"> <li>of distinguishing properties of matter and the relationship between visible and non-visible matter</li> <li>that the outcome of the combination of one or more substances is predictable based on the properties of the substances</li> </ul>
<b>PS2: Motion and Stability: Forces and Interactions</b>	<ul style="list-style-type: none"> <li>that objects are acted upon by forces that can cause predictable patterns of motion</li> <li>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</li> </ul>	<ul style="list-style-type: none"> <li>that an object's motion is a product of the net force acting on the object and can therefore cause predictable motion</li> <li>of how certain relationships among the interactions between objects are interconnected and can explain how the objects ultimately affect each other</li> </ul>	<ul style="list-style-type: none"> <li>of the relationship between net force and motion of an object in predicting future movement</li> <li>that the relationships between the interactions and the properties of objects are dependent upon systems in which the objects exist</li> </ul>

### Grade 5 Threshold Performance Level Descriptors (Physical Science)

*Students should be able to demonstrate knowledge:*

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

## Grade 5 Threshold Performance Level Descriptors (Physical Science)

*Students should be able to demonstrate knowledge:*

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

## Grade 5 Threshold Performance Level Descriptors (Life Science)

*Students should be able to demonstrate knowledge:*

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

## Grade 5 Threshold Performance Level Descriptors (Life Science)

*Students should be able to demonstrate knowledge:*

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

## Grade 5 Threshold Performance Level Descriptors (Life Science)

*Students should be able to demonstrate knowledge:*

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

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

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<b>ESS1: Earth's Place in the Universe</b>	<ul style="list-style-type: none"> <li>• that the Sun is an object in the sky and gives off light</li> <li>• that Earth is a rotating body in relative position to the Sun</li> <li>• that the Earth's rotation affects day and night</li> <li>• that there are observable patterns in Moon phases, shadows, and star patterns</li> <li>• that patterns of rock formations can contain fossils and can change due to Earth forces</li> </ul>	<ul style="list-style-type: none"> <li>• that distance affects relative size</li> <li>• 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</li> <li>• that fossil records can help identify rock layer formations because of changes caused by natural processes</li> </ul>	<ul style="list-style-type: none"> <li>• that relative distance affects brightness</li> <li>• that the 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</li> <li>• that a geological history can be determined by examining rock layers and fossil records</li> </ul>

DCI	Level 2	Level 3	Level 4
<b>ESS2: Earth's Systems</b>	<ul style="list-style-type: none"> <li>• that Earth's four major systems can interact with each other and that components of the systems can change</li> <li>• that maps can be used to locate Earth's features and processes</li> <li>• that Earth has oceans and areas of freshwater</li> <li>• that weather conditions in different areas change over time</li> <li>• that organisms affect the environment</li> </ul>	<ul style="list-style-type: none"> <li>• of how specific processes change components of Earth's four major systems and, in turn, have an effect on the systems themselves</li> <li>• that maps can be used to determine patterns of Earth's features and processes</li> <li>• of the distribution of water on Earth and its availability and accessibility</li> <li>• that patterns of weather form the basis of climate data</li> <li>• of how organisms affect the environment</li> </ul>	<ul style="list-style-type: none"> <li>• of patterns of processes affecting Earth's four major systems and how changes in those processes will likely affect the components of those systems</li> <li>• that the locations of Earth's features are related to geologic changes</li> <li>• that the water cycle affects the distribution of water on Earth</li> <li>• that climatic patterns can be used to predict future weather conditions of an area</li> <li>• that behavior of organisms in an environment can help predict changes to the physical characteristics of that environment</li> </ul>

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

*Students should be able to demonstrate knowledge:*

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

## Grade 5 SEP Threshold Performance Level Descriptors

*Students should be able to:*

SEP	Level 2	Level 3	Level 4
<p><b>Asking Questions (for Science) and Defining Problems (for engineering) (AQDP):</b></p> <p>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"> <li>identify or ask relevant questions that are testable and that can show cause-and-effect relationships in the natural or designed world</li> </ul>	<ul style="list-style-type: none"> <li>identify or ask relevant questions that can be investigated</li> <li>describe problems that can be solved</li> <li>predict reasonable outcomes</li> <li>clarify and redesign a solution to a problem</li> </ul>	<ul style="list-style-type: none"> <li>generate questions based on investigations incorporating variables to determine patterns while defining and solving a design problem</li> </ul>
<p><b>Developing and Using Models (DUM):</b></p> <p>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"> <li>describe or use a model to show the relationship among components in a phenomenon</li> </ul>	<ul style="list-style-type: none"> <li>develop or refine a model to minimize limitations, or test cause and effect relationships</li> </ul>	<ul style="list-style-type: none"> <li>evaluate and revise or develop models to show relationships in cause-and-effect systems</li> </ul>
<p><b>Planning and Carrying Out Investigations (PACI):</b></p> <p>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"> <li>plan an investigation and collect observational data using appropriate methods or tools that help identify outcomes from changing a variable</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>

## Grade 5 SEP Threshold Performance Level Descriptors

*Students should be able to:*

SEP	Level 2	Level 3	Level 4
<p><b>Analyzing and Interpreting Data (AID):</b></p> <p>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"> <li>organize relevant data to identify similarities or differences and describe how the data can be interpreted to make sense of phenomena</li> </ul>	<ul style="list-style-type: none"> <li>analyze and represent relevant data describing how the data can be interpreted to make sense of phenomena</li> </ul>	<ul style="list-style-type: none"> <li>evaluate and analyze data to refine a problem statement or make sense of phenomena</li> </ul>
<p><b>Using Mathematics and Computational Thinking (UMCT):</b></p> <p>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"> <li>identify ways to organize or analyze qualitative or quantitative data</li> </ul>	<ul style="list-style-type: none"> <li>collect and organize data to reveal patterns, determine whether qualitative or quantitative data would be more appropriate</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>

## Grade 5 SEP Threshold Performance Level Descriptors

*Students should be able to:*

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

### 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
<b>PS1: Matter and Its Interactions</b>	<ul style="list-style-type: none"> <li>that everything is made from atoms and that the states of matter have some unique characteristics</li> <li>that temperature and/or pressure have an effect on changes of state</li> <li>that chemical reactions create new substances while the mass does not change, and energy is involved</li> </ul>	<ul style="list-style-type: none"> <li>that substances are made from one or more types of atoms and that the particles in the states of matter have unique characteristics</li> <li>that atoms are regrouped and conserved during chemical processes, and energy is either released or stored</li> </ul>	<ul style="list-style-type: none"> <li>that substances can be made from two to thousands of atoms that can be combined in a variety of ways</li> <li>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</li> </ul>
<b>PS2: Motion and Stability: Forces and Interactions</b>	<ul style="list-style-type: none"> <li>that the movement of an object is the sum of its forces</li> <li>that forces among objects are either attractive or repulsive and are dependent upon the distance between the objects</li> </ul>	<ul style="list-style-type: none"> <li>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</li> <li>that the size of the electromagnetic force depends upon the magnitudes of the charges, currents, or magnetic strengths due to the fields created</li> </ul>	<ul style="list-style-type: none"> <li>of the effect of balanced versus unbalanced forces on the motion of objects</li> <li>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</li> </ul>

## Grade 8 Threshold Performance Level Descriptors (Physical Science)

*Students should be able to demonstrate knowledge:*

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

## Grade 8 Threshold Performance Level Descriptors (Life Science)

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<p><b>LS1: From Molecules to Organisms: Structures and Processes</b></p>	<ul style="list-style-type: none"> <li>• that cells contain special structures which may be specific to the type of cell in a living unicellular or multicellular organism</li> <li>• 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</li> <li>• that matter and energy cycle through plants, creating sugars which can be broken down or rearranged to release the energy</li> <li>• that sense receptors can send various signals to the brain</li> </ul>	<ul style="list-style-type: none"> <li>• 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</li> <li>• 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</li> <li>• 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</li> <li>• that nerves transmit sense receptor inputs to be processed in the brain, resulting in memories or responses</li> </ul>	<ul style="list-style-type: none"> <li>• of how parts of a cell function together in a manner similar to how systems interact in multicellular organisms</li> <li>• 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</li> <li>• of the relationship between photosynthesis and cellular respiration and of how an organism obtains energy to sustain life</li> <li>• of the different ways a sense receptor reacts to inputs and of the process by which the signal is processed</li> </ul>

## Grade 8 Threshold Performance Level Descriptors (Life Science)

*Students should be able to demonstrate knowledge:*

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

## Grade 8 Threshold Performance Level Descriptors (Life Science)

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<b>LS4: Biological Evolution: Unity and Diversity</b>	<ul style="list-style-type: none"> <li>• 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</li> <li>• that environmental conditions will drive trait commonality in species</li> </ul>	<ul style="list-style-type: none"> <li>• 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</li> <li>• of how environmental conditions can change a species over generations and of how distributions of traits reflect adaptation by natural selection</li> </ul>	<ul style="list-style-type: none"> <li>• of evolutionary history based on anatomical similarities and to predict predominance of certain traits in a population</li> <li>• to predict trait distribution in a species based on changing environmental conditions</li> </ul>

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

*Students should be able to demonstrate knowledge:*

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

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

*Students should be able to demonstrate knowledge:*

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

## Grade 8 SEP Threshold Performance Level Descriptors

*Students should be able to:*

SEP	Level 2	Level 3	Level 4
<p><b>Analyzing and Interpreting Data (AID):</b>            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"> <li>identify and/or interpret data, graphical displays, and/or concepts of statistics and/or their limitations to provide evidence for phenomena</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>
<p><b>Asking Questions (for science) and Defining Problems (for engineering) (AQDP):</b>            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"> <li>identify questions that arise from observations and models in order to clarify information and/or arguments, refine models, and/or determine relationships</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>

## Grade 8 SEP Threshold Performance Level Descriptors

*Students should be able to:*

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

## Grade 8 SEP Threshold Performance Level Descriptors

*Students should be able to:*

SEP	Level 2	Level 3	Level 4
<p><b>Obtaining, Evaluating, and Communicating Information (OEIC):</b> 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"> <li>• read and use information from multiple simple scientific sources to describe patterns, clarify claims, and/or assess accuracy</li> </ul>	<ul style="list-style-type: none"> <li>• integrate information from multiple, complex, qualitative sources to clarify claims, assess accuracy, and evaluate conclusions</li> </ul>	<ul style="list-style-type: none"> <li>• integrate information from multiple, complex, quantitative sources to describe patterns, clarify claims, assess accuracy, and evaluate conclusions</li> </ul>
<p><b>Planning and Carrying Out Investigations (PACI):</b> 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"> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• plan an investigation that includes the identification of variables and/or controls, or indicates how much data is sufficient to serve as evidence necessary to test a design solution, or evaluate an experimental design</li> </ul>	<ul style="list-style-type: none"> <li>• 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</li> </ul>
<p><b>Using Mathematics and Computational Thinking (UMCT):</b> 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"> <li>• identify qualitative and quantitative data and when the use of digital tools is warranted,</li> <li>• select appropriate mathematical representations, and</li> <li>• use algorithms to solve problems and/or address engineering questions</li> </ul>	<ul style="list-style-type: none"> <li>• decide whether to use qualitative or quantitative data,</li> <li>• use digital tools to analyze large data sets,</li> <li>• use mathematical representations, and</li> <li>• explain and/or evaluate algorithms or mathematical concepts for solving problems and/or addressing engineering questions</li> </ul>	<ul style="list-style-type: none"> <li>• explain when to use qualitative or quantitative data,</li> <li>• evaluate digital tools,</li> <li>• explain mathematical representations, and/or</li> <li>• create algorithms to solve problems and/or address engineering questions</li> </ul>

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

## Grade 11 Threshold Performance Level Descriptors (Physical Science)

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<b>PS2: Motion and Stability: Forces and Interactions</b>	<ul style="list-style-type: none"> <li>• of quantified acceleration and momentum</li> <li>• of types of fields and attractive/repulsive forces of gravitational and/or electric fields</li> <li>• that electrical energy can be stored or transmitted</li> </ul>	<ul style="list-style-type: none"> <li>• (quantified knowledge) of factors that affect Newton’s second law, single object momentum systems, and conservation of momentum</li> <li>• of how interactions happen at a distance due to fields</li> <li>• of electrical interactions at the atomic level</li> <li>• of the difference between magnetic and electric fields</li> </ul> <p style="text-align: center;"><b>OR</b></p> <ul style="list-style-type: none"> <li>• (quantified knowledge) of Coulomb’s law and Newton’s universal law of gravitation</li> <li>• of how electrical energy can be stored in a battery or transmitted by electric currents</li> </ul>	<ul style="list-style-type: none"> <li>• (quantified knowledge) of outside interactions that affect the momentum and acceleration of a single- or multiple-object system</li> <li>• of how to predict changes in electrical and gravitational forces</li> <li>• 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</li> </ul>

DCI	Level 2	Level 3	Level 4
<p><b>PS3: Energy</b></p>	<ul style="list-style-type: none"> <li>• of how different types of energy can be transferred</li> <li>• of systems in which energy is conserved and how the availability of energy restricts what is possible in a closed system</li> <li>• of the nature of the relationship between two objects interacting in a field using the energy prospective</li> <li>• of how energy can be converted to different forms</li> </ul>	<ul style="list-style-type: none"> <li>• of how energy manifests itself at the microscopic and macroscopic scale and how energy transfers in a system</li> <li>• (quantified knowledge) of how energy transfers in and out of a system</li> </ul> <p><b>OR</b></p> <ul style="list-style-type: none"> <li>• of possible and impossible events based on energy availability, and defined stable states</li> <li>• of how the distance between two objects affects the energy of a field</li> <li>• of how energy can be converted to less useful forms</li> <li>• of how solar energy can be captured and used for other processes, such as photosynthesis</li> </ul>	<ul style="list-style-type: none"> <li>• of the amount of various types of energy in a given situation and how microscopic changes affect macroscopic manifestations of energy</li> <li>• of how to evaluate physical changes in a system using the conservation of energy</li> <li>• of how to predict changes in energy in a field based on the position and nature of objects</li> <li>• of the importance of energy conservation and efficiency</li> </ul>

## Grade 11 Threshold Performance Level Descriptors (Physical Science)

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<p><b>PS4: Waves and Their Applications in Technologies for Information Transfer</b></p>	<ul style="list-style-type: none"> <li>• of how a wave travels through a medium, including understanding of examples of digitized information, and qualitative understanding of superposition principle</li> <li>• of the wave and particle models of electromagnetic radiation, the absorption of electromagnetic radiation, and the relationship between frequency and energy of light</li> <li>• of everyday experiences that involve waves and how wave signals are produced, transmitted, and captured</li> </ul>	<ul style="list-style-type: none"> <li>• (quantified knowledge) of the relationship among frequency, wavelength, and speed in a real-world phenomenon</li> </ul> <p style="text-align: center;"><b>OR</b></p> <ul style="list-style-type: none"> <li>• of the advantages and disadvantages of digitizing information</li> <li>• 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</li> <li>• of technologies used to produce, transmit, and/or capture signals and technologies used to store and interpret information</li> </ul>	<ul style="list-style-type: none"> <li>• 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</li> <li>• 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</li> <li>• of how technology can be used to store and/or interpret information</li> </ul>

**Grade 11 Threshold Performance Level Descriptors (Life Science)**

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<p><b>LS1: From Molecules to Organisms: Structures and Processes</b></p>	<ul style="list-style-type: none"> <li>• of how multicellular organisms utilize feedback mechanisms and have specialized cells that are organized and function according to the proteins coded by the DNA</li> <li>• of the role of cellular division (mitosis) in creating genetically identical cells that differentiate into complex multicellular organisms</li> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• 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</li> <li>• of how mitosis and differentiation produce and maintain complex organisms from a single cell</li> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• 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</li> <li>• of how genetic material from two variants of each chromosome pair is maintained as a single cell (fertilized egg) grows to a multicellular organism</li> <li>• of the interdependence of photosynthesis and cellular respiration and their role in the growth and maintenance of living systems</li> </ul>

## Grade 11 Threshold Performance Level Descriptors (Life Science)

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<b>LS2: Ecosystems: Interactions, Energy, and Dynamics</b>	<ul style="list-style-type: none"> <li>• of both living and non-living factors that contribute to the carrying capacity of the ecosystem</li> <li>• 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</li> <li>• 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</li> <li>• of how group behavior has evolved to increase individual and group survival</li> </ul>	<ul style="list-style-type: none"> <li>• of how carrying capacity is affected by challenges and/or availability of resources</li> <li>• 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</li> <li>• 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</li> <li>• of how group behaviors can increase the chances of survival for individuals and their genetic relatives</li> </ul>	<ul style="list-style-type: none"> <li>• of how carrying capacity affects the population size of a given species within an ecosystem</li> <li>• 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</li> <li>• 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</li> <li>• of how changes to the group or conditions can affect the survival of individuals and their genetic relatives</li> </ul>

**Grade 11 Threshold Performance Level Descriptors (Life Science)**

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<p><b>LS3: Heredity: Inheritance and Variation of Traits</b></p>	<ul style="list-style-type: none"> <li>• of how all cells have the same DNA containing genes that are the organisms' characteristics, but not all DNA codes for protein</li> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• that chromosomes contain genes that code for proteins and regions that do not code for proteins, and that different cells express different genes</li> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• of the mechanisms of gene regulation and different possible functions of segments of non-protein coding DNA</li> <li>• 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</li> </ul>

DCI	Level 2	Level 3	Level 4
<b>LS4: Biological Evolution: Unity and Diversity</b>	<ul style="list-style-type: none"> <li>• of the different types of evidence of evolution</li> <li>• of how natural selection allows inheritable advantageous traits to become more common if they increase chances of survival within populations</li> <li>• that natural selection selects for inheritable traits that provide a survival advantage for a particular environment</li> <li>• that changes to the environment may cause the selection of different traits leading to changes in the population known as adaptation</li> <li>• that the frequency of traits depends on natural selection forces that can change with a changing environment</li> <li>• of how biodiversity increases or decreases and how humans need resources and biodiversity, but are having adverse effects on biodiversity</li> </ul>	<ul style="list-style-type: none"> <li>• of how different sources of evidence for evolution can support each other</li> <li>• 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</li> <li>• that evolution occurs when there is genetic variation, competition, and selective reproduction of organisms with desirable genetic traits</li> <li>• that organisms with desirable traits will become more common, but as the environment changes, different traits may provide the selective advantages</li> <li>• that some populations may increase while others may go extinct</li> <li>• of specific results of human activities that affect the environment and biodiversity and reasons why preservation of biodiversity is desirable</li> </ul>	<ul style="list-style-type: none"> <li>• of how DNA sequences, amino acid sequences, and anatomical and embryological evidence support that evolution has occurred</li> <li>• of how natural selection leads to different levels of performance of the individual</li> <li>• that factors affecting natural selection work together creating changes in the diversity within populations and ecosystems</li> <li>• that changing environments cause changes in selection pressures that result in further adaptation or extinction</li> <li>• of ways that humans can maintain or increase biodiversity while meeting the needs of humanity and why this is beneficial to life on Earth</li> </ul>

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

*Students should be able to demonstrate knowledge:*

DCI	Level 2	Level 3	Level 4
<b>ESS1: Earth's Place in the Universe</b>	<ul style="list-style-type: none"> <li>• of the Big Bang, which allowed for the creation of galaxies and stars, where many elements are created</li> <li>• of identifying properties of orbits, factors that affect the orbit, and how the orbit affects the stellar body</li> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• that light spectra emitted from a star can give information about its life cycle, composition, and distance</li> <li>• of features of motion of orbital objects, what changes that motion, and the effects of changing the motion of the stellar body</li> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• of the life cycle of stars and explain how the characteristics of a star can support the Big Bang theory</li> <li>• of the laws explaining motions of orbiting objects, their changes, and the changes to the stellar bodies as a result of those changes</li> <li>• of why different areas of the Earth have rocks of different ages and the processes that are erasing the early rock history</li> </ul>

DCI	Level 2	Level 3	Level 4
<p><b>ESS2: Earth's Systems</b></p>	<ul style="list-style-type: none"> <li>• of how Earth has a series of interacting dynamic systems</li> <li>• that Earth's surface is in motion, and that motion can create physical features on the Earth's surface</li> <li>• of the properties of water that are essential to Earth's dynamics</li> <li>• of Earth's atmosphere and how it undergoes temperature changes</li> <li>• that dynamic and delicate feedbacks between the Earth's systems and biosphere exist</li> </ul>	<ul style="list-style-type: none"> <li>• of methods of investigation of Earth's dynamic systems and how the data can be used to describe the effects of these systems</li> <li>• that Earth's surface is in motion due to convection, creating physical features that have changed throughout history</li> <li>• of how the properties of water are essential to Earth's processes</li> <li>• 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</li> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• 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</li> <li>• of the theory of plate tectonics allowing for the prediction of future plate movements and interpretations of Earth's geologic history</li> <li>• of how the properties of water can be used to explain Earth's processes</li> <li>• of why Earth's atmosphere undergoes short-term and long-term temperature changes at the global scale</li> <li>• 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</li> </ul>

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

*Students should be able to demonstrate knowledge:*

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

## Grade 11 SEP Threshold Performance Level Descriptors

*Students should be able to:*

SEP	Level 2	Level 3	Level 4
<p><b>(Investigating)</b></p> <p><b>Asking Questions (for Science) and Defining Problems (for engineering) (AQDP):</b>                      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"> <li>ask relevant questions or define problems in different contexts, based on unexpected results, independent and dependent variables, models, theories, etc.</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>	<ul style="list-style-type: none"> <li>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</li> </ul>

SEP	Level 2	Level 3	Level 4
<p><b>(Sensemaking)</b></p> <p><b>Developing and Using Models (DUM):</b>  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"> <li>• use a model to generate data that test the model’s reliability and/or evaluates its merits and limitations</li> </ul>	<ul style="list-style-type: none"> <li>• 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</li> </ul>	<ul style="list-style-type: none"> <li>• develop or revise complex models that test and/or predict relationships/ phenomena based on the models’ merits and limitations</li> </ul>

## Grade 11 SEP Threshold Performance Level Descriptors

*Students should be able to:*

SEP	Level 2	Level 3	Level 4
<p><b>(Investigating)</b></p> <p><b>Planning and Carrying Out Investigations (PACI):</b> 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"> <li>identify ways to conduct an investigation (including making a directional hypothesis) or test a design solution through manipulating variables or acquiring data</li> </ul>	<ul style="list-style-type: none"> <li>plan and/or conduct an investigation (including making a directional hypothesis) or test a design solution through manipulating variables or acquiring data</li> </ul>	<ul style="list-style-type: none"> <li>revise and/or evaluate an investigation in which an independent variable is manipulated or an unsatisfactory performance is found</li> </ul>

SEP	Level 2	Level 3	Level 4
<p><b>(Sensemaking)</b></p> <p><b>Analyzing and Interpreting Data (AID):</b>            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"> <li>• identify the appropriate statistics and/or data, and/or their limitations, when providing evidence for claims, design solutions, or solving problems</li> </ul>	<ul style="list-style-type: none"> <li>• apply and/or analyze data and statistics to identify or solve scientific and engineering problems, or to make scientific claims</li> </ul>	<ul style="list-style-type: none"> <li>• evaluate the use of data and statistics and/or their limitations to solve problems, make claims, or design solutions</li> </ul>

## Grade 11 SEP Threshold Performance Level Descriptors

*Students should be able to:*

SEP	Level 2	Level 3	Level 4
<p><b>(Investigating)</b>  <b>Using Mathematics and Computational Thinking (UMCT):</b></p> <p>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"> <li>• apply/use mathematical concepts to describe conclusions that may require deciding when to use qualitative versus quantitative data</li> </ul>	<ul style="list-style-type: none"> <li>• apply/use mathematical computational representations to see if a model is viable, or decide if qualitative or quantitative data meet criteria for success</li> </ul>	<ul style="list-style-type: none"> <li>• through the use of evaluation of mathematical computations, create a model or justify the choice of qualitative versus quantitative data</li> </ul>

SEP	Level 2	Level 3	Level 4
<p><b>(Sensemaking)</b></p> <p><b>Constructing Explanations (for science) and Designing Solutions (for engineering) (CEDs):</b></p> <p>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"> <li>• identify and describe appropriate data and/or evidence for supporting claims, solving problems, constructing explanations, or designing solutions</li> </ul>	<ul style="list-style-type: none"> <li>• make or revise claims, explanations, or solutions by applying appropriate data and/or evidence</li> </ul>	<ul style="list-style-type: none"> <li>• evaluate, design, or construct claims, explanations, or solutions by applying appropriate data, evidence, and/or scientific theories and laws</li> </ul>

## Grade 11 SEP Threshold Performance Level Descriptors

*Students should be able to:*

SEP	Level 2	Level 3	Level 4
<p><b>(Critiquing)</b>  <b>Engaging in Argument from Evidence (EAE):</b>            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"> <li>identify and/or describe the main points of an argument or claim that is based on scientific evidence</li> </ul>	<ul style="list-style-type: none"> <li>evaluate and/or defend a claim or argument— or choose between competing arguments— related to currently accepted explanations or solutions</li> </ul>	<ul style="list-style-type: none"> <li>construct and/or critique an argument or claim by using scientific evidence</li> </ul>
<p><b>(Critiquing)</b>  <b>Obtaining, Evaluating, and Communicating Information (OEI):</b>            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"> <li>read and compare sources of information to describe patterns in evidence and/ or evidence for solving problems or answering scientific questions</li> </ul>	<ul style="list-style-type: none"> <li>integrate information from multiple sources to gather valid and reliable evidence for solving problems or answering scientific questions</li> </ul>	<ul style="list-style-type: none"> <li>evaluate information from multiple sources and determine the usefulness of evidence, ensuring it is valid and reliable, for solving problems or answering scientific questions</li> </ul>

## **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-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-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-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 choose 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-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	RPB	Median Time
518000_01	1	TE	DUM	LS2.A	S & SM	Life	Sensemaking	0.336	0.47	.26	69
518000_03	1	TE	OECI	LS2.A	S & SM	Life	Critiquing	0.445	0.43	.49	38
518000_06	1	MC	CEDS	LS2.A	C and E	Life	Sensemaking	0.490	0.44	.39	65
518004_01	1	TE	CEDS	PS3.B	E&M	Physical	Sensemaking	0.178	0.46	.41	79
518004_03	1	MC	EAE	PS3.B	E&M	Physical	Critiquing	-0.114	0.52	.47	44
518004_05	1	TE	PACI	PS3.B	SC	Physical	Investigating	0.521	0.40	.36	33
518054_01	1	TE	DUM	ESS1.B	PAT	Earth and Space	Sensemaking	1.026	0.26	.53	202
518054_03	1	MC	DUM	ESS1.B	PAT	Earth and Space	Sensemaking	1.324	0.26	.27	27
518054_06	1	MC	EAE	ESS1.B	C and E	Earth and Space	Critiquing	0.627	0.41	.32	67
518056_01	1	TE	DUM	LS2.A	S & SM	Life	Sensemaking	-1.285	0.70	.43	86
518056_02	1	TE	CEDS	LS2.A	S & SM	Life	Sensemaking	0.432	0.41	.48	47
519000_01a	1	TE	AQDP	LS2.D	C and E	Life	Investigating	0.530	0.41	.46	105
519000_05a	1	TE	UMCT	LS2.D	S, P, and Q	Life	Investigating	-0.699	0.67	.37	56
519000_07a	1	TE	EAE	LS2.D	C and E	Life	Critiquing	-1.160	0.69	.50	39
519013_02b	1	TE	AID	ESS2.D	PAT	Earth and Space	Sensemaking	0.221	0.47	.47	98
519013_04b	1	TE	AID	ESS2.D	PAT	Earth and Space	Sensemaking	-0.424	0.54	.55	59
519013_06b	2	TE	UMCT	ESS2.D	S & SM	Earth and Space	Investigating	0.904	0.66	.57	134
1905B007_01	1	TE	AID	LS3.A	SC	Life	Sensemaking	1.017	0.30	.61	140
1905B007_03	1	TE	UMCT	LS3.A	S, P, and Q	Life	Investigating	0.902	0.32	.58	93
1905B007_05	1	TE	AID	LS3.A	SC	Life	Sensemaking	0.553	0.38	.57	81
1905B007_08	4	CR	EAE	LS3.A	PAT	Life	Critiquing	0.393	1.52	.63	353
1905B007_10	1	TE	CEDS	LS1.C	S & SM	Life	Sensemaking	0.541	0.38	.65	37

Grade 5											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	RPB	Median Time
1905B009_01	1	TE	UMCT	LS4.B	S, P, and Q	Life	Investigating	0.350	0.43	.63	112
1905B009_02	1	TE	AID	LS4.B	S, P, and Q	Life	Sensemaking	-0.809	0.64	.67	46
1905B009_05	1	TE	AQDP	LS4.B	PAT	Life	Investigating	1.579	0.22	.23	75
1905B022_01	1	TE	AQDP	PS3.B	E&M	Physical	Investigating	1.136	0.30	.29	108
1905B022_03	1	TE	CEDS	PS3.B	E&M	Physical	Sensemaking	-0.644	0.62	.42	48
1905B022_07	1	TE	UMCT	PS3.B	E&M	Physical	Investigating	1.090	0.30	.56	88
1905B022_10	4	CR	EAE	PS3.B	E&M	Physical	Critiquing	0.952	1.38	.66	363
1905M005_01	1	TE	EAE	ESS2.D	S & SM	Earth and Space	Critiquing	1.014	0.31	.38	165
1905M005_03	1	TE	OECI	ESS2.D	PAT	Earth and Space	Critiquing	0.881	0.33	.54	93
1905M005_04	1	MC	UMCT	ESS2.D	PAT	Earth and Space	Investigating	-0.282	0.55	.20	59
1905M008_01	1	MC	EAE	ESS1.A	S, P, and Q	Earth and Space	Critiquing	-0.515	0.59	.55	107
1905M008_05	1	TE	AID	ESS1.A	PAT	Earth and Space	Sensemaking	0.111	0.47	.58	65
1905M008_06	1	MC	OECI	ESS1.A	C and E	Earth and Space	Critiquing	-0.794	0.64	.53	58
1905M018_03	1	MC	DUM	LS1.B	PAT	Life	Sensemaking	-0.026	0.50	.29	85
1905M018_04	1	MC	DUM	LS1.B	PAT	Life	Sensemaking	-0.188	0.53	.34	53
1905M018_05	1	MC	CEDS	LS1.B	PAT	Life	Sensemaking	-0.762	0.64	.39	52
1905M030_02	1	TE	CEDS	ESS3.A	PAT	Earth and Space	Sensemaking	-1.494	0.75	.44	91
1905M030_04	1	TE	EAE	ESS3.A	PAT	Earth and Space	Critiquing	-0.269	0.55	.55	71
1905M036_02	1	MC	EAE	PS4.B	S & SM	Physical	Critiquing	-0.055	0.50	.45	97
1905M036_03	1	MC	CEDS	PS4.B	C and E	Physical	Sensemaking	-0.644	0.62	.41	86
1905M036_04	1	MC	DUM	PS4.B	C and E	Physical	Sensemaking	-0.411	0.57	.47	62
1905M040_01	1	TE	PACI	PS2.A	PAT	Physical	Investigating	0.733	0.36	.47	155
1905M040_03	1	MC	AID	PS2.A	C and E	Physical	Sensemaking	-0.224	0.54	.46	79
1905M040_05	1	TE	AID	PS2.A	C and E	Physical	Sensemaking	0.775	0.35	.41	51

<b>Grade 5</b>											
<b>UIN</b>	<b>Points</b>	<b>Item Type</b>	<b>SEP</b>	<b>DCI</b>	<b>CCC</b>	<b>Domain</b>	<b>Practice</b>	<b>Rasch</b>	<b>Mean</b>	<b>RPB</b>	<b>Median Time</b>
1905M048_01	1	MC	AQDP	PS1.A	SF	Physical	Investigating	-0.020	0.50	.38	82
1905M048_02	1	TE	PACI	PS1.A	SF	Physical	Investigating	1.287	0.26	.38	42
1905M048_03	1	TE	AID	PS1.A	SF	Physical	Sensemaking	0.053	0.48	.46	58
519013_08b	4	CR	AID	ESS2.D	PAT	Earth and Space	Sensemaking	1.432	0.90	.57	478

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

Grade 8											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	RPB	Median Time
1908M014_01	1	TE	UMCT	ESS1.C	S, P, and Q	Earth and Space	Investigating	0.896	0.20	.31	94
1908M014_04	1	MC	AQDP	LS4.A	SF	Life	Investigating	0.427	0.24	.22	73
1908M014_05	1	TE	DUM	ESS1.C	S, P, and Q	Earth and Space	Sensemaking	-0.057	0.38	.40	54
818004_01b	1	TE	AID	PS3.A	E&M	Physical	Sensemaking	-1.955	0.76	.37	100
818004_02a	1	TE	EAE	PS3.A	E&M	Physical	Critiquing	0.587	0.27	.30	99
818004_03b	1	TE	AID	PS3.A	E&M	Physical	Sensemaking	-0.852	0.55	.17	74
818004_04a	1	TE	CEDS	PS3.B	S & SM	Physical	Sensemaking	2.058	0.09	.33	83
818015_01a	1	TE	DUM	ESS1.A	S & SM	Earth and Space	Sensemaking	0.527	0.30	.45	82
818015_02a	1	TE	UMCT	ESS1.B	S, P, and Q	Earth and Space	Investigating	0.097	0.36	.16	35
818015_03a	2	TE	DUM	ESS1.A	C and E	Earth and Space	Sensemaking	1.097	0.46	.21	68
818015_05b	4	CR	CEDS	ESS1.B	S & SM	Earth and Space	Sensemaking	0.268	1.10	.58	315
818041_02	1	TE	CEDS	LS4.A	PAT	Life	Sensemaking	0.849	0.24	.28	106
818041_03a	1	MC	CEDS	LS4.C	PAT	Life	Sensemaking	-0.278	0.45	.33	45
818041_04b	1	TE	EAE	LS4.B	C and E	Life	Critiquing	0.229	0.33	.45	72
818064	1	MC	AQDP	LS1.D	SF	Life	Investigating	-0.619	0.51	.47	59
818067_01	1	MC	PACI	LS1.B	SF	Life	Investigating	0.592	0.29	.06	102
818067_02	1	TE	EAE	LS1.B	SF	Life	Critiquing	1.485	0.14	.42	81
818165	1	TE	OECI	PS4.C	SC	Physical	Critiquing	1.043	0.18	.43	67
818197_02	1	MC	UMCT	LS3.A	C and E	Life	Investigating	-0.658	0.45	.38	40
818197_03	1	MC	UMCT	LS3.B	S, P, and Q	Life	Investigating	-0.379	0.41	.44	64
818251	1	MC	AQDP	PS1.B	S, P, and Q	Physical	Investigating	0.219	0.30	.34	87
818314	1	MC	OECI	PS4.C	SF	Physical	Critiquing	0.254	0.34	.23	80
818333	1	TE	PACI	PS3.C	E&M	Physical	Investigating	-0.058	0.35	.46	58

Grade 8											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	RPB	Median Time
1908B000_03	1	TE	CEDS	PS2.A	C and E	Physical	Sensemaking	1.003	0.20	.30	78
1908B000_04	1	TE	AID	PS2.B	C and E	Physical	Sensemaking	0.573	0.26	.40	70
1908B000_08	1	TE	EAE	PS3.A	S, P, and Q	Physical	Critiquing	0.624	0.26	.41	51
1908B000_11	4	CR	CEDS	PS2.A	C and E	Physical	Sensemaking	-0.027	1.41	.71	425
1908B000_12	1	TE	AID	PS2.A	SC	Physical	Sensemaking	0.589	0.26	.40	51
1908B001_01	1	TE	OECI	ESS3.C	C and E	Earth and Space	Critiquing	0.129	0.34	.33	113
1908B001_02	2	TE	UMCT	ESS3.A	SC	Earth and Space	Investigating	-1.102	1.20	.50	93
1908B007_01	1	TE	EAE	LS2.A	E&M	Life	Critiquing	-0.360	0.43	.52	78
1908B007_04	1	TE	CEDS	LS2.A	S & SM	Life	Sensemaking	-0.318	0.43	.41	52
1908B007_05	4	CR	EAE	LS2.A	C and E	Life	Critiquing	1.009	0.77	.62	437
1908B007_10	2	TE	PACI	LS2.A	C and E	Life	Investigating	0.195	0.66	.52	103
1908M003_02	1	TE	DUM	LS3.A	S, P, and Q	Life	Sensemaking	1.580	0.13	.16	89
1908M003_07	1	TE	DUM	LS3.B	PAT	Life	Sensemaking	1.413	0.15	.27	47
1908M003_08	1	MC	CEDS	LS3.B	PAT	Life	Sensemaking	0.397	0.30	.37	40
1908M004_02	1	MC	AQDP	LS2.B	C and E	Life	Investigating	-0.802	0.54	.54	81
1908M004_03	1	MC	PACI	LS1.C	E&M	Life	Investigating	0.787	0.23	.08	49
1908M004_04	1	TE	UMCT	LS1.B	S & SM	Life	Investigating	-0.318	0.43	.45	57
1908M004_05	1	MC	CEDS	LS2.A	S & SM	Life	Sensemaking	-0.064	0.38	.32	40
1908M005_01	1	MC	DUM	PS1.B	S & SM	Physical	Sensemaking	-0.679	0.51	.37	63
1908M005_02	1	TE	UMCT	PS1.A	S & SM	Physical	Investigating	1.080	0.19	.36	73
1908M005_03	1	TE	UMCT	PS1.B	E&M	Physical	Investigating	1.409	0.15	.49	71
1908M005_05	1	MC	AQDP	PS1.B	C and E	Physical	Investigating	-0.173	0.40	.38	44
1908M005_06	1	MC	EAE	PS1.B	PAT	Physical	Critiquing	-0.080	0.38	.19	35
1908M018_01	1	TE	DUM	ESS2.C	SC	Earth and Space	Sensemaking	0.990	0.20	.21	127

Grade 8											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	RPB	Median Time
1908M018_04	1	MC	AID	ESS2.D	C and E	Earth and Space	Sensemaking	-0.349	0.44	.29	120
1908M018_06	1	MC	CEDS	ESS2.C	PAT	Earth and Space	Sensemaking	-0.520	0.47	.26	82
1908M023_03	1	TE	OECI	PS4.C	PAT	Physical	Critiquing	0.592	0.26	.41	39
1908M023_07	1	TE	OECI	PS4.C	SF	Physical	Critiquing	1.390	0.15	.39	110
1908M026_01	1	TE	OECI	ESS3.D	C and E	Earth and Space	Critiquing	1.459	0.14	.41	151
1908M026_04	1	MC	AID	ESS3.D	PAT	Earth and Space	Sensemaking	-0.096	0.38	.32	60
1908M026_06	1	MC	EAE	ESS3.D	PAT	Earth and Space	Critiquing	-0.205	0.41	.31	61
1908M030_01	1	MC	EAE	LS1.A	S & SM	Life	Critiquing	.269	.32	.32	85
1908M030_02	1	TE	AID	LS1.B	C and E	Life	Sensemaking	.106	.34	.39	29
1908M030_05	1	MC	PACI	LS1.D	C and E	Life	Investigating	.362	.30	.37	44
1908M033_02	1	TE	UMCT	ESS2.B	PAT	Earth and Space	Investigating	.105	.35	.32	90
1908M033_03	1	TE	EAE	ESS3.B	PAT	Earth and Space	Critiquing	-.412	.45	.30	85
1908M033_04	1	MC	DUM	ESS2.B	SC	Earth and Space	Sensemaking	-.857	.54	.36	60

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

Grade 11											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	RPB	Median Time
1911B002_01A	1	TE	AID	ESS3.B	SC	Earth and Space	Sensemaking	0.251	0.44	.61	63
1911B002_05A	1	TE	EAE	ESS3.A	SC	Earth and Space	Critiquing	1.396	0.21	.42	72
1911M069_01	1	TE	AID	LS3.B	C and E	Life	Sensemaking	1.911	0.16	.42	131
1911M069_02	1	MC	AID	LS3.B	C and E	Life	Sensemaking	-0.005	0.49	.46	45
1911M069_03	1	MC	EAE	LS3.B	C and E	Life	Critiquing	-0.542	0.57	.47	43
1911M069_04	1	TE	CEDS	LS3.B	C and E	Life	Sensemaking	-0.430	0.54	.49	40
HS18006_01	1	MC	AID	LS3.B	S, P, and Q	Life	Sensemaking	-1.518	0.74	.43	68
HS18006_05	1	MC	EAE	LS3.B	SC	Life	Critiquing	0.978	0.29	.46	87
HS18006_07	1	TE	CEDS	LS3.B	C and E	Life	Sensemaking	-0.626	0.61	.38	48
HS18018_02	1	TE	AID	ESS3.D	SC	Earth and Space	Sensemaking	-0.998	0.64	.50	40
HS18018_03	1	TE	AID	ESS3.D	SC	Earth and Space	Sensemaking	1.041	0.24	.23	42
HS18018_05	1	TE	EAE	ESS3.D	SC	Earth and Space	Critiquing	0.829	0.33	.42	36
HS18018_08	1	TE	CEDS	ESS3.D	SC	Earth and Space	Sensemaking	0.786	0.32	.46	56
HS18047_02	1	TE	AID	LS1.A	SC	Life	Sensemaking	-0.475	0.52	.57	131
HS18047_03	1	MC	OECI	LS1.A	PAT	Life	Critiquing	-0.119	0.47	.40	60
HS18047_05	1	MC	AQDP	LS1.A	S, P, and Q	Life	Investigating	-0.212	0.50	.42	46
HS18047_07	1	TE	UMCT	LS1.A	SC	Life	Investigating	-1.185	0.66	.53	42
HS18051_02	1	MC	DUM	ESS2.D	E&M	Earth and Space	Sensemaking	0.830	0.30	.36	91
HS18051_04	1	TE	AID	ESS2.D	E&M	Earth and Space	Sensemaking	-0.027	0.46	.37	63
HS18051_05	1	TE	CEDS	ESS2.D	E&M	Earth and Space	Sensemaking	0.769	0.32	.37	44
HS18051_08	1	MC	CEDS	ESS2.D	E&M	Earth and Space	Sensemaking	0.482	0.36	.36	30
HS18082_03	1	MC	PACI	PS2.B	C and E	Physical	Investigating	0.303	0.39	.31	89
HS18082_06	1	TE	PACI	PS2.B	C and E	Physical	Investigating	0.806	0.31	.44	52

Grade 11											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	RPB	Median Time
HS18082_07	1	MC	UMCT	PS2.B	C and E	Physical	Investigating	0.225	0.36	.46	83
HS19003_01A	1	TE	DUM	PS4.C	C and E	Physical	Sensemaking	-0.417	0.52	.53	51
HS19003_03A	1	TE	UMCT	PS4.C	S & SM	Physical	Investigating	-0.448	0.52	.64	45
HS19003_05A	1	TE	AQDP	PS4.C	S & SM	Physical	Investigating	2.270	0.12	.23	44
HS19003_09A	1	TE	AID	PS4.C	PAT	Physical	Sensemaking	-0.690	0.58	.52	59
1911B002_08B	4	CR	EAE	ESS3.B	PAT	Earth and Space	Critiquing	0.139	1.57	.70	228
1911B002_10B	1	TE	DUM	ESS3.B	S & SM	Earth and Space	Sensemaking	0.904	0.28	.29	48
1911B009_01A	1	TE	OECI	LS2.D	S & SM	Life	Critiquing	-0.439	0.54	.55	70
1911B009_03A	1	TE	AID	LS2.D	PAT	Life	Sensemaking	0.479	0.37	.62	52
1911B009_05A	1	TE	EAE	LS2.D	S, P, and Q	Life	Critiquing	1.377	0.21	.46	57
1911B009_07A	4	CR	CEDS	LS2.D	S & SM	Life	Sensemaking	0.683	1.21	.61	248
1911B009_09A	1	TE	PACI	LS4.C	S & SM	Life	Investigating	1.331	0.22	.34	74
1911M002_01	1	TE	DUM	ESS2.A	PAT	Earth and Space	Sensemaking	1.627	0.18	.34	62
1911M002_04	1	MC	EAE	ESS2.A	SC	Earth and Space	Critiquing	0.909	0.28	.18	53
1911M002_05	1	MC	AQDP	ESS2.B	S, P, and Q	Earth and Space	Investigating	0.166	0.42	.32	27
1911M019_04	1	MC	EAE	PS1.A	SF	Physical	Critiquing	-0.549	0.56	.34	60
1911M019_08	1	TE	CEDS	PS1.A	SF	Physical	Sensemaking	1.344	0.22	.35	38
1911M019_10	1	MC	PACI	PS1.A	SF	Physical	Investigating	0.436	0.37	.40	44
1911M023_02	1	MC	AID	LS2.B	PAT	Life	Sensemaking	0.460	0.37	.54	84
1911M023_05	1	TE	AID	LS2.B	S, P, and Q	Life	Sensemaking	-0.600	0.57	.59	58
1911M023_06	1	MC	OECI	LS2.B	PAT	Life	Critiquing	0.420	0.37	.31	53
1911M023_07	1	MC	PACI	LS2.B	PAT	Life	Investigating	0.465	0.36	.23	38
1911M028_01	1	MC	UMCT	PS3.A	S & SM	Physical	Investigating	-0.617	0.57	.27	74
1911M028_03	1	MC	UMCT	PS3.A	S & SM	Physical	Investigating	1.373	0.23	.45	46

Grade 11											
UIN	Points	Item Type	SEP	DCI	CCC	Domain	Practice	Rasch	Mean	RPB	Median Time
1911M028_04	1	MC	UMCT	PS3.B	S & SM	Physical	Investigating	0.103	0.43	.23	36
1911M028_06	1	MC	UMCT	PS3.B	S & SM	Physical	Investigating	0.576	0.35	.25	68
1911M029_01	1	TE	DUM	PS3.B	E&M	Physical	Sensemaking	0.580	0.34	.32	74
1911M029_02	1	TE	EAE	PS3.B	SF	Physical	Critiquing	-1.265	0.69	.43	49
1911M029_03	1	MC	CEDS	PS3.A	SF	Physical	Sensemaking	-0.592	0.57	.29	64
1911M029_04	1	MC	DUM	PS3.A	SF	Physical	Sensemaking	1.034	0.27	.17	54
1911M079_02	1	MC	AID	ESS1.C	PAT	Earth and Space	Sensemaking	0.249	0.41	.24	113
1911M079_03	1	TE	DUM	ESS1.C	S, P, and Q	Earth and Space	Sensemaking	-0.368	0.52	.48	122
1911M079_04	1	TE	OECI	ESS1.C	S, P, and Q	Earth and Space	Critiquing	1.241	0.23	.46	134
1911M086_01	1	TE	EAE	ESS1.C	S & SM	Earth and Space	Critiquing	0.003	0.46	.55	93
1911M086_02	1	MC	AQDP	ESS1.C	S & SM	Earth and Space	Investigating	-1.265	0.69	.48	26
1911M150_01	1	MC	AQDP	LS4.C	C and E	Life	Investigating	-0.271	0.51	.40	49
1911M150_02	1	TE	CEDS	LS4.C	SF	Life	Sensemaking	0.185	0.42	.48	40
1911M150_03	1	TE	EAE	LS4.C	S & SM	Life	Critiquing	0.410	0.37	.57	77
HS18060_01	1	TE	OECI	PS1.C	E&M	Physical	Critiquing	1.747	0.17	.53	115
HS18060_03	1	MC	AID	PS1.C	E&M	Physical	Sensemaking	-0.032	0.46	.41	34
HS18060_04	1	MC	AID	PS1.C	E&M	Physical	Sensemaking	-0.192	0.49	.48	25
HS18060_06	1	TE	AID	PS1.C	E&M	Physical	Sensemaking	2.033	0.14	.32	32
HS19000_01a	1	TE	UMCT	PS3.B	S, P, and Q	Physical	Investigating	1.261	0.24	.54	114
HS19000_06b	1	TE	AQDP	PS3.B	C and E	Physical	Investigating	0.646	0.32	.11	39
HS19000_09a	1	TE	EAE	ESS3.B	S, P, and Q	Earth and Space	Critiquing	-0.406	0.53	.58	76
HS19003_07A	4	CR	EAE	PS4.C	C and E	Physical	Critiquing	0.706	1.28	.67	310

## 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	4	0.00	0.00	0.00	0.00	0.01	0.00	0.01
1	100	39	0.04	0.04	0.04	0.01	0.10	0.06	0.02
2	100	164	0.17	0.16	0.18	0.03	0.35	0.25	0.09
3	100	509	0.53	0.46	0.59	0.08	1.10	0.78	0.25
4	100	1,241	1.29	1.13	1.44	0.22	2.69	1.95	0.54
5	100	2,490	2.59	2.29	2.87	0.54	5.30	3.88	1.12
6	100	4,293	4.46	3.99	4.91	0.90	9.21	6.63	2.02
7	100	6,528	6.78	6.26	7.28	1.37	13.68	10.15	3.13
8	100	9,073	9.42	8.83	10.00	1.95	18.65	14.13	4.46
9	100	11,767	12.22	11.60	12.82	2.58	23.47	18.44	5.92
10	100	14,423	14.98	14.41	15.53	3.31	28.10	22.68	7.38
11	105	17,029	17.69	17.15	18.21	4.05	32.58	26.63	9.00
12	110	19,557	20.31	19.97	20.64	5.11	36.82	30.25	10.67
13	114	21,981	22.83	22.63	23.02	5.84	40.83	33.82	12.31
14	119	24,394	25.33	25.31	25.36	6.71	44.14	37.40	14.15
15	123	26,768	27.80	28.03	27.58	7.45	47.60	40.73	16.05
16	127	29,057	30.18	30.57	29.80	8.41	50.80	43.92	17.87
17	131	31,404	32.61	33.21	32.05	9.41	54.06	47.09	19.81
18	135	33,627	34.92	35.76	34.12	10.39	56.91	49.95	21.85
19	139	35,838	37.22	38.18	36.30	11.22	59.46	52.95	23.89
20	142	38,019	39.48	40.65	38.37	12.55	61.89	55.71	25.99
21	146	40,054	41.60	42.98	40.27	13.71	64.12	58.05	28.12
22	150	42,364	44.00	45.52	42.54	15.26	66.69	60.66	30.58
23	153	44,516	46.23	47.89	44.65	16.79	68.86	63.20	32.79
24	156	46,725	48.53	50.34	46.79	18.34	70.72	65.82	35.21
25	159	48,856	50.74	52.58	48.98	19.76	72.89	68.20	37.55
26	163	51,021	52.99	54.95	51.11	21.38	74.81	70.47	40.06
27	166	53,168	55.22	57.23	53.29	23.13	76.75	72.56	42.64
28	169	55,274	57.40	59.47	55.43	25.14	78.59	74.60	45.09
29	173	57,472	59.69	61.84	57.63	27.25	80.39	76.74	47.75
30	176	59,661	61.96	64.14	59.87	29.12	82.24	78.67	50.54
31	179	61,756	64.14	66.33	62.03	31.24	83.69	80.54	53.27
32	183	63,793	66.25	68.46	64.14	33.55	85.15	82.24	55.87
33	186	65,841	68.38	70.64	66.21	35.85	86.55	83.83	58.63
34	189	67,969	70.59	72.81	68.46	38.59	87.95	85.32	61.54

<b>Raw Score</b>	<b>Scale Score</b>	<b>All Cum. N</b>	<b>All Cum. %</b>	<b>Female Cum. %</b>	<b>Male Cum. %</b>	<b>Asian Cum. %</b>	<b>Black Cum. %</b>	<b>Hisp. Cum. %</b>	<b>White Cum. %</b>
35	193	70,031	72.73	75.00	70.55	41.27	89.41	86.88	64.21
36	196	71,731	74.50	76.71	72.37	43.45	90.41	88.10	66.60
37	200	73,940	76.79	78.93	74.74	46.57	91.77	89.55	69.62
38	203	75,877	78.80	80.80	76.88	49.70	92.71	90.90	72.23
39	207	77,687	80.68	82.59	78.85	52.71	93.80	91.97	74.73
40	211	79,566	82.63	84.35	80.99	55.96	94.78	93.13	77.30
41	214	81,332	84.47	86.06	82.93	59.15	95.60	94.09	79.85
42	218	83,039	86.24	87.69	84.85	62.47	96.26	94.96	82.31
43	222	84,752	88.02	89.41	86.68	66.14	96.85	95.79	84.76
44	226	86,359	89.69	90.97	88.46	69.63	97.37	96.58	87.01
45	230	87,834	91.22	92.28	90.20	73.12	97.87	97.21	89.09
46	235	89,210	92.65	93.59	91.74	76.36	98.42	97.71	91.06
47	243	90,507	94.00	94.82	93.20	79.75	98.82	98.17	92.90
48	244	91,681	95.22	95.91	94.54	83.00	99.14	98.61	94.46
49	249	92,764	96.34	96.86	95.84	86.29	99.36	98.98	95.93
50	254	93,655	97.27	97.63	96.92	89.36	99.57	99.27	97.06
51	259	94,391	98.03	98.31	97.76	92.22	99.71	99.46	97.93
52	265	94,984	98.65	98.82	98.48	94.33	99.84	99.62	98.62
53	272	95,460	99.14	99.28	99.00	96.39	99.92	99.79	99.11
54	279	95,779	99.47	99.52	99.42	97.71	99.96	99.90	99.45
55	287	96,018	99.72	99.73	99.70	98.75	99.98	99.94	99.72
56	297	96,151	99.86	99.88	99.84	99.40	99.99	99.96	99.86
57	300	96,239	99.95	99.96	99.94	99.78	100.00	99.99	99.94
58	300	96,277	99.99	99.99	99.99	99.92	100.00	100.00	99.99
59	300	96,284	100.00	100.00	100.00	99.97	100.00	100.00	100.00
60	300	96,288	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	2	0.00	0.00	0.00	0.00	0.01	0.00	0.00
1	100	10	0.01	0.01	0.01	0.00	0.04	0.01	0.00
2	100	38	0.04	0.04	0.04	0.00	0.09	0.05	0.02
3	100	132	0.13	0.12	0.14	0.01	0.35	0.16	0.06
4	100	316	0.31	0.27	0.35	0.04	0.79	0.43	0.13
5	100	795	0.78	0.65	0.92	0.16	1.76	1.11	0.35
6	100	1,655	1.63	1.35	1.91	0.31	3.50	2.39	0.73
7	104	3,045	3.00	2.49	3.50	0.62	6.10	4.49	1.38
8	110	5,188	5.12	4.30	5.90	1.00	10.32	7.66	2.39
9	115	8,097	7.99	6.78	9.15	1.69	15.62	11.96	3.90
10	120	11,646	11.49	9.98	12.93	2.41	21.78	17.17	5.89
11	125	15,815	15.60	13.76	17.37	3.43	28.47	23.08	8.52
12	129	20,275	20.00	17.97	21.95	4.70	35.32	29.39	11.42
13	133	24,754	24.42	22.32	26.44	6.07	41.82	35.60	14.53
14	136	29,219	28.82	26.77	30.80	7.40	47.91	41.83	17.77
15	140	33,473	33.02	31.07	34.90	8.76	53.61	47.38	21.16
16	143	37,625	37.12	35.34	38.82	10.38	58.99	52.61	24.58
17	146	41,477	40.92	39.46	42.32	12.26	63.46	57.31	27.92
18	150	45,041	44.43	43.19	45.63	14.16	67.48	61.34	31.30
19	152	48,289	47.64	46.79	48.45	15.98	71.11	64.70	34.62
20	155	51,273	50.58	49.94	51.21	17.67	74.06	67.90	37.68
21	158	54,109	53.38	53.09	53.66	19.78	76.41	70.80	40.76
22	161	56,911	56.14	56.07	56.23	22.12	78.77	73.58	43.75
23	163	59,536	58.73	58.82	58.66	24.21	80.91	76.09	46.72
24	166	62,027	61.19	61.45	60.96	26.15	82.72	78.22	49.77
25	168	64,335	63.46	63.84	63.13	28.30	84.30	80.15	52.60
26	171	66,525	65.63	66.17	65.13	30.62	85.68	81.95	55.29
27	173	68,700	67.77	68.44	67.15	33.05	87.12	83.70	57.91
28	176	70,746	69.79	70.66	68.98	35.48	88.40	85.31	60.39
29	178	72,774	71.79	72.75	70.89	37.95	89.71	86.81	62.92
30	181	74,601	73.59	74.63	72.62	40.64	90.61	88.08	65.22
31	183	76,407	75.37	76.53	74.28	43.00	91.50	89.23	67.68
32	185	78,155	77.10	78.33	75.94	45.54	92.35	90.34	69.96
33	187	79,835	78.76	80.01	77.57	48.11	93.24	91.29	72.18
34	190	81,486	80.38	81.71	79.13	50.54	94.04	92.26	74.41
35	192	82,937	81.82	83.20	80.51	52.81	94.66	93.02	76.43
36	194	84,386	83.24	84.72	81.84	55.17	95.39	93.82	78.36

Raw Score	Scale Score	All Cum. N	All Cum. %	Female Cum. %	Male Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	White Cum. %
37	197	85,562	84.40	85.91	82.98	57.24	95.81	94.39	80.00
38	200	87,007	85.83	87.35	84.39	60.04	96.33	95.06	81.95
39	201	88,300	87.11	88.62	85.67	62.55	96.91	95.63	83.69
40	203	89,391	88.18	89.65	86.79	64.82	97.27	96.17	85.13
41	206	90,495	89.27	90.70	87.92	67.24	97.61	96.68	86.59
42	208	91,562	90.32	91.75	88.97	69.78	97.95	97.18	87.92
43	210	92,540	91.29	92.65	90.00	72.14	98.31	97.57	89.16
44	213	93,501	92.24	93.50	91.04	74.44	98.56	97.96	90.41
45	215	94,361	93.08	94.24	91.98	76.75	98.78	98.21	91.55
46	218	95,138	93.85	94.94	92.81	78.74	98.98	98.51	92.54
47	220	95,831	94.53	95.57	93.54	80.67	99.13	98.73	93.41
48	223	96,519	95.21	96.14	94.33	82.51	99.23	98.95	94.35
49	225	97,146	95.83	96.69	95.02	84.42	99.34	99.12	95.13
50	228	97,723	96.40	97.20	95.64	86.17	99.45	99.28	95.88
51	231	98,293	96.96	97.68	96.28	88.11	99.51	99.43	96.59
52	233	98,822	97.49	98.11	96.89	89.93	99.61	99.57	97.20
53	236	99,236	97.89	98.46	97.36	91.28	99.68	99.64	97.73
54	239	99,639	98.29	98.78	97.83	92.80	99.74	99.73	98.20
55	242	99,968	98.62	99.05	98.21	93.93	99.79	99.82	98.57
56	246	100,228	98.87	99.23	98.53	94.81	99.87	99.86	98.89
57	249	100,498	99.14	99.41	98.88	95.82	99.91	99.91	99.18
58	253	100,697	99.34	99.55	99.13	96.72	99.93	99.92	99.40
59	256	100,867	99.50	99.68	99.33	97.52	99.95	99.94	99.56
60	260	101,011	99.64	99.80	99.50	98.20	99.96	99.95	99.71
61	265	101,115	99.75	99.87	99.63	98.68	99.97	99.97	99.79
62	269	101,207	99.84	99.91	99.77	99.17	99.97	99.98	99.87
63	274	101,266	99.90	99.94	99.86	99.42	99.99	99.98	99.92
64	280	101,297	99.93	99.96	99.90	99.57	100.00	99.99	99.94
65	286	101,333	99.96	99.98	99.94	99.77	100.00	100.00	99.97
66	292	101,347	99.98	99.99	99.96	99.87	100.00	100.00	99.98
67	300	101,359	99.99	99.99	99.98	99.93	100.00	100.00	99.99
68	300	101,367	100.00	100.00	99.99	99.98	100.00	100.00	100.00
69	300	101,368	100.00	100.00	99.99	100.00	100.00	100.00	100.00
70	300	101,370	100.00	100.00	100.00	100.00	100.00	100.00	100.00
71	300	101,371	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. %	Hispanic Cum. %	White Cum. %
1	100	10	0.01	0.01	0.01	0.00	0.02	0.01	0.01
2	100	27	0.03	0.02	0.04	0.01	0.07	0.02	0.03
3	100	84	0.09	0.07	0.11	0.03	0.23	0.12	0.04
4	100	241	0.26	0.19	0.33	0.08	0.57	0.39	0.13
5	100	528	0.57	0.40	0.73	0.17	1.16	0.89	0.28
6	100	1,105	1.19	0.94	1.44	0.33	2.35	1.84	0.63
7	100	2,065	2.22	1.77	2.67	0.57	4.40	3.44	1.16
8	100	3,495	3.76	3.10	4.41	0.83	7.22	5.75	2.12
9	100	5,360	5.77	4.84	6.68	1.22	10.96	8.87	3.27
10	100	7,601	8.18	6.98	9.35	1.95	15.14	12.49	4.78
11	100	10,204	10.98	9.54	12.39	2.84	20.05	16.49	6.61
12	100	12,974	13.96	12.28	15.60	3.58	25.41	20.83	8.56
13	101	15,614	16.80	15.04	18.52	4.54	30.00	25.01	10.48
14	106	18,232	19.62	17.86	21.34	5.46	34.45	29.16	12.40
15	111	20,637	22.21	20.59	23.79	6.41	38.43	32.88	14.26
16	115	22,979	24.73	23.18	26.24	7.40	42.31	36.26	16.16
17	120	25,230	27.15	25.58	28.69	8.26	45.78	39.52	18.10
18	124	27,374	29.46	27.98	30.90	9.04	48.71	42.67	20.02
19	128	29,468	31.71	30.38	33.02	10.06	51.62	45.69	21.87
20	132	31,460	33.86	32.63	35.06	11.23	54.65	48.29	23.72
21	135	33,495	36.05	34.97	37.11	12.46	57.33	50.98	25.65
22	139	35,421	38.12	37.25	38.98	13.43	59.71	53.52	27.63
23	143	37,343	40.19	39.43	40.94	14.78	62.25	55.83	29.58
24	146	39,245	42.23	41.64	42.83	15.92	64.53	58.31	31.49
25	149	41,074	44.20	43.88	44.53	17.17	66.58	60.44	33.50
26	153	42,909	46.18	46.06	46.30	18.26	68.81	62.53	35.53
27	158	44,682	48.08	48.13	48.06	19.55	70.67	64.55	37.54
28	159	46,485	50.02	50.27	49.81	20.99	72.50	66.60	39.58
29	162	48,308	51.99	52.35	51.65	22.43	74.19	68.60	41.71
30	165	49,999	53.81	54.35	53.30	23.78	75.95	70.47	43.66
31	168	51,684	55.62	56.22	55.05	25.15	77.64	72.22	45.66
32	171	53,450	57.52	58.27	56.81	26.75	79.10	74.02	47.83
33	174	55,127	59.32	60.18	58.51	28.09	80.45	75.85	49.86
34	177	56,800	61.12	62.16	60.14	29.71	81.88	77.33	52.04
35	180	58,502	62.96	64.10	61.86	31.19	83.39	78.78	54.30
36	183	60,152	64.73	65.93	63.58	32.80	84.71	80.38	56.37
37	186	61,691	66.39	67.63	65.20	34.39	85.82	81.75	58.38

Raw Score	Scale Score	All Cum. N	All Cum. %	Female Cum. %	Male Cum. %	Asian Cum. %	Black Cum. %	Hisp. Cum. %	White Cum. %
38	189	63,216	68.03	69.33	66.77	36.13	86.93	83.13	60.32
39	192	64,759	69.69	71.07	68.36	37.98	87.97	84.38	62.39
40	195	65,994	71.02	72.47	69.63	39.52	88.70	85.31	64.13
41	200	67,705	72.86	74.38	71.39	41.61	89.80	86.71	66.41
42	201	69,108	74.37	75.95	72.85	43.36	90.68	87.85	68.28
43	204	70,510	75.88	77.51	74.31	45.54	91.59	88.81	70.18
44	207	71,864	77.34	79.04	75.69	47.47	92.17	89.88	72.06
45	210	73,184	78.76	80.55	77.02	49.74	92.96	90.81	73.79
46	213	74,456	80.12	81.93	78.39	51.90	93.70	91.61	75.52
47	216	75,646	81.41	83.27	79.61	53.99	94.29	92.41	77.13
48	219	76,857	82.71	84.63	80.85	56.18	94.81	93.17	78.80
49	222	78,019	83.96	85.88	82.10	58.23	95.38	93.87	80.41
50	225	79,117	85.14	87.05	83.29	60.55	95.79	94.48	81.93
51	229	80,178	86.28	88.14	84.49	62.73	96.24	95.08	83.40
52	232	81,233	87.42	89.19	85.71	65.04	96.73	95.60	84.86
53	235	82,333	88.60	90.31	86.95	67.43	97.12	96.18	86.36
54	239	83,290	89.63	91.29	88.03	69.63	97.57	96.61	87.69
55	242	84,205	90.62	92.23	89.06	71.60	97.97	97.05	88.96
56	246	85,073	91.55	93.09	90.07	73.81	98.23	97.41	90.13
57	250	85,902	92.44	93.95	90.99	76.13	98.50	97.80	91.16
58	254	86,724	93.33	94.75	91.96	78.37	98.71	98.12	92.28
59	258	87,506	94.17	95.51	92.87	80.59	98.87	98.40	93.34
60	262	88,243	94.96	96.17	93.78	82.71	99.05	98.67	94.32
61	266	88,911	95.68	96.76	94.64	84.87	99.30	98.91	95.15
62	271	89,517	96.33	97.33	95.37	86.67	99.47	99.12	95.96
63	275	90,078	96.94	97.83	96.07	88.58	99.57	99.31	96.66
64	280	90,583	97.48	98.24	96.75	90.34	99.66	99.48	97.28
65	285	91,052	97.98	98.65	97.34	92.09	99.73	99.58	97.89
66	291	91,476	98.44	99.00	97.90	93.82	99.81	99.69	98.38
67	297	91,840	98.83	99.29	98.39	95.21	99.86	99.77	98.83
68	300	92,138	99.15	99.51	98.81	96.29	99.91	99.85	99.19
69	300	92,402	99.44	99.72	99.16	97.50	99.96	99.89	99.48
70	300	92,582	99.63	99.81	99.46	98.23	99.98	99.94	99.67
71	300	92,731	99.79	99.90	99.69	98.93	99.99	99.98	99.83
72	300	92,816	99.88	99.94	99.82	99.41	99.99	99.99	99.90
73	300	92,870	99.94	99.98	99.91	99.70	100.00	100.00	99.95
74	300	92,905	99.98	100.00	99.96	99.88	100.00	100.00	99.98

<b>Raw Score</b>	<b>Scale Score</b>	<b>All Cum. N</b>	<b>All Cum. %</b>	<b>Female Cum. %</b>	<b>Male Cum. %</b>	<b>Asian Cum. %</b>	<b>Black Cum. %</b>	<b>Hisp. Cum. %</b>	<b>White Cum. %</b>
75	300	92,918	99.99	100.00	99.99	99.93	100.00	100.00	100.00
76	300	92,923	100.00	100.00	100.00	99.98	100.00	100.00	100.00
77	300	92,925	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**

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
518000_01	0.336	1.28	1.43	.26	0.36	.18	0.47
518000_03	0.445	0.98	0.95	.49	1.04	.01	0.43
518000_06	0.490	1.11	1.20	.39	0.75	.09	0.44
518004_01	0.178	1.07	1.10	.41	0.84	.06	0.46
518004_03	-0.114	0.98	0.97	.47	1.04	.02	0.52
518004_05	0.521	1.13	1.21	.36	0.72	.07	0.40
518054_01	1.026	0.80	0.71	.53	1.35	.00	0.26
518054_03	1.324	1.15	1.54	.27	0.69	.06	0.26
518054_06	0.627	1.19	1.33	.33	0.58	.11	0.41
518056_01	-1.285	1.00	1.04	.43	0.99	.00	0.70
518056_02	0.432	0.97	0.96	.48	1.07	.00	0.41
519000_01a	0.530	1.01	1.00	.46	0.98	.01	0.41
519000_05a	-0.699	1.04	1.06	.37	0.92	.15	0.67
519000_07a	-1.160	0.92	0.84	.50	1.15	.00	0.69
519013_02b	0.221	0.99	0.97	.47	1.03	.00	0.47
519013_04b	-0.424	0.90	0.85	.55	1.24	.00	0.54
519013_06b	0.904	0.97	0.96	.57	1.06	.00	0.66
1905B007_01	1.017	0.77	0.68	.61	1.40	.00	0.30
1905B007_03	0.902	0.81	0.81	.59	1.33	.00	0.32
1905B007_05	0.553	0.84	0.81	.57	1.32	.00	0.38
1905B007_08	0.393	1.29	1.23	.63	0.77	.00	1.52
1905B007_10	0.541	0.75	0.68	.65	1.52	.00	0.38
1905B009_01	0.350	0.79	0.70	.63	1.48	.00	0.43
1905B009_02	-0.809	0.69	0.59	.67	1.62	.00	0.64
1905B009_05	1.579	1.19	1.52	.23	0.70	.05	0.22
1905B022_01	1.136	1.17	1.39	.29	0.67	.06	0.30
1905B022_03	-0.644	1.02	1.02	.42	0.95	.04	0.62
1905B022_07	1.090	0.86	0.72	.56	1.27	.00	0.30
1905B022_10	0.952	0.98	0.95	.66	1.03	.00	1.38
1905M005_01	1.014	1.07	1.15	.38	0.86	.03	0.31
1905M005_03	0.881	0.88	0.82	.54	1.23	.00	0.33
1905M005_04	-0.282	1.33	1.48	.20	0.22	.21	0.55
1905M008_01	-0.515	0.87	0.81	.55	1.28	.00	0.59
1905M008_05	0.111	0.85	0.81	.58	1.33	.00	0.47
1905M008_06	-0.794	0.87	0.81	.53	1.26	.00	0.64

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1905M018_03	-0.026	1.21	1.28	.29	0.50	.10	0.50
1905M018_04	-0.188	1.15	1.21	.34	0.65	.09	0.53
1905M018_05	-0.762	1.06	1.10	.39	0.87	.06	0.64

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

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1905M030_02	-1.494	0.91	0.86	.45	1.13	.00	0.75
1905M030_04	-0.269	0.87	0.84	.55	1.29	.00	0.55
1905M036_02	-0.055	1.02	1.01	.45	0.96	.02	0.50
1905M036_03	-0.644	1.05	1.03	.41	0.90	.08	0.62
1905M036_04	-0.411	0.97	0.95	.47	1.06	.00	0.57
1905M040_01	0.733	0.98	1.00	.47	1.03	.00	0.36
1905M040_03	-0.224	1.01	1.01	.46	0.99	.05	0.54
1905M040_05	0.775	1.05	1.06	.41	0.91	.02	0.35
1905M048_01	-0.020	1.10	1.14	.38	0.77	.04	0.50
1905M048_02	1.287	1.05	1.05	.38	0.93	.01	0.26
1905M048_03	0.053	1.00	1.01	.46	0.99	.00	0.48
519013_08b	1.432	1.13	1.06	.57	0.94	.00	0.90

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

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1908M014_01	0.896	0.98	1.01	.31	1.02	.00	0.20
1908M014_04	0.427	1.02	1.07	.22	0.94	.00	0.24
1908M014_05	-0.057	0.96	0.94	.41	1.11	.00	0.38
818004_01b	-1.955	0.90	0.82	.37	1.17	.00	0.76
818004_02a	0.587	1.07	1.09	.30	0.89	.02	0.27
818004_03b	-0.852	1.16	1.21	.17	0.40	.16	0.55
818004_04a	2.058	0.96	0.90	.33	1.03	.00	0.09
818015_01a	0.527	0.98	0.93	.45	1.05	.00	0.30
818015_02a	0.097	1.21	1.29	.16	0.50	.12	0.36
818015_03a	1.097	1.27	1.33	.21	0.64	.12	0.46
818015_05b	0.268	1.10	1.05	.58	0.99	.00	1.10
818041_02	0.849	1.09	1.24	.28	0.85	.04	0.24
818041_03a	-0.278	1.05	1.07	.33	0.84	.08	0.45
818041_04b	0.229	0.93	0.94	.45	1.13	.00	0.33
818064	-0.619	0.88	0.86	.47	1.42	.00	0.51
818067_01	0.592	1.35	1.54	.06	0.43	.12	0.29
818067_02	1.485	0.92	0.79	.42	1.09	.00	0.14
818165	1.043	0.86	0.82	.43	1.15	.00	0.18
818197_02	-0.658	0.99	0.98	.38	1.04	.00	0.45
818197_03	-0.379	0.92	0.90	.44	1.27	.00	0.41
818251	0.219	0.99	1.02	.34	1.01	.00	0.30
818314	0.254	1.15	1.16	.23	0.71	.05	0.34
818333	-0.058	0.89	0.86	.46	1.29	.00	0.35
1908B000_03	1.003	1.02	1.04	.30	0.97	.01	0.20
1908B000_04	0.573	0.96	0.94	.40	1.07	.00	0.26

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

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1908B000_08	0.624	0.96	0.98	.41	1.05	.00	0.26
1908B000_11	-0.027	0.81	0.77	.71	1.18	.00	1.41
1908B000_12	0.589	0.97	0.96	.40	1.05	.00	0.26
1908B001_01	0.129	1.03	1.07	.33	0.91	.02	0.34
1908B001_02	-1.102	0.88	0.88	.50	1.24	.00	1.20
1908B007_01	-0.360	0.84	0.81	.52	1.50	.00	0.43
1908B007_04	-0.318	0.95	0.94	.41	1.16	.00	0.43
1908B007_05	1.009	0.98	0.87	.62	1.08	.00	0.77
1908B007_10	0.195	0.92	0.92	.52	1.14	.00	0.66
1908M003_02	1.580	1.11	1.24	.16	0.90	.02	0.13
1908M003_07	1.413	1.02	1.14	.27	0.97	.01	0.15
1908M003_08	0.397	0.99	0.99	.37	1.02	.00	0.30
1908M004_02	-0.802	0.82	0.78	.54	1.65	.00	0.54
1908M004_03	0.787	1.21	1.46	.08	0.67	.07	0.23
1908M004_04	-0.318	0.92	0.91	.45	1.24	.00	0.43
1908M004_05	-0.064	1.05	1.09	.32	0.85	.05	0.38
1908M005_01	-0.679	0.99	0.97	.37	1.06	.03	0.51
1908M005_02	1.080	0.98	1.00	.36	1.02	.00	0.19
1908M005_03	1.409	0.88	0.71	.49	1.13	.00	0.15
1908M005_05	-0.173	0.98	0.97	.38	1.06	.00	0.40
1908M005_06	-0.080	1.16	1.21	.19	0.56	.10	0.38
1908M018_01	0.990	1.09	1.23	.21	0.87	.03	0.20
1908M018_04	-0.349	1.05	1.07	.29	0.83	.02	0.44
1908M018_06	-0.520	1.08	1.11	.26	0.71	.06	0.47
1908M023_03	0.592	0.95	1.00	.41	1.06	.00	0.26
1908M023_07	1.390	0.94	0.95	.39	1.05	.00	0.15
1908M026_01	1.459	0.92	0.78	.41	1.09	.00	0.14
1908M026_04	-0.096	1.03	1.03	.32	0.92	.00	0.38
1908M026_06	-0.205	1.05	1.07	.31	0.85	.04	0.41
1908M030_01	0.269	1.03	1.07	.32	0.92	.02	0.32
1908M030_02	0.106	0.97	0.99	.39	1.05	.00	0.34
1908M030_05	0.362	0.98	1.02	.38	1.02	.00	0.30
1908M033_02	0.105	1.04	1.05	.32	0.90	.02	0.35
1908M033_03	-0.412	1.05	1.07	.30	0.84	.01	0.45
1908M033_04	-0.857	0.99	0.99	.36	1.05	.00	0.54

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

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1911B002_01A	0.251	0.82	0.76	.61	1.42	.00	0.44
1911B002_05A	1.396	0.96	1.01	.42	1.04	.00	0.21
1911M069_01	1.911	0.97	0.82	.42	1.05	.00	0.16
1911M069_02	-0.005	0.99	0.97	.46	1.05	.03	0.49
1911M069_03	-0.542	0.94	0.91	.47	1.18	.00	0.57
1911M069_04	-0.430	0.92	0.90	.49	1.21	.00	0.54
HS18006_01	-1.518	0.89	0.88	.43	1.18	.00	0.74
HS18006_05	0.978	0.98	0.99	.46	1.03	.00	0.29
HS18006_07	-0.626	1.01	1.00	.38	0.98	.06	0.61
HS18018_02	-0.998	0.85	0.90	.50	1.30	.00	0.64
HS18018_03	1.041	1.13	1.13	.23	0.81	.02	0.24
HS18018_05	0.829	1.08	1.07	.42	0.88	.03	0.33
HS18018_08	0.786	0.98	0.91	.46	1.06	.00	0.32
HS18047_02	-0.475	0.84	0.79	.57	1.43	.00	0.52
HS18047_03	-0.119	1.03	1.05	.41	0.91	.05	0.47
HS18047_05	-0.212	1.01	1.02	.42	0.96	.03	0.50
HS18047_07	-1.185	0.86	0.78	.53	1.30	.00	0.66
HS18051_02	0.830	1.07	1.15	.36	0.86	.03	0.30
HS18051_04	-0.027	1.07	1.08	.37	0.83	.03	0.46
HS18051_05	0.769	1.07	1.09	.37	0.88	.03	0.32
HS18051_08	0.482	1.08	1.10	.36	0.83	.04	0.36
HS18082_03	0.303	1.14	1.21	.31	0.67	.08	0.39
HS18082_06	0.806	0.98	1.11	.44	0.99	.02	0.31
HS18082_07	0.225	0.94	0.95	.46	1.13	.00	0.36
HS19003_01A	-0.417	0.88	0.84	.53	1.32	.00	0.52
HS19003_03A	-0.448	0.76	0.70	.64	1.65	.00	0.52
HS19003_05A	2.270	1.13	1.32	.23	0.88	.02	0.12
HS19003_09A	-0.690	0.87	0.83	.52	1.33	.00	0.58
1911B002_08B	0.139	1.12	1.11	.70	0.97	.00	1.57
1911B002_10B	0.904	1.14	1.14	.29	0.79	.03	0.28
1911B009_01A	-0.439	0.86	0.81	.55	1.39	.00	0.54
1911B009_03A	0.479	0.79	0.74	.62	1.43	.00	0.37
1911B009_05A	1.377	0.92	0.85	.46	1.11	.00	0.21
1911B009_07A	0.683	1.18	1.14	.61	0.85	.00	1.21
1911B009_09A	1.331	1.02	1.26	.34	0.91	.02	0.22
1911M002_01	1.627	1.00	1.12	.34	0.98	.00	0.18
1911M002_04	0.909	1.23	1.56	.18	0.54	.10	0.28

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1911M002_05	0.166	1.12	1.17	.32	0.69	.07	0.42
1911M019_04	-0.549	1.09	1.10	.34	0.76	.11	0.56
1911M019_08	1.344	1.05	1.10	.35	0.93	.01	0.22

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

Item	Rasch	Infit	Outfit	Corr.	Discrim.	Lower Asym.	Item Mean
1911M019_10	0.436	1.02	1.08	.40	0.92	.04	0.37
1911M023_02	0.460	0.89	0.90	.54	1.22	.00	0.37
1911M023_05	-0.600	0.78	0.74	.59	1.56	.00	0.57
1911M023_06	0.420	1.13	1.16	.31	0.72	.05	0.37
1911M023_07	0.465	1.22	1.32	.23	0.52	.09	0.36
1911M028_01	-0.617	1.15	1.19	.27	0.59	.09	0.57
1911M028_03	1.373	0.94	1.01	.45	1.05	.00	0.23
1911M028_04	0.103	1.22	1.29	.23	0.44	.13	0.43
1911M028_06	0.576	1.19	1.33	.25	0.58	.10	0.35
1911M029_01	0.580	1.11	1.19	.32	0.76	.05	0.34
1911M029_02	-1.265	0.92	0.91	.43	1.15	.00	0.69
1911M029_03	-0.592	1.13	1.20	.29	0.64	.12	0.57
1911M029_04	1.034	1.25	1.66	.17	0.53	.10	0.27
1911M079_02	0.249	1.22	1.33	.24	0.47	.10	0.41
1911M079_03	-0.368	0.92	0.92	.48	1.19	.00	0.52
1911M079_04	1.241	0.93	0.92	.46	1.09	.00	0.23
1911M086_01	0.003	0.87	0.83	.55	1.34	.00	0.46
1911M086_02	-1.265	0.87	0.81	.48	1.26	.00	0.69
1911M150_01	-0.271	1.04	1.04	.40	0.89	.08	0.51
1911M150_02	0.185	0.94	0.95	.48	1.12	.00	0.42
1911M150_03	0.410	0.85	0.78	.57	1.35	.00	0.37
HS18060_01	1.747	0.82	0.65	.53	1.20	.00	0.17
HS18060_03	-0.032	1.03	1.07	.41	0.90	.05	0.46
HS18060_04	-0.192	0.95	0.93	.48	1.15	.00	0.49
HS18060_06	2.033	1.00	1.16	.32	0.98	.00	0.14
HS19000_01a	1.261	0.85	0.78	.54	1.20	.00	0.24
HS19000_06b	0.646	1.34	1.50	.11	0.33	.12	0.32
HS19000_09a	-0.406	0.81	0.77	.58	1.50	.00	0.53
HS19003_07A	0.706	1.05	1.01	.67	0.98	.00	1.28

## 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	-5.251	1.834	-61.346	100	15	100	115
1	-4.027	1.014	-9.371	100	15	100	115
2	-3.306	0.727	21.246	100	15	100	115
3	-2.872	0.601	39.675	100	15	100	115
4	-2.556	0.528	53.094	100	15	100	115
5	-2.304	0.478	63.795	100	15	100	115
6	-2.093	0.442	72.755	100	15	100	115
7	-1.911	0.414	80.483	100	15	100	115
8	-1.748	0.392	87.405	100	15	100	115
9	-1.601	0.374	93.647	100	15	100	115
10	-1.467	0.360	99.337	100	15	100	115
11	-1.342	0.347	104.645	105	15	100	120
12	-1.225	0.337	109.613	110	14	100	124
13	-1.115	0.328	114.284	114	14	100	128
14	-1.010	0.320	118.743	119	14	105	133
15	-0.910	0.313	122.990	123	13	110	136
16	-0.814	0.307	127.066	127	13	114	140
17	-0.721	0.302	131.015	131	13	118	144
18	-0.631	0.298	134.837	135	13	122	148
19	-0.543	0.294	138.574	139	12	127	151
20	-0.458	0.291	142.183	142	12	130	154
21	-0.374	0.288	145.750	146	12	134	158
22	-0.292	0.286	149.232	150	12	138	162
23	-0.211	0.284	152.672	153	12	141	165
24	-0.130	0.282	156.111	156	12	144	168
25	-0.051	0.281	159.466	159	12	147	171
26	0.027	0.280	162.778	163	12	151	175
27	0.106	0.279	166.133	166	12	154	178
28	0.184	0.279	169.445	169	12	157	181
29	0.262	0.279	172.757	173	12	161	185
30	0.339	0.279	176.027	176	12	164	188
31	0.418	0.280	179.382	179	12	167	191
32	0.496	0.280	182.694	183	12	171	195
33	0.575	0.281	186.048	186	12	174	198

Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS CSEM	Lower SS	Upper SS
34	0.654	0.283	189.403	189	12	177	201
35	0.735	0.284	192.843	193	12	181	205
36	0.816	0.286	196.282	196	12	184	208
37	0.898	0.288	199.764	200	12	188	212
38	0.981	0.290	203.289	203	12	191	215
39	1.066	0.293	206.898	207	12	195	219
40	1.153	0.296	210.593	211	13	198	224
41	1.241	0.299	214.329	214	13	201	227
42	1.332	0.303	218.194	218	13	205	231
43	1.425	0.307	222.143	222	13	209	235
44	1.520	0.311	226.177	226	13	213	239
45	1.618	0.316	230.338	230	13	217	243
46	1.720	0.322	234.670	235	14	221	249
47	1.826	0.328	239.171	243	14	229	257
48	1.936	0.335	243.842	244	14	230	258
49	2.051	0.344	248.725	249	15	234	264
50	2.172	0.353	253.863	254	15	239	269
51	2.301	0.365	259.341	259	15	244	274
52	2.440	0.380	265.244	265	16	249	281
53	2.591	0.399	271.656	272	17	255	289
54	2.760	0.424	278.832	279	18	261	297
55	2.954	0.457	287.070	287	19	268	300
56	3.184	0.505	296.837	297	21	276	300
57	3.475	0.578	309.194	300	25	275	300
58	3.879	0.705	326.349	300	25	275	300
59	4.567	0.996	355.564	300	25	275	300
60	5.764	1.823	406.394	300	25	275	300

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

**Table I.2: Grade 8 – Operational**

**NJSLA–S Grade 8**

<b>Raw Score</b>	<b>Theta</b>	<b>Standard Error</b>	<b>Unrounded Scale Score</b>	<b>Scale Score (SS)</b>	<b>SS SCME</b>	<b>Lower SS</b>	<b>Upper SS</b>
0	-5.452	1.833	-21.681	100	17	100	117
1	-4.229	1.013	24.524	100	17	100	117
2	-3.510	0.725	51.688	100	17	100	117
3	-3.078	0.599	68.009	100	17	100	117
4	-2.765	0.525	79.834	100	17	100	117
5	-2.517	0.474	89.204	100	17	100	117
6	-2.310	0.437	97.024	100	17	100	117
7	-2.131	0.409	103.787	104	15	100	119
8	-1.973	0.386	109.756	110	15	100	125
9	-1.831	0.368	115.121	115	14	101	129
10	-1.702	0.352	119.994	120	13	107	133
11	-1.582	0.339	124.528	125	13	112	138
12	-1.471	0.328	128.722	129	12	117	141
13	-1.367	0.318	132.651	133	12	121	145
14	-1.269	0.309	136.353	136	12	124	148
15	-1.176	0.301	139.867	140	11	129	151
16	-1.088	0.294	143.191	143	11	132	154
17	-1.003	0.288	146.403	146	11	135	157
18	-0.922	0.282	149.463	150	11	139	161
19	-0.844	0.277	152.410	152	10	142	162
20	-0.768	0.273	155.281	155	10	145	165
21	-0.694	0.269	158.077	158	10	148	168
22	-0.623	0.265	160.759	161	10	151	171
23	-0.553	0.262	163.404	163	10	153	173
24	-0.485	0.260	165.973	166	10	156	176
25	-0.419	0.257	168.466	168	10	158	178
26	-0.353	0.255	170.960	171	10	161	181
27	-0.289	0.253	173.378	173	10	163	183
28	-0.225	0.251	175.795	176	9	167	185
29	-0.162	0.250	178.176	178	9	169	187
30	-0.100	0.249	180.518	181	9	172	190
31	-0.039	0.248	182.823	183	9	174	192
32	0.023	0.247	185.165	185	9	176	194
33	0.083	0.246	187.432	187	9	178	196
34	0.144	0.246	189.736	190	9	181	199
35	0.204	0.246	192.003	192	9	183	201

Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS SCME	Lower SS	Upper SS
36	0.265	0.246	194.308	194	9	185	203
37	0.325	0.246	196.575	197	9	188	206
38	0.386	0.246	198.879	200	9	191	209
39	0.446	0.247	201.146	201	9	192	210
40	0.507	0.247	203.450	203	9	194	212
41	0.569	0.248	205.793	206	9	197	215
42	0.630	0.249	208.097	208	9	199	217
43	0.693	0.250	210.478	210	9	201	219
44	0.755	0.252	212.820	213	10	203	223
45	0.819	0.253	215.238	215	10	205	225
46	0.884	0.255	217.694	218	10	208	228
47	0.950	0.258	220.187	220	10	210	230
48	1.017	0.260	222.718	223	10	213	233
49	1.085	0.263	225.287	225	10	215	235
50	1.155	0.266	227.932	228	10	218	238
51	1.227	0.270	230.652	231	10	221	241
52	1.301	0.274	233.448	233	10	223	243
53	1.377	0.279	236.319	236	11	225	247
54	1.456	0.284	239.304	239	11	228	250
55	1.538	0.289	242.402	242	11	231	253
56	1.624	0.296	245.651	246	11	235	257
57	1.713	0.303	249.013	249	11	238	260
58	1.808	0.311	252.602	253	12	241	265
59	1.907	0.320	256.343	256	12	244	268
60	2.012	0.330	260.309	260	12	248	272
61	2.125	0.341	264.579	265	13	252	278
62	2.246	0.355	269.150	269	13	256	282
63	2.377	0.370	274.099	274	14	260	288
64	2.521	0.389	279.539	280	15	265	295
65	2.681	0.412	285.584	286	16	270	300
66	2.862	0.440	292.422	292	17	275	300
67	3.072	0.477	300.356	300	18	282	300
68	3.322	0.527	309.801	300	18	282	300
69	3.637	0.601	321.702	300	18	282	300
70	4.071	0.727	338.099	300	18	282	300
71	4.792	1.014	365.338	300	18	282	300
72	6.015	1.833	411.543	300	18	282	300

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	-5.461	1.832	-113.541	100	17	100	117
1	-4.241	1.012	-49.102	100	17	100	117
2	-3.524	0.723	-11.230	100	17	100	117
3	-3.096	0.597	11.376	100	17	100	117
4	-2.785	0.522	27.803	100	17	100	117
5	-2.539	0.472	40.796	100	17	100	117
6	-2.334	0.435	51.624	100	17	100	117
7	-2.157	0.407	60.973	100	17	100	117
8	-2.001	0.384	69.213	100	17	100	117
9	-1.861	0.365	76.608	100	17	100	117
10	-1.733	0.350	83.368	100	17	100	117
11	-1.616	0.336	89.548	100	17	100	117
12	-1.506	0.325	95.358	100	17	100	117
13	-1.404	0.315	100.746	101	17	100	118
14	-1.308	0.306	105.816	106	16	100	122
15	-1.217	0.298	110.623	111	16	100	127
16	-1.131	0.291	115.165	115	15	100	130
17	-1.048	0.284	119.549	120	15	105	135
18	-0.969	0.278	123.722	124	15	109	139
19	-0.893	0.273	127.736	128	14	114	142
20	-0.820	0.268	131.592	132	14	118	146
21	-0.749	0.264	135.342	135	14	121	149
22	-0.680	0.260	138.987	139	14	125	153
23	-0.613	0.257	142.526	143	14	129	157
24	-0.548	0.253	145.959	146	13	133	159
25	-0.485	0.251	149.286	149	13	136	162
26	-0.422	0.248	152.614	153	13	140	166
27	-0.362	0.246	155.783	158	13	145	171
28	-0.302	0.244	158.952	159	13	146	172
29	-0.243	0.242	162.069	162	13	149	175
30	-0.185	0.240	165.132	165	13	152	178
31	-0.128	0.239	168.143	168	13	155	181
32	-0.071	0.238	171.153	171	13	158	184
33	-0.015	0.237	174.111	174	13	161	187
34	0.041	0.236	177.069	177	12	165	189

Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS CSEM	Lower SS	Upper SS
35	0.097	0.236	180.027	180	12	168	192
36	0.152	0.235	182.932	183	12	171	195
37	0.208	0.235	185.890	186	12	174	198
38	0.263	0.235	188.795	189	12	177	201
39	0.318	0.235	191.700	192	12	180	204
40	0.374	0.236	194.658	195	12	183	207
41	0.430	0.236	197.616	200	12	188	212
42	0.486	0.237	200.574	201	13	188	214
43	0.542	0.238	203.531	204	13	191	217
44	0.599	0.239	206.542	207	13	194	220
45	0.656	0.240	209.553	210	13	197	223
46	0.715	0.242	212.669	213	13	200	226
47	0.773	0.243	215.733	216	13	203	229
48	0.833	0.245	218.902	219	13	206	232
49	0.893	0.247	222.071	222	13	209	235
50	0.955	0.249	225.346	225	13	212	238
51	1.017	0.251	228.620	229	13	216	242
52	1.081	0.253	232.001	232	13	219	245
53	1.145	0.256	235.381	235	14	221	249
54	1.212	0.259	238.920	239	14	225	253
55	1.279	0.262	242.459	242	14	228	256
56	1.349	0.265	246.156	246	14	232	260
57	1.420	0.269	249.907	250	14	236	264
58	1.493	0.273	253.762	254	14	240	268
59	1.569	0.277	257.777	258	15	243	273
60	1.647	0.282	261.896	262	15	247	277
61	1.728	0.287	266.175	266	15	251	281
62	1.812	0.293	270.612	271	15	256	286
63	1.901	0.300	275.312	275	16	259	291
64	1.993	0.308	280.172	280	16	264	296
65	2.090	0.316	285.295	285	17	268	300
66	2.193	0.326	290.736	291	17	274	300
67	2.303	0.337	296.546	297	18	279	300
68	2.421	0.350	302.778	300	18	282	300
69	2.549	0.366	309.539	300	18	282	300
70	2.689	0.384	316.934	300	18	282	300
71	2.845	0.406	325.174	300	18	282	300
72	3.021	0.435	334.470	300	18	282	300

Raw Score	Theta	Standard Error	Unrounded Scale Score	Scale Score (SS)	SS CSEM	Lower SS	Upper SS
73	3.226	0.472	345.298	300	18	282	300
74	3.471	0.522	358.238	300	18	282	300
75	3.782	0.596	374.665	300	18	282	300
76	4.209	0.723	397.219	300	18	282	300
77	4.925	1.011	435.037	300	18	282	300
78	6.145	1.832	499.476	300	18	282	300

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	-4.108	1.855	-6.891	-1.326	Below
1	-2.829	1.053	-4.409	-1.250	Below
2	-2.027	0.781	-3.199	-0.856	Below
3	-1.510	0.668	-2.512	-0.508	Below
4	-1.107	0.606	-2.016	-0.198	Below
5	-0.764	0.568	-1.616	0.088	Below
6	-0.456	0.544	-1.272	0.360	Below
7	-0.168	0.530	-0.963	0.627	Below
8	0.108	0.522	-0.675	0.891	Below
9	0.379	0.519	-0.400	1.158	Near/Met
10	0.649	0.521	-0.133	1.431	Near/Met
11	0.923	0.525	0.136	1.711	Near/Met
12	1.202	0.532	0.404	2.000	Near/Met
13	1.489	0.540	0.679	2.299	Near/Met
14	1.787	0.551	0.961	2.614	Above
15	2.100	0.571	1.244	2.957	Above
16	2.448	0.614	1.527	3.369	Above
17	2.878	0.711	1.812	3.945	Above
18	3.557	0.983	2.083	5.032	Above
19	4.728	1.810	2.013	7.443	Above

**Table J.2: Grade 5 Life Science Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.310	1.847	-7.081	-1.540	Below
1	-3.052	1.038	-4.609	-1.495	Below
2	-2.281	0.760	-3.421	-1.141	Below
3	-1.799	0.641	-2.761	-0.838	Below
4	-1.433	0.573	-2.293	-0.574	Below
5	-1.132	0.528	-1.924	-0.340	Below
6	-0.869	0.498	-1.616	-0.122	Below
7	-0.632	0.477	-1.348	0.084	Below
8	-0.412	0.462	-1.105	0.281	Below
9	-0.204	0.452	-0.882	0.474	Below
10	-0.002	0.446	-0.671	0.667	Below
11	0.196	0.444	-0.470	0.862	Below
12	0.393	0.445	-0.275	1.061	Near/Met
13	0.592	0.449	-0.082	1.266	Near/Met
14	0.796	0.456	0.112	1.480	Near/Met
15	1.010	0.468	0.308	1.712	Near/Met
16	1.237	0.487	0.507	1.968	Near/Met
17	1.487	0.515	0.715	2.260	Near/Met
18	1.772	0.557	0.937	2.608	Above
19	2.118	0.624	1.182	3.054	Above
20	2.579	0.745	1.462	3.697	Above
21	3.326	1.027	1.786	4.867	Above
22	4.568	1.841	1.807	7.330	Above

**Table J.3: Grade 5 Physical Science Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-3.958	1.843	-6.723	-1.194	Below
1	-2.708	1.032	-4.256	-1.160	Below
2	-1.947	0.755	-3.080	-0.815	Below
3	-1.469	0.640	-2.429	-0.509	Below
4	-1.102	0.577	-1.968	-0.237	Below
5	-0.792	0.539	-1.601	0.017	Below
6	-0.515	0.516	-1.289	0.259	Below
7	-0.257	0.502	-1.010	0.496	Below
8	-0.008	0.496	-0.752	0.736	Below
9	0.237	0.496	-0.507	0.981	Near/Met
10	0.484	0.500	-0.266	1.234	Near/Met
11	0.739	0.509	-0.025	1.503	Near/Met
12	1.005	0.523	0.221	1.790	Near/Met
13	1.287	0.541	0.476	2.099	Near/Met
14	1.593	0.565	0.746	2.441	Near/Met
15	1.930	0.598	1.033	2.827	Above
16	2.316	0.648	1.344	3.288	Above
17	2.793	0.745	1.676	3.911	Above
18	3.518	1.002	2.015	5.021	Above
19	4.707	1.811	1.991	7.424	Above

**Table J.4: Grade 5 Sensemaking Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.659	1.843	-7.424	-4.659	Below
1	-3.411	1.031	-4.958	-3.411	Below
2	-2.654	0.751	-3.781	-2.654	Below
3	-2.184	0.632	-3.132	-2.184	Below
4	-1.830	0.563	-2.675	-1.830	Below
5	-1.539	0.518	-2.316	-1.539	Below
6	-1.287	0.487	-2.018	-1.287	Below
7	-1.061	0.464	-1.757	-1.061	Below
8	-0.854	0.447	-1.525	-0.854	Below
9	-0.660	0.435	-1.313	-0.660	Below
10	-0.475	0.426	-1.114	-0.475	Below
11	-0.296	0.419	-0.925	-0.296	Below
12	-0.122	0.415	-0.745	-0.122	Below
13	0.049	0.414	-0.572	0.049	Below
14	0.220	0.414	-0.401	0.220	Below
15	0.392	0.416	-0.232	0.392	Near/Met
16	0.567	0.420	-0.063	0.567	Near/Met
17	0.745	0.426	0.106	0.745	Near/Met
18	0.930	0.433	0.281	0.930	Near/Met
19	1.121	0.443	0.457	1.121	Near/Met
20	1.323	0.454	0.642	1.323	Near/Met
21	1.535	0.468	0.833	1.535	Near/Met
22	1.761	0.484	1.035	1.761	Above
23	2.005	0.505	1.248	2.005	Above
24	2.275	0.536	1.471	2.275	Above
25	2.589	0.591	1.703	2.589	Above
26	2.997	0.699	1.949	2.997	Above
27	3.662	0.977	2.197	3.662	Above
28	4.826	1.806	2.117	4.826	Above

**Table J.5: Grade 5 Critiquing Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-3.946	1.847	-6.717	-3.946	Below
1	-2.689	1.037	-4.245	-2.689	Below
2	-1.921	0.759	-3.060	-1.921	Below
3	-1.439	0.642	-2.402	-1.439	Below
4	-1.071	0.577	-1.937	-1.071	Below
5	-0.761	0.539	-1.570	-0.761	Below
6	-0.483	0.517	-1.259	-0.483	Below
7	-0.223	0.505	-0.981	-0.223	Below
8	0.029	0.500	-0.721	0.029	Below
9	0.280	0.502	-0.473	0.280	Near/Met
10	0.534	0.508	-0.228	0.534	Near/Met
11	0.798	0.519	0.020	0.798	Near/Met
12	1.076	0.537	0.271	1.076	Near/Met
13	1.376	0.562	0.533	1.376	Near/Met
14	1.711	0.598	0.814	1.711	Above
15	2.099	0.651	1.123	2.099	Above
16	2.580	0.746	1.461	2.580	Above
17	3.299	0.993	1.810	3.299	Above
18	4.460	1.792	1.772	4.460	Above

**Table J.6: Grade 5 Investigating Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-3.428	1.856	-6.212	-3.428	Below
1	-2.145	1.056	-3.729	-2.145	Below
2	-1.333	0.788	-2.515	-1.333	Below
3	-0.803	0.679	-1.822	-0.803	Below
4	-0.383	0.622	-1.316	-0.383	Below
5	-0.019	0.589	-0.903	-0.019	Below
6	0.317	0.572	-0.541	0.317	Near/Met
7	0.640	0.567	-0.211	0.640	Near/Met
8	0.963	0.572	0.105	0.963	Near/Met
9	1.298	0.588	0.416	1.298	Near/Met
10	1.662	0.620	0.732	1.662	Near/Met
11	2.080	0.678	1.063	2.080	Above
12	2.607	0.787	1.427	2.607	Above
13	3.416	1.055	1.834	3.416	Above
14	4.698	1.855	1.916	4.698	Above

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

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.402	1.845	-7.170	-1.635	Below
1	-3.149	1.036	-4.703	-1.595	Below
2	-2.381	0.759	-3.520	-1.243	Below
3	-1.897	0.643	-2.862	-0.933	Below
4	-1.527	0.579	-2.396	-0.659	Below
5	-1.217	0.537	-2.023	-0.412	Below
6	-0.944	0.510	-1.709	-0.179	Below
7	-0.694	0.490	-1.429	0.041	Below
8	-0.462	0.475	-1.175	0.251	Below
9	-0.241	0.464	-0.937	0.455	Near/Met
10	-0.030	0.456	-0.714	0.654	Near/Met
11	0.175	0.449	-0.499	0.849	Near/Met
12	0.375	0.445	-0.293	1.043	Near/Met
13	0.572	0.444	-0.094	1.238	Near/Met
14	0.771	0.448	0.099	1.443	Near/Met
15	0.976	0.459	0.288	1.665	Near/Met
16	1.195	0.478	0.478	1.912	Above
17	1.438	0.510	0.673	2.203	Above
18	1.722	0.559	0.884	2.561	Above
19	2.077	0.636	1.123	3.031	Above
20	2.558	0.763	1.414	3.703	Above
21	3.342	1.048	1.770	4.914	Above
22	4.618	1.855	1.836	7.401	Above

**Table J.8: Grade 8 Life Science Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.341	1.840	-7.101	-1.581	Below
1	-3.099	1.026	-4.638	-1.560	Below
2	-2.352	0.745	-3.470	-1.235	Below
3	-1.891	0.624	-2.827	-0.955	Below
4	-1.547	0.554	-2.378	-0.716	Below
5	-1.266	0.509	-2.030	-0.503	Below
6	-1.025	0.476	-1.739	-0.311	Below
7	-0.809	0.453	-1.489	-0.130	Below
8	-0.613	0.435	-1.266	0.040	Below
9	-0.430	0.422	-1.063	0.203	Below
10	-0.256	0.412	-0.874	0.362	Below
11	-0.089	0.405	-0.697	0.519	Near/Met
12	0.073	0.402	-0.530	0.676	Near/Met
13	0.234	0.400	-0.366	0.834	Near/Met
14	0.394	0.401	-0.208	0.996	Near/Met
15	0.556	0.405	-0.052	1.164	Near/Met
16	0.722	0.411	0.106	1.339	Near/Met
17	0.894	0.419	0.266	1.523	Near/Met
18	1.074	0.430	0.429	1.719	Above
19	1.266	0.445	0.599	1.934	Above
20	1.472	0.464	0.776	2.168	Above
21	1.698	0.488	0.966	2.430	Above
22	1.952	0.521	1.171	2.734	Above
23	2.247	0.567	1.397	3.098	Above
24	2.606	0.636	1.652	3.560	Above
25	3.084	0.756	1.950	4.218	Above
26	3.849	1.035	2.297	5.402	Above
27	5.103	1.845	2.336	7.871	Above

**Table J.9: Grade 8 Physical Science Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.263	1.859	-7.052	-1.475	Below
1	-2.977	1.054	-4.558	-1.396	Below
2	-2.183	0.770	-3.338	-1.028	Below
3	-1.692	0.643	-2.657	-0.728	Below
4	-1.329	0.566	-2.178	-0.480	Below
5	-1.038	0.516	-1.812	-0.264	Below
6	-0.790	0.481	-1.512	-0.069	Below
7	-0.571	0.458	-1.258	0.116	Below
8	-0.368	0.443	-1.033	0.297	Below
9	-0.176	0.435	-0.829	0.477	Near/Met
10	0.011	0.432	-0.637	0.659	Near/Met
11	0.198	0.432	-0.450	0.846	Near/Met
12	0.386	0.436	-0.268	1.040	Near/Met
13	0.579	0.443	-0.086	1.244	Near/Met
14	0.779	0.452	0.101	1.457	Near/Met
15	0.990	0.465	0.293	1.688	Near/Met
16	1.213	0.481	0.492	1.935	Above
17	1.455	0.503	0.701	2.210	Above
18	1.723	0.534	0.922	2.524	Above
19	2.030	0.577	1.165	2.896	Above
20	2.400	0.644	1.434	3.366	Above
21	2.887	0.762	1.744	4.030	Above
22	3.661	1.039	2.103	5.220	Above
23	4.921	1.848	2.149	7.693	Above

**Table J.10: Grade 8 Sensemaking Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.714	1.846	-7.483	-1.945	Below
1	-3.458	1.035	-5.011	-1.906	Below
2	-2.695	0.754	-3.826	-1.564	Below
3	-2.223	0.631	-3.170	-1.277	Below
4	-1.872	0.558	-2.709	-1.035	Below
5	-1.588	0.509	-2.352	-0.825	Below
6	-1.348	0.473	-2.058	-0.639	Below
7	-1.138	0.446	-1.807	-0.469	Below
8	-0.949	0.424	-1.585	-0.313	Below
9	-0.776	0.408	-1.388	-0.164	Below
10	-0.615	0.395	-1.208	-0.023	Below
11	-0.463	0.385	-1.041	0.115	Below
12	-0.318	0.378	-0.885	0.249	Below
13	-0.177	0.373	-0.737	0.383	Below
14	-0.039	0.370	-0.594	0.516	Near/Met
15	0.097	0.368	-0.455	0.649	Near/Met
16	0.232	0.367	-0.319	0.783	Near/Met
17	0.367	0.367	-0.184	0.918	Near/Met
18	0.502	0.369	-0.052	1.056	Near/Met
19	0.639	0.371	0.083	1.196	Near/Met
20	0.779	0.376	0.215	1.343	Near/Met
21	0.922	0.382	0.349	1.495	Near/Met
22	1.071	0.391	0.485	1.658	Above
23	1.229	0.404	0.623	1.835	Above
24	1.399	0.421	0.768	2.031	Above
25	1.585	0.442	0.922	2.248	Above
26	1.792	0.470	1.087	2.497	Above
27	2.031	0.507	1.271	2.792	Above
28	2.313	0.558	1.476	3.150	Above
29	2.664	0.631	1.718	3.611	Above
30	3.136	0.754	2.005	4.267	Above
31	3.899	1.035	2.347	5.452	Above
32	5.154	1.846	2.385	7.923	Above

**Table J.11: Grade 8 Critiquing Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-3.666	1.844	-6.432	-0.900	Below
1	-2.414	1.034	-3.965	-0.863	Below
2	-1.652	0.754	-2.783	-0.521	Below
3	-1.177	0.636	-2.131	-0.223	Below
4	-0.817	0.569	-1.671	0.037	Below
5	-0.519	0.526	-1.308	0.270	Below
6	-0.257	0.499	-1.006	0.492	Near/Met
7	-0.017	0.482	-0.740	0.706	Near/Met
8	0.211	0.473	-0.499	0.921	Near/Met
9	0.433	0.471	-0.274	1.140	Near/Met
10	0.656	0.474	-0.055	1.367	Near/Met
11	0.885	0.483	0.161	1.610	Near/Met
12	1.125	0.497	0.380	1.871	Above
13	1.383	0.518	0.606	2.160	Above
14	1.666	0.548	0.844	2.488	Above
15	1.988	0.590	1.103	2.873	Above
16	2.373	0.656	1.389	3.357	Above
17	2.875	0.772	1.717	4.033	Above
18	3.663	1.046	2.094	5.232	Above
19	4.932	1.851	2.156	7.709	Above

**Table J.12: Grade 8 Investigating Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.379	1.842	-7.142	-1.616	Below
1	-3.131	1.032	-4.679	-1.583	Below
2	-2.371	0.754	-3.502	-1.240	Below
3	-1.895	0.637	-2.851	-0.940	Below
4	-1.532	0.572	-2.390	-0.674	Below
5	-1.230	0.531	-2.027	-0.434	Below
6	-0.963	0.503	-1.718	-0.209	Below
7	-0.720	0.485	-1.448	0.008	Below
8	-0.491	0.472	-1.199	0.217	Below
9	-0.272	0.465	-0.970	0.426	Near/Met
10	-0.058	0.461	-0.750	0.634	Near/Met
11	0.155	0.462	-0.538	0.848	Near/Met
12	0.369	0.466	-0.330	1.068	Near/Met
13	0.590	0.474	-0.121	1.301	Near/Met
14	0.820	0.487	0.090	1.551	Near/Met
15	1.066	0.506	0.307	1.825	Near/Met
16	1.336	0.534	0.535	2.137	Above
17	1.643	0.576	0.779	2.507	Above
18	2.011	0.642	1.048	2.974	Above
19	2.494	0.759	1.356	3.633	Above
20	3.263	1.036	1.709	4.817	Above
21	4.518	1.846	1.749	7.287	Above

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

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.094	1.848	-6.866	-1.322	Below
1	-2.835	1.038	-4.392	-1.278	Below
2	-2.067	0.756	-3.201	-0.933	Below
3	-1.592	0.633	-2.542	-0.643	Below
4	-1.240	0.559	-2.079	-0.402	Below
5	-0.956	0.509	-1.720	-0.193	Below
6	-0.716	0.472	-1.424	-0.008	Below
7	-0.506	0.445	-1.174	0.162	Below
8	-0.317	0.426	-0.956	0.322	Below
9	-0.141	0.413	-0.761	0.479	Below
10	0.025	0.405	-0.583	0.633	Near/Met
11	0.188	0.401	-0.414	0.790	Near/Met
12	0.349	0.402	-0.254	0.952	Near/Met
13	0.512	0.407	-0.098	1.123	Near/Met
14	0.681	0.415	0.059	1.304	Near/Met
15	0.858	0.427	0.218	1.499	Near/Met
16	1.047	0.442	0.384	1.710	Near/Met
17	1.251	0.462	0.558	1.944	Above
18	1.475	0.486	0.746	2.204	Above
19	1.727	0.519	0.949	2.506	Above
20	2.019	0.565	1.172	2.867	Above
21	2.376	0.634	1.425	3.327	Above
22	2.849	0.753	1.720	3.979	Above
23	3.609	1.033	2.060	5.159	Above
24	4.859	1.844	2.093	7.625	Above

**Table J.14: Grade 11 Life Science Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.549	1.844	-7.315	-1.783	Below
1	-3.298	1.033	-4.848	-1.749	Below
2	-2.538	0.753	-3.668	-1.409	Below
3	-2.065	0.633	-3.015	-1.116	Below
4	-1.710	0.564	-2.556	-0.864	Below
5	-1.418	0.519	-2.197	-0.640	Below
6	-1.166	0.487	-1.897	-0.436	Below
7	-0.940	0.464	-1.636	-0.244	Below
8	-0.732	0.448	-1.404	-0.060	Below
9	-0.537	0.435	-1.190	0.116	Below
10	-0.352	0.426	-0.991	0.287	Below
11	-0.173	0.421	-0.805	0.459	Below
12	0.002	0.417	-0.624	0.628	Near/Met
13	0.176	0.416	-0.448	0.800	Near/Met
14	0.350	0.418	-0.277	0.977	Near/Met
15	0.526	0.421	-0.106	1.158	Near/Met
16	0.705	0.427	0.065	1.346	Near/Met
17	0.890	0.435	0.238	1.543	Near/Met
18	1.084	0.446	0.415	1.753	Near/Met
19	1.289	0.460	0.599	1.979	Above
20	1.510	0.481	0.789	2.232	Above
21	1.755	0.510	0.990	2.520	Above
22	2.036	0.554	1.205	2.867	Above
23	2.379	0.622	1.446	3.312	Above
24	2.837	0.743	1.723	3.952	Above
25	3.582	1.026	2.043	5.121	Above
26	4.822	1.840	2.062	7.582	Above

**Table J.15: Grade 11 Physical Science Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.353	1.843	-7.118	-1.589	Below
1	-3.105	1.030	-4.650	-1.560	Below
2	-2.351	0.749	-3.475	-1.228	Below
3	-1.885	0.628	-2.827	-0.943	Below
4	-1.537	0.557	-2.373	-0.702	Below
5	-1.254	0.510	-2.019	-0.489	Below
6	-1.011	0.477	-1.727	-0.296	Below
7	-0.795	0.453	-1.475	-0.116	Below
8	-0.598	0.435	-1.251	0.055	Below
9	-0.414	0.423	-1.049	0.221	Below
10	-0.240	0.414	-0.861	0.381	Below
11	-0.071	0.408	-0.683	0.541	Near/Met
12	0.094	0.405	-0.514	0.702	Near/Met
13	0.258	0.404	-0.348	0.864	Near/Met
14	0.422	0.406	-0.187	1.031	Near/Met
15	0.587	0.409	-0.027	1.201	Near/Met
16	0.756	0.413	0.137	1.376	Near/Met
17	0.930	0.420	0.300	1.560	Near/Met
18	1.109	0.428	0.467	1.751	Above
19	1.297	0.439	0.639	1.956	Above
20	1.495	0.452	0.817	2.173	Above
21	1.706	0.469	1.003	2.410	Above
22	1.937	0.492	1.199	2.675	Above
23	2.193	0.523	1.409	2.978	Above
24	2.489	0.567	1.639	3.340	Above
25	2.848	0.635	1.896	3.801	Above
26	3.323	0.754	2.192	4.454	Above
27	4.084	1.033	2.535	5.634	Above
28	5.336	1.844	2.570	8.102	Above

**Table J.16: Grade 11 Sensemaking Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.728	1.839	-7.487	-1.970	Below
1	-3.489	1.025	-5.027	-1.952	Below
2	-2.745	0.742	-3.858	-1.632	Below
3	-2.290	0.619	-3.219	-1.362	Below
4	-1.952	0.548	-2.774	-1.130	Below
5	-1.678	0.501	-2.430	-0.927	Below
6	-1.445	0.467	-2.146	-0.745	Below
7	-1.238	0.442	-1.901	-0.575	Below
8	-1.052	0.422	-1.685	-0.419	Below
9	-0.881	0.407	-1.492	-0.271	Below
10	-0.720	0.395	-1.313	-0.128	Below
11	-0.569	0.385	-1.147	0.009	Below
12	-0.424	0.377	-0.990	0.142	Below
13	-0.284	0.371	-0.841	0.273	Below
14	-0.148	0.367	-0.699	0.403	Below
15	-0.015	0.363	-0.560	0.530	Near/Met
16	0.117	0.361	-0.425	0.659	Near/Met
17	0.247	0.361	-0.295	0.789	Near/Met
18	0.377	0.361	-0.165	0.919	Near/Met
19	0.507	0.362	-0.036	1.050	Near/Met
20	0.639	0.364	0.093	1.185	Near/Met
21	0.772	0.367	0.222	1.323	Near/Met
22	0.908	0.371	0.352	1.465	Near/Met
23	1.048	0.377	0.483	1.614	Above
24	1.192	0.384	0.616	1.768	Above
25	1.343	0.393	0.754	1.933	Above
26	1.501	0.404	0.895	2.107	Above
27	1.670	0.419	1.042	2.299	Above
28	1.853	0.438	1.196	2.510	Above
29	2.055	0.462	1.362	2.748	Above
30	2.284	0.496	1.540	3.028	Above
31	2.553	0.544	1.737	3.369	Above
32	2.886	0.615	1.964	3.809	Above
33	3.337	0.738	2.230	4.444	Above
34	4.076	1.023	2.542	5.611	Above
35	5.312	1.839	2.554	8.071	Above

**Table J.17: Grade 11 Critiquing Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.009	1.846	-6.778	-1.240	Below
1	-2.755	1.034	-4.306	-1.204	Below
2	-1.997	0.749	-3.121	-0.874	Below
3	-1.534	0.623	-2.469	-0.600	Below
4	-1.195	0.547	-2.016	-0.375	Below
5	-0.924	0.496	-1.668	-0.180	Below
6	-0.696	0.460	-1.386	-0.006	Below
7	-0.496	0.435	-1.149	0.157	Below
8	-0.315	0.417	-0.941	0.311	Below
9	-0.147	0.405	-0.755	0.461	Below
10	0.015	0.399	-0.584	0.614	Near/Met
11	0.173	0.398	-0.424	0.770	Near/Met
12	0.332	0.400	-0.268	0.932	Near/Met
13	0.495	0.407	-0.116	1.106	Near/Met
14	0.664	0.417	0.039	1.290	Near/Met
15	0.844	0.431	0.198	1.491	Near/Met
16	1.036	0.448	0.364	1.708	Near/Met
17	1.246	0.469	0.543	1.950	Above
18	1.478	0.494	0.737	2.219	Above
19	1.738	0.527	0.948	2.529	Above
20	2.038	0.572	1.180	2.896	Above
21	2.402	0.639	1.444	3.361	Above
22	2.881	0.756	1.747	4.015	Above
23	3.644	1.034	2.093	5.195	Above
24	4.896	1.844	2.130	7.662	Above

**Table J.18: Grade 11 Investigating Score Table**

<b>Raw Score</b>	<b>Theta</b>	<b>CSEM</b>	<b>Lower</b>	<b>Upper</b>	<b>Level</b>
0	-4.160	1.852	-6.938	-1.382	Below
1	-2.889	1.047	-4.460	-1.319	Below
2	-2.099	0.773	-3.259	-0.940	Below
3	-1.594	0.659	-2.583	-0.606	Below
4	-1.205	0.595	-2.098	-0.313	Below
5	-0.876	0.554	-1.707	-0.045	Below
6	-0.585	0.528	-1.377	0.207	Below
7	-0.316	0.511	-1.083	0.451	Below
8	-0.061	0.500	-0.811	0.689	Near/Met
9	0.186	0.495	-0.557	0.929	Near/Met
10	0.431	0.495	-0.312	1.174	Near/Met
11	0.679	0.501	-0.073	1.431	Near/Met
12	0.934	0.512	0.166	1.702	Near/Met
13	1.205	0.529	0.412	1.999	Near/Met
14	1.499	0.557	0.664	2.335	Above
15	1.830	0.598	0.933	2.727	Above
16	2.224	0.663	1.230	3.219	Above
17	2.736	0.778	1.569	3.903	Above
18	3.535	1.052	1.957	5.113	Above
19	4.815	1.855	2.033	7.598	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	96,288	57.83	30.94	11.22	59.59	28.86	11.55	54.26	36.77	8.96
Male	49,121	55.53	31.80	12.67	57.87	29.53	12.60	52.69	37.64	9.67
Female	47,156	60.24	30.04	9.72	61.38	28.17	10.45	55.91	35.86	8.23
Am. Indian	150	60.67	24.67	14.67	61.33	26.67	12.00	52.00	43.33	4.67
Asian	10,437	27.38	41.20	31.42	27.56	40.86	31.58	24.74	51.27	23.99
Black	13,848	78.42	18.72	2.85	78.91	17.59	3.50	74.83	22.46	2.71
Hispanic	31,155	73.72	22.16	4.11	76.28	19.32	4.40	69.83	26.62	3.56
Pacific Islander	195	51.28	36.41	12.31	49.74	36.92	13.33	46.67	44.62	8.72
White	37,680	46.28	39.44	14.28	48.21	37.18	14.61	42.76	45.90	11.34
EL–Yes	8,533	90.54	8.77	0.69	92.82	6.52	0.67	89.01	10.51	0.48
EL–No	87,752	54.65	33.10	12.25	56.35	31.04	12.61	50.88	39.33	9.79
EconDis–Yes	32,449	77.58	19.33	3.09	79.49	17.13	3.38	73.16	24.06	2.78
EconDis–No	63,836	47.79	36.85	15.36	49.47	34.83	15.70	44.66	43.23	12.11
SWD–Yes	19,887	77.06	18.24	4.70	77.61	17.32	5.07	74.52	21.77	3.71
SWD–No	76,399	52.83	34.25	12.92	54.89	31.87	13.24	48.99	40.68	10.33
CBT	77,237	53.27	33.96	12.77	54.87	31.92	13.21	49.42	40.40	10.18
PBT	47	65.96	27.66	6.38	65.96	29.79	4.26	72.34	21.28	6.38
TTS	16,877	74.50	20.00	5.50	77.13	17.57	5.30	72.09	23.48	4.43
SP	1,204	90.95	8.64	0.42	91.69	7.31	1.00	88.79	10.63	0.58
SP TTS	561	93.58	6.42	0.00	94.30	5.53	0.18	90.91	8.91	0.18
Human Reader	308	87.99	10.06	1.95	88.31	10.71	0.97	82.79	15.26	1.95

**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	96,288	55.24	37.33	7.43	60.02	29.20	10.78	54.87	32.95	12.18
Male	49,121	53.87	37.78	8.35	57.27	30.35	12.39	54.03	33.24	12.73
Female	47,156	56.67	36.85	6.48	62.89	28.00	9.10	55.74	32.65	11.61
Am. Indian	150	52.67	44.67	2.67	56.67	31.33	12.00	62.00	28.67	9.33
Asian	10,437	24.83	53.92	21.25	28.49	41.20	30.31	24.17	43.07	32.76
Black	13,848	75.76	22.19	2.05	80.13	16.98	2.89	75.63	20.74	3.63
Hispanic	31,155	70.87	26.41	2.72	76.58	19.51	3.91	71.12	24.00	4.88
Pacific Islander	195	47.18	43.59	9.23	51.79	33.85	14.36	48.21	38.97	12.82
White	37,680	43.85	46.80	9.35	48.46	37.90	13.64	43.01	41.65	15.34
EL–Yes	8,533	88.90	10.69	0.41	93.14	6.28	0.57	89.53	9.69	0.77
EL–No	87,752	51.96	39.92	8.12	56.80	31.43	11.77	51.50	35.21	13.29
EconDis–Yes	32,449	74.34	23.48	2.18	80.08	17.08	2.84	74.89	21.39	3.73
EconDis–No	63,836	45.53	44.37	10.11	49.83	35.36	14.81	44.69	38.83	16.48
SWD–Yes	19,887	75.26	21.60	3.14	78.32	17.21	4.47	74.96	20.14	4.90
SWD–No	76,399	50.03	41.42	8.55	55.26	32.32	12.42	49.64	36.29	14.07
CBT	77,237	50.55	40.96	8.49	55.39	32.30	12.31	49.93	36.15	13.92
PBT	47	65.96	31.91	2.13	74.47	19.15	6.38	72.34	21.28	6.38
TTS	16,877	72.64	23.83	3.53	77.03	17.91	5.05	73.19	21.16	5.65
SP	1,204	87.87	11.96	0.17	93.85	5.65	0.50	88.54	10.55	0.91
SP TTS	561	88.95	11.05	0.00	94.30	5.53	0.18	90.73	8.56	0.71
Human Reader	308	82.79	16.23	0.97	88.31	10.06	1.62	86.36	12.99	0.65

**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	101,371	62.00	33.73	4.27	70.16	22.88	6.95	66.82	26.77	6.42
Male	51,908	62.26	32.94	4.81	70.31	22.21	7.48	64.84	27.12	8.05
Female	49,405	61.75	34.55	3.69	70.03	23.58	6.39	68.91	26.39	4.71
Am. Indian	149	58.39	38.93	2.68	65.77	26.17	8.05	63.76	26.17	10.07
Asian	10,766	30.93	54.88	14.19	38.48	39.44	22.08	32.60	46.06	21.34
Black	15,042	81.05	18.14	0.81	87.71	10.83	1.46	86.05	12.62	1.33
Hispanic	31,831	77.31	21.58	1.12	84.16	13.75	2.10	82.57	15.65	1.77
Pacific Islander	191	56.54	39.27	4.19	62.83	26.70	10.47	61.78	30.89	7.33
White	40,913	51.74	43.07	5.18	61.68	29.79	8.52	57.05	35.18	7.77
EL–Yes	5,908	93.35	6.55	0.10	96.28	3.42	0.30	95.38	4.37	0.25
EL–No	95,454	60.06	35.42	4.52	68.55	24.09	7.36	65.05	28.16	6.80
EconDis–Yes	31,966	79.72	19.49	0.79	86.49	11.90	1.61	84.63	13.95	1.42
EconDis–No	69,392	53.84	40.30	5.87	62.64	27.95	9.41	58.60	32.68	8.72
SWD–Yes	20,268	80.13	18.29	1.58	85.67	11.55	2.78	84.30	13.31	2.39
SWD–No	81,097	57.47	37.59	4.94	66.29	25.72	7.99	62.45	30.13	7.42
CBT	84,156	58.65	36.54	4.81	67.12	25.06	7.81	63.39	29.26	7.35
PBT	56	83.93	16.07	0.00	82.14	17.86	0.00	82.14	17.86	0.00
TTS	15,049	76.24	21.95	1.81	83.52	13.40	3.08	81.76	16.13	2.11
SP	1,494	93.84	6.16	0.00	95.58	4.08	0.33	96.52	3.41	0.07
SP TTS	416	93.03	6.97	0.00	97.12	2.64	0.24	96.39	3.61	0.00
Human Reader	173	94.22	5.78	0.00	96.53	2.89	0.58	98.27	1.73	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	101,371	63.37	32.06	4.57	71.23	23.33	5.44	62.73	29.77	7.50
Male	51,908	63.20	31.36	5.43	69.73	23.70	6.57	63.39	28.43	8.19
Female	49,405	63.57	32.77	3.67	72.82	22.92	4.26	62.06	31.16	6.78
Am. Indian	149	60.40	36.24	3.36	68.46	26.85	4.70	57.72	31.54	10.74
Asian	10,766	31.24	52.47	16.29	38.97	42.98	18.06	30.34	46.15	23.52
Black	15,042	82.64	16.58	0.78	89.12	9.85	1.03	81.04	17.08	1.88
Hispanic	31,831	78.47	20.49	1.04	86.07	12.57	1.36	77.30	20.24	2.46
Pacific Islander	191	56.54	38.74	4.71	68.59	24.61	6.81	58.12	29.84	12.04
White	40,913	53.55	40.98	5.47	62.09	31.20	6.71	53.64	37.35	9.01
EL–Yes	5,908	94.08	5.82	0.10	97.02	2.79	0.19	90.93	8.73	0.34
EL–No	95,454	61.47	33.68	4.85	69.63	24.60	5.77	60.98	31.07	7.95
EconDis–Yes	31,966	80.86	18.37	0.77	88.21	10.77	1.02	79.86	18.26	1.87
EconDis–No	69,392	55.31	38.37	6.32	63.41	29.11	7.48	54.83	35.07	10.10
SWD–Yes	20,268	81.09	17.27	1.63	87.22	10.78	2.00	80.01	16.93	3.06
SWD–No	81,097	58.94	35.75	5.31	67.23	26.47	6.30	58.41	32.98	8.61
CBT	84,156	59.99	34.86	5.15	68.07	25.72	6.21	59.36	32.19	8.45
PBT	56	82.14	17.86	0.00	87.50	12.50	0.00	82.14	17.86	0.00
TTS	15,049	77.73	20.25	2.01	85.14	12.95	1.91	77.59	19.17	3.24
SP	1,494	95.72	4.22	0.07	97.72	2.28	0.00	90.29	9.37	0.33
SP TTS	416	94.23	5.77	0.00	97.60	2.40	0.00	90.63	9.38	0.00
Human Reader	173	94.22	5.78	0.00	97.11	2.89	0.00	93.64	5.78	0.58

**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	92,925	53.85	34.56	11.59	56.21	28.92	14.87	56.33	29.86	13.81
Male	47,026	54.22	32.93	12.85	55.85	28.30	15.85	54.94	28.87	16.20
Female	45,829	53.49	36.23	10.28	56.60	29.55	13.85	57.78	30.87	11.34
Am. Indian	121	59.50	36.36	4.13	63.64	26.45	9.92	61.16	28.93	9.92
Asian	10,025	24.86	44.36	30.78	26.98	35.61	37.41	26.31	36.18	37.51
Black	12,436	75.16	21.61	3.23	76.56	18.65	4.79	77.05	19.17	3.78
Hispanic	27,349	69.83	25.87	4.30	71.87	21.83	6.30	71.91	22.98	5.11
Pacific Islander	202	48.51	40.59	10.89	49.01	32.18	18.81	52.97	31.68	15.35
White	40,898	44.15	41.70	14.14	47.04	35.08	17.88	47.33	36.04	16.63
EL–Yes	4,830	93.02	6.56	0.41	94.74	4.82	0.43	92.34	7.27	0.39
EL–No	88,095	51.70	36.09	12.20	54.10	30.24	15.66	54.36	31.10	14.54
EconDis–Yes	24,929	72.09	23.98	3.92	73.81	20.42	5.76	74.17	21.24	4.59
EconDis–No	67,996	47.16	38.43	14.40	49.76	32.03	18.21	49.79	33.02	17.18
SWD–Yes	18,225	74.11	20.61	5.28	75.88	17.32	6.80	75.03	18.78	6.19
SWD–No	74,697	48.91	37.96	13.13	51.41	31.75	16.84	51.77	32.57	15.66
CBT	85,010	51.85	35.85	12.31	54.28	29.97	15.75	54.42	30.92	14.67
PBT	62	75.81	19.35	4.84	72.58	16.13	11.29	74.19	17.74	8.06
TTS	6,525	71.80	23.59	4.61	73.56	20.06	6.38	73.76	20.80	5.44
SP	1,149	93.12	6.61	0.26	93.73	5.92	0.35	92.34	7.57	0.09
SP TTS	108	87.96	10.19	1.85	92.59	6.48	0.93	92.59	6.48	0.93
Human Reader	56	98.21	1.79	0.00	98.21	1.79	0.00	96.43	3.57	0.00

**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	92,925	55.43	35.03	9.54	54.79	28.80	16.40	56.00	30.07	13.93
Male	47,026	54.64	33.67	11.68	53.76	28.04	18.19	56.20	28.81	14.99
Female	45,829	56.26	36.41	7.33	55.87	29.58	14.55	55.82	31.36	12.83
Am. Indian	121	62.81	28.93	8.26	61.98	28.93	9.09	59.50	31.40	9.09
Asian	10,025	26.12	45.49	28.39	26.02	34.18	39.79	25.76	37.34	36.91
Black	12,436	74.81	22.96	2.24	75.51	19.11	5.39	77.76	18.22	4.02
Hispanic	27,349	70.10	26.77	3.13	70.88	22.18	6.94	71.99	22.79	5.23
Pacific Islander	202	50.00	36.63	13.37	47.03	32.18	20.79	49.01	38.12	12.87
White	40,898	47.19	41.59	11.22	45.13	34.77	20.10	46.45	36.62	16.93
EL–Yes	4,830	87.91	11.90	0.19	93.87	5.69	0.43	94.53	5.09	0.37
EL–No	88,095	53.65	36.29	10.05	52.65	30.07	17.28	53.89	31.44	14.67
EconDis–Yes	24,929	72.40	24.86	2.74	72.78	20.85	6.37	74.26	20.92	4.83
EconDis–No	67,996	49.21	38.75	12.03	48.20	31.72	20.08	49.31	33.42	17.27
SWD–Yes	18,225	73.70	22.11	4.19	74.40	18.00	7.60	76.01	17.87	6.12
SWD–No	74,697	50.98	38.18	10.85	50.01	31.44	18.55	51.12	33.05	15.84
CBT	85,010	53.64	36.21	10.15	52.78	29.81	17.41	54.04	31.17	14.79
PBT	62	74.19	20.97	4.84	72.58	17.74	9.68	74.19	17.74	8.06
TTS	6,525	72.02	24.49	3.49	72.93	20.43	6.64	73.64	20.78	5.58
SP	1,149	87.12	12.71	0.17	93.65	6.09	0.26	94.34	5.48	0.17
SP TTS	108	89.81	10.19	0.00	90.74	8.33	0.93	92.59	6.48	0.93
Human Reader	56	98.21	1.79	0.00	100.00	0.00	0.00	98.21	1.79	0.00

## APPENDIX L: Executive Summary of the NJSLA–S Alignment Evaluation Study

### 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 school year; thus, the 2021 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). Items within the NJSLA-S each represent 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). 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.

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

### **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 their 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 & Space Science portion of the PLDs, due primarily to vagueness in the PLD wording. Overall, for all three forms evaluated, the PLDs 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 & 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 to display strong evidence of alignment and to be engaging. 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 & 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. We commend 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 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.