

AASA
Arizona ELA & Mathematics Assessments

2023
Technical Report

Submitted to the
Arizona Department of Education
July 2024



TABLE OF CONTENTS

Chapter 1: Introduction	10
1.1. Assessment Overview	10
1.2. Participation	11
1.3. Purpose and Intended Use of Test Scores	11
1.4. Educator Involvement	11
Chapter 2: Test Design	13
2.1. Arizona Academic Standards	13
2.2. Item Specifications	14
2.3. Test Blueprint	14
2.4. Depth of Knowledge	16
2.5. Item Types	16
2.6. Test Designs	18
2.6.1. ELA	19
2.6.2. Mathematics	20
Chapter 3: Test Development	21
3.1. Content Development and Management Tool	22
3.2. Item Bank Analysis	22
3.3. Passage and Item Development	23
3.4. Item Review	25
3.5. Form Construction	25
3.5.1. Preparation for Item Selection	26
3.5.2. Item Selection and Positioning	26
3.5.3. Sampling Plan	26
3.6. Data Review	26
3.7. Accommodated Forms	28
Chapter 4: Test Administration	29
4.1. Test Units	29
4.2. Administration Materials	30
4.3. Administration Training	30
4.4. Sample Tests	31
4.5. Accommodations	32
4.6. Universal Test Administration Conditions	33
4.7. Universal Test Tools	34
4.8. Pearson Customer Support	35
4.9. Test Security	35
Chapter 5: Scoring and Reporting	38
5.1. Human Scoring of Open-Ended Items	38
5.1.1. Scorer Recruitment	38
5.1.2. Training	39
5.1.2.1. Writing	39
5.1.2.2. Mathematics and Reading	39
5.1.3. Quality Control	40
5.1.4. Security	41
5.2. Automated Scoring of ELA Writing Prompts	41
5.2.1. Calibration of IEA	41
5.2.2. Smart Routing	42
5.2.3. Quality Control	43
5.3. Reporting	44

Chapter 6: Classical Item Analysis.....	50
6.1. Data.....	50
6.2. Descriptive Statistics.....	51
6.3. Classical Item Analysis.....	52
6.4. Distractor Analysis.....	53
Chapter 7: Calibration, Equating, and Scaling.....	55
7.1. Calibration Sample.....	55
7.2. Calibration Methods.....	55
7.3. Calibration Results.....	56
7.4. Equating.....	57
7.5. Scaling Methods.....	58
7.6. IRT Assumptions	58
7.6.1. Unidimensionality	59
7.6.2. Local Item Independence	60
7.6.3. Item Fit.....	60
Chapter 8: Test Results	62
Chapter 9: Reliability and Validity	69
9.1. Reliability.....	69
9.1.1. Internal Consistency	69
9.1.2. Inter-rater Reliability.....	72
9.2. Differential Item Functioning	74
9.3. Correlations Among Reporting Categories	76
9.4. Validity Evidence.....	79
9.4.1. Evidence Based on Test Content.....	80
9.4.2. Evidence Based on Response Processes.....	81
9.4.3. Evidence Based on Internal Structure	81
9.4.4. Evidence Based on Performance Standards	82
9.4.5. Evidence Based on Relation to Other Variables	82
9.4.6. Summary	83
Chapter 10: Classification into Performance Levels	84
10.1. Standard Setting.....	84
10.2. Classification Consistency and Accuracy	85
10.3. MOWR Policy	87
References.....	88
Appendix A: Item-Level CTT Statistics.....	91
Appendix B: Item-Level IRT Statistics.....	115
Appendix C: Administration Results	185
Appendix D: ACT Grade 8 Linking Study	198

LIST OF TABLES

Table 1.1. Schedule of Major Events.....	12
Table 2.1. AASA ELA Blueprint, Grades 3–8	15
Table 2.2. AASA Mathematics Blueprint, Grades 3–5	15
Table 2.3. AASA Mathematics Blueprint, Grades 6–7	15
Table 2.4. AASA Mathematics Blueprint, Grade 8	15
Table 2.5. DOK Levels.....	16
Table 2.6. Percentage of Points by DOK Level.....	16
Table 2.7. Item Types	17
Table 2.8. AASA Test Design—ELA	18
Table 2.9. AASA Test Design—Mathematics.....	18
Table 3.1. Number of Newly Developed Items	23
Table 3.2. Passage Lexile Measures and Word Count.....	23
Table 3.3. Item Statistical Flagging Criteria.....	27
Table 3.4. Data Review Results: Number of Field Tested Items	27
Table 4.1. Estimated Testing Time by Test Unit	30
Table 4.2. Administration Materials	30
Table 4.3. Administration Trainings.....	31
Table 4.4. Number of Items on the AASA Sample Tests	32
Table 4.5. Available Accommodations.....	32
Table 4.6. Universal Test Tools.....	34
Table 5.1. Scoring Qualification Standards	39
Table 6.1. Number of Students in the Calibration Sample by Subgroup—ELA	50
Table 6.2. Number of Students in the Calibration Sample by Subgroup—Mathematics	51
Table 6.3. Classical Test Analysis Statistics.....	51
Table 6.4. Classical Item Analysis Summary	52
Table 6.5. Distractor Analysis Summary: Point-Biserial Correlations for Correct Options	53
Table 6.6. Distractor Analysis Summary: Point-Biserial Correlations for Incorrect Options.....	54
Table 7.1. IRT Statistics Summary.....	56
Table 7.2. Summary of Anchor Items.....	58
Table 7.3. Eigenvalues from PCA	59
Table 7.4. Q3 Statistics.....	60
Table 7.5. IRT Item Fit Summary Statistics	61
Table 8.1. Overall Test Results.....	62
Table 8.2. Performance Distributions by Reporting Category: Percentage of Students at each Level of Mastery—ELA.....	63
Table 8.3. Performance Distributions by Reporting Category: Percentage of Students at each Level of Mastery—Mathematics	63
Table 8.4. Test Results by Accommodation—ELA	64
Table 8.5. Test Results by Accommodation—Mathematics.....	65
Table 8.6. Scale Score Distribution by Performance Level—ELA	67
Table 8.7. Scale Score Distribution by Performance Level—Mathematics	67
Table 9.1. Coefficient Alpha and SEM by Total and Reporting Category Score—ELA, Form 1	70
Table 9.2. Coefficient Alpha and SEM by Total and Reporting Category Score—ELA, Form 2	70
Table 9.3. Coefficient Alpha and SEM by Total and Reporting Category Score—Mathematics.....	71
Table 9.4. Inter-rater Reliability Statistics	73
Table 9.5. DIF Flag Categories.....	75
Table 9.6. Number of Items Exhibiting Strong DIF	76
Table 9.7. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—ELA Form 1	77

Table 9.8. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—ELA Form 2	78
Table 9.9. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—Mathematics Grades 3–5	78
Table 9.10. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—Mathematics Grades 6 and 7	79
Table 9.11. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—Mathematics Grade 8	79
Table 9.12. Correlation between AASA ELA and Mathematics Scale Scores	82
Table 10.1. Performance Level Cut Scores	84
Table 10.2. CSEM at Performance Level Cuts	85
Table 10.3. Classification Consistency for the <i>Proficient</i> Cut	86
Table 10.4. Classification Accuracy for the <i>Proficient</i> Cut	86
Table 10.5. Classification Consistency and Accuracy Results	87
Table A.1. Item-Level CTT Statistics, ELA Grade 3	91
Table A.2. Item-Level CTT Statistics, ELA Grade 4	92
Table A.3. Item-Level CTT Statistics, ELA Grade 5	93
Table A.4. Item-Level CTT Statistics, ELA Grade 6	94
Table A.5. Item-Level CTT Statistics, ELA Grade 7	95
Table A.6. Item-Level CTT Statistics, ELA Grade 8	97
Table A.7. Item-Level CTT Statistics, Mathematics Grade 3	98
Table A.8. Item-Level CTT Statistics, Mathematics Grade 4	99
Table A.9. Item-Level CTT Statistics, Mathematics Grade 5	100
Table A.10. Item-Level CTT Statistics, Mathematics Grade 6	101
Table A.11. Item-Level CTT Statistics, Mathematics Grade 7	102
Table A.12. Item-Level CTT Statistics, Mathematics Grade 8	103
Table A.13. Distractor Analysis of Multiple-Choice Items, ELA Grade 3	105
Table A.14. Distractor Analysis of Multiple-Choice Items, ELA Grade 4	106
Table A.15. Distractor Analysis of Multiple-Choice Items, ELA Grade 5	107
Table A.16. Distractor Analysis of Multiple-Choice Items, ELA Grade 6	108
Table A.17. Distractor Analysis of Multiple-Choice Items, ELA Grade 7	109
Table A.18. Distractor Analysis of Multiple-Choice Items, ELA Grade 8	110
Table A.19. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 3	111
Table A.20. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 4	111
Table A.21. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 5	112
Table A.22. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 6	112
Table A.23. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 7	113
Table A.24. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 8	114
Table B.1. Item-Level IRT Statistics, ELA Grade 3	115
Table B.2. Item-Level IRT Statistics, ELA Grade 4	116
Table B.3. Item-Level IRT Statistics, ELA Grade 5	117
Table B.4. Item-Level IRT Statistics, ELA Grade 6	118
Table B.5. Item-Level IRT Statistics, ELA Grade 7	120
Table B.6. Item-Level IRT Statistics, ELA Grade 8	121
Table B.7. Item-Level IRT Statistics, Mathematics Grade 3	122
Table B.8. Item-Level IRT Statistics, Mathematics Grade 4	123
Table B.9. Item-Level IRT Statistics, Mathematics Grade 5	124
Table B.10. Item-Level IRT Statistics, Mathematics Grade 6	125
Table B.11. Item-Level IRT Statistics, Mathematics Grade 7	126
Table B.12. Item-Level IRT Statistics, Mathematics Grade 8	127
Table B.13. Raw-to-Scale Score Conversion, ELA Grade 3	129
Table B.14. Raw-to-Scale Score Conversion, ELA Grade 4	131

Table B.15. Raw-to-Scale Score Conversion, ELA Grade 5	133
Table B.16. Raw-to-Scale Score Conversion, ELA Grade 6	136
Table B.17. Raw-to-Scale Score Conversion, ELA Grade 7	138
Table B.18. Raw-to-Scale Score Conversion, ELA Grade 8	140
Table B.19. Raw-to-Scale Score Conversion, Mathematics Grade 3	143
Table B.20. Raw-to-Scale Score Conversion, Mathematics Grade 4	144
Table B.21. Raw-to-Scale Score Conversion, Mathematics Grade 5	145
Table B.22. Raw-to-Scale Score Conversion, Mathematics Grade 6	146
Table B.23. Raw-to-Scale Score Conversion, Mathematics Grade 7	147
Table B.24. Raw-to-Scale Score Conversion, Mathematics Grade 8	148
Table C.1. Test Results by Subgroup, ELA Grade 3	185
Table C.2. Test Results by Subgroup, ELA Grade 4	186
Table C.3. Test Results by Subgroup, ELA Grade 5	186
Table C.4. Test Results by Subgroup, ELA Grade 6	187
Table C.5. Test Results by Subgroup, ELA Grade 7	187
Table C.6. Test Results by Subgroup, ELA Grade 8	188
Table C.7. Test Results by Subgroup, Mathematics Grade 3	188
Table C.8. Test Results by Subgroup, Mathematics Grade 4	189
Table C.9. Test Results by Subgroup, Mathematics Grade 5	189
Table C.10. Test Results by Subgroup, Mathematics Grade 6	190
Table C.11. Test Results by Subgroup, Mathematics Grade 7	190
Table C.12. Test Results by Subgroup, Mathematics Grade 8	191
Table D.1. Summary Statistics of ACT Aspire Scores to Create Short-to-Long Concordance—Mathematics	199
Table D.2. Summary Statistics of ACT Aspire Scores to Create Short-to-Long Concordance—Reading.....	199
Table D.3. Summary Statistics of ACT Aspire Scores to Create Short-to-Long Concordance—ELA	199
Table D.4. Summary Statistics of AASA Scores to Create Long-to-Short Concordance—Mathematics	199
Table D.5. Summary Statistics of AASA Scores to Create Long-to-Short Concordance—ELA	199
Table D.6. Item Statistics for the Common Items—Mathematics	200
Table D.7. Item Statistics for the Common Items—Reading/ELA	200
Table D.8. Final Concordance Table—Mathematics.....	201
Table D.9. Final Concordance Table—ELA	203

LIST OF FIGURES

Figure 3.1. Item Development Process	21
Figure 3.2. Text Complexity Worksheet Example	24
Figure 4.1. Test Security Agreement.....	36
Figure 5.1. Dynamic Model Development and Deployment.....	42
Figure 5.2. Smart Routing	43
Figure 5.3. Sample Reports—Confidential Student Score Report, Grade 3	45
Figure 5.4. Sample Reports—Confidential Student Score Report, Grade 8	47
Figure 5.5. Sample Reports—Confidential Roster Report with Summary	49
Figure B.1. Item-Person Map, ELA Grade 3, Form 1	149
Figure B.2. Item-Person Map, ELA Grade 4, Form 1	149
Figure B.3. Item-Person Map, ELA Grade 5, Form 1	150
Figure B.4. Item-Person Map, ELA Grade 6, Form 1	150
Figure B.5. Item-Person Map, ELA Grade 7, Form 1	151
Figure B.6. Item-Person Map, ELA Grade 8, Form 1	151
Figure B.7. Item-Person Map, ELA Grade 3, Form 2	152
Figure B.8. Item-Person Map, ELA Grade 4, Form 2	152
Figure B.9. Item-Person Map, ELA Grade 5, Form 2	153
Figure B.10. Item-Person Map, ELA Grade 6, Form 2	153

Figure B.11. Item-Person Map, ELA Grade 7, Form 2	154
Figure B.12. Item-Person Map, ELA Grade 8, Form 2	154
Figure B.13. Item-Person Map, Mathematics Grade 3	155
Figure B.14. Item-Person Map, Mathematics Grade 4	155
Figure B.15. Item-Person Map, Mathematics Grade 5	156
Figure B.16. Item-Person Map, Mathematics Grade 6	156
Figure B.17. Item-Person Map, Mathematics Grade 7	157
Figure B.18. Item-Person Map, Mathematics Grade 8	157
Figure B.19. TCC, ELA Grade 3, Form 1	158
Figure B.20. CSEM, ELA Grade 3, Form 1	158
Figure B.21. TCC, ELA Grade 4, Form 1	159
Figure B.22. CSEM, ELA Grade 4, Form 1	159
Figure B.23. TCC, ELA Grade 5, Form 1	160
Figure B.24. CSEM, ELA Grade 5, Form 1	160
Figure B.25. TCC, ELA Grade 6, Form 1	161
Figure B.26. CSEM, ELA Grade 6, Form 1	161
Figure B.27. TCC, ELA Grade 7, Form 1	162
Figure B.28. CSEM, ELA Grade 7, Form 1	162
Figure B.29. TCC, ELA Grade 8, Form 1	163
Figure B.30. CSEM, ELA Grade 8, Form 1	163
Figure B.31. TCC, ELA Grade 3, Form 2	164
Figure B.32. CSEM, ELA Grade 3, Form 2	164
Figure B.33. TCC, ELA Grade 4, Form 2	165
Figure B.34. CSEM, ELA Grade 4, Form 2	165
Figure B.35. TCC, ELA Grade 5, Form 2	166
Figure B.36. CSEM, ELA Grade 5, Form 2	166
Figure B.37. TCC, ELA Grade 6, Form 2	167
Figure B.38. CSEM, ELA Grade 6, Form 2	167
Figure B.39. TCC, ELA Grade 7, Form 2	168
Figure B.40. CSEM, ELA Grade 7, Form 2	168
Figure B.41. TCC, ELA Grade 8, Form 2	169
Figure B.42. CSEM, ELA Grade 8, Form 2	169
Figure B.43. TCC, Mathematics Grade 3	170
Figure B.44. CSEM, Mathematics Grade 3	170
Figure B.45. TCC, Mathematics Grade 4	171
Figure B.46. CSEM, Mathematics Grade 4	171
Figure B.47. TCC, Mathematics Grade 5	172
Figure B.48. CSEM, Mathematics Grade 5	172
Figure B.49. TCC, Mathematics Grade 6	173
Figure B.50. CSEM, Mathematics Grade 6	173
Figure B.51. TCC, Mathematics Grade 7	174
Figure B.52. CSEM, Mathematics Grade 7	174
Figure B.53. TCC, Mathematics Grade 8	175
Figure B.54. CSEM, Mathematics Grade 8	175
Figure B.55. Scree Plot, ELA Grade 3, Form 1	176
Figure B.56. Scree Plot, ELA Grade 4, Form 1	176
Figure B.57. Scree Plot, ELA Grade 5, Form 1	177
Figure B.58. Scree Plot, ELA Grade 6, Form 1	177
Figure B.59. Scree Plot, ELA Grade 7, Form 1	178
Figure B.60. Scree Plot, ELA Grade 8, Form 1	178
Figure B.61. Scree Plot, ELA Grade 3, Form 2	179
Figure B.62. Scree Plot, ELA Grade 4, Form 2	179

Figure B.63. Scree Plot, ELA Grade 5, Form 2	180
Figure B.64. Scree Plot, ELA Grade 6, Form 2	180
Figure B.65. Scree Plot, ELA Grade 7, Form 2	181
Figure B.66. Scree Plot, ELA Grade 8, Form 2	181
Figure B.67. Scree Plot, Mathematics Grade 3	182
Figure B.68. Scree Plot, Mathematics Grade 4	182
Figure B.69. Scree Plot, Mathematics Grade 5	183
Figure B.70. Scree Plot, Mathematics Grade 6	183
Figure B.71. Scree Plot, Mathematics Grade 7	184
Figure B.72. Scree Plot, Mathematics Grade 8	184
Figure C.1. Total Scale Score Distribution, ELA Grade 3	192
Figure C.2. Total Scale Score Distribution, ELA Grade 4	192
Figure C.3. Total Scale Score Distribution, ELA Grade 5	193
Figure C.4. Total Scale Score Distribution, ELA Grade 6	193
Figure C.5. Total Scale Score Distribution, ELA Grade 7	194
Figure C.6. Total Scale Score Distribution, ELA Grade 8	194
Figure C.7. Total Scale Score Distribution, Mathematics Grade 3	195
Figure C.8. Total Scale Score Distribution, Mathematics Grade 4	195
Figure C.9. Total Scale Score Distribution, Mathematics Grade 5	196
Figure C.10. Total Scale Score Distribution, Mathematics Grade 6	196
Figure C.11. Total Scale Score Distribution, Mathematics Grade 7	197
Figure C.12. Total Scale Score Distribution, Mathematics Grade 8	197

Chapter 1: INTRODUCTION

This technical report documents the design, development, administration, technical processes, and results of the Spring 2023 administration of Arizona’s Academic Standards Assessment (AASA) in English language arts (ELA) and mathematics in Grades 3–8 to support test users in evaluating the intended purposes, uses, and interpretations of the test scores. The technical information herein is intended for use by those who evaluate tests, interpret scores, or use test results in making educational decisions. It is assumed that the reader has technical knowledge of test construction and measurement procedures, as stated in the *Standards for Educational and Psychological Testing* (AERA et al., 2014).

1.1. Assessment Overview

AASA is the statewide achievement test for Arizona students in ELA and mathematics in Grades 3–8 aligned with the Arizona Academic Standards as described in state and federal law (State Law ARS 15-741; Federal Law: 34 CFR 200.2 *Participation in Assessments*). It is a summative, criterion-referenced assessment designed to promote increasingly higher academic outcomes for students to prepare them for a broader array of post-secondary outcomes. It is available as a computer-based test (CBT) or paper-based test (PBT), with CBT as the default administration mode.

In November 2014, the State Board of Education adopted Arizona’s Measurement of Educational Readiness to Inform Teaching (AzMERIT) to measure student mastery of the Arizona academic standards and progress toward college and career readiness, with the first administration in Spring 2015. The current Arizona Academic Standards were adopted by the Arizona State Board of Education in December 2016. Beginning in 2019–2020, AzMERIT was renamed AzM2. Beginning in 2021–2022, AzM2 was renamed to AASA. The assessment is still aligned to the same 2016 academic content standards and has the same cut scores.

A Writing standalone field test (SAFT) was administered in 2022 to all students in Grades 3–8 to build Arizona’s item bank for extended writing items. Oral Reading Fluency (ORF) field test items were also embedded on the Grade 3 operational AASA test in Spring 2022 to enhance coverage of the Grade 3 ELA standards. They were field tested again in Spring 2023 to further explore their functioning and performance. In line with this work, AASA writing rubrics and guides were made available in August 2022. The rubrics can be used in classrooms to score students’ work to prepare them for the AASA Writing test unit. The writing guides are included in the test environment with each prompt as a reminder to students of key pieces from the rubric to include in their essays. They can be used in classrooms with assignments or to help students complete classroom or district essays throughout the school year.

Beginning in Spring 2023, an ACT predicted score was included for Grade 8 students on the Confidential Student Score Reports for both ELA and mathematics. Test scores from the AASA were linked to the ACT scale to obtain the ACT predicted score range that indicates the score a student would likely receive if they were to take the ACT test, which will help students understand their predicted college readiness and plan future course work. Students who score at or above the ACT score indicated on the Confidential Student Score Report are more likely to be successful in college courses taken by first-year students. Appendix D presents the study results.

1.2. Participation

Students in Grades 3–8 participate in the spring administration of the AASA test. The state and federal laws mandate that all public school students participate in the assessments that measure student achievement of grade-level content standards. Students with significant cognitive disabilities whose Individualized Education Program (IEP) designates them as eligible for an alternate assessment, the Multi-State Alternate Assessment (MSAA), should not be administered the AASA assessment.

1.3. Purpose and Intended Use of Test Scores

The primary intended score interpretation of AASA is that AASA test scores provide reliable and valid information about important knowledge and skills in grade-level numeracy and literacy that students are attaining. Furthermore, while ultimate use of the test scores is determined by Arizona educators and other stakeholders, the primary intended uses of the AASA test scores include the following:

- Schools and districts use the AASA assessment and its results to (a) monitor trends in student performance and (b) design professional development for teachers.
- Teachers use the AASA assessment and its results to integrate assessment with their instructional planning.
- Parents/guardians use the AASA assessment and its results to get information about (a) what their child knows and can do and (b) their child’s progress from year to year.

1.4. Educator Involvement

This section addresses the involvement of Arizona educators in test development as indicated by Standard 4.8 of the *Standards for Educational and Psychological Testing* (AERA et al., 2014). Arizona educators were involved in many steps of the process, as shown in Table 1.1 that presents the major events regarding the development, administration, and reporting of the Spring 2023 AASA assessments.

Arizona educators participated in meetings and provided feedback on assets developed for field testing. These meetings were held virtually and included educators from across the state. The committee meetings included a passage review that enabled educators to review ELA passages for content, grade-level appropriateness, and bias and sensitivity; a content and bias item review that enabled educators to review items for content, standard alignment, grade-level appropriateness, and bias and sensitivity; and a bias and sensitivity community review that enabled community members, including past and present Arizona educators, to evaluate items for bias and sensitivity concerns.

Table 1.1. Schedule of Major Events

Event	Date(s)
ELA Passage Review	February 4, 2022
ELA Community Passage Review	February 8, 2022
Content and Bias Item Review	June 27 – July 1, 2022
Bias and Sensitivity Community Review	July 18–22, 2022
Technical Advisory Committee (TAC) Meeting	October 12–13, 2022
Administration Training	December 12, 2022 – April 21, 2023
AASA Additional Order Window for Test Materials	March 24 – April 21, 2023
Spring 2023 AASA Test Administration Window (CBT)	April 3–28, 2023
Spring 2023 AASA Test Administration Window (PBT)	April 3–12, 2023
Spring 2023 AASA Test Administration Window (ELA Writing)	April 3–14, 2023
Release of Grade 3 Electronic Score Reports	May 15, 2023
Release of Grades 4–8 Electronic Score Reports	May 25, 2023
Release of Grades 3–8 Paper Reports to Districts	June 15, 2023
Data Review	July 10–14, 2023

Chapter 2: TEST DESIGN

This chapter provides information regarding test design as indicated by Standards 1.11, 4.0, 4.1, 4.12, 12.4, and 12.8 (AERA et al., 2014). AASA is designed to be administered online, with paper accommodated forms available as needed. The needs of the student are also addressed through other supports, such as assessment features built into the online platform and accommodations such as using assistive technology, a scribe, and/or sign language (see Chapter 4 for more information). ELA includes 42 operational items consisting of multiple-choice and technology-enhanced item types, along with an open-response writing prompt at each grade level. In grade 3 ELA only, 3 short Oral Reading Fluency (ORF) passages are included. Mathematics consists of 53–55 operational multiple-choice and technology-enhanced items. Field test items are also embedded on each assessment that do not count toward students' scores.

Accessibility was the foundation of the AASA test design to make sure all students have access to the content based on the college- and career-ready Arizona Academic Standards, which begins with rigorous curriculum, instructional resources, and training for teachers. Principles of Universal Design are adhered to throughout the item and test creation process to accommodate the needs and abilities of all learners. AASA is available to be administered in online settings including group, small group, or one-on-one settings. AASA is also available in appropriate accommodations including ASL, Braille, Large Print, or Regular Print format.

2.1. Arizona Academic Standards

In 2016, the State Board of Education adopted new academic content standards in ELA and mathematics that reflect high expectations of all Arizona students and strive to ensure that high school graduates are college- and career-ready. The Arizona Academic Standards define the knowledge, understanding, and skills that need to be taught and learned so all students are ready to succeed in credit-bearing, college-entry courses and/or in the workplace.

The ELA standards describe the reading, writing, language, speaking, and listening skills that students should acquire from Grades K–12, and the mathematics standards describe expectations for learning in Grades K–8 and the first three high school courses (Algebra I, Geometry, Algebra II; Mathematics 1, 2, 3), plus specific standards that could be included in a fourth high school credit mathematics course. The standards are located on the Arizona Department of Education (ADE) website at <https://www.azed.gov/standards-practices>.

The standards work together in a clear progression from Grades K–12. Each standard builds on the standard that came before and toward the standard that comes in the next grade level. They are the foundation to guide the construction and evaluation of programs in Arizona K–12 schools and the broader Arizona community. The Arizona Academic Standards are

- focused in coherent progressions across Grades K–12;
- aligned with college and workforce expectations;
- inclusive of rigorous content and applications of knowledge through higher-order thinking;
- research and evidence based;
- broad in nature, allowing for the widest possible range of student learning; and
- designed as an integrated approach to literacy (ELA).

2.2. Item Specifications

AASA item specifications are available for each grade and content area on the ADE website at <https://www.azed.gov/assessment/aasa>. These item specifications, refined by Pearson and ADE content experts, are used to guide the item development process by defining the content limit, model tasks, and response types for a specific standard. During each level of review, items are compared to the item specifications to ensure their alignment to the standard, grade-level appropriateness, and adherence to the content limits set forth in the item specifications.

The item specifications were developed using a vertical alignment for each standard, wherein the suggested task demands and cognitive complexity of items build upon those of the previous grade level, just as the standards themselves do. The item specifications also provide models for item writers that include item samples that target different Depth of Knowledge (DOK) and difficulty levels. These item models annotate the information to communicate the intent of the standard and DOK and clarify how to manipulate the item difficulty while keeping the cognitive demands the same for the writer. The item specifications document includes the following:

- **Content Limits.** This section delineates the specific content measured by the standard and the extent to which the content is different across grade levels. For example, in Grade 3 Mathematics, fraction denominators are limited to 2, 3, 4, 6, and 8.
- **Acceptable Response Mechanisms.** This section identifies the various ways in which students may respond to a prompt (e.g., multiple choice, graphic response, equation response, matching, multi-select).
- **Task Demands.** In this section, the standards are broken down into specific task demands aligned to the standard. In addition, each task demand is assigned a common item format relevant to that particular task demand.

2.3. Test Blueprint

The test blueprint, in concert with the item specifications, defines the content and structure of the test. Table 2.1 – Table 2.4 present a summary of the blueprints based on the 2016 standards for Grades 3–8 in ELA and mathematics. External, public-facing blueprints are available on the ADE website at <https://www.azed.gov/assessment/aasa>. More detailed blueprints are used internally by ADE and the vendor. The blueprint defines the standards to be assessed for each test form, the number of items per standard, the number of item types, the number of points per item type, and the total number of items and points per test form. Inherent in the number of points per test is the relative weighting associated with the standards and the reporting categories being assessed.

Table 2.1. AASA ELA Blueprint, Grades 3–8

Reporting Category	Grades 3–5		Grades 6–8	
	Min.	Max.	Min.	Max.
Reading Standards for Literature	26%	35%	24%	31%
Reading Standards for Informational Text	26%	35%	30%	38%
Reading for Informational Text	26%	22%	30%	25%
Listening Comprehension	0%	13%	0%	13%
Writing and Language	26%	38%	30%	38%
Writing	13%	19%	17%	19%
Language	13%	19%	13%	19%

Note. Listening standards are only assessed on the online assessment.

Table 2.2. AASA Mathematics Blueprint, Grades 3–5

Reporting Category	Grade 3		Grade 4		Grade 5	
	Min.	Max.	Min.	Max.	Min.	Max.
Operations and Algebraic Thinking and Numbers and Operations in Base Ten	49%	53%	46%	54%	38%	42%
Operations and Algebraic Thinking	38%	42%	22%	26%	4%	8%
Numbers in Base Ten	9%	13%	24%	28%	31%	35%
Numbers and Operations – Fractions	18%	22%	29%	33%	31%	35%
Measurement and Data and Geometry	26%	30%	15%	19%	24%	28%
Measurement and Data	26%	28%	9%	13%	18%	20%
Geometry	1%	4%	4%	7%	7%	11%

Table 2.3. AASA Mathematics Blueprint, Grades 6–7

Reporting Category	Grade 6		Grade 7	
	Min.	Max.	Min.	Max.
Ratios and Proportions	19%	23%	19%	23%
The Number System	28%	32%	19%	23%
Expressions and Equations	29%	33%	23%	27%
Geometry and Statistics and Probability	15%	19%	27%	35%
Geometry	6%	15%	15%	19%
Statistics and Probability	6%	11%	12%	16%

Table 2.4. AASA Mathematics Blueprint, Grade 8

Reporting Category	Grade 8	
	Min.	Max.
Functions	21%	25%
Expressions & Equations	29%	33%
Geometry	17%	21%
Statistics and Probability and The Number System	19%	27%
Statistics and Probability	4%	8%
The Number System	15%	19%

2.4. Depth of Knowledge

All items are aligned according to DOK, the cognitive complexity of the item, and the cognitive demands on the student. DOK refers to the level of rigor or sophistication of the task in an item designed to reflect the complexity of the Arizona Academic Standards. Table 2.5 presents a description of the DOK levels as provided in the item specifications, and Table 2.6 presents the percentage of points by DOK level as provided in the blueprints.

Table 2.5. DOK Levels

DOK Level	ELA	Mathematics
Level 1: Recall	Focuses on basic tasks such as correcting grammatical and spelling errors, defining terms, and locating details or facts in texts.	Focuses on the recall of information, such as definitions, terms, and simple procedures.
Level 2: Skill/Concept	Requires a greater degree of engagement and cognitive processing than DOK 1 items. DOK 2 items may require students to show relationships or identify examples, use context to identify meaning, identify structures or features of texts, or distinguish between facts and opinions.	Requires students to make decisions, solve problems, or recognize patterns. In general, DOK 2 items require a greater degree of engagement and cognitive processing than DOK 1 items.
Level 3: Strategic Thinking	Features higher-order cognitive tasks that assess students' capacities to read complex texts and think abstractly and focuses on critical thinking, developing, and assessing logical arguments, making inferences, and citing evidence to support claims or conclusions.	Features higher-order cognitive tasks that assess students' capacities to approach abstract or complex problems.
Level 4: Extended Thinking (Writing only)	Requires creativity, extensive planning, and/or sophisticated reasoning in the composition and organization of written essays.	N/A

Table 2.6. Percentage of Points by DOK Level

DOK Level	ELA	Mathematics
Level 1	10–20%	10–20%
Level 2	50–60%	60–70%
Level 3	15–25%	12–30%
Level 4	16–19% (Writing)	N/A

2.5. Item Types

The AASA assessments include traditional multiple-choice items and technology-enhanced items (TEIs), as shown in Table 2.7. Examples of each item type may be found in the AASA sample tests accessed through TestNav (see Section 0 for more information).

TEIs require students to interact with test content to select, construct, and/or support their responses and are better able to assess a deeper level of understanding. For paper-based assessments, TEIs are modified or replaced with another item type that assesses the same standards so they can be scanned and scored electronically or hand scored. For example, gap match/gap match table, match – table grid, and short-constructed response items may be replaced with another item type that assesses the same standard and can be scanned and scored electronically. Inline choice items are modified so the student fills in a circle to indicate the correct word or phrase, and hot text items are modified so the student fills in a circle to indicate a selection.

Table 2.7. Item Types

Item Type	Description
Multiple-Choice (MC)	The student selects only one correct answer from among a number of options.
Multiple-Select (MS)	The student selects all of the correct answers from among a number of options.
Evidence-Based Selected Response (EBSR) (ELA only)	<ul style="list-style-type: none"> MC/MS Format: The student answers a Part A multiple-choice item based on a passage and then provides evidence in support of that answer by completing another Part B multiple-choice item or a Part B multi-select item. MC/TEI Format: The student answers a Part A multiple-choice item based on a passage and then provides evidence in support of that answer by completing a Part B technology-enhanced item.
Bar Graph (mathematics only)	The student drags bars vertically or horizontally along numerical values. Individual bars, histograms, and clusters are supported.
Equation Editor (mathematics only)	The student uses a palette of buttons to enter a numerical response or to create mathematical expressions.
Fraction Model (mathematics only)	The student divides a shape (circle or rectangle) into varying numbers of segments by clicking a ‘Fewer’ or ‘More’ button and selects those segments to shade those segments with a solid color.
Point Graph (mathematics only)	The student plots points, line segments, continuous lines, and/or polygons. Point graph items can use one or multiple graph interactions (composite graphs).
Shape Transformation (mathematics only)	The student chooses one of four variants of a single shape, drags it onto a four-quadrant grid, and positions it on the grid.
Inline Choice (IC)	The student selects a single text option from a drop-down menu within a table or inline text, similar to a fill-in-the-blank item. The item may contain multiple blanks.
Gap Match (GM)	Certain numbers, words, phrases, or sentences may be designated “draggable” in this item type. The student can click on the option, hold down the mouse button, and drag it to a graphic or other format.
Gap Match Table (GMT)	Same as the gap match item except the drop zone is in a table format.
Match – Table Grid (MTG)	The student selects radio buttons or checks boxes in cells to indicate if information from a column header matches information from a row.
Hot Text (HT) (ELA only)	The student selects one or more areas called hot spots on an image. For ELA, excerpted sentences from the text are presented in this item type. Certain words, phrases, or sentences are highlighted to indicate that the text is selectable (“hot”). The student can then click on an option to select it.
Hot Spot (mathematics only)	The student selects one or more areas called hot spots on an image. An example for mathematics is selecting a point on a number line. The student can click on an option to select it.
Short Constructed Response (SCR) (ELA only)	The student uses the keyboard to enter a response into a text field. These items can usually be answered in a sentence or two.
Writing Prompt (ELA only)	These items may require the student to use features of an online word processor. The student can perform various tasks within the online word processor such as bold text, use bullet points, underline, etc.

2.6. Test Designs

Table 2.8 and Table 2.9 present the test designs for the ELA and mathematics assessments. As shown in the tables, the AASA test consists of the following test units:

- ELA Oral Reading Fluency (ORF) test unit (Grade 3 only)
- ELA Writing test unit
- ELA Reading/Language Test Unit 1 and Test Unit 2
- Math Test Unit 1 and Test Unit 2

Each grade-level ELA and mathematics test form includes the same operational items but a different set of embedded field test items. The ELA assessments consist of three test units (Writing, Reading/Language Test Unit 1, Reading/Language Test Unit 2), with a fourth ORF unit for Grade 3 only. The mathematics assessments consist of two test units (Math Test Unit 1 and Math Test Unit 2). The tables indicate the number of operational and field test items included on the test form for each unit. Given the nature of passage-based item sets in Reading/Language, field test items are confined to their associated set in only one unit of the test.

Table 2.8. AASA Test Design—ELA

Grade	#Forms	#Passages	Overall			#Items by Test Unit											
			#Items			Writing			Reading/Language Test Unit 1			Reading/Language Test Unit 2			Oral Reading Fluency (ORF)		
			OP	FT	Total	OP	FT	Total	OP	FT	Total	OP	FT	Total	OP	FT	Total
3	21	8	42	11	53	1	–	1	16	8	24	25	–	25	0	3	3
4	21	8	42	8	50	1	–	1	16	8	24	25	–	25	–	–	–
5	21	7	42	8	50	1	–	1	20	8	28	21	–	21	–	–	–
6	21	7	42	8	50	1	–	1	18	8	26	23	–	23	–	–	–
7	21	7	42	8	50	1	–	1	18	8	26	23	–	23	–	–	–
8	18	7	42	8	50	1	–	1	18	8	26	23	–	23	–	–	–

Note. Each writing prompt is worth 10 points. The test design for ELA is based on the number of items, and the total points per operational form vary from 52–56 points. For Grade 3, the ORF passages are worth 2 points each. The #Passages are specific to the two Reading test units.

Table 2.9. AASA Test Design—Mathematics

Grade	#Forms	#Items								
		Overall			Test Unit 1			Test Unit 2		
		Total	OP	FT	Total	OP	FT	Total	OP	FT
3	11	53	45	8	27	23	4	26	22	4
4	11	53	45	8	27	23	4	26	22	4
5	11	53	45	8	27	23	4	26	22	4
6	11	55	47	8	27	23	4	28	24	4
7	11	55	47	8	27	23	4	28	24	4
8	11	55	47	8	27	23	4	28	24	4

Note. Each operational item is worth 1 point. Grades 3–5 have 45 points possible, and Grades 6–8 have 47 points possible.

2.6.1. ELA

The ELA ORF test unit consists of three short passages that students read aloud to measure oral reading fluency. Students have one minute to read each passage, and they receive a score of 0, 1, or 2. Word counts for ORF passages range between 250 and 400 words. Students are presented with three passages, each with a different difficulty level of low, medium, and high. These levels are based on Lexile ranges (600L – 650L, 650L – 700L, and 700L – 750L). Each student receives a combination of fiction and nonfiction genres. The ELA test has a Writing part and a Reading Part 1 and Part 2 for all grade levels. Writing consists of one writing prompt, which is an extended text/essay response. The Reading/Language is a long test, so it is split into two units. Each unit includes both reading and language items.

The ELA passages represent a variety of genres and topics. Pearson’s content experts develop informational texts from multiple content areas, such as history, science, and technical subjects. Literary texts represent authentic pieces from multiple genres, including stories, poetry, and drama. The ratio of informational to literary texts increases at each grade band, with a greater percentage of informational texts in the upper grades. The AASA uses both single passages and passage sets in which students are asked to synthesize information across texts. The number of items associated with each varies depending on the actual set and what standards are assessed.

The AASA ELA assessment is designed to reflect the importance of using evidence and reading complex texts outlined in the Arizona Academic Standards. It includes extended-writing tasks that provide students with meaningful contexts in which to construct their responses. Each writing prompt presents students with various stimuli (at least 2–3 per task) that serve as a springboard for an informed piece of writing. Students are given research articles, charts and graphs, and narratives to serve as the basis for their written responses. Students can then use this information, along with their own reasoning, to formulate an essay that is a clear and coherent expression of their own thinking while being grounded in research and evidence.

Each student is administered a single informative/explanatory or opinion/argumentative writing essay. While each student will only see one type of writing essay, both types are administered operationally at every grade level each year. Informative/explanatory writing is focused on conveying information accurately and seeks to enlighten the reader about processes or procedures, phenomena, states of affairs, and terminology. To produce this kind of writing, students draw from what they already know and from primary and secondary sources and develop a main idea and a primary focus as they relate facts, details, and examples.

Opinion (Grades 3–5) and argumentative (Grades 6–11) prompts ask students to analyze primary and secondary sources, make sound judgments, and present their opinions or arguments in a coherent manner that weaves personal opinions with evidence from the texts. The stimuli present opposing points of view about a topic so that students have enough information to take a stand. The stimuli are followed by a prompt that asks students to write an opinion or argumentative essay. The students must synthesize information across the passages to write the essay and cite specific details to support the ideas they present. For example, the prompt might require students to describe the steps in a process or describe problems that need to be solved.

The reading level of the stimulus does not exceed the easy Lexile range for the grade level to enable the students to attend to the content of the passages and not struggle with unfamiliar language and non-content-related vocabulary. Moreover, this helps ensure that students are assessed on their writing skills and not their reading abilities.

2.6.2. Mathematics

The mathematics items are created to address key components of the Arizona mathematics standards that check a student's conceptual understanding of mathematics and their procedural skills. The standards for Math Practices are embedded within all AASA items. The items are written in accordance with the item specifications to address key components of the standards and assess a range of important skills and performance levels based on the performance level descriptors (PLDs) that provide a standard level description of the level of knowledge and skills required at each performance level of the assessment. This provides an opportunity for students at all performance levels to show their understanding of the mathematics standards in the assessment. Each item is also aligned to a DOK level and the overall percentage of points by DOK level, as outlined in Table 2.6.

Equation editor items are an item type unique to mathematics. For an equation editor response, students type with a keyboard or use a palette of buttons to enter a response that could be a number, an expression, or an equation. The response may contain scaffolding where students are given part of a solution and fill in the missing parts. Two types of palettes are used in equation editor items that provides quick access to mathematical operators and symbols. For numerical responses, an abbreviated palette is given that contains the digits 0–9, a decimal point, a negative sign, a button to add a fraction, and a button to add a mixed number button. For expression or equation responses, the palette contains everything from the abbreviated palette plus additional mathematical operators and symbols depending on the grade level.

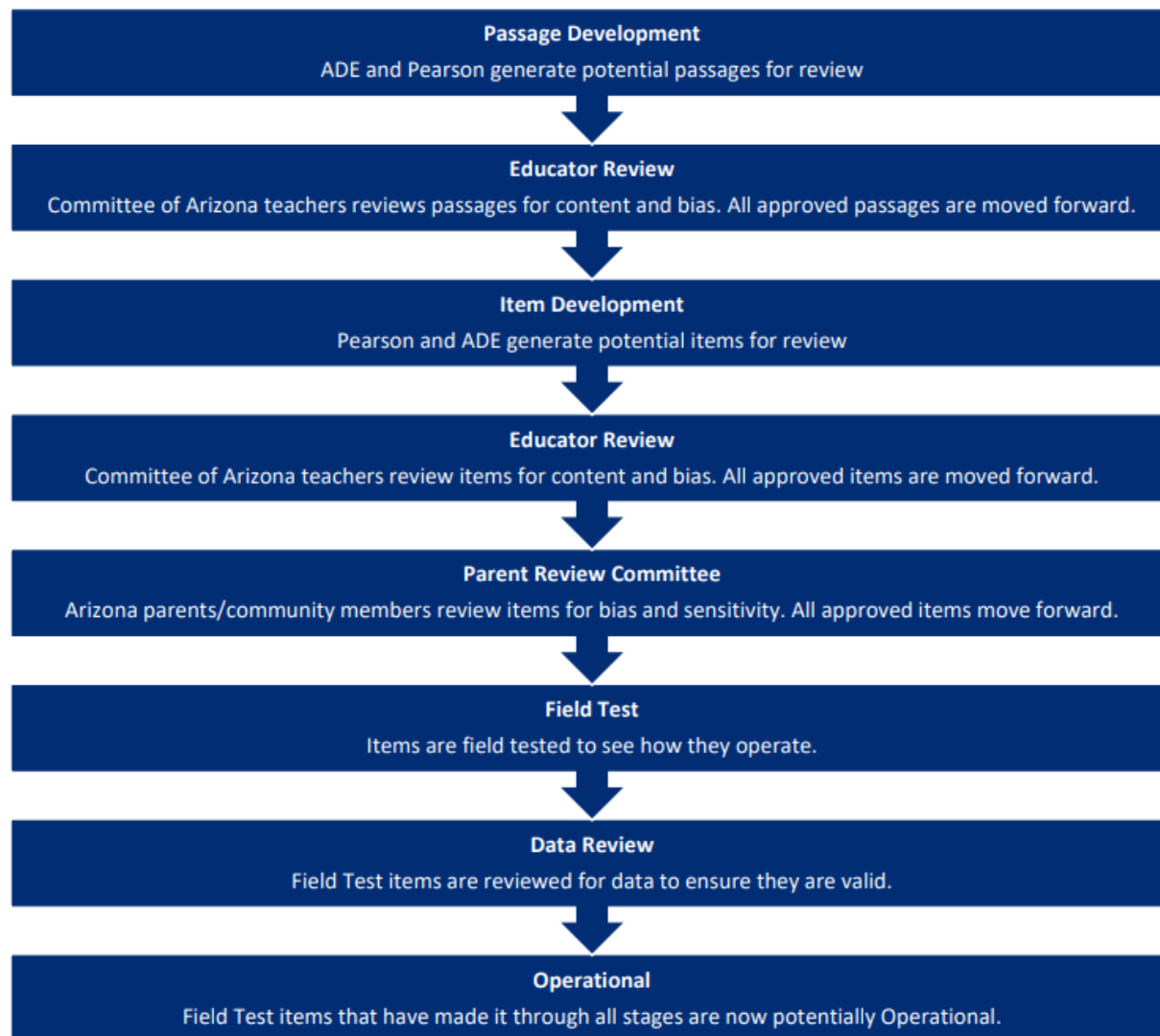
Calculators are not allowed for the mathematics assessments in Grades 3–6. For the Grades 7 and 8 assessments, where calculator use is allowed for some item types, the items are grouped into two units administered separately to students: calculator and no calculator. The construct of the items dictates in which section they are to be assessed.

Arizona has determined content emphasis in the standards at the cluster level for each grade and course. Major clusters are considered as groups of related standards that require greater emphasis than some of the others due to the depth of the ideas and the time it takes to master these groups of related standards. Supporting clusters are considered as groups of related standards that support standards within the major cluster in and across grade levels. Supporting clusters also encompass pre-requisite and extension of grade-level content. Arizona suggests instructional time encompass a range of at least 65%–75% for major clusters and a range of 25%–35% for supporting cluster instruction. Content emphasis can be found at the beginning of all grade-level standards documents at <https://www.azed.gov/standards-practices/k-12standards/mathematics-standards>. The major and supporting clusters align with the test blueprints for AASA mathematics.

Chapter 3: TEST DEVELOPMENT

This chapter addresses Standards 1.11, 3.2, 3.6, 4.0, 4.1, 4.4, 4.6, 4.7, 4.8, 4.10, 4.12, 7.0, 7.2, 12.4, and 12.8 (AERA et al., 2014) regarding item development and test construction. ADE and Pearson worked together to construct the AASA tests based on the steps depicted in Figure 3.1.

Figure 3.1. Item Development Process



Items used to develop the Spring 2023 operational test forms were drawn from the item pool of Arizona-owned items and writing prompts custom-developed to align to the Arizona Academic Standards. The item development process is iterative, allowing for multiple opportunities for review of the items by various stakeholders including ADE and external passage and item content and bias review participants. Newly developed items are then field tested during the spring administration, followed by a data analysis and data review process with Arizona stakeholders. Items that pass data review are added to the operational item bank.

This multistage development and review process provides ample opportunity to evaluate items for their accessibility, appropriateness, and adherence to the principles of Universal Design. In this way, accessibility serves as a primary area of consideration throughout the item development process. This focus on accessibility is critical in developing an assessment that allows for the widest range of student participation as educators seek to provide access to the general education curriculum and foster higher expectations for students.

3.1. Content Development and Management Tool

The item pool, as well as content development and test construction processes, are managed within Pearson's Assessment Banking and Building solutions for Interoperable assessments tool (ABBI) that acts as a content development and management tool, item bank, and publication system supporting both paper-pencil and online publication. The item development workflow is designed to move items and assets from inception through a series of content, fairness, graphic, and other reviews to final publication. The system captures the outcomes at each review and maintains previous versions of each item. As items travel through the review process, every version of each asset is archived, along with each comment received in any review. Reviewers have immediate access to all older versions, providing version control throughout development.

ABBI allows remote internet access by item writers and reviewers while ensuring security with individualized passwords for all users, limited access for external users, and strong encryption of all information. Forms are also built in ABBI. After items are used, ABBI stores the resulting statistics, including exposure statistics and classical and item response theory (IRT) statistics.

The item development process is predicated on a high level of interaction between test developers at Pearson and ADE, as well as with Arizona educators and stakeholders. Pearson's ABBI manages item content throughout the entire lifecycle of an item. It also manages item content beyond the operational life of the item, including items identified for use in sample tests or other training materials. ABBI provides on-demand reports of the content and item bank status. Each item is directed through a sequence of reviews and approvals by Pearson and ADE before it is identified for field test or operational administration.

3.2. Item Bank Analysis

Pearson conducted an item bank analysis at the start of the test development cycle to identify gaps that were then used to determine the priorities for new item development. For ELA, the gap analysis examined the Arizona-owned items in the bank eligible for operational use. A comparison to the blueprint requirements revealed the standards underrepresented in the bank as the focus for new development. For mathematics, the gap analysis identified areas of need for standards coverage as the focus for new item development.

An item development plan was created based on the item bank analysis that outlines the number of items needed to be developed by item type, standard, and DOK. Table 3.1 presents the number of newly developed items that varied by grade and content area depending on the needs of the bank. Standards that were underrepresented in the item bank, or were represented by items with poorly performing statistics, were identified as candidates for item development. Blueprint requirements were also used to determine which standards most needed new item development.

Table 3.1. Number of Newly Developed Items

Content Area	Grade	#Items for FT
ELA	3	133
	4	138
	5	130
	6	130
	7	130
	8	144
	Total	805
Math	3	60
	4	57
	5	56
	6	58
	7	60
	8	61
	Total	352

3.3. Passage and Item Development

Item development for ELA began with the development of reading passages. All new reading passages are commissioned by professional writers who are current or retired educators, while some legacy passages are permissioned. To ensure that all passages align to the correct grade level and provide sufficient complexity for close analytical reading, test developers adhere to detailed passage specifications. The passage specifications call for a close examination of both quantitative measures, such as word counts and Lexile readabilities as shown in Table 3.2, and qualitative measures such as passage structure and levels of meaning, all of which are defined as important measures of text complexity. For example, content experts use passage text complexity worksheets based on the passage specifications to analyze each passage in-depth, as illustrated in Figure 3.2. Table 3.2 also presents the Lexile measures and word count for passages used in the Grade 3 ORF test.

Table 3.2. Passage Lexile Measures and Word Count

Grade	Lexile Range	Word Count Range	ORF Lexile Range	ORF Word Count
3	420–820	100–700	600–750	250–400
4	740–1010	100–900	–	–
5	740–1010	200–1,000	–	–
6	925–1185	200–1,100	–	–
7	925–1185	300–1,100	–	–
8	925–1185	350–1,200	–	–

Note. ORF = Oral Reading Fluency

Figure 3.2. Text Complexity Worksheet Example

UIN:	Word Count:	
Title:		
Genre:	Sub-Genre:	
Quantitative Measures	Flesch-Kincaid:	Lexile:
Qualitative Considerations		
Identify the theme and/or central message and describe how it is adequately developed. (<i>Theme and central message should be similar or the same across paired texts.</i>)		
Briefly describe how the characters are adequately developed, including how they respond to an event or how they change throughout the story.		
Describe the overall structure of a text and how it contributes to the development of the theme, setting, or plot.		
Briefly describe additional plot elements (<i>setting, climax, rising and falling action</i>) that demonstrate clear plot development and how they are similar and/or different across multiple texts. (<i>Paired text only.</i>)		
Explain how you, the author, develops the points of view from which each text is narrated.		
Compare/contrast the differences between the texts when considering genre, theme, and topic.		
Identify one higher level words used in the passage(s) and identify its text support for understanding meaning.		
List grade-level appropriate examples of literary devices used throughout the passage (<i>e.g., metaphor, onomatopoeia, flashback, foreshadowing, voice, irony, symbolism</i>).		
Identify a phrase from the text that has a figurative or connotative meaning and describe the text support.		
Holistically, this text should be considered: ACCESSIBLE MODERATELY COMPLEX HIGHLY COMPLEX for grade _____.		

The next step of item development for ELA and the first step for mathematics was training item writers and introducing them to project requirements. Writers relied on existing item specifications and the Arizona Academic Standards to guide item development. Items were submitted in batches and revised as needed based on feedback from Pearson, with open communication throughout the writing process. Queries were addressed in a timely manner to facilitate a deeper understanding of the Arizona standards and ADE expectations.

Throughout all steps, Pearson responded to ADE feedback, revised, and resubmitted for approval as needed. An integral part of this process was a review by Pearson research librarians who verified accuracy of information and by Pearson copyeditors who reviewed for clarity and correct use of grammar, punctuation, and spelling. All asset creators and reviewers also apply the principles of Universal Design to meet the goal of maximizing accessibility and minimizing construct-irrelevant demands for all items. To meet these goals, text complexity was controlled, graphics were designed to be clear, and subject matter that might affect the student's performance was monitored. Pearson also paid close attention to respecting the diverse cultures of the American Indian tribes in Arizona, particularly to the presentation of topics related to animals.

3.4. Item Review

ADE pre-review was the first of several external reviews of the newly developed passages and items. Educators and community members also had opportunities to participate in review committees. Content and bias review allowed educators to apply their familiarity with Arizona students and the Arizona Academic Standards to provide feedback on the accuracy and appropriateness of the item and stimulus content. A bias and sensitivity community review also allowed parents and other community stakeholders to review assets.

Prior to beginning review, committee members received training from Pearson assessment specialists and were provided resources, including a checklist, to guide the review process. All feedback was recorded in ABBI. The overall goals for both committees were to confirm alignment to the standards, ensure that assets had no bias or sensitivity issues, and revise the assets as needed to be appropriate for Arizona students. An additional benefit of these interactions was that Pearson gained insight to help guide future item development.

ADE and Pearson engaged in a reconciliation process to review committee feedback. Pearson revised assets based on ADE guidance and made the newly edited versions available for ADE review. With ADE approval, the assets went through a final editorial review at Pearson to confirm that they met style expectations and that no errors had unintentionally been introduced.

3.5. Form Construction

Once the newly developed items were ready for field testing, the next step was to construct the test forms, beginning with selecting and positioning the items.

3.5.1. Preparation for Item Selection

Parameters based on the test construction blueprint for each grade were loaded into ABBI by Pearson psychometricians and verified by Pearson assessment specialists. Different test map views were configured based on the specific needs of various users, including the Pearson assessment specialists, ADE and Pearson psychometricians, and Pearson publishing teams. Test maps for each stage were maintained throughout all steps of production. Pearson updated the test maps when any replacements or changes to items or item metadata were made.

Pearson psychometricians had previously loaded statistics from the Spring 2022 administration, and Pearson assessment specialists had updated the ABBI item status used to indicate eligibility for operational or field test selection based on the results from data review. Item statistics included, but were not limited to, classical difficulty (p -value) and item response theory difficulty (Rasch), item discrimination (point-biserial correlation by total score and by reporting category score), the Rasch model fit indices (infit/outfit), differential item functioning (DIF) flags as a measure of possible bias, coefficient alpha, kappa, and distractor analysis.

3.5.2. Item Selection and Positioning

The overriding goal in selecting items for the forms was adhering to the blueprint requirements. Additional criteria for item selection included item positioning and both content and statistical considerations. For each grade, a Pearson assessment specialist did an initial pull of operational items using the tools embedded in ABBI to verify blueprint alignment and acceptable statistics according to the test construction specifications. A different assessment specialist reviewed the form and provided feedback, identifying issues such as clueing. After issues were resolved, a Pearson psychometrician reviewed the form and provided feedback based on statistical considerations. This process repeated until the form met psychometric approval.

The form is also reviewed by the ADE content and psychometrics teams who work with Pearson throughout the process, including final item selection for each form (including the paper and braille versions) and ensuring the psychometric thresholds. Revisions were made based on ADE feedback, and ADE provided the final approval. Once the operational forms were approved, Pearson selected the field test items, with ADE reviewing the field test selections and Pearson revising as needed.

3.5.3. Sampling Plan

Grades 3–7 ELA had 21 forms, whereas Grade 8 ELA had 18 forms. All grades for mathematics had 11 forms. The operational items were the same on all forms within a grade. The test forms were randomly assigned at a student level within a testing group, created by a district, by TestNav, Pearson’s online test delivery platform. Only one paper-pencil version was available per grade.

3.6. Data Review

Field tested items were flagged based on the criteria in Table 3.3. During data review, committee members reviewed the flagged items and their item statistics to determine whether they were eligible for the operational item pool. Two different committees meet for data review. One committee group focused solely on the items flagged for DIF, while another group reviewed the items flagged by the remaining statistics (e.g., item difficulty, point biserial, distractor analysis and Rasch values). The DIF committee looks at the possibility of bias in each item flagged for DIF.

The meeting began with a training session that introduced the item review process, including an overview of the item statistics and how they should be used to evaluate items. Decisions about an item's quality cannot be made on statistics alone; the item itself and the content it measures should also be considered. Thus, the groups also reviewed the content of the items and how the items functioned according to the statistics before making a consensus decision about whether the item should be accepted or rejected for operational use. Revisions were recommended for the rejected items if applicable.

Table 3.3. Item Statistical Flagging Criteria

Statistic	Criterion	Possible Indication
<i>P</i> -value	< 0.2 or > 0.9	Very difficult or easy item
Point-biserial correlation	< 0.25	Poorly discriminating item
Distractor point-biserial correlation (MC only)	> 0.05	Possible miskey*
Omit rate	> 2%	Skipped item
Rasch difficulty	< -3 or > 3	Easy or difficult item
Item fit statistics	< 0.6 or > 1.4	Poor fit
Score point percentage (multi-point items only)	< 1%**	Very few students got a certain score
Differential item functioning (DIF)	B, C	Item could be biased toward a certain student demographic group

*Possible miskey because the key should have a positive point-biserial correlation

**I.e., there should be at least 1% of students at each score point (multi-point items only)

Table 3.4 presents the data review results based on the Spring 2023 data. Committee members made these decisions based on the item content, using the item statistics to guide their discussion. Accepted items were added to the operational item pool for future use. Because the data review committee only reviewed the flagged field tested items, this table does not reflect the total number of field tested items because many did not have any statistical issues or they had fatal statistical issues (e.g., negative point-biserial) that removed them from the item pool.

Table 3.4. Data Review Results: Number of Field Tested Items

Content Area	Grade	#Accepted	#Accepted w/Edits	#Rejected
ELA	3	126	0	7
	4	128	0	10
	5	119	0	11
	6	109	0	21
	7	106	0	24
	8	121	0	23
Math	3	61	1	2
	4	53	0	0
	5	57	2	1
	6	53	6	2
	7	59	2	0
	8	60	4	0

3.7. Accommodated Forms

Each grade and content area had one form of the paper-pencil Special Paper Version (SPV). The Pearson content team worked with ADE to produce paper-equivalent versions of the items used on the online test form. Upon approval of the item set, the Pearson publishing team worked with ADE to determine an approved paper-based test template for each grade. There were three rounds of review between ADE and Pearson before the document was approved to print. A final PDF printer proof was provided to ADE.

Upon approval of the paper-pencil form, Pearson began work on the Large Print and Braille forms. The Large Print forms are enlarged versions of the paper-pencil test forms. The publishing team enlarged the entire test book file to reach an 18-point font equivalent. The final Large Print printer proof file was posted for ADE's review and approval.

The Inkprint Braille version of the test was modified based on the Braille modification document to reflect any item omissions or modifications on the Student Braille Test Book. Pearson Braille Services reviewed all forms presented for Braille to determine if forms were well-suited for Braille testers. Any recommended modifications were reviewed in conjunction with ADE to arrive at final decisions. ADE then reviewed the Inkprint Test Book, the Student Braille Test Book proof, the Braille Test Administration Directions, and the Braille memo before production of the Braille material commenced.

Each grade and content area also had one form created for ASL testers. After approval by ADE of the online test form, Pearson ASL team began work for ASL translation. The Pearson ASL team created scripts to be used for filming of the ASL translation by professional ASL signers. Video sessions for ASL Filming were attended by the Pearson ASL team as well as Pearson content for any questions that arose during translation. ADE had final approval of any modifications necessary for successful ASL filming. All ASL videos and test forms were reviewed and approved by ADE before final production.

Chapter 4: TEST ADMINISTRATION

This chapter describes how the AASA assessments were administered, including the procedures used to ensure that the test administration was conducted in a secure and standardized manner, as indicated by Standards 1.10, 3.1, 3.9, 3.10, 4.2, 4.5, 4.15, 4.16, 4.21, 6.1, 6.2, 6.3, 6.4, 6.5, 6.6, 6.7, 7.0, and 7.8 (AERA et al., 2014). The AASA assessment is administered online via TestNav, Pearson’s online testing platform that students use to access the assessment, with accommodated forms available as needed. PearsonAccess^{next} (PAN) is the student test management portal that test administrators use to manage student tests and registrations and order materials if needed.

District Test Coordinators (DTCs), School Test Coordinators (STCs), and Test Administrators (TAs) received online training and the supporting documents to ensure fidelity of implementation and the validity of the assessment results and to help prevent, detect, and respond to irregularities in academic testing and maintain testing integrity practices for technology-based assessments. For example, TAs were instructed to use the *Test Administration Directions* (TAD) for the online and paper administrations, as well as for the Special Paper Version (SPV) tests and entering student responses into TestNav.

When all TAs use the same well-defined administration procedures and are provided the same training, manuals, and supporting documents, administration is optimally standardized and poised to be fair to all students. DTCs were responsible for supporting the TAs in understanding and following the administration procedures. Comprehensive test coordinator training and materials targeted to their role and responsibility ensure that they are appropriately prepared to support the test administrators.

4.1. Test Units

Table 4.1 presents the estimated time to complete each test unit. A test unit must be completed prior to starting the next one. All ELA Writing and Reading test units must be administered to receive an ELA score, and both mathematics test units must be administered to receive a mathematics score. The ELA Writing test must be administered on a separate day than the ELA Reading and mathematics units. ELA Reading and mathematics test units could be administered in any order, with no more than two test units plus the Grade 3 ORF unit in a single day. If two test units were administered on the same day, there must be a significant break between them. ADE requires that a test unit be submitted within the day that it is started. Any test that is not complete at the end of the testing day is marked complete and submitted for scoring by Pearson.

As part of the operational test administration, Grade 3 students also participated in the ORF test unit that was field tested in Spring 2022 and again in Spring 2023; thus, the items were not included in scoring. Each student read three separate passages, with a time limit of one minute per passage. The ORF online test unit was to be administered in small groups, with no more than six students testing simultaneously in a classroom or a computer lab environment. For paper-based testing, ORF was administered one-on-one on a speaker telephone.

Table 4.1. Estimated Testing Time by Test Unit

Unit	Testing Time
ELA Writing	60–90 minutes
ELA Reading Test Unit 1	45–75 minutes
ELA Reading Test Unit 2	45–75 minutes
Grade 3 Oral Reading Fluency (ORF)	15 minutes
Math Test Unit 1	60–85 minutes
Math Test Unit 2	60–85 minutes

Note. The testing time is the same for the CBT and PBT administrations.

4.2. Administration Materials

Table 4.2 describes the materials provided to support the standardized administration of the AASA assessments and ensure fair testing for all students. The TAD and *Test Coordinator’s Manual* (TCM) were produced in collaboration with ADE. The Pearson program team drafted each manual using the previous year’s version as a template. The manuals were then composed in desktop publishing software and sent for an editorial review. After a review of all comments and edits by the program team, the file was delivered for ADE review. There were multiple rounds of review between ADE and Pearson before the document was approved to print. ADE was provided with a final web-ready 508 compliant version in addition to the final printer’s proof. Hard copies were sent automatically to all participating schools, and a limited number were available for additional order during the additional order window. The materials are available on the ADE website at <https://www.azed.gov/assessment/aasa>.

Table 4.2. Administration Materials

Material	Description
<i>Test Administration Directions</i> (TAD)	Provides an overview of the AASA test administration, including the user roles in PAN and the test administration schedule, and directions about what to do before, during, and after testing. Provided for both the CBT and PBT assessments.
<i>Test Coordinator’s Manual</i> (TCM)	Indicates the responsibilities of the DTCs before, during, and after testing and explains the procedures for test administration. DTCs must review the TCM and the TAD well in advance of training STCs and TAs and before administering the tests. DTCs are responsible for ensuring the appropriate and correct administration of the AASA in all schools within the district or under the same charter.
<i>PAN User’s Guide</i>	Explains how to navigate PAN and the tasks related to the AASA test administration.
<i>Arizona Accommodation Manual</i>	Lists the current accommodations, accessibility features, and tools available on Arizona’s achievement assessments.

4.3. Administration Training

Mandatory test administration training was provided by ADE and Pearson and delivered through Pearson’s online Training Management System (TMS) that contained the training modules summarized in Table 4.3 that were required for DTCs, STCs, TAs, and other school staff involved in testing or test results.

The online training modules were available prior to the beginning of the testing window and throughout the testing window. The training modules addressed the specific responsibilities of the DTC and provided important information from the three documents TAs are required to use (i.e., the TAD, TCM, and *PAN User's Guide*). These training modules are updated for each test administration in correspondence with the updates to the required documents. Each module requires approximately 30–45 minutes to complete. DTCs are required to view the training modules in sequence and to successfully complete a final quiz after viewing all modules. DTCs must obtain a score of 80% or higher on the final quiz to be certified to access the secure test administration materials. DTCs are allowed multiple attempts to obtain a score of 80% or higher on the final quiz.

Table 4.3. Administration Trainings

Training	Description
AASA Training for Test Coordinators	This training covered the AASA test administration for Grades 3–8, including an overview of the test administration, websites and resources, and responsibilities before, during, and after testing.
Accommodations	This training covered the test accommodations. This was required for all DTCs but could be shared with staff members.
Achievement Test Administration Responsibilities	This training covered the test administration of AASA and AzSCI for all employees who administered, proctored, or were in contact with test materials. The purpose of this training was to provide guidance on consistent test administration across the state, increase the number of valid student tests, reduce test improprieties, and limit staff exposure to accusations of testing violations and discipline.
Test Security and Ethics	This training covered policies and practices to ensure the security and confidentiality of testing materials and the reliability and validity of test score interpretation. This training module was required for all employees who administered, proctored, or came in contact with testing materials.
PearsonAccess ^{next} (PAN)	This training covered PAN and was required for DTCs, STCs, and other testing staff who assisted with registering students or managing test sessions in PAN.
Technology Training	This training outlined the critical steps necessary to prepare the network, testing devices, and other technology related items required for a successful test administration.

4.4. Sample Tests

In addition to the module training, TAs are instructed to become familiar with the online system by accessing sample items. Sample tests are available in TestNav year-round to help TAs and students become familiar with the AASA item types. The sample tests were created following Pearson's standard item and test development process, including item content and bias review by Arizona educators and community members. The sample tests reflect the AASA test specifications and blueprints and had 1–25 items on each test, as shown in Table 4.4. Because the sample tests do not include an item for each of the aligned Arizona Academic Standards and do not provide scores for students, they should NOT be used to evaluate a student's performance level. Students access the test as a guest, so no personal information needs to be provided.

There is a sample test for each grade and content area, and every eligible item type is represented, and an accompanying scoring guide identifies standard and DOK alignment. The portal and scoring guides are both available on ADE website at <https://www.azed.gov/assessment/aasa>. Scoring guides for the sample tests are also available.

Table 4.4. Number of Items on the AASA Sample Tests

Grade	ELA	Writing	ORF	Mathematics
3	24	1	3	25
4	24	1	–	25
5	24	1	–	25
6	24	1	–	25
7	24	1	–	25
8	24	1	–	25

4.5. Accommodations

Accommodations are specific practices and procedures that provide students with equitable access during the assessment. They are made to provide a student equal access to learning and equal opportunity to demonstrate what is known and are intended to reduce or even eliminate the effects of a student's disability. Accommodations can be changes in the presentation, response, setting, and timing/scheduling of educational activities. There should be a direct connection between a student's disability, special education need, or language need and the accommodation(s) provided to the student during educational activities, including assessment.

Students should receive the same accommodations for classroom instruction, classroom assessments, district assessments, and state assessments. No accommodations should be provided during assessments that are not also provided during instruction. However, not all accommodations appropriate for instruction are appropriate for use during a standardized state assessment. Table 4.5 presents the accommodations available to students while testing on Arizona assessments.

Table 4.5. Available Accommodations

Accommodation	Description
Abacus	Students may use an abacus without restrictions for any mathematics test (for students taking the Braille test only). Students may use an abacus without restrictions for any mathematics test or a talking calculator for students taking Unit 1 of the Grades 7 or 8 mathematics test.
Adult Scribe	A student who requires one-on-one adult assistance during daily instruction may orally dictate or use gestures to indicate a selected response for multiple-choice items only while an adult enters this in the test. The adult may not ask or answer any questions during the session or influence student responses in any way.
American Sign Language (ASL)	ASL requires the use of a different test form that must be indicated in PearsonAccess ^{next} (PAN). The ASL test form must be requested using the Additional Accommodations online request form.
Braille test booklet	Braille tests must be requested using the special paper version (SPV) test online request form. Requires adult transcription: An adult must transfer the student's response exactly as written into the TestNav system.

Accommodation	Description
Large print test booklet	Large Print tests must be requested using the special paper version (SPV) test online request form. The 504 plan or IEP must clearly state the font size used for instruction and the type of materials teachers enlarge for the student. Requires adult Transcription: An adult must transfer the student's response exactly as written into the TestNav system.
Paper test booklet	A student who cannot access the computer for classroom work due to injury, illness, or vision impairments may need a paper test in lieu of taking the test with peers on the computer. Requires adult transcription: An adult must transfer the student's response exactly as written into the TestNav system.
Sign test content	Any student who requires signing of content during daily instruction may have any of the content of writing, mathematics, and science signed.
Simplified test administration directions	The test administrator may provide verbal directions in simplified English for the scripted directions from the <i>Test Administration Directions</i> manual. This must take place in a setting that does not disturb other students.
Translated test administration directions	Exact oral translation, in the student's native language, of the scripted directions from the <i>Test Administration Directions</i> manual are permitted. No test content or directions embedded within the test may be translated.
Translation dictionary	During testing, students may use the word-for-word published paper translation dictionary that is used regularly for classroom instruction. Students with a visual impairment may use an electronic dictionary with other features turned off.

4.6. Universal Test Administration Conditions

The following Universal Test Administration Conditions are testing situations and conditions that may be offered to any student to provide a comfortable and distraction-free testing environment. They do not require an accommodations request. While some of the items listed as Universal Test Administration Conditions might be included in an IEP or 504 plan as an accommodation, for achievement testing purposes these are not considered testing accommodations and are available to any student who needs them.

- Testing in a small group, 1:1, or in a separate location on campus or in a study carrel
- Being seated in a specific location within the testing room or at special furniture
- Having the test administered by a familiar test administrator
- Using a special pencil or pencil grip
- Using a placeholder
- Read-aloud (text-to-speech or human reader) content of the ELA writing, mathematics, and science assessments
- Using devices that allow the student to see the test: glasses, contacts, magnification, and special lighting
- Using different contrast settings or color overlays
- Using devices that allow the student to hear the test directions: hearing aids and amplification
- Wearing noise buffers after the scripted directions from the *Test Administration Directions* manual have been read
- Signing the scripted directions from the *Test Administration Directions* manual
- Repeating the scripted directions from the *Test Administration Directions* manual
- Having assistance with logging into an online test
- Reading the test quietly to themselves as long as other students are not disrupted

- A phone or electronic device needed for medical care is permitted. The phone needs to stay close to the Test Administrator or proctor as well as the student and should be monitored to assure the device is only being used for medical purposes during testing
- Individual students may take a stretch break (1 or 2 minutes) during the test session (students may not talk, use electronic devices, go to lunch, or leave the testing room)
 - Paper test booklet and scratch paper must be collected
 - Students must sign out of TestNav without submitting the test. The test administrator will need to resume the student's test session using PAN.
- Students may use the restroom (only one student at a time)
 - The TA must collect the student's paper test booklet and scratch paper.
 - Students must sign out of TestNav without submitting the test. The test administrator will need to resume the student's test session using PAN.
- The use of scratch paper (plain, lined, or graph; school provided). Scratch paper must be securely shredded at the conclusion of testing
- Each testing session must be completed in the same school day in which it was started. The AASA and AzSCI are untimed. Do not start a test unit unless there is sufficient time to complete the test in the same school day.
- Students cannot leave for lunch during a test session. Test units should be scheduled in a way that provides the student more than adequate time to complete the test.

4.7. Universal Test Tools

The Universal Test Tools provided in Table 4.6 are available to all students taking the AASA assessment and cannot be disabled.

Table 4.6. Universal Test Tools

Universal Test Tool	Description
Alternate Mouse Pointer	There are six alternate mouse pointers available for students in TestNav. Alternate options include a medium, large, or extra-large sized white pointer, and extra-large sized black, green, or yellow pointer.
Answer Masking	Allows student to electronically cover and reveal individual answer choices.
Answer Eliminator	Cross out answer options for multiple-choice and multi-select items.
Area Boundaries	Allows student to click anywhere on the selected response text or button for multiple choice items.
Bookmark for Review	Mark an item for review so that it can be easily found later.
Contrast	Allows the student to change the background and text color based on need or preference. The Contrast setting will not change images or artwork. The options are white background with black text; cream background with black text; light blue background with black text; black background with white text; light magenta background with black text; and blue background with yellow text.
Expand/Collapse Passage	Expand a passage for easier readability. Expanded passages can also be collapsed.
Highlighter	Highlight text in a passage or item.
Line Reader	An adjustable box allows the student to focus on one line or a few lines at a time. The box can be adjusted to increase or decrease the number of lines shown. The Line Reader and Magnifier tools may be used simultaneously.
Magnifier	Allows the student to make part of the screen larger. When in use, the magnifier can be moved around the screen as needed.
Pause and Restart	Students may sign out of TestNav. Before the student can resume testing, the Test Administrator will need to resume the student's session in TestNav.

Universal Test Tool	Description
Notes/Comments	Allows student to open an on-screen notepad and take notes or make comments. Notes carry over within a passage set. In non-passage items, notes are attached to the specific test item on which they are entered.
Review Test	Allows student to review the test before submitting it.
System Settings	Adjust audio (volume) during the test.
Text-to-Speech	Text-to-Speech for content of writing, mathematics, and science.
Tutorial	Learn and practice using TestNav tools and responding to each item type.
Writing Tools	Editing tools (cut, copy, and paste) and basic text formatting tools (bold, underline, and italic) for extended response items.
Zoom In/Zoom Out	Enlarge the font and images in the test up to 200%. Undo zoom in and return the font and images in the test to original size.

4.8. Pearson Customer Support

To provide support to schools before, during, and after testing, Pearson provides tiered technical support Monday – Friday from 7:00 a.m. to 7:00 p.m. CST. DTCs, STCs, and TAs can contact the customer support line with questions pertaining to the TestNav and PAN system and test administration procedures. The toll-free support number, e-mail address, and chat link are disseminated to the field through the AASA system and related communications.

4.9. Test Security

All test coordinators, test administrators, and proctors must be trained in proper test security procedures, must sign an Achievement Tests Staff Security Agreement form (as shown in Figure 4.1), and must adhere to test security procedures. Test materials should be secured prior to, and at the conclusion of, all testing sessions. Test Administrators and proctors may not assist students in answering test items and may not translate, reword, or explain any test content. No test content may ever be discussed before, during, or after test administration. It is unethical and shall be viewed as a violation of test security for any person to:

- Log into TestNav as a student unless assisting student with log in procedures
- Share their username/password for PAN
- Capture images of any part of the test via any electronic device
- Duplicate in any way any part of the test
- Examine, read, or review the content of any portion of the test
- Disclose, or allow to be disclosed, the content of any portion of the test before, during, or after test administration
- Discuss any test item before, during, or after test administration
- Allow students access to test content prior to testing
- Provide any reference sheets to students during the mathematics test administration or graphic organizers during the Writing test administration
- Allow students to share information during test administration
- Read any parts of the test to students, except as indicated in TAD or as part of an approved accommodation
- Influence students' responses by making any kind of gestures (e.g., pointing to items, holding up fingers to signify item numbers or answer options) while students are taking the test
- Instruct students to go back and reread/redo responses after they have finished their test since this instruction may only be given before the students take the test

- Review students' responses
- Change students' answer choices
- Read or review students' scratch paper
- Participate in, direct, aid, counsel, assist in, encourage, or fail to report any violations of these test administration security procedures

Figure 4.1. Test Security Agreement



Assessments Achievement

Achievement Tests (AASA, AzSCI, ACT Aspire, and ACT) School Year 2022-2023 Staff Test Security Agreement

I acknowledge that all Achievement Tests are secure tests and agree to the following conditions of use to ensure the security of the test. For this document, Achievement Tests refers to AASA, AzSCI, ACT Aspire, and ACT.

1. I shall take necessary precautions to safeguard test materials.
 - a. I shall sign an *Achievement Tests Staff Security Agreement* for School Year 2022-2023.
 - b. Access to test materials, including online tests, is restricted. I shall not attempt to gain access to test materials beyond that which is granted to me by my school/district test coordinator, superintendent, or charter representative.
 - c. If test materials are distributed to me, I shall keep them under lock and key except during actual test times. This includes any student data sheets or student information sheets provided to me.
 - d. I shall not permit students to remove test material from the testing room except under the supervision of staff.
 - e. I shall not examine, read, or review the Achievement Tests.
 - i. I shall not disclose, nor allow to be disclosed, the content of the test.
 - ii. I shall not discuss any test item at any time.
 - iii. I shall not examine, read, or review any student responses.
 - iv. I shall not log into any student online test.
 - f. I shall not erase or change any student responses or any marks (including stray marks) on a scorable test booklet or answer document.
 - g. If test materials are distributed to me, I shall return all test materials to the school/district test coordinator immediately upon the completion of testing.
 - h. I shall not use any test materials for instruction before or after test administration. I shall follow *Test Preparation and Administration Practices*, the guidelines approved by the State Board of Education in January 2003 and updated in December 2007.
 - i. I shall not provide prohibited or inappropriate resources to students during testing, including but not limited to graphic organizers, reference sheets, and calculators, except for tests and test sections where calculators are allowed.
2. I understand that the district superintendent or charter representative will develop, distribute, and enforce disciplinary procedures for the violation of test security by staff.

Individuals who will administer or proctor Achievement Tests for school year 2022-2023 must also agree to the following conditions to ensure the correct administration of the tests.

3. I shall participate in training activities prior to administering the tests.
4. I shall review the appropriate Test Administration Directions prior to administering the test.
5. I shall follow all instructions in the appropriate Test Administration Directions including reading the directions to students exactly as scripted.

By signing my name to this document, I am assuring my district/charter and the Arizona Department of Education that I will abide by the above conditions and that anyone I supervise, who will have access to the Achievement Tests, will also sign a Test Security Agreement.

Signed By: _____ Date: _____

Printed Name: _____

Title: _____ School: _____

Please return signed copy as per instructions from your school/district test coordinator.

In addition to test security procedures required of all educators involved in the testing process, TestNav has built-in security features for the test content and personal data that relies on multiple levels of protection, including restricted user access, encryption of data in transit and at rest, systems monitoring for abnormal behavior, application, server, and network security testing, and qualified, verified and trusted support personnel.

Pearson uses Advanced Encryption Standard (AES) encryption for data at rest and Hypertext Transfer Protocol Secure (HTTPS) to provide encryption and data-in-motion security for online testing by creating a secure channel on the network with the Secure Socket Layer (SSL) /Transport Layer Security (TLS) protocols. Test content can only be viewed through a valid test registration and login, all of which are logged within the platform's audit trail system and cannot be deleted.

TestNav also locks down the student's desktop during testing to prevent students from accessing outside resources that could be used for cheating, such as email, instant messaging, or internet browsing. TestNav will stop students' tests if another background application attempts to interfere with or take "focus" away from the secure testing environment. These types of interruption cannot be blocked during testing and therefore could present additional opportunities for students to access unauthorized resources. However, TestNav also has a blocklist feature that prevents students from starting their test if certain applications that pose a threat to disrupt testing are running at the time TestNav is launched. In these situations, the student and/or proctor are prompted to shut down the offending application before attempting to start TestNav again.

Chapter 5: SCORING AND REPORTING

This chapter describes the human-scoring procedures used by the Pearson Performance Scoring Center (PSC) to score the AASA writing, reading, and mathematics open-ended items, as well as the automated scoring procedures for the writing prompts. This section addresses Standards 2.7, 4.18, 4.19, 4.20, 6.8, and 6.9 (AERA et al., 2014) regarding the scoring of the assessments.

The AASA machine-scored items were scored with maximum likelihood estimation (MLE) scoring, with an attemptedness rule that a student needed to answer one item in each operational unit. Both ELA and mathematics have their own scale score ranges. Students received a scale score in each content area, and student performance was reported as one of four performance levels: Level 1: *Minimally Proficient*, Level 2: *Partially Proficient*, Level 3: *Proficient*, and Level 4: *Highly Proficient*.

Student performance on reporting categories is reported as one of three levels of mastery: *Below Mastery*, *At/Near Mastery*, or *Above Mastery*. Students who score *Below Mastery* demonstrate performance in the reporting category that was clearly below *Proficient*. Students who score *At/Near Mastery* demonstrate performance in the reporting category that was exactly at or immediately above/below *Proficient*. Students who score *Above Mastery* demonstrate performance in the reporting category that was clearly *Proficient* or higher.

5.1. Human Scoring of Open-Ended Items

The AASA assessments contain open-ended items that prompt students to write a short answer or extended response (i.e., a paragraph or multi-paragraph essay) that require scoring by professionally trained scorers. These items were the writing prompts on the ELA Writing test, short constructed-response items on the ELA Reading test, and the paper-equivalent of the technology-enhanced (TE) items on the ELA Reading and mathematics assessments. Writing was scored via a distributed scoring model (i.e., scorers were trained in a self-paced model), whereas Reading and mathematics were scored using a synchronous model (i.e., scorers were trained by instructors). Human scoring was conducted in Pearson's scoring platform known as the Electronic Performance Evaluation Network (ePEN2).

5.1.1. Scorer Recruitment

Scorers are recruited by Pearson, with scorers who have extensive experience scoring this type of rubric on previous projects being given priority. Scorers receive performance ratings based on internal quality metrics of inter-rater reliability and validity. Those who have achieved a high-performance rating on previous writing, reading, and mathematics responses are recruited for the AASA assessment. Upon being hired, scorers sign a confidentiality agreement in which they pledge to keep all information and student responses confidential.

Scoring supervisors are chosen based on demonstrated expertise in the scoring process, including strong organizational abilities and training, practical skills, leadership abilities, and sensitivity to interpersonal communication requirements. Supervisors also possess the essential capability of helping scorers understand the AASA scoring requirements. Supervisors provide continuous feedback to the scorers through the validity and calibration process and monitor the quality of their assigned scorers. All scoring, including the scorers and supervisors, is supervised by a content specialist who is responsible for training and leading the entirety of the project.

5.1.2. Training

Scorers and scoring supervisors were trained to learn the rubric and score responses according to the AASA scoring guidelines. At the beginning of the scoring project, all scoring supervisors and scorers completed project-specific training consisting of a review of the rubric and prompts for the items being scored and a review of the anchor responses selected and approved by ADE for each prompt. Training for the ELA Writing prompts differed than the training for the Reading and mathematics open-ended items. Writing established training materials that could be inserted into modules for self-paced training, whereas training materials for Reading and mathematics were created as the students completed testing. This could be accomplished because the Reading and mathematics open-ended items were only 0,1 score point items.

5.1.2.1. Writing

The training for ELA Writing was conducted in a distributed environment using online modules designed to take scorers through the background of the assessment and the rubric and anchor sets for each item. A module is an online set of training materials that can be delivered to scorers individually at their own pace. These modules are embedded into the ePEN2 system and are set up so as not to allow scorers to advance in their training until all proceeding modules are complete and correct.

Scoring supervisors and scorers were both required to take one set of practice papers and two sets of qualification papers once they completed the item-specific modules. They must have passed one of the two qualification sets for the items they were assigned before they could score on the project based on the criteria in Table 5.1. Their scores were compared to the “true score” approved by ADE for each training response. Once the scorer completed the item-specific training and had qualified, they were allowed to score live responses for that item or set of items. Different scoring rubrics are used for the different item types and are posted on the ADE website at <https://www.azed.gov/assessment/aasa>.

Table 5.1. Scoring Qualification Standards

Reporting Category	Score Points	Qualification %:Perfect/Adjacent Agreement	#Sets
Writing Multi-trait	1–4	60/90 for each trait at least once across the two sets	2

5.1.2.2. Mathematics and Reading

Prior to scorer training for reading and mathematics, scoring directors reviewed items/passages and rubrics and selected actual student responses to review and discuss at rangefinding sessions with ADE staff. The rangefinding sessions allowed Pearson and ADE to discuss any questions regarding possible correct answers and assign final scores for the student responses. These scored student responses from the rangefinding sessions used to create an anchor and practice set for reading and “prototype” items for mathematics used as initial training items for an item type that included an anchor set and practice set. The sets were shared with ADE and adjusted as needed for final approval.

Training was conducted in the train-score-train-score model live via online conferencing where scoring directors trained scorers on the content for a single item and worked with the team to score that item before moving to train the second item. There were two separate ELA teams and two separate mathematics teams, each led by a scoring director. Mathematics scoring directors began content training on a prototype item, reviewing the prompt, rubric, and the anchor set for the item. The team then took and discussed a practice set to test their knowledge of rubric application before moving into live scoring. Subsequent similar items were trained with bridge sets. For such items, the scoring director would prepare the team by covering the prompt, rubric, and bridge set. Reading scoring directors began content training on every item, reviewing the prompt and passages, the rubric, and anchor set. The team then took and discussed a practice set to test their knowledge of rubric application before moving into live scoring.

5.1.3. Quality Control

A variety of reports are produced throughout the scoring process to monitor the progress of the project, the reliability of scores assigned, and individual scorers' work:

- Daily and Cumulative Interrater Reliability Reports by item and scorer that indicate how many times scorers were in exact agreement or assigned adjacent scores. The reliability is computed and is monitored daily and cumulatively for the project.
- Daily and Cumulative Validity Reports by item and scorer that indicate how many times scorers were in exact agreement or assigned adjacent scores to responses deemed True Scores. The validity is computed and monitored daily and cumulatively for the project.
- Daily and Cumulative Frequency Distributions that show how many times each score point has been assigned to the item being scored. The frequency distributions are produced daily and cumulatively for the entire scoring project. This report allows scoring supervisors and directors to see whether scorers tend to score consistently high or low.

The most immediate method of monitoring a scorer's performance is through backreading by scoring supervisors and directors. If a scoring supervisor discovers that a scorer is consistently assigning scores other than those the scoring supervisor would assign, they can send a message to that scorer using the backreading function and through the ePEN2 instant messaging system.

With the help of the individual scorer reliability metrics and through backreading, the scoring staff can closely monitor each scorer's performance. Scorers are also monitored using the scorer exception process for validity and scoring rate. A scorer must meet and maintain the quality metrics established for AASA in the designated area to continue scoring the project. If a scorer fails to maintain the established validity perfect agreement and perfect plus adjacent agreement percentage, they will receive a targeted calibration set consisting of 10 anchor-type responses similar to a qualification set. If the scorer fails to pass the calibration set, they will be locked out of scoring and dismissed from the project.

Scorers with low inter-rater reliability or a lower- or higher-than-desired scoring rate are closely monitored in backreading and through reports. If, in the opinion of the scoring director and content specialist, these scorers are still performing below acceptable standards after receiving sufficient feedback and being given every reasonable opportunity to improve, they are manually locked out of the system and dismissed from the project.

5.1.4. Security

To ensure that test security is never compromised, the following safeguards are employed:

- Scorers and scoring staff personnel must sign a non-disclosure and confidentiality form in which they agree not to use or divulge any information concerning the tests.
- All contact with the press is handled through ADE.
- ePEN2 is accessed via a secure website with login credentials required for each user. Only Pearson project support staff can issue user IDs to scorers to access ePEN2.

5.2. Automated Scoring of ELA Writing Prompts

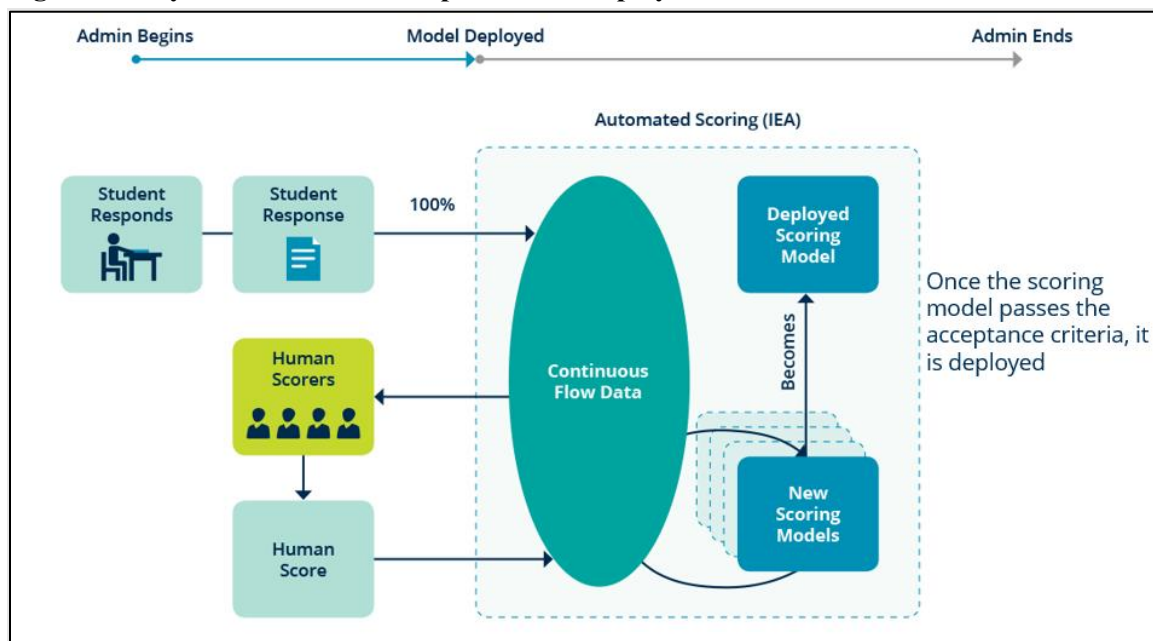
Pearson’s automated scoring engine, the Intelligent Essay Assessor (IEA), is the default option for scoring the AASA ELA writing prompts. For the operational writing prompts in Spring 2023, the automated scoring engine was calibrated based on previously tested and human-scored field test responses. During the scoring window, human-scored student responses were used to continue improving and validating the scoring models via Continuous Flow. All the ELA prompts were scored at least in part by IEA in the spring. For 10% of responses, a second reliability score was assigned by human scorers to provide data for evaluating the consistency of scoring, which is done by evaluating scoring agreement.

IEA is trained by humans anytime a new writing prompt is introduced and follows the Continuous Flow process that incorporates human scoring to ensure the highest-quality scores. Responses flow between the engine and human scorers so the engine can learn from humans in real time and challenging responses can be instantly routed to human scorers (known as Smart Routing). When the engine is less confident in scoring a response, the response is marked with a low confidence flag that automatically routes it to human scorers. Human scoring is applied to responses that are scored while IEA is being trained, as well as to the Smart Routing responses. When multiple scores are assigned for a given response, the IEA score is reported operationally if it is a high confidence score. If the IEA score is low confidence, the human score is assigned.

5.2.1. Calibration of IEA

With Continuous Flow, human scorers begin the scoring process and IEA learns from them. This process can begin with previously tested and human scored field test items or during the operational scoring window. For the 2023 Writing prompts, IEA used a combination of human-scored field test and operational data to calibrate the automated scoring engine. The field test data were used to build initial models. Some prompts had enough data to build a full scoring model that passed the criteria described in Section 5.2.3, while other prompts required additional data to meet the criteria using the Continuous Flow process to supplement the models with human-scored operation data until it met all the quality criteria. Figure 5.1 presents scoring model development and deployment in the Continuous Flow scoring approach.

Figure 5.1. Dynamic Model Development and Deployment



The early performance of human scoring was monitored based on the following characteristics to verify that an appropriate set of data was available for training IEA:

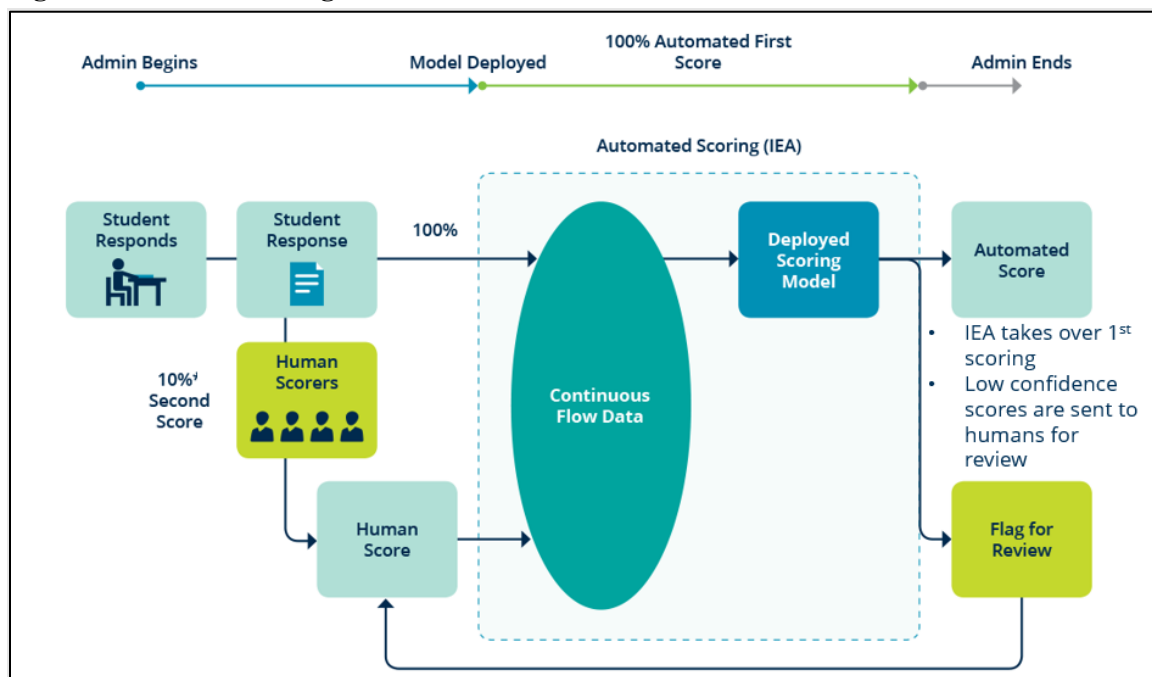
- Exact agreement between human scorers (with the goal of at least 65%)
- Exact agreement between human scores conditioned on score point (with the goal of at least 50%)
- The number of responses at each score point
- The number of responses with two human scores assigned (IEA via Continuous Flow “ordered” additional scoring of responses during the sampling period as needed)

Although the desired characteristics of the training data were easily achieved for some prompts, they were more challenging to achieve for others. For some prompts, a subset of scores were reset and clarifying directions were provided to scorers to improve human-human agreement. A healthy percentage of responses were also backread during the sampling period. These scores in addition to the double human scores were all part of the data used to train IEA.

5.2.2. Smart Routing

As illustrated in Figure 5.2, once IEA is trained, it takes over first scoring with human scorers providing the 10% second score for reliability. Smart Routing refers to the practice of using automated scoring results to detect responses that are likely to be challenging to score and applying automated routing rules to obtain one or more additional human scores on those responses. Smart Routing can be applied prompt-by-prompt to the extent needed to meet scoring quality criteria for automated scoring. When the engine is less confident in scoring a response, the response is marked with a low confidence flag that automatically routes it for human scorers.

Figure 5.2. Smart Routing



5.2.3. Quality Control

IEA performance on the writing prompts was evaluated based on IEA-human exact agreement and compared to agreement based on responses that were double-scored by humans. The following industry-standard measures were computed between pairs of human scores and between IEA and humans to evaluate scoring performance:

- Pearson correlation between IEA-human should be at least 0.70 and within 0.1 of human-human.
- Quadratic-weighted kappa between IEA-human should be at least 0.70 and within 0.1 of human-human.
- Standardized mean difference between IEA-human should be less than $|0.15|$.
- With Smart Routing applied as needed, exact agreement between IEA-human should meet the inter-rater reliability requirement of at least 65% and be within 5.25% of human-human exact agreement. If the IEA-human agreement is within 5.25% of the human-human agreement, IEA can be deployed operationally. This is the primary criterion for evaluating IEA.

In addition to the overall comparison, the following performance thresholds were targeted in the test data set: (1) at least 65% overall IEA-human agreement and (2) 50% IEA-human agreement by score point (conditioned on the human score).

5.3. Reporting

The following AASA reports were available in PAN at <https://az.pearsonaccessnext.com>. PDF versions of the reports and district-wide electronic student data files were also available for downloading. District-level user roles provided access to all school-level reports and district-level reports, including all Confidential Student Score Reports for students who tested in the district. School-level user roles provided access to all school-level reports and all Confidential Student Score Reports for students who tested in the school. A Family Guide for interpreting reports was also available for download. Figure 5.3 and Figure 5.5 present sample reports, including an example from Grade 8 to show the ACT predicted score.

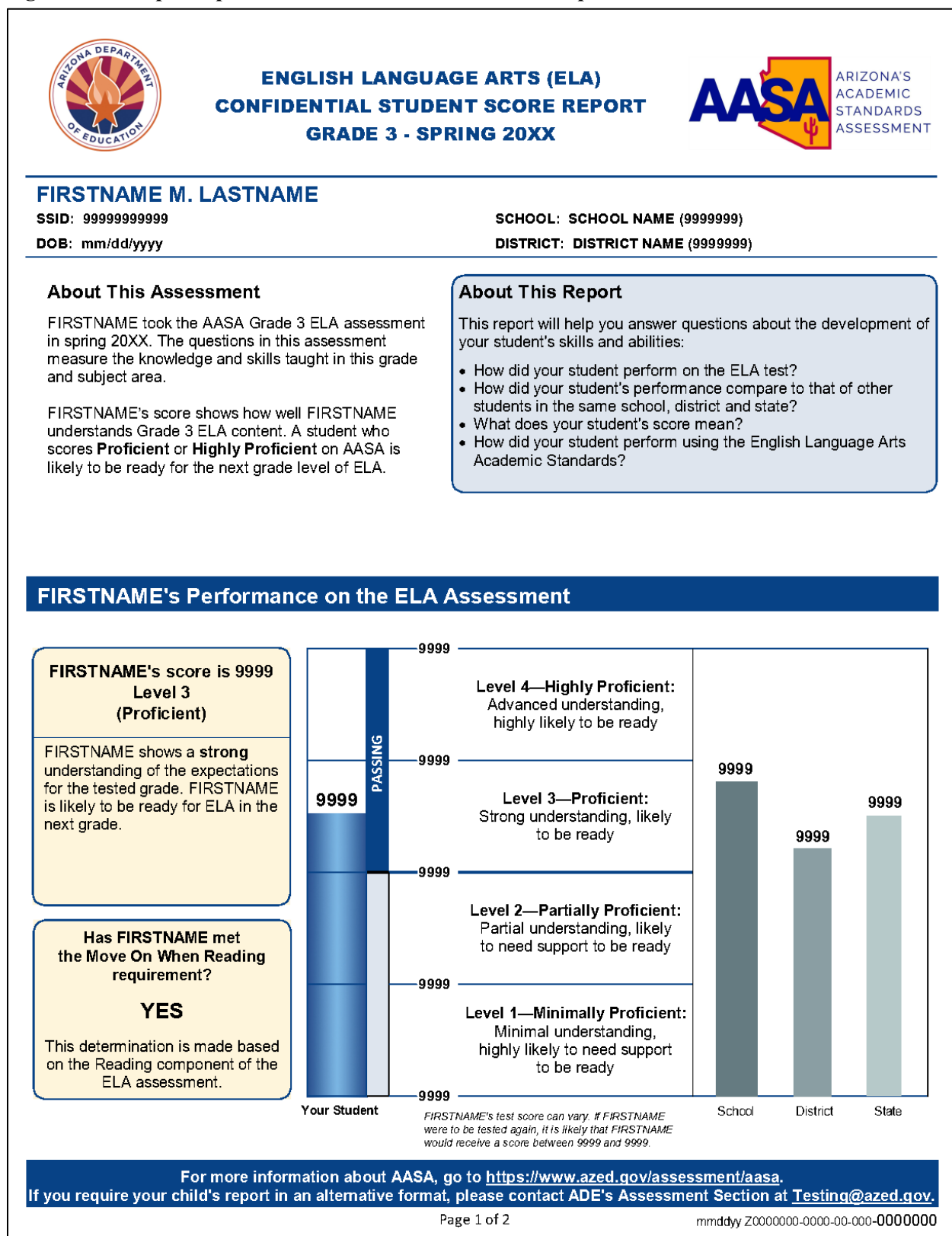
- District-level
 - District Confidential Roster Report with Summary (district-level, student roster by grade and content area)
 - Student Data File
- School-level
 - Confidential Student Score Report (individual student report by grade and content area)
 - Informe del Estudiante (individual student report in Spanish)
 - Confidential Roster Report with Summary (school-level, student roster by grade and content area)

AASA reports have been designed with the user’s comprehension in mind. The goal of these reports is to deliver accurate assessment data and ensure that it is correctly interpreted and understood. Similar colors are used for groups of similar elements, such as performance levels, throughout the design to guide the user to compare like elements and avoid comparison of dissimilar elements. All score report data are based on the total number of students whose tests have been scored. All score report data in PAN, except for individual students’ score reports, can be disaggregated into testing groups if they were set up by the school during the specified timeframe. The Confidential Student Score Report (individual student report) includes the average scale scores for the school, district, and state to allow for visual comparison. Two copies of the printed Confidential Student Score Report and Family Report Guide were also provided. Printed reports are packed by the school and shipped to participating districts.

Additionally, beginning in Spring 2023, an ACT predicted score is included for Grade 8 students on the Confidential Student Score Reports based on a study by ACT that linked AASA test scores to the ACT scale to obtain the ACT predicted score range (i.e., the score a student would likely receive if they were to be tested again). Students who score at or above the ACT score are more likely to be successful in college courses taken by first-year students. Appendix D presents the study results.

The AASA score reports are also available in the Parent Portal, which is an optional resource for schools and districts to use that allows families to securely access and view their student’s online individual student report. After creating a user account, families enter the student’s information, including the student’s claim code, to retrieve the AASA Student Report. The claim codes file (in CSV format) is available for request in PAN for authorized district and school users. The *Parent Portal Access Guide* is also available to families and includes the steps that should be followed to access their student’s information on the Parent Portal.

Figure 5.3. Sample Reports—Confidential Student Score Report, Grade 3



Legend: Reporting Categories



= Below Mastery



= At/Near Mastery



= Above Mastery

ELA Reporting Categories

Reading for Information



FIRSTNAME performed **above mastery** in Reading for Information.

What was assessed?

Students find the main idea and the supporting details of a text. They connect events, ideas, steps, sentences, paragraphs, and illustrations to one another. They find similarities and differences between two texts on the same topic.

What do these results mean?

Your student almost always finds connections between concepts, ideas, or events; uses the text and pictures to make conclusions to ask and answer questions; and finds the similarities and differences between important ideas and key details in two texts on the same topic.

Reading for Literature



FIRSTNAME performed **at or near mastery** in Reading for Literature.

What was assessed?

Students ask and answer questions about a text. They tell how characters and their actions affect a story. They explain how pictures help tell a story. They read two texts by one author and tell the similarities and differences. They find the central message of a story.

What do these results mean?

Your student can often find similarities and differences between the settings or plots of stories written by the same author; tell how one part of a story affects another part; use key details to retell a story and find the main idea; and tell the point of view in a story.

Writing and Language



FIRSTNAME performed **below mastery** in Writing and Language.

What was assessed?

Students write to give information or state opinions. They write on a topic giving supporting details or facts. They use correct capitalization, punctuation, and spelling. They use sentences, a glossary, or a dictionary to figure out the meaning of new words.

What do these results mean?

Your student may have trouble organizing writing for a purpose (like to give information or give opinions); using clues in a text to understand the meaning of new words; spelling commonly used words correctly; and writing simple sentences with correct capitalization and punctuation.

The Writing and Language portion of the ELA assessment requires that each student complete an essay. The essay is evaluated on three criteria.

Writing Essay Performance

Statement of Purpose, Focus & Organization

Your student earned 3 out of 4 possible points. Your student's essay mainly stays on topic. The opinion is clearly stated and mostly focused. Context supporting the opinion fits the purpose. The response is organized and has few mistakes. There is some variety of transitions used. There is a clear progression of ideas within the essay. There is a clear beginning and end.

Evidence & Elaboration

Your student earned 2 out of 4 possible points. Your student's essay includes details, facts, and sources that somewhat support its opinion. This evidence is unevenly integrated into the response. The words used are sometimes inappropriate for audience and purpose.

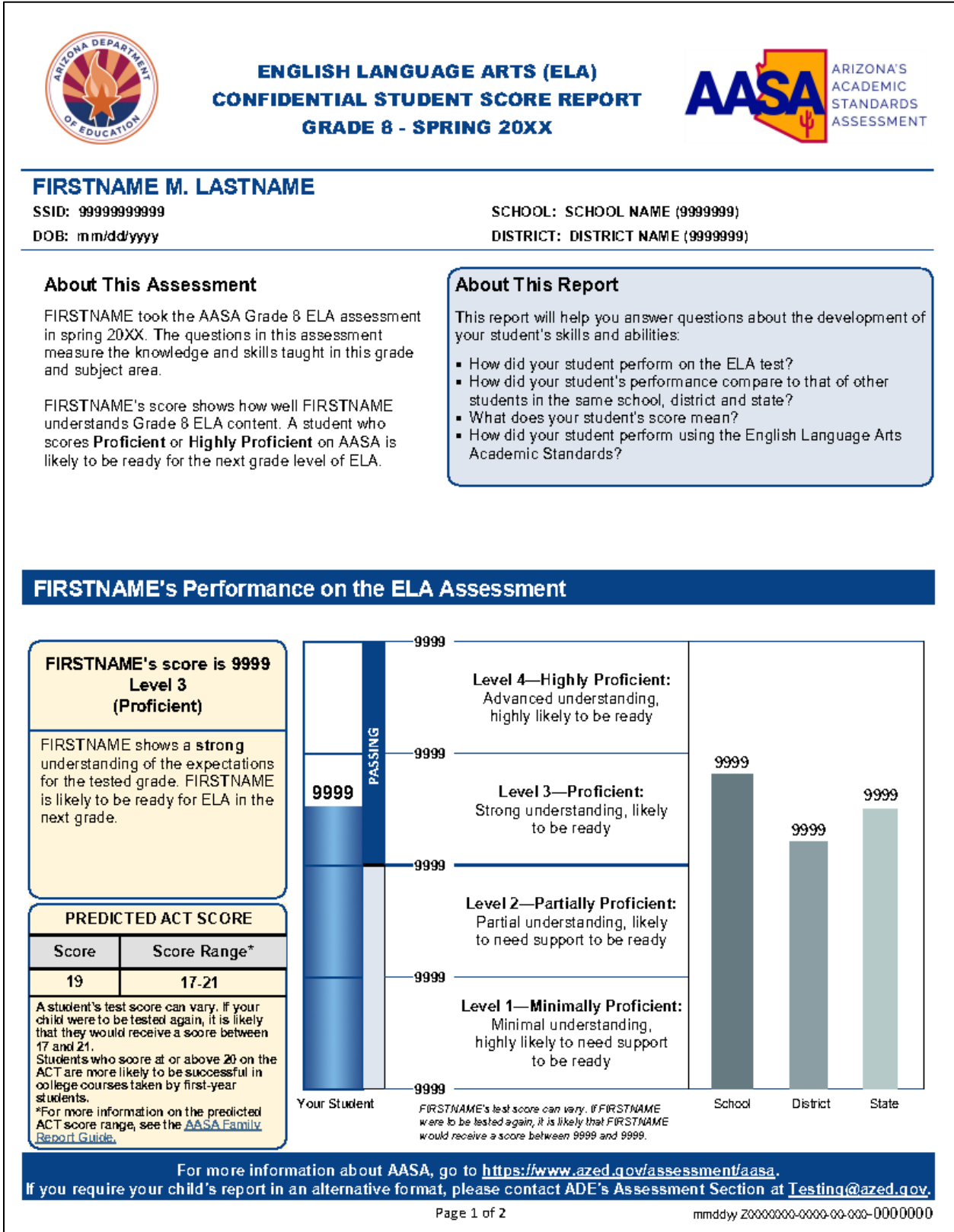
Conventions & Editing

Your student earned 2 out of 2 possible points. Your student's essay shows a strong understanding of sentence structure and language conventions. There are few mistakes in punctuation, capitalization, and spelling present in the response.

For more information about AASA, go to <https://www.azed.gov/assessment/aasa>.

If you require your child's report in an alternative format, please contact ADE's Assessment Section at Testing@azed.gov.

Figure 5.4. Sample Reports—Confidential Student Score Report, Grade 8



Legend: Reporting Categories



= Below Mastery



= At/Near Mastery



= Above Mastery

ELA Reporting Categories

Reading for Information



FIRSTNAME performed above mastery in Reading for Information.

What was assessed?

Students explain how reasoning and evidence shape and support the main idea of a text. They examine how a text makes connections between different individuals, ideas, or events. They show how an author of a text responds to evidence that does not support his or her point of view.

What do these results mean?

Your student gives an objective summary of a text; uses evidence from a text to make and support conclusions; explains how an author addresses a conflicting viewpoint; determines if information in a text is needed; evaluates how presentation (like text or audio) affects information.

Reading for Literature



FIRSTNAME performed at or near mastery in Reading for Literature.

What was assessed?

Students find the main idea of a text and examine how it is developed. They determine how specific words and phrases can change the meaning and tone of a text. They analyze how a character's point of view affects a text. They recognize the influence of other literature on a text.

What do these results mean?

Your student often uses supporting details to explain the theme or main idea; shows how a story moves forward; describes the effect of point of view on a text; recognizes the influences of other literature on a text; compares the structure of two or more texts.

Writing and Language



FIRSTNAME performed below mastery in Writing and Language.

What was assessed?

Students write to inform or make an argument. They use evidence and clear reasoning to support their writing. Their evidence comes from many different sources. They determine the meaning of new words and figurative language. They spell correctly and use correct grammar.

What do these results mean?

Your student may have trouble stating a claim clearly and providing supporting details to make an argument when writing; using citations correctly when doing research; using verb tenses and punctuation correctly; using other words or word parts to figure out the meaning of new words.

The Writing and Language portion of the ELA assessment requires that each student complete an essay. The essay is evaluated on three criteria.

Writing Essay Performance

Statement of Purpose, Focus & Organization

Your student earned 3 out of 4 possible points. In general, your student's essay stays on topic and is focused. The main idea of the topic is given context and addresses the audience and purpose for writing. The response is organized and develops connections between ideas. It uses transitions and has an introduction and conclusion.

Evidence & Elaboration

Your student earned 2 out of 4 possible points. Your student's essay includes some support or evidence for the main idea. It uses some facts and details from other sources but does not use citations regularly. The response does not expand on ideas or make clear connections between ideas. It uses simple and sometimes inappropriate language for the audience and purpose.

Conventions & Editing

Your student earned 2 out of 2 possible points. Your student's essay shows an understanding of sentence formation and other conventions. The response may have some mistakes, but they are not repeated often in the text. It uses correct punctuation, capitalization, and spelling.

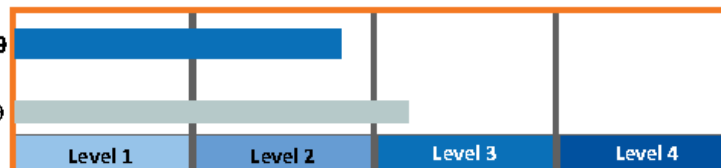
FIRSTNAME's ELA Assessment Progress

This chart displays your student's performance in ELA assessments over time. It reports the proficiency level for the most recently completed tests in ELA (if available). You can use this information to determine your student's progress in ELA.

<CurrYear> 9999

<CurrYear-1> N/A

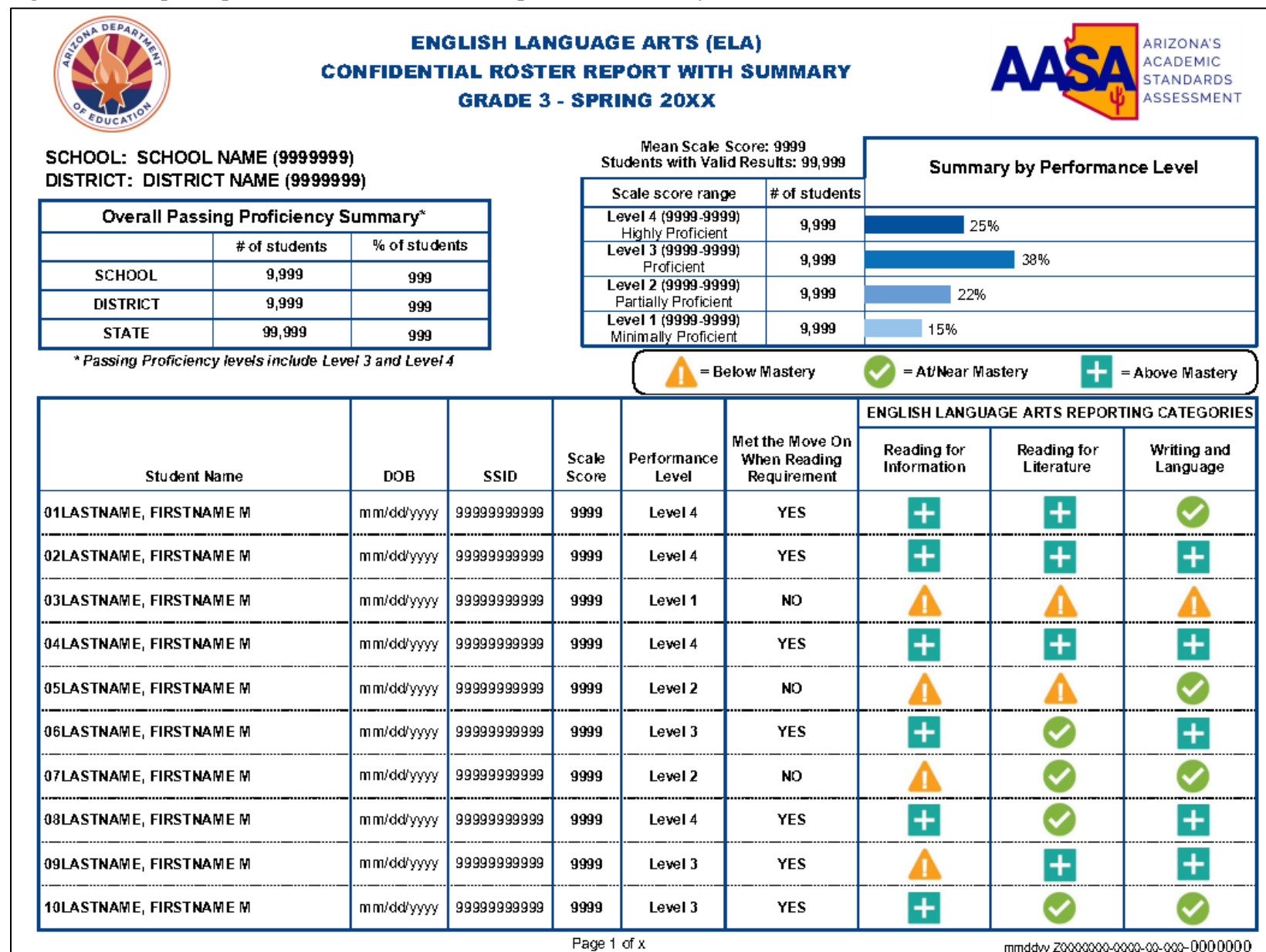
<CurrYear-2> 9999



For more information about AASA, go to <https://www.azed.gov/assessment/aasa>.

If you require your child's report in an alternative format, please contact ADE's Assessment Section at Testing@azed.gov.

Figure 5.5. Sample Reports—Confidential Roster Report with Summary



Chapter 6: CLASSICAL ITEM ANALYSIS

This chapter presents classical statistics for the data used for calibration, equating, and scaling of the Spring 2023 AASA assessments as indicated by Standards 1.8, 1.10, 2.5, 2.19, 3.6, 4.14, and 7.4 (AERA et al., 2014).

Each grade in ELA had two core online forms with different embedded field test sets. The core online forms differed by only a writing prompt. For ELA Grades 3–7, the first 11 forms had one writing prompt (referred to as Form 1), and the next 10 forms had another writing prompt (referred to as Form 2). For ELA Grade 8, the first nine forms had one writing prompt (Form 1), and the next nine forms had another writing prompt (Form 2). Where appropriate, statistics are reported for both ELA core online forms. Mathematics only had one core online form with different embedded field test sets for each grade, with 11 online forms total.

6.1. Data

The classical item analysis was conducted based on the calibration samples described in Section 7.1. Table 6.1 and Table 6.2 present demographic information of the students included in the calibration sample by gender, ethnicity (Hispanic or Not-Hispanic), race, and special education, English learner (EL), and low socioeconomic status (SES). Because only a few students took the accommodated forms, these students were not included in the item analysis. Students who did not complete the test were also excluded.

Table 6.1. Number of Students in the Calibration Sample by Subgroup—ELA

Subgroup	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
All	77,029	76,809	77,090	76,856	77,812	81,290
Male	39,022	38,868	39,189	38,912	39,656	41,779
Female	38,007	37,941	37,901	37,944	38,156	39,511
Hispanic	37,718	37,232	37,024	37,165	37,926	39,648
Non-Hispanic	39,311	39,577	40,066	39,691	39,886	41,642
American Indian	4,065	4,105	4,121	4,267	4,399	4,618
Asian	2,190	2,145	2,166	2,015	1,997	2,061
Black or African American	5,559	5,645	5,487	5,510	5,552	5,636
Multi-racial	4,815	4,701	4,676	4,430	4,387	4,382
Native Hawaiian or Other Pacific Islander	456	438	418	473	424	429
White	59,129	58,917	59,397	59,181	59,895	62,892
Missing	815	858	825	980	1,158	1,272
Special Education	10,460	10,845	10,628	10,162	9,639	9,586
English Learner (EL)	9,107	8,839	8,266	7,270	7,479	7,483
Low Socioeconomic Status (SES)	32,034	31,927	31,445	31,364	30,944	31,461

Table 6.2. Number of Students in the Calibration Sample by Subgroup—Mathematics

Subgroup	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8
All	78,131	77,590	77,610	77,507	78,511	82,075
Male	39,742	39,354	39,486	39,271	40,021	42,183
Female	38,389	38,236	38,124	38,236	38,490	39,892
Hispanic	38,271	37,611	37,265	37,505	38,296	40,070
Non-Hispanic	39,860	39,979	40,345	40,002	40,215	42,005
American Indian	4,177	4,200	4,191	4,338	4,490	4,697
Asian	2,212	2,166	2,177	2,029	2,005	2,079
Black or African American	5,720	5,732	5,545	5,556	5,617	5,690
Multi-racial	4,885	4,744	4,706	4,471	4,429	4,413
Native Hawaiian or Other Pacific Islander	464	439	421	480	426	434
White	59,810	59,414	59,718	59,602	60,352	63,428
Missing	863	895	852	1,031	1,192	1,334
Special Education	10,784	11,077	10,744	10,303	9,798	9,714
English Learner (EL)	9,304	8,948	8,339	7,354	7,577	7,580
Low Socioeconomic Status (SES)	32,523	32,256	31,643	31,620	31,203	31,686

6.2. Descriptive Statistics

Table 6.3 presents the descriptive statistics on total raw scores for the spring AASA assessment, including the number of students included in the classical analysis, the number of operational items on the assessment, the maximum possible raw score, the mean raw score, the standard deviation (SD) of the raw score, and the minimum/maximum obtained raw score.

Table 6.3. Classical Test Analysis Statistics

Content Area	Grade	#Students	#Items	Max. Possible Raw Score	Mean Raw Score	SD Raw Score	Min. Raw Score	Max. Raw Score
ELA, Form 1	3	38,882	44	56	27.37	11.22	2	56
	4	40,148	44	57	30.15	11.72	2	57
	5	40,456	44	55	28.36	11.42	3	55
	6	40,364	44	55	28.76	11.20	3	55
	7	40,842	44	56	30.86	10.88	5	56
	8	41,039	44	55	28.34	10.89	4	54
ELA, Form 2	3	38,147	44	56	27.92	11.25	2	56
	4	36,661	44	57	29.62	11.32	4	57
	5	36,634	44	55	27.95	11.41	2	55
	6	36,492	44	55	29.03	11.07	4	54
	7	36,970	44	56	30.99	10.94	4	56
	8	40,251	44	55	28.27	10.96	3	54
Mathematics	3	78,131	45	45	25.02	11.96	0	45
	4	77,590	45	45	23.29	11.99	0	45
	5	77,610	45	45	20.34	11.55	0	45
	6	77,507	47	47	19.51	11.72	0	47
	7	78,511	47	47	19.24	11.70	0	47
	8	82,075	47	47	18.96	10.32	0	47

6.3. Classical Item Analysis

Classical item analysis was conducted to show how the items performed for each grade-level assessment. Item difficulty is measured by the *p*-value bounded by 0.0 and 1.0 that indicates how easy or hard an item is for students. The *p*-value for 1-point items is based on the proportion of students who answered an item correctly and is derived by dividing the number of students who got the item correct by the total number of students who answered it. For multiple-point items, the *p*-value is the average item score (i.e., the sum of student scores on an item divided by the total number of students who responded to the item) divided by the number of possible score points on the item. A high *p*-value indicates that an item is easy (high proportion of students answered it correctly), whereas a low *p*-value indicates that an item is difficult. For example, a *p*-value of 0.79 indicates that 79% of students answered the item correctly. Easy and hard items are both necessary to include on an assessment to balance the test difficulty. The AASA assessment targets *p*-values in the range of 0.20 to 0.90.

Item discrimination is represented by the point-biserial correlation bounded by -1.0 and 1.0 that indicates how well an item discriminates, or distinguishes, between low-performing and high-performing students. The point-biserial correlation is based on the relationship between student performance on a specific item and performance on the entire test based on their test score. Students who do well on a test are expected to select the right answer to any given item, and students who do poorly are expected to select the wrong answer. This means that for a highly discriminating item, students who get the item correct will have a higher average test score than students who get the item incorrect. An item with a high positive point-biserial correlation discriminates between low-performing and high-performing students better than an item with a point-biserial correlation near zero. A negative point-biserial correlation indicates that lower-performing students did better on that item than higher-performing students. The AASA assessment targets point-biserial correlations of 0.25 or higher.

Table 6.4 presents a summary of the classical item analysis, and Appendix A presents the statistics for each item. If the classical item statistics for the operational items were outside of the item selection criteria presented in Table 3.3, the items will be reviewed during test construction of the next testing cycle for possible replacement in future administrations.

Table 6.4. Classical Item Analysis Summary

Content Area	Grade	#Items	Mean <i>P</i> -Value	Mean Point-Biserial
ELA, Form 1	3	44	0.49	0.48
	4	44	0.53	0.48
	5	44	0.50	0.48
	6	44	0.52	0.47
	7	44	0.55	0.45
	8	44	0.49	0.46
ELA, Form 2	3	44	0.49	0.48
	4	44	0.52	0.47
	5	44	0.50	0.48
	6	44	0.52	0.47
	7	44	0.55	0.45
	8	44	0.49	0.47

Content Area	Grade	#Items	Mean <i>P</i> -Value	Mean Point-Biserial
Mathematics	3	45	0.56	0.56
	4	45	0.52	0.56
	5	45	0.45	0.53
	6	47	0.42	0.52
	7	47	0.41	0.54
	8	47	0.40	0.47

6.4. Distractor Analysis

Table 6.5 and Table 6.6 present the point-biserial correlations associated with a correct option and the incorrect options at various percentiles. As expected, the point-biserial correlation for a correct option was around 0.20 or higher for most items, whereas the point-biserial correlation for incorrect options was negative or very close to zero. The results show that students with higher proficiency tended to choose a correct option, and students with lower proficiency tended to choose an incorrect option. This indicates that the distractors appear to perform appropriately.

Table 6.5. Distractor Analysis Summary: Point-Biserial Correlations for Correct Options

Content Area	Grade	#MC Items	Min.	P25	P50	P75	Max.
ELA, Form 1	3	26	0.26	0.39	0.47	0.54	0.59
	4	28	0.17	0.36	0.42	0.51	0.58
	5	21	0.29	0.36	0.41	0.46	0.59
	6	25	0.29	0.37	0.44	0.47	0.59
	7	30	0.11	0.36	0.40	0.46	0.57
	8	26	0.27	0.38	0.43	0.48	0.64
ELA, Form 2	3	26	0.26	0.39	0.47	0.54	0.59
	4	28	0.17	0.36	0.42	0.51	0.58
	5	21	0.29	0.36	0.41	0.46	0.59
	6	25	0.29	0.37	0.44	0.47	0.59
	7	30	0.11	0.36	0.40	0.46	0.57
	8	26	0.27	0.38	0.43	0.48	0.64
Mathematics	3	12	0.33	0.42	0.46	0.56	0.62
	4	12	0.35	0.41	0.44	0.54	0.61
	5	10	0.33	0.36	0.38	0.48	0.53
	6	13	0.26	0.30	0.45	0.47	0.57
	7	17	0.25	0.36	0.40	0.49	0.59
	8	27	0.29	0.36	0.40	0.49	0.63

Note. Min.= minimum, P25 = 25th percentile, P50 = 50th percentile (median), P75 = 75th percentile, Max. = maximum. This analysis is conducted for MC items only.

Table 6.6. Distractor Analysis Summary: Point-Biserial Correlations for Incorrect Options

Content Area	Grade	#MC Items	Min.	P25	P50	P75	Max.
ELA, Form 1	3	26	-0.35	-0.27	-0.23	-0.17	-0.03
	4	28	-0.36	-0.26	-0.21	-0.16	0.00
	5	21	-0.38	-0.29	-0.23	-0.14	0.07
	6	25	-0.39	-0.28	-0.22	-0.15	0.05
	7	30	-0.34	-0.26	-0.20	-0.14	0.06
	8	26	-0.36	-0.25	-0.20	-0.15	-0.05
ELA, Form 2	3	26	-0.35	-0.27	-0.23	-0.17	-0.03
	4	28	-0.36	-0.26	-0.21	-0.16	0.00
	5	21	-0.38	-0.29	-0.23	-0.14	0.07
	6	25	-0.39	-0.28	-0.22	-0.15	0.05
	7	30	-0.34	-0.26	-0.20	-0.14	0.06
	8	26	-0.36	-0.25	-0.20	-0.15	-0.05
Mathematics	3	12	-0.42	-0.30	-0.24	-0.16	-0.03
	4	12	-0.42	-0.28	-0.23	-0.15	-0.08
	5	10	-0.37	-0.26	-0.21	-0.16	0.08
	6	13	-0.32	-0.28	-0.16	-0.10	-0.01
	7	17	-0.35	-0.24	-0.19	-0.15	0.00
	8	27	-0.31	-0.22	-0.19	-0.13	0.00

Note. Min.= minimum, P25 = 25th percentile, P50 = 50th percentile (median), P75 = 75th percentile, Max. = maximum. This analysis is conducted for MC items only.

A distractor analysis was also conducted for each multiple-choice item, as presented in Appendix A. The response distribution for an item across all possible choices (e.g., a correct option and distractors) was calculated. The point-biserial correlation associated with each response option was calculated as well. Typically, a negative point-biserial correlation is sought for distractors because less-proficient students should be more likely to choose an incorrect option.

Chapter 7: CALIBRATION, EQUATING, AND SCALING

This chapter describes the calibration, equating, and scaling procedures that took place for the Spring 2023 AASA assessments, addressing Standards 1.10, 5.1, 5.2, 5.3, 7.2, 7.4, and 12.9 (AERA et al., 2014).

7.1. Calibration Sample

To ensure valid calibration results, several data cleaning steps occurred upon receipt of raw data from the scanning and scoring processes. These steps allowed for calibration to be conducted on valid student responses. The cleaning process removed the following records from the calibration datasets for each grade level:

- Records with invalidated tests that are marked Do Not Report (DNR) in PearsonAccess^{next} (PAN)
- Records that indicate the student took an accommodated form
- Records with non-valid attempts noted by less than one response
- Duplicate records (e.g., students indicated as taking the test more than once)
- Records in which a student was enrolled in an exclusionary school list from ADE

7.2. Calibration Methods

Item response theory (IRT) models were used in the item calibration. All tests were calibrated separately by grade. If there was more than one operational form, all operational forms were calibrated concurrently. All calibration activities were replicated with two psychometricians independently as a quality control measure. The calibration results were also reviewed independently by a senior-level psychometrician at Pearson.

The Rasch model (Rasch, 1960) was used for 1-point items and the partial-credit model (Masters, 1982) was used for multiple-point items for calibration. Parameter estimation for items was implemented using Winsteps 4.8.1.0 (Linacre, 2022b) that uses joint maximum likelihood estimation (JMLE) as described by Wright and Masters (1982).

The Rasch model estimates item difficulty and student ability on the same scale. Under the Rasch model, the probability that student j with ability θ answers item i with difficulty of b correctly is as follows:

$$P_i(\theta_j) = \frac{\exp(\theta_j - b_i)}{1 + \exp(\theta_j - b_i)}$$

The partial-credit model is an extension of the Rasch model for items in which students may receive partial credit. Thus, the partial-credit model reduces to the Rasch model when items have only two response categories (i.e., 0 or 1). According to the partial-credit model, the probability that student j scores x on item i , which has a maximum possible score of m ($k = m+1$ possible response categories), is expressed as follows:

$$P_{ix}(\theta_j) = \frac{\exp \sum_{l=0}^x (\theta_j - D_{il})}{\sum_{k=0}^{m_i} [\exp \sum_{l=0}^k (\theta_j - D_{il})]}$$

where $x = 0, 1, \dots, m_i$, D_{il} is a step difficulty for score l and by definition,

$$\sum_{l=0}^0 (\theta_j - D_{il}) = 0$$

The step difficulty D_{il} can be decomposed such that

$$D_{il} = b_i + h_{il}$$

where b_i is an overall difficulty for item i , and h_{il} is a threshold for score l (Embretson & Reise, 2000; Linacre, 2022a). This parameterization allows b_i in the partial-credit model to be comparable to b_i in the Rasch model.

7.3. Calibration Results

All items converged during calibration using typical procedures for Winsteps software. Standard error of estimates for the Rasch difficulty measures indicated that the parameters were well-estimated. Table 7.1 presents a summary of the IRT statistics, and Appendix B presents the item-level IRT statistics resulting from the calibration of the spring AASA assessments.

Table 7.1. IRT Statistics Summary

Content Area	Grade	#Items	Mean Rasch
ELA, Form 1	3	44	0.05
	4	44	0.29
	5	44	0.13
	6	44	0.14
	7	44	0.04
	8	44	0.17
ELA, Form 2	3	44	0.04
	4	44	0.33
	5	44	0.16
	6	44	0.12
	7	44	0.03
	8	44	0.18
Mathematics	3	45	0.22
	4	45	-0.01
	5	45	0.11
	6	47	0.07
	7	47	0.21
	8	47	-0.12

An item-person map shows the distribution of item difficulty and the distribution of student ability in one graph, as they are on the same scale. This graph is useful for Rasch models to evaluate the extent to which the item difficulty and student ability distributions are aligned because they assume the probability of a correct answer is affected only by a student's ability and the item difficulty. Figure B.1 – Figure B.18 in Appendix B present the item difficulty distribution on the lefthand side and the student ability distribution on the right. Each marker in the item difficulty distribution is an item, and the item difficulty values are rounded with an increment of 0.20 before they are plotted. Horizontal dotted lines represent the three performance level cuts (i.e., *Partially Proficient*, *Proficient*, and *Highly Proficient*) for the total test.

In addition to the item-person map, two more graphs are presented to summarize the characteristics of each operational assessment. The test characteristic curve (TCC) shows an expected total raw score across different student abilities, whereas the CSEM curve presents an amount of standard error across different student abilities. The CSEM has an inverse relationship with the test information function (TIF) as follows:

$$SE(\theta) = \frac{1}{TI(\theta)}$$

where $SE(\theta)$ is the CSEM, and $TI(\theta)$ is the TIF (Embretson & Reise, 2000). Because the CSEM can be interpreted on the ability scale, the CSEM curve is presented over the TIF curve in this technical report.

7.4. Equating

The Spring 2023 AASA tests were equated and placed on the operational AASA scale using a non-equivalent groups anchor item (NEAT) design. A set of anchor items was selected from the existing item bank. The anchor items were selected such that they contributed approximately 30% of the total score points and their content representation was as similar as possible to the blueprint. The location of all anchor items stayed within three positions from where they were in the previous year.

A fixed anchor parameter equating was implemented within Winsteps to place the tests on the operational reporting scale. This was implemented by constraining the parameter estimates in the existing item bank for the anchor items to equal the final parameter estimates obtained in the original AASA calibration analyses. The displacement statistic, which estimates the difference between the fixed parameter and the estimate had the item parameter not been constrained, was evaluated for each anchor item.

Items with a displacement statistic greater than 0.30 or less than -0.30 were reiteratively removed from the anchor set. The criterion of 0.30 has been used to flag displaced anchor items under a common item, non-equivalent group equating design for many state programs (Miller et al., 2004). If more than one anchor item was flagged, the item with the largest magnitude of displacement value was dropped from the anchor set. The displacement values of the remaining anchor items were then re-estimated by implementing the fixed anchor parameter equating with the remaining anchor items. This process was repeated until all the anchor items had displacement values of a magnitude smaller than 0.30 and greater than -0.30.

Table 7.2 presents the number of items for the initial anchor set of each grade and the number of items dropped from each initial anchor set.

Table 7.2. Summary of Anchor Items

Content Area	Grade	#Items in Initial Anchor Set	#Items Dropped from Anchor
ELA	3	17	1
	4	22	0
	5	22	1
	6	27	2
	7	17	0
	8	22	1
Mathematics	3	17	1
	4	17	1
	5	15	1
	6	15	1
	7	17	0
	8	19	0

7.5. Scaling Methods

The AASA reporting scale was established in 2015 when the first administration took place (known as the AzMERIT statewide achievement assessment at that time). These tests were placed on a vertical scale for the total score as a result of a previous study (American Institutes for Research, 2015, Appendix J). Scaling constants for the total score were determined such that the vertically scaled theta score, based on the total test, was transformed by solving the following equation:

$$Scale\ Score = VS_A \times \theta + VS_B$$

where VS_A and VS_B are scaling constants on the vertical scale that are used to transform θ , which are the performance level cuts on the theta (ability) scale, into scale scores. For reporting, θ is truncated at -3.5 and 3.5 for the lower and upper ends, respectively.

The AASA reporting scale ranged from 2395 to 2658 across grades for ELA and from 3395 to 3776 across grades for mathematics. In addition to a total score, a subscore was also calculated for each reporting category by grade using the same formula. The scaling constants were applied to a theta score based on items associated with a reporting category to transform it to a scale score. Table B.13 – Table B.24 in Appendix B presents the raw-to-scale score conversion tables for each content area and grade.

7.6. IRT Assumptions

It is important to evaluate how the Rasch models fit the data because reported scale scores are derived from theta estimated under the IRT models. Three major assumptions are investigated: (1) unidimensionality, (2) local item independence, and (3) item fit.

7.6.1. Unidimensionality

An assumption under the Rasch models is unidimensionality, that there is exactly one latent variable (e.g., mathematics proficiency) that an instrument intends to measure. This is a more traditional and strict definition of the unidimensionality assumption. On the other hand, essential unidimensionality, in which there is one dominant latent variable with some minor latent variable(s), is a more practically applicable assumption (Stout, 1990).

Principal component analysis (PCA) is a statistical technique widely applied to investigate the dimensionality of data (Jackson, 1993; Velicer & Jackson, 1990). Many decision rules have been proposed to determine the number of dimensions using PCA results. Horn's (1965) parallel analysis is a Monte Carlo simulation technique used to determine the number of factors to retain from a PCA. Parallel analysis compares the observed eigenvalues from a correlation matrix to be analyzed with those obtained from uncorrelated normal variables (Ledesma & Valero-Mora, 2007). In other words, expected eigenvalues are obtained by simulating normal, random samples that "parallel" the observed data in terms of sample size and number of variables. Numerous studies have shown parallel analysis to be an effective and appropriate method to determine the number of factors underlying a construct (Glorfeld, 1995; Humphreys & Montanelli, 1975; Zwick & Velicer, 1986), including the least variability and sensitivity to different factors.

PCA was conducted for the operational form in each content area and grade. Table 7.3 presents the first 10 eigenvalues from PCA for each operational form. Because the same blueprint was used to construct the operational forms, only one set of eigenvalues from the parallel analysis is presented. The graphical presentations of eigenvalues (i.e., scree plot) are presented in Figure B.55 – Figure B.72 in Appendix B. The PCA results with the parallel analysis criterion show only one significant dimension for each grade, which supports unidimensionality.

Table 7.3. Eigenvalues from PCA

Content Area	Grade	1	2	3	4	5	6	7	8	9	10
ELA, Form 1	3	15.94	1.40	1.21	1.06	0.96	0.90	0.88	0.85	0.84	0.82
	4	15.64	1.71	1.17	1.09	1.04	0.95	0.89	0.87	0.85	0.85
	5	15.52	1.39	1.16	1.04	0.99	0.92	0.91	0.88	0.86	0.83
	6	15.69	1.29	1.28	1.06	0.97	0.92	0.91	0.87	0.86	0.83
	7	13.86	1.59	1.34	1.01	0.99	0.95	0.92	0.90	0.88	0.86
	8	15.01	1.67	1.22	1.05	0.97	0.95	0.95	0.90	0.87	0.84
ELA, Form 2	3	15.96	1.46	1.15	1.04	0.93	0.91	0.88	0.84	0.82	0.79
	4	15.38	1.70	1.18	1.11	1.05	0.93	0.92	0.88	0.84	0.84
	5	15.78	1.27	1.15	1.04	0.97	0.89	0.88	0.87	0.84	0.82
	6	15.46	1.32	1.26	1.05	0.98	0.94	0.91	0.88	0.87	0.82
	7	13.96	1.56	1.33	1.00	0.99	0.95	0.92	0.91	0.88	0.87
	8	15.28	1.66	1.25	1.04	0.98	0.96	0.95	0.91	0.87	0.84
Mathematics	3	22.57	1.21	1.16	1.07	0.87	0.81	0.79	0.75	0.73	0.70
	4	22.31	1.67	1.21	1.11	0.96	0.80	0.78	0.74	0.72	0.71
	5	20.31	1.73	1.26	1.06	0.96	0.91	0.85	0.83	0.82	0.79
	6	20.38	1.49	1.39	1.11	0.99	0.94	0.91	0.90	0.82	0.78
	7	22.58	1.70	1.18	0.96	0.90	0.88	0.85	0.82	0.80	0.77
	8	17.13	1.68	1.36	1.12	1.05	1.00	0.95	0.92	0.89	0.86

7.6.2. Local Item Independence

Local item independence is another assumption under the Rasch models that assumes any item pair is uncorrelated, conditioned on the latent trait an instrument is intended to measure (e.g., mathematics proficiency). A violation of local item dependence would impact parameter estimation under the Rasch models because JMLE performed by Winsteps (Linacre, 2022b) relies on uncorrelated item pairs. Winsteps produces raw score residual correlations for pairs of items on a test, which are analogous to Yen's Q3 statistics (Yen, 1984). For an item pair with the residual correlation greater than 0.70, only one item is needed on the test (Linacre, 2022a).

Table 7.4 summarizes the distribution of the residual correlations. Most residual correlations are slightly negative or slightly positive and only five (out of more than 900 per grade) are greater than 0.70, which indicates that the local item independence assumption holds for the AASA tests.

Table 7.4. Q3 Statistics

Content Area	Grade	#Item Pairs	Mean	SD	Min.	P10	P25	P50	P75	P90	Max.	#Items Exceeding 0.70
ELA	3	1,081	-0.02	0.05	-0.09	-0.05	-0.04	-0.02	-0.01	0.01	0.68	0
	4	1,081	-0.02	0.05	-0.10	-0.06	-0.04	-0.02	-0.01	0.01	0.81	1
	5	1,081	-0.02	0.05	-0.09	-0.05	-0.04	-0.03	-0.01	0.01	0.79	1
	6	1,081	-0.02	0.05	-0.09	-0.06	-0.04	-0.02	-0.01	0.01	0.76	2
	7	1,081	-0.02	0.05	-0.12	-0.06	-0.04	-0.02	-0.01	0.01	0.77	1
	8	1,081	-0.02	0.04	-0.11	-0.05	-0.04	-0.02	-0.01	0.01	0.63	0
Math	3	990	-0.02	0.04	-0.10	-0.06	-0.04	-0.03	-0.01	0.02	0.31	0
	4	990	-0.02	0.04	-0.11	-0.06	-0.04	-0.03	-0.01	0.02	0.43	0
	5	990	-0.02	0.05	-0.11	-0.07	-0.05	-0.03	0.00	0.02	0.33	0
	6	1,081	-0.02	0.04	-0.14	-0.07	-0.05	-0.02	0.00	0.02	0.40	0
	7	1,081	-0.02	0.05	-0.13	-0.07	-0.05	-0.03	0.00	0.03	0.50	0
	8	1,081	-0.02	0.04	-0.13	-0.06	-0.04	-0.03	-0.01	0.02	0.29	0

Note. SD = standard deviation, min. = minimum, P10 = 10th percentile, P25 = 25th percentile, P50 = 50th percentile, P75 = 75th percentile, P90 = 90th percentile, max. = maximum

7.6.3. Item Fit

Item fit was monitored using weighted mean-square (MNSQ) that indicates the degree of accuracy and predictability with which the data fit the model (Linacre, 2022b). In Winsteps and Rasch literature, weighted mean-square is also referred to as infit MNSQ. The infit MNSQ is sensitive to unexpected responses at or near the item's calibrated level. Items were flagged for misfit using a set of conservative criteria. For infit MNSQ, values less than 0.60 or greater than 1.40 were flagged, in accordance with Wright and Linacre's (1994) recommendation.

Table 7.5 presents a summary of the item fit statistics, and Table B.1 – Table B.12 in Appendix B present the statistics for each item. Items flagged by Winsteps' infit statistics are reviewed during test construction for possible replacement in future administrations.

Table 7.5. IRT Item Fit Summary Statistics

Content Area	Grade	#Items	#Flagged Items by Infit	%Flagged
ELA	3	47	0	0
	4	47	1	2
	5	47	0	0
	6	47	0	0
	7	47	0	0
	8	47	0	0
Mathematics	3	45	0	0
	4	45	0	0
	5	45	0	0
	6	47	0	0
	7	47	0	0
	8	47	0	0

Chapter 8: TEST RESULTS

This chapter presents the test results of the Spring 2023 AASA administration, addressing Standards 1.8, 2.11, 2.15, 3.1, 3.3, 3.6, 3.15, 5.3, 7.4, 12.17, and 12.18 (AERA et al., 2014). The results, summarized below, are based on the population data contained within the final electronic data files (note that the data in this chapter are different from the calibration sample). The results in this section of the technical report may differ slightly from the final testing results presented on the ADE website due to small differences in the application of exclusion rules. Official results typically use more detailed school-level information than is used to conduct research analyses. Please note that the results in the following tables are presented as evidence of reliability and validity of the test scores and should not be used for state accountability purposes.

- Table 8.1 presents the test results for all students by content area and grade, including the mean and standard deviation (SD) of the scale scores and the percentage of students in the overall performance levels. Overall performance levels are determined based on the performance levels for the total score.
- Table 8.2 and Table 8.3 present the percentage of students in each level of mastery by reporting category.
- Appendix C presents the test results by demographics. Histograms of the scale score distribution for the total score are also presented.
- Table 8.4 and Table 8.5 present the mean and standard deviation of the scale scores and the performance level distributions by accommodation for students who used the available accommodations. These tables only include the accommodations captured in the student data file (i.e., accommodations used by students during the Spring 2023 administration).
- Table 8.6 and Table 8.7 present the frequency distribution statistics for total scale score by performance level. Results indicate that average scale scores increase when moving from lower to higher performance levels across all grades and content areas.

Table 8.1. Overall Test Results

Content Area	Grade	N	SS Mean	SS SD	%Level 1	%Level 2	%Level 3	%Level 4
ELA	3	80,814	2498.68	35.51	47.4	11.7	27.1	13.8
	4	80,659	2519.03	33.92	41.1	14.2	30.5	14.3
	5	80,917	2528.74	35.21	40.2	22.6	28.3	8.9
	6	81,369	2543.19	34.76	35.7	22.4	35.5	6.5
	7	82,061	2552.90	31.66	38.9	19.6	32.8	8.7
	8	85,232	2557.60	34.17	41.4	22.2	26.5	9.8
Mathematics	3	81,986	3517.02	48.72	31.2	25.5	28.3	14.9
	4	81,480	3543.68	49.02	39.3	21.9	26.2	12.6
	5	81,451	3578.26	44.04	36.5	27.3	25.0	11.2
	6	82,066	3606.26	41.62	48.7	20.9	20.8	9.6
	7	82,799	3625.99	46.01	54.1	15.8	15.3	14.8
	8	86,031	3653.61	37.08	53.9	19.1	16.6	10.4

Note. SS = scale score, SD = standard deviation, Level 1 = *Minimally Proficient*, Level 2 = *Partially Proficient*, Level 3 = *Proficient*, Level 4 = *Highly Proficient*

Table 8.2. Performance Distributions by Reporting Category: Percentage of Students at each Level of Mastery—ELA

Grade	Reporting Category	N	%Level 1	%Level 2	%Level 3
3	Reading for Information	80,814	44.6	26.4	29.0
	Reading for Literature	80,814	46.7	28.7	24.6
	Writing and Language	80,814	41.7	29.7	28.6
4	Reading for Information	80,659	40.4	30.8	28.9
	Reading for Literature	80,659	39.1	31.3	29.6
	Writing and Language	80,659	38.8	26.6	34.6
5	Reading for Information	80,917	49.4	28.0	22.6
	Reading for Literature	80,917	47.1	32.6	20.3
	Writing and Language	80,917	48.4	27.2	24.4
6	Reading for Information	81,369	44.6	27.8	27.6
	Reading for Literature	81,369	38.5	36.5	24.9
	Writing and Language	81,369	41.2	34.5	24.3
7	Reading for Information	82,061	43.1	29.6	27.2
	Reading for Literature	82,061	43.6	33.4	23.0
	Writing and Language	82,061	42.1	29.9	28.0
8	Reading for Information	85,232	49.4	26.7	23.9
	Reading for Literature	85,232	53.8	22.7	23.5
	Writing and Language	85,232	45.7	28.8	25.5

Note. Level 1 = *Below Mastery*, Level 2 = *At or Around Mastery*, Level 3 = *Above Mastery*

Table 8.3. Performance Distributions by Reporting Category: Percentage of Students at each Level of Mastery—Mathematics

Grade	Reporting Category	N	%Level 1	%Level 2	%Level 3
3	Operations, Algebraic Thinking, and Numbers in Base Ten	81,986	46.3	19.2	34.5
	Numbers and Operations – Fractions	81,986	48.1	23.2	28.8
	Measurement, Data, and Geometry	81,986	46.5	32.0	21.5
4	Operations, Algebraic Thinking, and Numbers in Base Ten	81,480	52.1	19.2	28.7
	Numbers and Operations – Fractions	81,480	49.2	20.8	30.0
	Measurement, Data, and Geometry	81,480	38.3	43.4	18.4
5	Operations, Algebraic Thinking, and Numbers in Base Ten	81,451	49.7	21.3	29.1
	Numbers and Operations – Fractions	81,451	54.6	25.5	19.8
	Measurement, Data, and Geometry	81,451	50.3	28.8	20.9
6	Ratio and Proportional Relationships	82,066	53.2	26.5	20.3
	The Number System	82,066	58.6	19.3	22.1
	Expressions and Equations	82,066	59.4	22.7	18.0
	Geometry, Statistics and Probability	82,066	53.5	29.0	17.6
7	Ratio and Proportional Relationships	82,799	53.5	28.9	17.6
	The Number System	82,799	58.4	23.7	17.9
	Expressions & Equations	82,799	60.6	21.1	18.4
	Geometry, Statistics and Probability	82,799	57.2	27.2	15.6

Grade	Reporting Category	N	%Level 1	%Level 2	%Level 3
8	Expressions and Equations	86,031	58.5	22.0	19.5
	Functions	86,031	57.2	25.6	17.3
	Geometry	86,031	54.2	33.5	12.3
	Statistics and Probability and The Number System	86,031	53.0	33.4	13.7

Note. Level 1 = *Below Mastery*, Level 2 = *At or Around Mastery*, Level 3 = *Above Mastery*

Table 8.4. Test Results by Accommodation—ELA

Grade	Accommodation	N	SS Mean	SS SD	%Level 1	%Level 2	%Level 3	%Level 4
3	Adult Transcription	35	2460.34	17.02	97.1	2.9	0.0	0.0
	American Sign Language	15	2456.33	20.14	93.3	0.0	6.7	0.0
	Assistive Technology	19	2462.00	25.86	89.5	5.3	5.3	0.0
	Braille Test Booklet	3	*	*	*	*	*	*
	Large Print Test Booklet	8	*	*	*	*	*	*
	Read Aloud Content	109	2470.18	26.64	84.4	5.5	8.3	1.8
	Sign Test Content	8	*	*	*	*	*	*
	Simplified Directions	465	2463.70	23.60	88.8	5.6	4.7	0.9
	Translate Directions	2	*	*	*	*	*	*
	Translation Dictionary	3	*	*	*	*	*	*
4	Adult Transcription	15	2481.13	22.46	86.7	6.7	6.7	0.0
	American Sign Language	16	2479.31	25.88	93.8	0.0	0.0	6.3
	Assistive Technology	20	2479.65	23.87	85.0	5.0	10.0	0.0
	Braille Test Booklet	5	*	*	*	*	*	*
	Large Print Test Booklet	11	2492.27	35.64	54.5	18.2	27.3	0.0
	Read Aloud Content	48	2485.88	24.82	81.3	8.3	10.4	0.0
	Sign Test Content	10	*	*	*	*	*	*
	Simplified Directions	516	2486.33	22.56	84.3	7.8	7.8	0.2
	Translate Directions	117	2487.36	20.36	86.3	6.8	6.8	0.0
	Translation Dictionary	181	2486.01	21.14	87.8	6.1	6.1	0.0
5	Adult Transcription	14	2502.57	39.52	85.7	0.0	7.1	7.1
	American Sign Language	19	2484.11	10.60	100.0	0.0	0.0	0.0
	Assistive Technology	16	2492.38	23.58	87.5	6.3	6.3	0.0
	Braille Test Booklet	8	*	*	*	*	*	*
	Large Print Test Booklet	10	*	*	*	*	*	*
	Read Aloud Content	64	2496.45	25.10	87.5	7.8	4.7	0.0
	Sign Test Content	0	*	*	*	*	*	*
	Simplified Directions	441	2495.58	23.93	82.5	13.6	3.6	0.2
	Translate Directions	86	2493.53	24.42	82.6	12.8	4.7	0.0
	Translation Dictionary	132	2493.57	25.09	84.8	11.4	2.3	1.5
6	Adult Transcription	10	*	*	*	*	*	*
	American Sign Language	20	2498.20	23.85	90.0	10.0	0.0	0.0
	Assistive Technology	16	2521.88	37.94	62.5	18.8	12.5	6.3
	Braille Test Booklet	5	*	*	*	*	*	*
	Large Print Test Booklet	10	*	*	*	*	*	*
	Read Aloud Content	51	2508.84	28.26	74.5	21.6	3.9	0.0
	Sign Test Content	1	*	*	*	*	*	*
	Simplified Directions	397	2507.40	25.85	81.1	11.6	7.3	0.0
	Translate Directions	74	2503.01	27.69	86.5	5.4	8.1	0.0

Grade	Accommodation	N	SS Mean	SS SD	%Level 1	%Level 2	%Level 3	%Level 4
7	Translation Dictionary	115	2504.82	27.09	82.6	10.4	7.0	0.0
	Adult Transcription	13	2530.85	28.12	84.6	7.7	0.0	7.7
	American Sign Language	37	2514.11	13.98	94.6	5.4	0.0	0.0
	Assistive Technology	10	*	*	*	*	*	*
	Braille Test Booklet	2	*	*	*	*	*	*
	Large Print Test Booklet	7	*	*	*	*	*	*
	Read Aloud Content	42	2534.93	22.26	66.7	19.0	14.3	0.0
	Sign Test Content	3	*	*	*	*	*	*
	Simplified Directions	318	2523.75	21.88	81.1	12.3	6.3	0.3
	Translate Directions	76	2515.41	19.57	92.1	6.6	1.3	0.0
8	Translation Dictionary	177	2517.76	17.83	91.0	7.9	1.1	0.0
	Adult Transcription	6	*	*	*	*	*	*
	American Sign Language	24	2518.92	22.55	87.5	12.5	0.0	0.0
	Assistive Technology	18	2546.50	27.69	61.1	16.7	22.2	0.0
	Braille Test Booklet	3	*	*	*	*	*	*
	Large Print Test Booklet	9	*	*	*	*	*	*
	Read Aloud Content	30	2533.67	19.22	86.7	10.0	3.3	0.0
	Sign Test Content	1	*	*	*	*	*	*
	Simplified Directions	331	2524.86	23.16	85.8	10.3	3.3	0.6
	Translate Directions	80	2519.73	21.77	88.8	10.0	1.3	0.0

SS = scale score, SD = standard deviation, Level 1 = *Minimally Proficient*, Level 2 = *Partially Proficient*, Level 3 = *Proficient*, Level 4 = *Highly Proficient*. Statistics for subgroups with less than 11 students are omitted in compliance with FERPA regulations. Read aloud is for Writing only.

Table 8.5. Test Results by Accommodation—Mathematics

Grade	Accommodation	N	SS Mean	SS SD	%Level 1	%Level 2	%Level 3	%Level 4
3	Adult Transcription	33	3461.58	31.29	87.9	9.1	3.0	0.0
	American Sign Language	15	3457.33	32.68	93.3	0.0	6.7	0.0
	Assistive Technology	20	3458.65	33.15	95.0	0.0	0.0	5.0
	Braille Test Booklet	4	*	*	*	*	*	*
	Large Print Test Booklet	9	*	*	*	*	*	*
	Read Aloud Content	124	3485.48	41.98	56.5	26.6	16.1	0.8
	Sign Test Content	7	*	*	*	*	*	*
	Simplified Directions	448	3466.28	41.46	75.2	16.5	7.4	0.9
	Translate Directions	81	3465.48	40.23	74.1	17.3	8.6	0.0
	Translation Dictionary	102	3467.32	41.06	71.6	18.6	9.8	0.0
4	Adult Transcription	16	3493.94	37.08	81.3	12.5	6.3	0.0
	American Sign Language	17	3473.88	48.76	88.2	0.0	5.9	5.9
	Assistive Technology	20	3488.60	39.18	80.0	15.0	5.0	0.0
	Braille Test Booklet	5	*	*	*	*	*	*
	Large Print Test Booklet	12	3498.67	55.49	75.0	0.0	25.0	0.0
	Read Aloud Content	62	3507.68	35.18	72.6	17.7	9.7	0.0
	Sign Test Content	11	3488.27	36.19	81.8	9.1	9.1	0.0
	Simplified Directions	460	3499.63	39.91	76.5	15.9	7.0	0.7
	Translate Directions	80	3496.53	36.17	81.3	11.3	7.5	0.0
	Translation Dictionary	119	3501.83	35.80	78.2	12.6	9.2	0.0

Grade	Accommodation	N	SS Mean	SS SD	%Level 1	%Level 2	%Level 3	%Level 4
5	Adult Transcription	6	*	*	*	*	*	*
	American Sign Language	19	3529.47	36.55	84.2	10.5	5.3	0.0
	Assistive Technology	16	3537.31	28.34	81.3	18.8	0.0	0.0
	Braille Test Booklet	8	*	*	*	*	*	*
	Large Print Test Booklet	10	*	*	*	*	*	*
	Read Aloud Content	66	3545.71	31.89	62.1	31.8	6.1	0.0
	Sign Test Content	0	*	*	*	*	*	*
	Simplified Directions	418	3544.77	34.08	69.6	20.8	8.1	1.4
	Translate Directions	77	3547.06	34.59	67.5	23.4	7.8	1.3
	Translation Dictionary	110	3545.35	35.57	68.2	23.6	6.4	1.8
6	Adult Transcription	9	*	*	*	*	*	*
	American Sign Language	21	3558.57	31.27	95.2	0.0	4.8	0.0
	Assistive Technology	13	3580.85	48.68	84.6	0.0	7.7	7.7
	Braille Test Booklet	6	*	*	*	*	*	*
	Large Print Test Booklet	11	3578.91	49.45	54.5	18.2	27.3	0.0
	Read Aloud Content	46	3576.02	34.21	80.4	15.2	4.3	0.0
	Sign Test Content	1	*	*	*	*	*	*
	Simplified Directions	388	3572.89	28.70	82.5	13.1	4.1	0.3
	Translate Directions	76	3572.13	26.85	81.6	15.8	2.6	0.0
7	Translation Dictionary	105	3576.88	32.14	81.0	10.5	7.6	1.0
	Adult Transcription	5	*	*	*	*	*	*
	American Sign Language	38	3582.24	19.27	97.4	2.6	0.0	0.0
	Assistive Technology	5	*	*	*	*	*	*
	Braille Test Booklet	3	*	*	*	*	*	*
	Large Print Test Booklet	7	*	*	*	*	*	*
	Read Aloud Content	29	3596.72	27.67	86.2	10.3	3.4	0.0
	Sign Test Content	2	*	*	*	*	*	*
	Simplified Directions	293	3588.75	32.02	90.1	4.4	3.8	1.7
8	Translate Directions	66	3584.03	25.03	93.9	4.5	1.5	0.0
	Translation Dictionary	128	3587.61	26.98	90.6	5.5	3.9	0.0
	Adult Transcription	5	*	*	*	*	*	*
	American Sign Language	24	3625.96	19.10	91.7	4.2	4.2	0.0
	Assistive Technology	2	*	*	*	*	*	*
	Braille Test Booklet	3	*	*	*	*	*	*
	Large Print Test Booklet	10	*	*	*	*	*	*
	Read Aloud Content	14	3626.86	27.09	85.7	7.1	7.1	0.0
	Sign Test Content	1	*	*	*	*	*	*
	Simplified Directions	307	3625.51	19.50	89.9	7.8	2.0	0.3
	Translate Directions	65	3624.62	16.07	90.8	7.7	1.5	0.0

SS = scale score, SD = standard deviation, Level 1 = *Minimally Proficient*, Level 2 = *Partially Proficient*, Level 3 = *Proficient*, Level 4 = *Highly Proficient*. Statistics for subgroups with less than 11 students are omitted in compliance with FERPA regulations.

Table 8.6. Scale Score Distribution by Performance Level—ELA

Grade	Performance Level	N	Average Scale Score	%	Cumulative %
3	Level 1	38,268	2467.55	47.4	47.4
	Level 2	9,475	2501.76	11.7	59.1
	Level 3	21,917	2522.66	27.1	86.2
	Level 4	11,154	2555.77	13.8	100.0
4	Level 1	33,144	2486.34	41.1	41.1
	Level 2	11,421	2515.73	14.2	55.3
	Level 3	24,598	2538.53	30.5	85.8
	Level 4	11,496	2574.83	14.3	100.0
5	Level 1	32,511	2493.63	40.2	40.2
	Level 2	18,322	2530.74	22.6	62.8
	Level 3	22,884	2557.19	28.3	91.1
	Level 4	7,200	2591.75	8.9	100.0
6	Level 1	29,027	2505.72	35.7	35.7
	Level 2	18,192	2541.36	22.4	58.0
	Level 3	28,847	2569.95	35.5	93.5
	Level 4	5,303	2609.08	6.5	100.0
7	Level 1	31,935	2521.20	38.9	38.9
	Level 2	16,085	2551.18	19.6	58.5
	Level 3	26,941	2575.87	32.8	91.4
	Level 4	7,100	2612.18	8.7	100.0
8	Level 1	35,273	2524.32	41.4	41.4
	Level 2	18,952	2560.10	22.2	63.6
	Level 3	22,624	2585.48	26.5	90.2
	Level 4	8,383	2616.79	9.8	100.0

Note. 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table 8.7. Scale Score Distribution by Performance Level—Mathematics

Grade	Performance Level	N	Average Scale Score	%	Cumulative %
3	Level 1	25,610	3458.43	31.2	31.2
	Level 2	20,924	3512.74	25.5	56.8
	Level 3	23,204	3548.47	28.3	85.1
	Level 4	12,248	3587.28	14.9	100.0
4	Level 1	31,996	3494.19	39.3	39.3
	Level 2	17,838	3544.81	21.9	61.2
	Level 3	21,339	3578.75	26.2	87.4
	Level 4	10,307	3622.76	12.7	100.0
5	Level 1	29,725	3532.44	36.5	36.5
	Level 2	22,205	3576.96	27.3	63.8
	Level 3	20,398	3611.90	25.0	88.8
	Level 4	9,123	3655.46	11.2	100.0
6	Level 1	39,985	3571.54	48.7	48.7
	Level 2	17,146	3614.31	20.9	69.6
	Level 3	17,035	3643.18	20.8	90.4
	Level 4	7,900	3684.93	9.6	100.0

Grade	Performance Level	N	Average Scale Score	%	Cumulative %
7	Level 1	44,776	3590.87	54.1	54.1
	Level 2	13,112	3639.33	15.8	69.9
	Level 3	12,668	3662.93	15.3	85.2
	Level 4	12,243	3701.97	14.8	100.0
8	Level 1	46,383	3626.40	53.9	53.9
	Level 2	16,420	3659.66	19.1	73.0
	Level 3	14,317	3686.96	16.6	89.6
	Level 4	8,911	3730.56	10.4	100.0

Note. 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Chapter 9: RELIABILITY AND VALIDITY

This chapter provides evidence supporting the reliability and validity of scores on the Spring 2023 AASA assessments, addressing Standards 1.8, 1.9, 1.21, 2.3, 2.7, 2.8, 2.11, 2.15, 2.19, 3.1, 3.3, 3.6, 3.15, and 7.4 (AERA et al., 2014).

9.1. Reliability

The *Standards for Educational and Psychological Testing* (AERA et al., 2014) refer to reliability as the “consistency of scores across replications of a testing procedure” (p. 33). A reliable test produces stable scores, meaning that very similar score distributions would result if the test were administered repeatedly under similar conditions to the same students without memory or fatigue affecting the scores. The level of reliability/precision of scores has implications for validity in that the scores must be consistent and precise enough to be useful for intended purposes. If scores are to be meaningful, tests should produce stable scores if the same group of students were to take the same test repeatedly without any fatigue or memory of the test. The range of certainty around the score should also be small enough to support educational decisions.

9.1.1. Internal Consistency

Reliability was evaluated based on the internal consistency for all tests. For test reliability, coefficient alpha, which is based on classical test theory (CTT), is a frequently used measure of internal consistency. Coefficient alpha is computed as follows:

$$\alpha = \frac{k}{k-1} \left(1 - \frac{\sum \sigma_i^2}{\sigma_X^2} \right)$$

where k is the number of items, σ_X^2 is the variance of the total score, and σ_i^2 is the variance of item i (Crocker & Algina, 1986; Cronbach, 1951).

Typically, a test score is obtained from a single observation of performance and represents an estimate of the trait being measured. As an estimate, an observed test score contains some measurement error and does not perfectly reflect an individual’s true score. The degree of measurement error in a test score can be estimated using a statistic called the standard error of measurement (SEM), which is calculated as follows:

$$SEM = \sigma_X \sqrt{1-r}$$

where σ_X is a standard deviation of total score X , and r is a reliability coefficient, such as the coefficient alpha (Crocker & Algina, 1986).

Table 9.1, Table 9.2, and Table 9.3 present coefficient alphas and SEMs (computed based on the calibration sample) for the total and reporting category scores. These results suggest that the AASA assessments produce reliable scores.

Table 9.1. Coefficient Alpha and SEM by Total and Reporting Category Score—ELA, Form 1

Grade	Reporting Category	N	#Items	Coefficient Alpha	SEM
3	Total	38,645	44	0.92	3.12
	Reading for Information	38,670	19	0.83	1.87
	Reading for Literature	38,648	16	0.80	1.79
	Writing and Language	38,645	9	0.81	1.65
4	Total	40,115	44	0.92	3.32
	Reading for Information	40,107	18	0.77	1.98
	Reading for Literature	40,114	17	0.83	1.86
	Writing and Language	40,115	9	0.84	1.74
5	Total	40,445	44	0.92	3.18
	Reading for Information	40,442	19	0.83	1.88
	Reading for Literature	40,277	16	0.79	1.75
	Writing and Language	40,445	9	0.83	1.73
6	Total	40,337	44	0.92	3.16
	Reading for Information	40,341	20	0.82	1.95
	Reading for Literature	40,350	15	0.78	1.66
	Writing and Language	40,337	9	0.82	1.73
7	Total	40,812	44	0.91	3.32
	Reading for Information	40,820	20	0.81	1.93
	Reading for Literature	40,683	14	0.74	1.82
	Writing and Language	40,812	10	0.78	1.90
8	Total	40,998	44	0.92	3.16
	Reading for Information	40,987	20	0.80	1.94
	Reading for Literature	40,906	15	0.79	1.69
	Writing and Language	40,998	9	0.82	1.71

Table 9.2. Coefficient Alpha and SEM by Total and Reporting Category Score—ELA, Form 2

Grade	Reporting Category	N	#Items	Coefficient Alpha	SEM
3	Total	37,319	44	0.92	3.12
	Reading for Information	37,357	19	0.82	1.87
	Reading for Literature	37,330	16	0.80	1.79
	Writing and Language	37,319	9	0.81	1.64
4	Total	36,636	44	0.92	3.26
	Reading for Information	36,637	18	0.77	1.98
	Reading for Literature	36,642	17	0.83	1.86
	Writing and Language	36,636	9	0.81	1.69
5	Total	36,613	44	0.92	3.15
	Reading for Information	36,611	19	0.83	1.88
	Reading for Literature	36,461	16	0.79	1.75
	Writing and Language	36,613	9	0.82	1.70
6	Total	36,461	44	0.92	3.17
	Reading for Information	36,465	20	0.82	1.95
	Reading for Literature	36,470	15	0.78	1.67
	Writing and Language	36,461	9	0.81	1.74

Grade	Reporting Category	N	#Items	Coefficient Alpha	SEM
7	Total	36,937	44	0.91	3.32
	Reading for Information	36,945	20	0.81	1.93
	Reading for Literature	36,846	14	0.74	1.82
	Writing and Language	36,937	10	0.78	1.91
8	Total	40,217	44	0.92	3.13
	Reading for Information	40,204	20	0.80	1.94
	Reading for Literature	40,141	15	0.80	1.68
	Writing and Language	40,217	9	0.83	1.66

Table 9.3. Coefficient Alpha and SEM by Total and Reporting Category Score—Mathematics

Grade	Reporting Category	N	#Items	Coefficient Alpha	SEM
3	Total	77,974	45	0.95	2.67
	Operations, Algebraic Thinking, and Numbers in Base Ten	78,029	23	0.93	1.83
	Numbers and Operations – Fractions	77,993	9	0.79	1.22
	Measurement, Data, and Geometry	77,974	13	0.80	1.47
4	Total	77,503	45	0.95	2.70
	Operations, Algebraic Thinking, and Numbers in Base Ten	77,515	23	0.91	1.94
	Numbers and Operations – Fractions	77,503	14	0.88	1.45
	Measurement, Data, and Geometry	77,443	8	0.73	1.14
5	Total	77,495	45	0.94	2.77
	Operations, Algebraic Thinking, and Numbers in Base Ten	77,495	18	0.90	1.69
	Numbers and Operations – Fractions	77,507	15	0.82	1.64
	Measurement, Data, and Geometry	77,528	12	0.80	1.41
6	Total	77,446	47	0.94	2.86
	Ratio and Proportional Relationships	77,395	10	0.83	1.26
	The Number System	77,290	14	0.85	1.52
	Expressions and Equations	77,433	15	0.84	1.61
	Geometry, Statistics and Probability	77,446	8	0.60	1.24
7	Total	78,270	47	0.95	2.72
	Ratio and Proportional Relationships	78,145	10	0.85	1.22
	The Number System	78,450	10	0.83	1.15
	Expressions and Equations	78,442	12	0.82	1.41
	Geometry, Statistics and Probability	78,270	15	0.80	1.59
8	Total	82,005	47	0.92	2.90
	Functions	81,619	11	0.73	1.44
	Expressions and Equations	82,015	15	0.84	1.57
	Geometry	82,005	9	0.68	1.22
	Statistics and Probability and The Number System	81,977	12	0.72	1.51

In contrast to the CTT-based SEM, an IRT-based SEM (i.e., CSEM) varies across an ability continuum. The CSEM should be lower around important performance level cuts (e.g., *Proficient*), which indicates higher measurement precision. The CSEM tends to be higher for upper and lower ends of the ability continuum because there are usually fewer items that measure those difficulty levels. Figure B.19 – Figure B.54 in Appendix B present the TCC and CSEM curves of the assessments. As expected, the CSEMs around the performance level cuts were the lowest.

9.1.2. Inter-rater Reliability

For the hand-scored ELA writing prompts, the consistency with which two raters assign scores to student responses is determined by inter-rater agreement, also referred to as rater agreement, which indicates the level of agreement between two scores assigned to student responses. It is the measure of how often scorers agree with each other. Rater agreement is calculated between the human-scored and IEA-scored prompts, and rater agreement statistics include the percentage of exact and adjacent scores for each item that received two scores. For 10% of responses, a second “reliability” score was assigned by a second scorer.

The expectation is an inter-rater agreement of 65% or higher between the first and second scores. When IEA provided a high confidence score, the second reliability score was from a human rater. For the subset of responses where IEA provided a low confidence score, the first and second score were both from human raters. Pearson scoring staff used inter-rater agreement indices as one factor in determining the needs for continuing training and intervention on both individual and group levels.

Two other statistical indices are also used to measure reliability in the hand-scoring process: Cohen’s kappa and intraclass correlation. The quadratic weighted kappa (Cohen, 1968) allows rater disagreements to be weighted differentially (e.g., magnitude of a 1-point difference in ratings versus a 2-point difference) and is calculated with the weighted differences included, which are defined by the following formulas:

$$w_{ij} = \frac{(|i - j|)^2}{(k - 1)^2}$$

$$\kappa_w = 1 - \frac{\sum w_{ij} O_{ij}}{\sum w_{ij} E_{ij}}$$

where $|i - j|$ is the number of categories by which raters disagree, k is the total number of score categories, and w_{ij} is the weighted level of disagreement. E_{ij} is the expected matrix, and O_{ij} is the observed matrix. The quadratic weighted kappa ranges from -1.0 to 1.0, with higher, more positive values indicative of greater rater agreement.

The intraclass correlation is defined by Shrout and Fleiss (1979) as “the correlation between one measurement (either a single rating or a mean of ratings) on a target and another measurement obtained on that target” (p. 422). In the context of the AASA assessments, the “target” was the student response and each measurement was obtained by a rater randomly assigned to that response. Therefore, $ICC(1,1)$ was used to estimate the intraclass correlation. $ICC(1,1)$ is estimated as follows (Shrout & Fleiss, 1979):

$$ICC(1,1) = \frac{BMS - WMS}{BMS + (k - 1)WMS}$$

where *BMS* is the between-targets mean square, *WMS* is the within-targets mean square, and *k* is the number of raters rating each target. Table 9.4 presents the quadratic weighted kappa and intraclass correlation by reporting category. Items with a kappa statistic lower than 0.20, considered as slight rater agreement (Landis & Koch, 1977) and of which there were none, were flagged for potential replacement in future administrations.

Table 9.4. Inter-rater Reliability Statistics

Grade	OE Item	Trait	Score Range	N	Quadratic Kappa	ICC	%Exact Agreement	%Adjacent Agreement
3	WR 1	Statement of Purpose, Focus & Organization	1–4	3,825	0.87	0.87	0.85	0.15
		Evidence & Elaboration	1–4	3,825	0.87	0.87	0.87	0.13
		Conventions & Editing	0–2	3,825	0.95	0.95	0.94	0.06
	WR 2	Statement of Purpose, Focus & Organization	1–4	3,097	0.73	0.73	0.71	0.29
		Evidence & Elaboration	1–4	3,097	0.71	0.70	0.72	0.28
		Conventions & Editing	0–2	3,097	0.98	0.98	0.98	0.02
4	WR 1	Statement of Purpose, Focus & Organization	1–4	4,018	0.87	0.87	0.81	0.18
		Evidence & Elaboration	1–4	4,018	0.87	0.87	0.81	0.18
		Conventions & Editing	0–2	4,018	0.93	0.93	0.94	0.06
	WR 2	Statement of Purpose, Focus & Organization	1–4	3,673	0.69	0.69	0.68	0.31
		Evidence & Elaboration	1–4	3,673	0.73	0.73	0.77	0.23
		Conventions & Editing	0–2	3,673	0.86	0.86	0.86	0.13
5	WR 1	Statement of Purpose, Focus & Organization	1–4	4,035	0.84	0.84	0.80	0.20
		Evidence & Elaboration	1–4	4,035	0.82	0.82	0.78	0.22
		Conventions & Editing	0–2	4,035	0.88	0.88	0.90	0.10
	WR 2	Statement of Purpose, Focus & Organization	1–4	3,752	0.88	0.88	0.83	0.16
		Evidence & Elaboration	1–4	3,752	0.87	0.87	0.84	0.16
		Conventions & Editing	0–2	3,752	0.91	0.91	0.92	0.08
6	WR 1	Statement of Purpose, Focus & Organization	1–4	4,033	0.83	0.83	0.73	0.27
		Evidence & Elaboration	1–4	4,033	0.81	0.81	0.74	0.26
		Conventions & Editing	0–2	4,033	0.90	0.90	0.90	0.10
	WR 2	Statement of Purpose, Focus & Organization	1–4	3,648	0.78	0.78	0.65	0.34
		Evidence & Elaboration	1–4	3,648	0.78	0.78	0.69	0.31
		Conventions & Editing	0–2	3,648	0.88	0.88	0.91	0.09
	RI	Human-Scored Reading Item	0–1	7,588	0.93	0.93	0.97	0.03
7	WR 1	Statement of Purpose, Focus & Organization	1–4	4,083	0.81	0.81	0.70	0.30
		Evidence & Elaboration	1–4	4,083	0.77	0.77	0.69	0.31
		Conventions & Editing	0–2	4,083	0.87	0.87	0.88	0.12
	WR 2	Statement of Purpose, Focus & Organization	1–4	3,698	0.76	0.76	0.66	0.34
		Evidence & Elaboration	1–4	3,698	0.77	0.77	0.68	0.32
		Conventions & Editing	0–2	3,698	0.86	0.86	0.88	0.12
8	WR 1	Statement of Purpose, Focus & Organization	1–4	4,100	0.81	0.81	0.71	0.29
		Evidence & Elaboration	1–4	4,100	0.83	0.83	0.74	0.25
		Conventions & Editing	0–2	4,100	0.87	0.87	0.89	0.11
	WR 2	Statement of Purpose, Focus & Organization	1–4	4,030	0.76	0.76	0.65	0.34
		Evidence & Elaboration	1–4	4,030	0.80	0.80	0.69	0.31
		Conventions & Editing	0–2	4,030	0.86	0.86	0.90	0.10

Note. OE = open-ended, ICC = intraclass correlation

9.2. Differential Item Functioning

Because test scores can have many sources of variation, the test developers' task is to create assessments that measure the intended abilities and skills without introducing extraneous elements or construct-irrelevant variance. When tests measure something other than what they are intended to measure, test scores will reflect these unintended skills and knowledge, as well as what is purportedly assessed by the test. If this occurs, these tests can be called biased (Angoff, 1993; Camilli & Shepard, 1994; Green, 1975; Zumbo, 1999). One of the factors that may render test scores biased is differing cultural and socioeconomic experiences.

Analysis of DIF is a statistical method to detect potential bias of an item. DIF is defined as a difference between groups (e.g., male and female) in the probability of answering an item correctly. DIF analyses are conditioned on the ability that the assessment is intended to measure (e.g., mathematics proficiency). DIF is an indicator that the item might exhibit bias for one group over the other, not that it actually does. If DIF exists on an item, a committee composed of subject experts reviews the item to determine whether it actually shows bias.

The Mantel-Haenszel (MH) method (Holland & Thayer, 1988; Mantel & Haenszel, 1959) was used to investigate DIF on one-point items. The MH method is frequently used and efficient in terms of statistical power (Clauser & Mazor, 1998). The Mantel-Haenszel chi-square statistic is computed as follows:

$$MH - \chi^2 = \frac{(\sum_k F_k - \sum_k E(F_k))^2}{\sum_k Var(F_k)}$$

where F_k is the sum of scores for the focal group at the k th level of the matching variable (Zwick et al., 1993). The MH statistic is sensitive to N such that larger sample sizes increase the value of chi-square.

In addition to the MH chi-square statistic, the MH delta statistic (ΔMH) was computed. Educational Testing Service (ETS) first developed the ΔMH DIF statistic. To compute the ΔMH DIF, the MH alpha (the odds ratio) is first computed:

$$\sigma_{MH} = \frac{\sum_{k=1}^K N_{r1k} N_{f0k} / N_k}{\sum_{k=1}^K N_{f1k} N_{r0k} / N_k}$$

where N_{r1k} is the number of correct responses in the reference group at ability level k , N_{f0k} is the number of incorrect responses in the focal group at ability level k , N_k is the total number of responses, N_{f1k} is the number of correct responses in the focal group at ability level k , and N_{r0k} is the number of incorrect responses in the reference group at ability level k . The ΔMH DIF is computed as follows:

$$\Delta MH \text{ DIF} = -2.35 \ln(\alpha_{MH})$$

Positive values of $\Delta MH\ DIF$ indicate items that favor the focal group, whereas negative values indicate items that favor the reference group. The MH chi-square statistic and the $\Delta MH\ DIF$ were used in combination to identify both the operational and field test items that exhibit strong, weak, or no DIF for single-point items.

The standardized mean difference (SMD) is another DIF method applied to multiple-point items (Dorans & Schmitt, 1991; Zwick et al., 1993). The SMD is an effect size index of DIF that compares the mean scores of the reference and focal groups for an item, adjusting for the distribution of the reference and focal groups on the conditioned variable, which for the analyses is the raw score. The SMD is computed as follows:

$$SMD = \sum_k P_{F_k} (m_{F_k} - m_{R_k})$$

where P_{F_k} is the proportion of the focal group at the k th level of the matching variable, m_{F_k} is the mean score on the item for the focal group at the k th level of the matching variable, and m_{R_k} is the mean score on the item for the reference group at the k th level of the matching variable (Zwick et al., 1993). A negative SMD value indicates an item in which the focal group has a lower mean than the reference group, conditioned on the matching variable (e.g., science proficiency), whereas a positive SMD value indicates an item for which the reference group has a lower mean than the focal group, conditioned on the matching variable.

Table 9.5 presents the summary of DIF classification criteria for both the MH method and SMD. An alpha level of 0.05 was used for all MH and SMD statistics.

Table 9.5. DIF Flag Categories

Category	Description	MH Criterion	SMD Criterion
A	No DIF	MH chi-square not significantly different from 0 ($p < 0.05$) or $ \Delta MH\ DIF < 1.0$	MH chi-square not significantly different from 0 ($p < 0.05$) or $ SMD \leq 0.17$
B	Weak DIF	MH chi-square significantly different from 0 ($p < 0.05$) and $1.0 \leq \Delta MH\ DIF < 1.5$	MH chi-square significantly different from 0 ($p < 0.05$) and $0.17 < SMD \leq 0.25$
C	Strong DIF	MH chi-square significantly higher than 1 ($p < 0.05$) and $ \Delta MH\ DIF \geq 1.5$	MH chi-square significantly different from 0 ($p < 0.05$) and $ SMD > 0.25$

The DIF analysis was conducted for 10 different group pairs:

1. Female vs. Male
2. Hispanic vs. Non-Hispanic
3. American Indian vs. White
4. Asian vs. White
5. Black or African American vs. White
6. Native Hawaiian or Other Pacific Islander vs. White
7. Multi-racial vs. White

8. Students with Disability vs. Students without Disability
9. Economically Disadvantaged vs. Not Economically Disadvantaged
10. English Learner vs. English as a First Language

Table 9.6 presents the number of operational items exhibiting strong DIF between any two groups. Any items that display strong DIF are flagged for possible replacement in the future administration, as strong DIF is one of the holistic item replacement evaluation criteria used for item selection. DIF results with a sample size of less than 200 per group should not be considered statistically reliable (Clauser & Mazor, 1998; Mazor et al., 1992).

Table 9.6. Number of Items Exhibiting Strong DIF

Content Area	Grade	#Items	#Items with Strong DIF
ELA, Form 1	3	44	0
	4	44	0
	5	44	1
	6	44	1
	7	44	1
	8	44	1
ELA, Form 2	3	44	0
	4	44	0
	5	44	1
	6	44	1
	7	44	1
	8	44	2
Mathematics	3	45	1
	4	45	0
	5	45	0
	6	47	0
	7	47	1
	8	47	2

9.3. Correlations Among Reporting Categories

Correlations were examined between the total raw score and the reporting category raw scores. The data used to calculate the correlations were based on the calibration sample described in Chapter 7. Disattenuated correlations between were also computed, calculated based on the following formula:

$$r_{T_{xy}} = \frac{r_{xy}}{\sqrt{r_x r_y}}$$

where $r_{T_{xy}}$ is a corrected correlation for attenuation between scores x and y , r_{xy} is an observed correlation between the scores x and y , and r_x and r_y are reliabilities for x and y , respectively. Coefficient alphas (presented in Table 9.1, Table 9.2, and Table 9.3) were used to calculate the corrected correlation coefficients for attenuation. The disattenuated correlations could be greater than 1.00.

Table 9.7 – Table 9.11 present the test correlations and disattenuated correlations between the total raw score and the reporting category raw scores. The numbers in the lower diagonal of the table are the disattenuated correlations.

Table 9.7. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—ELA Form 1

Grade	Score	Total	Reading for Information	Reading for Literature	Writing and Language
3	Total	1.00	0.93	0.92	0.89
	Reading for Information	1.06	1.00	0.81	0.73
	Reading for Literature	1.07	0.99	1.00	0.72
	Writing and Language	1.03	0.89	0.89	1.00
4	Total	1.00	0.91	0.92	0.89
	Reading for Information	1.08	1.00	0.77	0.71
	Reading for Literature	1.05	0.96	1.00	0.71
	Writing and Language	1.01	0.88	0.85	1.00
5	Total	1.00	0.93	0.91	0.90
	Reading for Information	1.06	1.00	0.78	0.74
	Reading for Literature	1.07	0.96	1.00	0.72
	Writing and Language	1.03	0.89	0.89	1.00
6	Total	1.00	0.93	0.90	0.90
	Reading for Information	1.07	1.00	0.78	0.75
	Reading for Literature	1.06	0.98	1.00	0.72
	Writing and Language	1.04	0.91	0.90	1.00
7	Total	1.00	0.92	0.89	0.89
	Reading for Information	1.07	1.00	0.75	0.71
	Reading for Literature	1.08	0.97	1.00	0.68
	Writing and Language	1.06	0.89	0.90	1.00
8	Total	1.00	0.92	0.90	0.89
	Reading for Information	1.07	1.00	0.77	0.71
	Reading for Literature	1.06	0.97	1.00	0.69
	Writing and Language	1.02	0.88	0.86	1.00

Table 9.8. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—ELA Form 2

Grade	Score	Total	Reading for Information	Reading for Literature	Writing and Language
3	Total	1.00	0.93	0.92	0.88
	Reading for Information	1.07	1.00	0.81	0.72
	Reading for Literature	1.07	1.00	1.00	0.72
	Writing and Language	1.02	0.88	0.89	1.00
4	Total	1.00	0.91	0.92	0.88
	Reading for Information	1.08	1.00	0.77	0.70
	Reading for Literature	1.05	0.96	1.00	0.70
	Writing and Language	1.02	0.89	0.85	1.00
5	Total	1.00	0.93	0.91	0.91
	Reading for Information	1.06	1.00	0.78	0.76
	Reading for Literature	1.07	0.96	1.00	0.73
	Writing and Language	1.05	0.92	0.91	1.00
6	Total	1.00	0.93	0.90	0.90
	Reading for Information	1.07	1.00	0.78	0.74
	Reading for Literature	1.06	0.98	1.00	0.72
	Writing and Language	1.04	0.91	0.91	1.00
7	Total	1.00	0.92	0.89	0.89
	Reading for Information	1.07	1.00	0.76	0.71
	Reading for Literature	1.08	0.98	1.00	0.68
	Writing and Language	1.06	0.89	0.90	1.00
8	Total	1.00	0.93	0.90	0.89
	Reading for Information	1.08	1.00	0.77	0.72
	Reading for Literature	1.05	0.96	1.00	0.70
	Writing and Language	1.02	0.88	0.86	1.00

Table 9.9. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—Mathematics Grades 3–5

Grade	Score	Total	Operations, Algebraic Thinking, and Numbers in Base Ten	Numbers and Operations – Fractions	Measurement, Data, and Geometry
3	Total	1.00	0.97	0.87	0.93
	Operations, Algebraic Thinking, and Numbers in Base Ten	1.03	1.00	0.77	0.85
	Numbers and Operations – Fractions	1.00	0.90	1.00	0.76
	Measurement, Data, and Geometry	1.07	0.99	0.96	1.00
4	Total	1.00	0.97	0.94	0.86
	Operations, Algebraic Thinking, and Numbers in Base Ten	1.04	1.00	0.84	0.78
	Numbers and Operations – Fractions	1.03	0.94	1.00	0.76
	Measurement, Data, and Geometry	1.03	0.96	0.95	1.00
5	Total	1.00	0.95	0.92	0.91
	Operations, Algebraic Thinking, and Numbers in Base Ten	1.03	1.00	0.80	0.80
	Numbers & Operations – Fractions	1.05	0.93	1.00	0.76
	Measurement, Data, and Geometry	1.05	0.94	0.94	1.00

Table 9.10. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—Mathematics Grades 6 and 7

Grade	Score	Total	Ratio and Proportional Relationships	The Number System	Expressions and Equations	Geometry, Statistics and Probability
6	Total	1.00	0.92	0.94	0.93	0.76
	Ratio and Proportional Relationships	1.04	1.00	0.83	0.80	0.62
	The Number System	1.05	0.99	1.00	0.82	0.64
	Expressions and Equations	1.05	0.96	0.97	1.00	0.64
	Geometry, Statistics and Probability	1.01	0.88	0.90	0.90	1.00
7	Total	1.00	0.92	0.92	0.93	0.91
	Ratio and Proportional Relationships	1.02	1.00	0.81	0.82	0.77
	The Number System	1.04	0.96	1.00	0.82	0.79
	Expressions and Equations	1.05	0.98	0.99	1.00	0.77
	Geometry, Statistics and Probability	1.04	0.93	0.97	0.95	1.00

Table 9.11. Correlations and Disattenuated Correlations between Total and Reporting Category Raw Score—Mathematics Grade 8

Grade	Score	Total	Expressions and Equations	Functions	Geometry	Statistics and Probability and The Number System
8	Total	1.00	0.93	0.88	0.82	0.87
	Expressions and Equations	1.06	1.00	0.76	0.69	0.74
	Functions	1.07	0.97	1.00	0.65	0.69
	Geometry	1.04	0.91	0.92	1.00	0.64
	Statistics and Probability and The Number System	1.07	0.95	0.95	0.91	1.00

9.4. Validity Evidence

According to the *Standards for Educational and Psychological Testing* (AERA et al., 2014), “Validity refers to the degree to which evidence and theory support the interpretations of test scores entailed for proposed uses of tests. Validity is, therefore, the most fundamental consideration in developing and evaluating tests” (p. 11). The purpose of test score validation is not to validate the test itself but to validate interpretations of the test scores for a particular purpose or use.

A validity argument should begin with clear statements regarding the purpose of a test and intended interpretations and uses of the test results. The purpose of the AASA tests is to assess the ELA and mathematics proficiency of students based on the Arizona Academic Standards. The objective of the proceeding sections is to highlight validity evidence for each aspect and to guide interested readers where to look for the evidence. Different aspects of validity evidence, which are in line with the *Standards* (AERA et al., 2014), are considered throughout this technical report. Providing validity evidence is an ongoing activity for any assessment as it matures.

9.4.1. Evidence Based on Test Content

Validity evidence based on test content refers to the extent to which a test is aligned with the construct the assessment is intended to measure (AERA et al., 2014, p. 14). AASA measures a student's level of ELA and mathematics proficiency based on the skills specified in the Arizona Academic Standards. Although the validity of AASA test score interpretations is evaluated along several dimensions as a criterion-referenced system of tests, the meaning of test scores is critically evaluated by the degree to which test content is aligned with the standards. The AASA ELA and Mathematics assessments are rigorously examined in accordance with the guidelines in the *Standards for Educational and Psychological Testing* (AERA et al., 2014). The Elementary and Secondary Education Act (ESEA) legislation also describes the evidence that is necessary to validate assessment scores for their intended purposes based on these standards.

Alignment of content standards is achieved through a rigorous, iterative test development process that proceeds from the content standards and begins with the item specifications and test blueprints, the core documents that ensure that the assessments are aligned to the Arizona Academic Standards. The item specifications define the content limit, model tasks, and response types for a specific standard, and the test blueprint defines the standards to be assessed for each test form, the number of items per standard, the number of item types, the number of points per item type, and the total number of items and points per test form. In addition to ensuring that test items are aligned with their intended content standards, each assessment is intended to measure a representative sample of the knowledge and skills identified in the standards. Thus, the blueprints also represent a policy document specifying the relative importance of content strands and standards in addition to meeting important measurement goals.

Once the item specifications and blueprints are established, item and test development can begin. It was a rigorous and iterative process involving the Pearson content team and ADE to ensure that the AASA assessments meet the test blueprints and other content criteria and psychometric targets, as described in Chapter 3. Beyond the test blueprint, ADE and Pearson attempted to include items measuring different levels of rigor to cover the Arizona Academic Standards as much as possible.

Alignment of test forms to the test blueprints is a thoughtful, careful task that involves collaboration among assessment specialists, psychometricians, and ADE. Developing test forms is challenging because test blueprints can be highly complex, specifying not only the range of items and points for each reporting category and standard, but also cross-cutting criteria such as distribution across item types, DOK, writing genre, etc. In addition to meeting complex blueprint requirements, test developers worked to meet psychometric goals so that accommodated test forms measure equivalently across the range of student ability.

9.4.2. Evidence Based on Response Processes

Evidence based on response processes refers to the cognitive process engaged in by students when answering test items, or the “evidence concerning the fit between the construct and the detailed nature of performance or response actually engaged in by examinees” (AERA et al., 2014, p. 15). A standalone field test was administered in Spring 2022 for the ELA Writing test to increase the number of eligible writing prompts in the item bank for operational use in future administrations. New items were field tested in Grade 3 for to assess students Oral Reading Fluency (ORF), which is one of the reading foundation standards, in Spring 2022 and again in Spring 2023 to further explore their functioning and performance. The ORF items were designed to align with low, medium, and high levels of difficulty (based on Lexiles) and gauge students’ ability to read aloud words.

As presented in Chapter 3, all newly developed items also go through a rigorous item review process, including content, bias, and sensitivity committees with Arizona educators, parents, and community members. Reviewers evaluated the item for its alignment to the Arizona Science Standards, grade appropriateness, editorial completeness and accuracy, and the presence of any content that could be biased or sensitive in nature. Only the items accepted by the committees were considered eligible to be field tested.

9.4.3. Evidence Based on Internal Structure

Validity evidence based on internal structure refers to the extent to which an item or a component of a test ties to the assessment it is intended to measure (AERA et al., 2014, p. 16). AASA is designed to measure students’ overall ELA and mathematics proficiency based on the Arizona Academic Standards, which are composed of various reporting categories for each content area. AASA items across all reporting categories were calibrated concurrently under the unidimensional Rasch models (Masters, 1982; Rasch, 1960) as indicated in Chapter 7. To evaluate the unidimensionality assumption of the Rasch models, PCA was conducted for each operational form. The results of PCA analysis with the parallel analysis criterion (Horn, 1965), presented in Table 7.3, indicated that there is one dominant dimension for both ELA and mathematics and the remaining components are non-significant.

Another assumption under the Rasch models is local item independence. The local item independence assumption is typically evaluated using Q3 statistics (Yen, 1984); Winsteps (Linacre, 2022b) produces raw score residual correlations for pairs of items on a test, which are analogous to the Q3 statistics. A distribution of the residual correlations by form, presented in Table 7.4, showed that most statistics are either slightly negative or slightly positive, which indicates that the item independence assumption generally holds for the AASA tests.

In addition to the total scale score, the scale score for each reporting category is reported individually. The scale scores for the reporting categories are generated by including the items associated with each reporting category and using the item parameter estimates from the concurrent calibration across all reporting categories. Details about scaling methods are described in Section 7.5. Correlations between the total score and reporting category score presented in Section 9.3 show that they are at least moderately correlated to each other, if not highly correlated, as expected.

A point-biserial correlation, as an indicator of interrelationship between an item and a construct that it is intended to measure, is calculated as a correlation between an item raw score and a total raw score. The point-biserial correlations should be higher than or equal to 0.25, as any item with a lower correlation is flagged during item selection. It is one of the psychometric criteria considered for item selection. The point-biserial correlation was calculated for distractors of multiple-choice items as well. Table 6.5 and Table 6.6 show that all the multiple-choice items have negative point-biserial correlations, except for a few distractors with a slightly positive correlation close to zero. The results indicate that the distractors work as expected.

Differential item functioning (DIF) analysis is a statistical method to detect potential bias of an item for (or against) a manifest group (e.g., female). DIF is defined as a difference between groups (e.g., male and female) in the probability of getting an item correct, given the same level of ability within the construct that an assessment is intended to measure. Details on DIF analysis are presented in Section 9.2. Items showing strong DIF are flagged for possible replacement in future administrations.

9.4.4. Evidence Based on Performance Standards

Validity evidence concerning performance standards refers to the extent to which passing scores are aligned to performance standards (Kane, 1994). Performance level descriptors (PLDs) highlight the knowledge, skills, and processes students possess at different performance levels (Egan et al., 2012). The PLDs are the foundation of standard setting meetings. The PLDs for AASA, provided on the ADE website at <https://www.azed.gov/assessment/aasa>, were drafted prior to the 2015 standard setting workshop and included educator input. ADE considered any need for clarification or revision that arose throughout the standard setting process prior to publishing the final versions (American Institutes for Research, 2015). See Section 10.1 for more details on standard setting.

9.4.5. Evidence Based on Relation to Other Variables

Validity evidence concerning a relation to other variables refers to the extent to which test scores are related to other external measures (AERA et al., 2014, p. 16). Because both the ELA and mathematics AASA assessments are administered to all eligible Arizona students, scores on the tests are expected to be positively correlated. Table 9.12 presents the correlation between AASA ELA and mathematics scale scores from the Spring 2023 administration. The correlations range from 0.73 to 0.80.

Table 9.12. Correlation between AASA ELA and Mathematics Scale Scores

Grade	N	Correlation
3	80,649	0.79
4	80,469	0.79
5	80,725	0.76
6	81,104	0.77
7	81,688	0.80
8	84,723	0.73

9.4.6. Summary

Overall, the validity evidence supports the use of AASA scores. The PCA revealed unidimensionality of AASA, which supports the use of unidimensional Rasch models. The AASA ELA and mathematics scores were also positively correlated. Test score validation is not a quantifiable property but an ongoing process, beginning at initial conceptualization and continuing throughout the entire assessment process. Additional evidence should and will be added to the AASA technical report in the future, as appropriate.

Chapter 10: CLASSIFICATION INTO PERFORMANCE LEVELS

This chapter provides information regarding classification of students into performance levels for the Spring 2023 AASA assessments, addressing Standards 1.8, 1.9, 2.13, 2.14, 2.16, 5.5, 5.21, 5.22, 5.23, and 7.4 (AERA et al., 2014).

Scores from the AASA tests are used to classify students into one of four performance levels: *Minimally Proficient*, *Partially Proficient*, *Proficient*, and *Highly Proficient*. This section provides information regarding classification of students into these four categories, including the consistency and accuracy with which students who took the Spring 2023 AASA assessment were assigned to the performance levels.

10.1. Standard Setting

Arizona educators made recommendations for cut scores for each performance level on the AASA assessments during the standard setting workshop conducted from July 13–16, 2015, following the first operational administration of the AASA in Spring 2015 (known as the AzMERIT assessments at that time) using the bookmark standard setting procedure. The State Board of Education adopted the panelist-recommended performance standards on August 14, 2015. See the standard setting report for a detailed account of the workshop process and outcomes (American Institutes for Research, 2015).

Table 10.1 presents the final scale score ranges for the AASA performance levels, and Table 10.2 presents the scale score and associated CSEM at the performance level cuts. The CSEM is very similar across all grades and content areas within each cut.

Table 10.1. Performance Level Cut Scores

Content Area	Grade	<i>Minimally Proficient</i>	<i>Partially Proficient</i>	<i>Proficient</i>	<i>Highly Proficient</i>
ELA	3	2395–2496	2497–2508	2509–2540	2541–2605
	4	2400–2509	2510–2522	2523–2558	2559–2610
	5	2419–2519	2520–2542	2543–2577	2578–2629
	6	2431–2531	2532–2552	2553–2596	2597–2641
	7	2438–2542	2543–2560	2561–2599	2600–2648
	8	2448–2550	2551–2571	2572–2603	2604–2658
Mathematics	3	3395–3494	3495–3530	3531–3572	3573–3605
	4	3435–3529	3530–3561	3562–3605	3606–3645
	5	3478–3562	3563–3594	3595–3634	3635–3688
	6	3512–3601	3602–3628	3629–3662	3663–3722
	7	3529–3628	3629–3651	3652–3679	3680–3739
	8	3566–3649	3650–3672	3673–3704	3705–3776

Note. The scale score cut for Move on When Reading (MOWR) in Grade 3 is 2446.

Table 10.2. CSEM at Performance Level Cuts

Content Area	Grade	<i>Partially Proficient</i> Cut		<i>Proficient</i> Cut		<i>Highly Proficient</i> Cut	
		Scale Score	CSEM	Scale Score	CSEM	Scale Score	CSEM
ELA	3	2497	9	2509	9	2541	11
	4	2510	9	2523	9	2559	11
	5	2520	9	2543	9	2578	11
	6	2532	9	2553	9	2597	12
	7	2543	9	2561	9	2600	12
	8	2551	9	2572	9	2604	11
Mathematics	3	3495	10	3531	10	3573	14
	4	3530	10	3562	10	3606	14
	5	3563	10	3595	10	3635	12
	6	3602	9	3629	9	3663	11
	7	3629	10	3652	10	3680	11
	8	3650	10	3673	9	3705	11

Performance classifications for reporting categories are determined by student performance on the reporting categories compared to the respective *Proficient* performance standard. For each reporting category, a mid-range band is established by extending one CSEM below and above the *Proficient* performance standard scale score cut. If a student's scale score for a reporting category is fallen into the mid-range band, the student performance is classified as *At/Near Mastery* for the reporting category. On the other hand, if a student's scale score is above or below the mid-range band, the student performance is classified as *Above Mastery* or *Below Mastery*, respectively.

10.2. Classification Consistency and Accuracy

Classification consistency is the agreement between students' performance level classification from two independent administrations of the same test (or two parallel forms of the test). Classification accuracy refers to the agreement between the actual classifications using observed cut scores and true classifications based on known true cut scores (Livingston & Lewis, 1995).

In conjunction with internal consistency, classification consistency is an important type of reliability and is particularly relevant to high-stakes decisions, such as passing or not passing the AASA tests. As a form of reliability, classification consistency represents how reliably students can be classified into performance levels. For tests such as AASA, classification consistency is most important for students whose ability is near the *Proficient* cut score. Students whose ability is far above or far below the value established for *Proficient* are unlikely to be misclassified because repeated administration of the test will nearly always result in the same classification. Students whose true scores are close to the cut score are a more serious concern. These students' true scores will likely lie within the SEM of the cut score. For this reason, the measurement error at the cut scores should be considered when evaluating the classification consistency of a test.

Classification consistency and accuracy were estimated using the total scale score for the *Proficient* cut based on procedures described by Livingston and Lewis (1995). Classification consistency is calculated as the proportion of students in the diagonal in Table 10.3 (i.e., students classified consistently between two parallel forms, listed in bold). Similarly, classification accuracy is calculated as the proportion of students in the diagonal in Table 10.4 (i.e., students classified the same between observed scores and true scores, listed in bold).

Table 10.3. Classification Consistency for the *Proficient* Cut

		Expected Performance on Parallel Form	
		Not Proficient	Proficient
Observed Performance on Actual Form	Not Proficient	Consistent Classification	Inconsistent Classification
	Proficient	Inconsistent Classification	Consistent Classification

Table 10.4. Classification Accuracy for the *Proficient* Cut

		Expected Performance on Test	
		Not Proficient	Proficient
Observed Performance on Test	Not Proficient	Accurate Classification	False Negative
	Proficient	False Positive	Accurate Classification

Cohen's kappa (κ) coefficient (Cohen, 1960) is another way of expressing overall consistency. This statistic assesses the proportion of consistent classification expected beyond chance and is therefore most often lower than the unadjusted value of overall consistency. Cohen's kappa is calculated as follows:

$$\kappa = \frac{P - P_c}{1 - P_c}$$

where P_c is the probability of consistent classification by chance, and P is the probability of consistent classification (unadjusted by chance). Students can be misclassified in one of two ways. Students who are truly not *Proficient* but were classified as being *Proficient*, based on the assessment, are false positives. Similarly, students who are truly *Proficient* but were classified as being not *Proficient* are false negatives.

Table 10.5 presents the classification consistency and accuracy results, generated by BB-class (Brennan, 2004). These results are for classifying students into four performance levels using the total score on the assessment for students in the calibration sample. Included in the table are the sample size (N), classification consistency (Consistency), classification inconsistency (Inconsistency), probability of consistent classification by chance (Chance), Cohen's Kappa (κ), classification accuracy (Accuracy), false positive (False Positive), and false negative (False Negative). Inconsistency is defined as one minus Consistency.

Table 10.5. Classification Consistency and Accuracy Results

Content Area	Grade	N	Consistency	Inconsistency	Chance	κ	Accuracy	False Positive	False Negative
ELA, Form 1	3	38,703	0.73	0.27	0.33	0.60	0.80	0.11	0.09
	4	40,146	0.73	0.27	0.30	0.61	0.80	0.11	0.09
	5	40,455	0.72	0.28	0.31	0.60	0.80	0.11	0.09
	6	40,363	0.73	0.27	0.31	0.61	0.80	0.11	0.09
	7	40,841	0.72	0.28	0.31	0.59	0.79	0.11	0.09
	8	41,039	0.73	0.27	0.31	0.61	0.80	0.11	0.09
ELA, Form 2	3	37,390	0.74	0.26	0.34	0.60	0.80	0.10	0.09
	4	36,660	0.73	0.27	0.31	0.61	0.80	0.11	0.09
	5	36,628	0.72	0.28	0.30	0.60	0.80	0.11	0.09
	6	36,489	0.73	0.27	0.31	0.61	0.80	0.11	0.09
	7	36,970	0.72	0.28	0.31	0.59	0.79	0.11	0.09
	8	40,248	0.73	0.27	0.31	0.61	0.80	0.11	0.09
Mathematics	3	78,131	0.74	0.26	0.27	0.65	0.82	0.10	0.09
	4	77,590	0.77	0.23	0.29	0.67	0.83	0.09	0.08
	5	77,610	0.76	0.24	0.29	0.67	0.83	0.09	0.08
	6	77,507	0.78	0.22	0.34	0.67	0.84	0.09	0.07
	7	78,511	0.80	0.20	0.38	0.68	0.85	0.08	0.07
	8	82,075	0.78	0.22	0.38	0.65	0.84	0.09	0.07

10.3. MOWR Policy

Arizona’s Move On When Reading (MOWR) policy is designed to provide students with evidence-based, effective reading instruction in Grades K–3 to position them for success as they progress through school, college, and career. The heart of the legislation emphasizes early identification and immediate intervention for struggling readers. Grade 3 students must meet the MOWR cut score of 2446 on the AASA ELA Reading portion, as established by the State Board of Education, to be promoted to Grade 4, with some exemptions. Students who are retained receive an extra year of specialized support so they are ready to enter Grade 4 as strong readers. For more information, refer to the ADE website at <https://www.azed.gov/mowr/>.

REFERENCES

- American Educational Research Association (AERA), American Psychological Association (APA), & National Council on Measurement in Education (NCME). (2014). *Standards for Educational and Psychological Testing*. AERA.
- American Institutes for Research (AIR). (2015). *Recommending AzMERIT performance standards: English language arts grades 3–11, math grades 3–8, Algebra I, Geometry, and Algebra II*. https://www.azed.gov/sites/default/files/2016/12/spring-2015-azmerit-standard-setting_091415.pdf?id=5846d5b4aadebe0cf0337f5e
- Angoff, W. (1993). Perspective on differential item functioning methodology. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 3–24). Lawrence Erlbaum Associates.
- Brennan, R. L. (2004). *BB-CLASS: A computer program that uses the beta-binomial model for classification consistency and accuracy [computer software] (Version 1.0)*. University of Iowa.
- Brennan, R. L., & Prediger, D. J. (1981). Coefficient kappa: Some uses, misuses, and alternatives. *Educational and Psychological Measurement*, 41, 687–699.
- Camilli, G., & Shepard, L. A. (1994). *Methods for identifying biased test items*. Sage.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20, 37–46. <http://dx.doi.org/10.1177/001316446002000104>
- Cohen, J. (1968). Weighted kappa: Nominal scale agreement provision for scaled disagreement or partial credit. *Psychological Bulletin*, 70, 213.
- Clauser, B. E., & Mazor, K. M. (1998). Using statistical procedures to identify differentially functioning test items. *Educational Measurement: Issues and Practice*, 17, 31–44.
- Crocker, L., & Algina, J. (1986). *Introduction to classical and modern test theory*. Holt, Rinehart, and Winston.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 12, 671–684.
- Dorans, N. J., & Schmitt, A. P. (1991). *Constructed response and differential item functioning: A pragmatic approach*. ETS Research Report 91-47. Educational Testing Service.
- Egan, K. A., Schneider, C., & Ferrara, S. (2012). Performance level descriptors: History, practice, and a proposed work. In G. J. Cizek (Ed.), *Setting performance standards: Foundations, methods, and innovations* (2nd ed., pp. 79–106). Routledge.
- Embretson, S. E., & Reise, S. P. (2000). *Item response theory for psychologists*. Erlbaum.

- Glorfeld, L. W. (1995). An improvement on Horn's parallel analysis methodology for selecting the correct number of factors to retain. *Educational and Psychological Measurement*, 55, 377–393.
- Green, D. R. (1975, December). *Procedures for assessing bias in achievement tests*. Presented at the National Institute of Education Conference on Test Bias, Annapolis, MD.
- Holland, P. W., & Thayer, D. T. (1988). Differential item functioning and the Mantel-Haenszel procedure. In H. Wainer & H. I. Braun (Eds.), *Test validity* (pp. 129–145). Lawrence Erlbaum Associates.
- Horn, J. L. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika*, 30, 179–185.
- Humphreys, L. G., & Montanelli, R. G. (1975). An investigation of the parallel analysis criterion for determining the number of common factors. *Multivariate Behavioral Research*, 10, 193–206.
- Jackson, D. A. (1993). Stopping rules in principal components analysis: A comparison of heuristical and statistical approaches. *Ecology*, 74(8), 2204–2214.
- Kane, M. T. (1994). Validating interpretive arguments for licensure and certification examinations. *Evaluation & the Health Professions*, 17, 133–159.
- Lane, S., Raymond, M. R., & Haladyna, T. M. (Eds.). (2015). *Handbook of test development*. Routledge.
- Ledesma, R. D., & Valero-Mora, P. (2007). Determining the number of factors to retain in EFA: An easy-to-use computer program for carrying out parallel analysis. *Practical Assessment, Research, and Evaluation*, 12, 2.
- Linacre, J. M. (2022a). *Winsteps® Rasch measurement computer program user's guide, Version 4.8.1.0*. Winsteps.com.
- Linacre, J. M. (2022b). *Winsteps® (Version 4.8.1.0)* [Computer Software]. <http://www.winsteps.com/>
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32, 179–197.
- Mantel, N., & Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute*, 22, 719–748.
- Masters, G. N. (1982). A Rasch model for partial credit scoring. *Psychometrika*, 47, 149–174.
- Mazor, K. M., Clauser, B. E., & Hambleton, R. K. (1992). The effect of sample size on the functioning of the Mantel-Haenszel statistic. *Educational and Psychological Measurement*, 52(2), 443–451. <https://doi.org/10.1177/0013164492052002020>

- Miller, E. G., Ourania, R., & Twing, J. S. (2004). Evaluation of the 0.3 logits screening criterion in common item equating. *Journal of Applied Measurement*, 5, 172–177.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Danmarks Paedagogiske Institut.
- Stout, W. F. (1990). A new item response theory modelling approach and applications to unidimensionality assessment and ability estimation. *Psychometrika*, 55, 293–325.
- Velicer, W. F., & Jackson, D. N. (1990). Component analysis versus common factor analysis: Some issues in selecting an appropriate procedure. *Multivariate Behavioral Research*, 25(1), 1–28.
- Wright, B. D., & Linacre, J. M. (1994). Reasonable mean-square fit values. *Rasch Measurement Transactions*, 8, 370.
- Wright, B. D., & Masters, G. N. (1982). *Rating scale analysis*. 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(2), 125–145.
- Zumbo, B. D. (1999). *A handbook on the theory and methods of differential item functioning (DIF): Logistic regression modeling as a unitary framework for binary and Likert-type (ordinal) item scores*. Directorate of Human Resources Research and Evaluation, Department of National Defense.
- Zwick, R., & Velicer, W. F. (1986). Comparison of five rules for determining the number of components to retain. *Psychological Bulletin*, 99, 432–442.
- Zwick, R., Donoghue, J. R., & Grima, A. (1993). Assessment of differential item functioning for performance tasks. *Journal of Educational Measurement*, 26, 44–66.

Appendix A: ITEM-LEVEL CTT STATISTICS

This appendix includes the following item-level CTT results:

- Table A.1 – Table A.12 present the item-level CTT statistics for each content area and grade, including item type, maximum number of points possible, number of students (N), *p*-value, and the point-biserial correlation between an item and total raw score.
- Table A.13 – Table A.24 present the item-level distractor analysis for multiple-choice items, including the percentage of students who selected the correct and incorrect response options, the point-biserial correlation associated with each option, and the overall omission rate for the item.

Table A.1. Item-Level CTT Statistics, ELA Grade 3

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	OE	4	38,703	0.22	0.73
2	OE	4	38,703	0.20	0.67
3	OE	2	38,703	0.60	0.70
4	OE	4	37,390	0.28	0.72
5	OE	4	37,390	0.27	0.68
6	OE	2	37,390	0.64	0.72
7	MC	1	76,977	0.77	0.54
8	MC	1	76,962	0.45	0.40
9	MC	1	75,394	0.22	0.31
10	MX	2	76,971	0.53	0.50
11	MC	1	76,975	0.69	0.54
12	MC	1	76,939	0.62	0.59
13	MC	1	76,931	0.40	0.36
14	MC	1	76,930	0.67	0.59
15	MC	1	76,919	0.60	0.52
16	MC	1	74,936	0.32	0.41
17	MX	1	76,928	0.24	0.36
18	MC	1	76,905	0.34	0.26
19	XI	1	76,914	0.57	0.48
20	MC	1	76,886	0.67	0.57
21	MX	2	76,887	0.54	0.43
22	MX	2	76,851	0.68	0.54
23	MC	1	76,982	0.48	0.49
24	MC	1	76,973	0.80	0.51
25	MC	1	76,974	0.62	0.49
26	MX	1	76,987	0.41	0.55
27	MC	1	76,973	0.64	0.45
28	MC	1	76,959	0.38	0.42
29	MC	1	76,960	0.56	0.54
30	MX	1	76,916	0.37	0.41
31	MX	1	76,971	0.34	0.47
32	MC	1	76,962	0.53	0.36
33	MC	1	76,962	0.47	0.43

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
34	MC	1	76,951	0.47	0.36
35	MC	1	76,948	0.50	0.54
36	MC	1	76,936	0.42	0.38
37	MC	1	76,917	0.35	0.40
38	MC	1	76,912	0.40	0.28
39	MC	1	76,918	0.58	0.46
40	MC	1	76,910	0.64	0.52
41	MX	1	76,812	0.28	0.56
42	XI	1	76,920	0.38	0.60
43	MC	1	75,488	0.52	0.55
44	MC	1	76,912	0.41	0.48
45	MC	1	76,917	0.56	0.39
46	MX	2	76,923	0.57	0.66
47	MX	2	76,898	0.38	0.31

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.2. Item-Level CTT Statistics, ELA Grade 4

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	OE	4	40,146	0.33	0.76
2	OE	4	40,146	0.31	0.73
3	OE	2	40,146	0.66	0.71
4	OE	4	36,660	0.25	0.61
5	OE	4	36,660	0.23	0.68
6	OE	2	36,660	0.60	0.71
7	MC	1	76,784	0.77	0.54
8	MC	1	76,783	0.65	0.51
9	MC	1	75,989	0.58	0.56
10	MC	1	76,772	0.57	0.42
11	MC	1	76,778	0.69	0.51
12	MC	1	76,771	0.74	0.48
13	MC	1	76,764	0.50	0.37
14	MC	1	76,757	0.63	0.41
15	MC	1	76,754	0.45	0.31
16	MX	2	76,756	0.36	0.42
17	MC	1	76,755	0.31	0.17
18	MC	1	76,747	0.37	0.33
19	MC	1	76,742	0.49	0.39
20	MC	1	76,711	0.46	0.47
21	MX	2	76,722	0.56	0.53
22	MX	2	76,677	0.50	0.50
23	MC	1	76,785	0.61	0.35
24	MC	1	76,097	0.44	0.49
25	MC	1	76,774	0.46	0.32
26	MC	1	76,783	0.79	0.43

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
27	MC	1	76,772	0.62	0.42
28	MX	2	76,774	0.49	0.59
29	MC	1	76,764	0.46	0.47
30	MC	1	76,767	0.47	0.34
31	MC	1	76,762	0.38	0.44
32	MC	1	76,771	0.49	0.51
33	MC	1	76,760	0.70	0.58
34	MC	1	76,122	0.49	0.62
35	MC	1	76,758	0.59	0.48
36	MC	1	76,743	0.44	0.34
37	MC	1	76,757	0.75	0.51
38	MC	1	76,028	0.34	0.46
39	MC	1	76,747	0.72	0.55
40	MC	1	76,736	0.39	0.40
41	MC	1	75,792	0.33	0.47
42	MX	1	76,739	0.24	0.51
43	MC	1	75,784	0.36	0.57
44	MC	1	76,759	0.62	0.40
45	MC	1	76,752	0.75	0.51
46	MX	2	76,757	0.52	0.52
47	MX	2	76,754	0.72	0.51

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.3. Item-Level CTT Statistics, ELA Grade 5

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	OE	4	40,455	0.32	0.72
2	OE	4	40,455	0.30	0.72
3	OE	2	40,455	0.74	0.67
4	OE	4	36,628	0.24	0.72
5	OE	4	36,628	0.22	0.73
6	OE	2	36,628	0.73	0.66
7	MC	1	77,079	0.51	0.36
8	MC	1	76,855	0.41	0.52
9	MX	1	77,075	0.26	0.37
10	MC	1	77,062	0.40	0.42
11	XI	1	76,995	0.43	0.45
12	MX	1	77,063	0.28	0.45
13	MC	1	76,774	0.60	0.52
14	XI	1	76,975	0.45	0.51
15	MC	1	77,059	0.57	0.41
16	MC	1	77,049	0.66	0.33
17	MX	1	77,056	0.56	0.61
18	MX	1	77,050	0.40	0.42
19	MC	1	77,051	0.42	0.36

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
20	MC	1	77,053	0.67	0.56
21	MC	1	77,049	0.57	0.59
22	MC	1	76,722	0.30	0.37
23	MC	1	77,049	0.59	0.46
24	MC	1	77,039	0.59	0.42
25	MX	2	77,044	0.54	0.56
26	MX	2	77,043	0.68	0.52
27	MC	1	77,080	0.57	0.41
28	MC	1	76,912	0.48	0.50
29	MC	1	77,071	0.52	0.33
30	MC	1	77,066	0.42	0.38
31	XI	1	77,065	0.24	0.37
32	MC	1	77,074	0.60	0.38
33	MC	1	77,071	0.49	0.36
34	MC	1	77,069	0.46	0.29
35	MC	1	77,070	0.62	0.57
36	MC	1	77,069	0.48	0.35
37	MX	1	77,063	0.52	0.47
38	MC	1	77,064	0.77	0.54
39	MC	1	76,745	0.54	0.54
40	MC	1	76,880	0.50	0.62
41	MX	1	77,054	0.25	0.49
42	MC	1	77,059	0.64	0.56
43	MC	1	76,867	0.46	0.58
44	MC	1	77,060	0.56	0.45
45	MC	1	77,060	0.58	0.44
46	MX	2	77,066	0.62	0.53
47	MX	2	77,065	0.53	0.61

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.4. Item-Level CTT Statistics, ELA Grade 6

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	OE	4	40,363	0.37	0.76
2	OE	4	40,363	0.32	0.73
3	OE	2	40,363	0.74	0.70
4	OE	4	36,489	0.42	0.74
5	OE	4	36,489	0.35	0.72
6	OE	2	36,489	0.79	0.67
7	MC	1	76,839	0.78	0.47
8	MC	1	76,840	0.65	0.46
9	MC	1	76,834	0.70	0.40
10	MC	1	76,579	0.43	0.45
11	MX	1	76,831	0.47	0.48
12	MC	1	76,829	0.67	0.42

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
13	MC	1	76,823	0.47	0.36
14	MC	1	76,486	0.41	0.41
15	MC	1	76,822	0.64	0.49
16	MC	1	76,822	0.37	0.33
17	MC	1	76,825	0.74	0.52
18	MX	1	76,811	0.29	0.54
19	MC	1	76,448	0.40	0.57
20	XI	1	76,812	0.21	0.27
21	MC	1	76,810	0.57	0.36
22	MC	1	76,797	0.89	0.44
23	MX	2	76,795	0.43	0.43
24	MX	2	76,784	0.43	0.50
25	MC	1	76,843	0.55	0.37
26	MX	1	76,839	0.24	0.34
27	MC	1	76,820	0.44	0.36
28	OE	1	76,262	0.57	0.57
29	MC	1	76,830	0.49	0.38
30	MC	1	76,832	0.56	0.34
31	MC	1	76,825	0.53	0.47
32	MC	1	76,829	0.66	0.51
33	MC	1	76,828	0.47	0.37
34	MC	1	76,824	0.50	0.52
35	MC	1	76,824	0.56	0.44
36	MC	1	76,824	0.74	0.54
37	MC	1	76,821	0.48	0.46
38	MX	1	76,812	0.43	0.57
39	MC	1	76,805	0.44	0.38
40	MC	1	76,810	0.29	0.42
41	MX	1	76,812	0.51	0.66
42	MC	1	76,807	0.63	0.47
43	MX	1	76,800	0.43	0.54
44	MC	1	76,810	0.36	0.29
45	MC	1	76,808	0.75	0.59
46	MX	2	76,817	0.51	0.50
47	MX	2	76,802	0.55	0.60

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.5. Item-Level CTT Statistics, ELA Grade 7

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	OE	4	40,841	0.40	0.72
2	OE	4	40,841	0.34	0.70
3	OE	2	40,841	0.76	0.67
4	OE	4	36,970	0.41	0.73
5	OE	4	36,970	0.37	0.71

Appendix A: Item-Level CTT Statistics

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
6	OE	2	36,970	0.79	0.66
7	MC	1	77,784	0.65	0.47
8	MC	1	77,449	0.73	0.57
9	MC	1	77,791	0.48	0.34
10	MX	1	77,771	0.58	0.42
11	MC	1	77,519	0.49	0.54
12	MC	1	77,771	0.57	0.51
13	MC	1	77,779	0.50	0.29
14	MC	1	77,763	0.56	0.42
15	MC	1	77,393	0.45	0.54
16	MC	1	77,758	0.38	0.30
17	MC	1	77,766	0.61	0.40
18	MX	2	77,758	0.39	0.47
19	MC	1	77,760	0.43	0.31
20	MC	1	77,751	0.48	0.45
21	MC	1	77,759	0.56	0.39
22	MC	1	77,742	0.36	0.11
23	MX	2	77,741	0.55	0.54
24	MX	2	77,724	0.68	0.52
25	MC	1	77,776	0.78	0.47
26	MC	1	77,787	0.77	0.39
27	MC	1	77,784	0.82	0.47
28	MC	1	77,785	0.72	0.42
29	MC	1	77,783	0.61	0.40
30	MC	1	77,786	0.46	0.31
31	MC	1	77,789	0.52	0.38
32	MC	1	77,786	0.56	0.37
33	MC	1	77,784	0.60	0.46
34	MC	1	77,534	0.41	0.54
35	MC	1	77,780	0.48	0.36
36	MC	1	77,530	0.53	0.61
37	MC	1	77,767	0.47	0.45
38	MC	1	77,762	0.61	0.57
39	MC	1	77,754	0.54	0.41
40	MC	1	77,755	0.57	0.39
41	MC	1	77,756	0.63	0.54
42	MC	1	77,767	0.71	0.53
43	MC	1	77,760	0.52	0.41
44	MC	1	77,766	0.37	0.33
45	MC	1	77,762	0.49	0.36
46	MX	2	77,764	0.46	0.47
47	MX	2	77,750	0.50	0.38

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.6. Item-Level CTT Statistics, ELA Grade 8

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	OE	4	41,039	0.43	0.71
2	OE	4	41,039	0.27	0.66
3	OE	2	41,039	0.80	0.62
4	OE	4	40,248	0.37	0.75
5	OE	4	40,248	0.29	0.71
6	OE	2	40,248	0.78	0.66
7	MC	1	81,268	0.47	0.43
8	MX	1	81,279	0.37	0.54
9	MC	1	81,266	0.62	0.43
10	MC	1	81,262	0.51	0.45
11	MX	1	81,265	0.31	0.58
12	MC	1	81,261	0.60	0.43
13	MC	1	81,258	0.52	0.34
14	MC	1	81,256	0.43	0.40
15	MC	1	81,248	0.36	0.28
16	MC	1	81,249	0.53	0.47
17	MC	1	81,245	0.43	0.38
18	MC	1	81,245	0.44	0.38
19	MC	1	81,240	0.40	0.39
20	MC	1	81,245	0.42	0.41
21	XI	1	81,245	0.30	0.43
22	MC	1	81,239	0.71	0.44
23	MX	2	81,244	0.63	0.41
24	MX	2	81,226	0.70	0.59
25	MC	1	81,273	0.78	0.50
26	MX	1	81,260	0.28	0.42
27	MC	1	81,246	0.47	0.30
28	MC	1	81,243	0.44	0.50
29	MC	1	81,245	0.58	0.46
30	MC	1	81,235	0.42	0.32
31	MC	1	81,242	0.34	0.27
32	MC	1	81,249	0.73	0.49
33	MC	1	81,238	0.46	0.45
34	MX	1	81,240	0.52	0.50
35	MC	1	81,234	0.67	0.64
36	MC	1	81,035	0.46	0.57
37	MC	1	81,229	0.55	0.59
38	MC	1	81,050	0.43	0.48
39	MC	1	81,227	0.54	0.56
40	MX	1	81,223	0.35	0.50
41	MX	1	81,227	0.30	0.41
42	MC	1	81,227	0.49	0.31
43	MC	1	81,029	0.26	0.44
44	XI	1	81,194	0.18	0.33
45	MC	1	81,221	0.83	0.48

Item Number	Item Type	Max. Points	N	P-Value	Point-Biserial
46	MX	2	81,227	0.60	0.50
47	MX	2	81,218	0.67	0.63

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.7. Item-Level CTT Statistics, Mathematics Grade 3

Item Number	Item Type	Max. Points	N	P-Value	Point-Biserial
1	MC	1	78,097	0.86	0.33
2	XI	1	78,089	0.43	0.57
3	MC	1	78,089	0.38	0.47
4	XI	1	78,022	0.72	0.47
5	XI	1	78,023	0.65	0.62
6	XI	1	78,031	0.69	0.55
7	XI	1	77,932	0.68	0.51
8	MX	1	77,979	0.35	0.47
9	XI	1	78,058	0.56	0.56
10	XI	1	77,957	0.50	0.65
11	XI	1	77,690	0.39	0.50
12	XI	1	77,851	0.30	0.50
13	XI	1	77,916	0.36	0.59
14	XI	1	77,880	0.43	0.62
15	MC	1	77,983	0.42	0.41
16	MC	1	77,925	0.41	0.49
17	MC	1	77,941	0.48	0.40
18	XI	1	77,904	0.59	0.51
19	XI	1	77,930	0.61	0.64
20	XI	1	77,887	0.69	0.67
21	MC	1	77,944	0.72	0.62
22	XI	1	77,887	0.72	0.64
23	XI	1	77,866	0.69	0.65
24	MC	1	78,099	0.77	0.45
25	MC	1	78,075	0.33	0.61
26	XI	1	78,062	0.69	0.67
27	MC	1	78,047	0.63	0.46
28	XI	1	78,007	0.58	0.67
29	MC	1	78,058	0.56	0.58
30	MC	1	78,058	0.76	0.55
31	XI	1	78,008	0.49	0.39
32	XI	1	78,007	0.32	0.56
33	MC	1	78,055	0.43	0.62
34	XI	1	78,031	0.57	0.66
35	XI	1	77,999	0.47	0.66
36	XI	1	77,990	0.77	0.64
37	MC	1	78,023	0.71	0.58
38	MC	1	78,007	0.43	0.44
39	MC	1	78,034	0.85	0.43

Item Number	Item Type	Max. Points	N	P-Value	Point-Biserial
40	XI	1	78,003	0.46	0.63
41	XI	1	77,993	0.48	0.62
42	XI	1	77,982	0.68	0.64
43	XI	1	77,976	0.46	0.66
44	MC	1	78,029	0.59	0.68
45	XI	1	77,974	0.43	0.59

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.8. Item-Level CTT Statistics, Mathematics Grade 4

Item Number	Item Type	Max. Points	N	P-Value	Point-Biserial
1	XI	1	77,352	0.40	0.67
2	XI	1	77,528	0.46	0.60
3	XI	1	77,493	0.73	0.60
4	XI	1	77,477	0.26	0.54
5	MX	1	77,511	0.60	0.46
6	XI	1	77,497	0.60	0.52
7	MC	1	77,534	0.41	0.37
8	XI	1	77,533	0.68	0.53
9	MC	1	77,500	0.56	0.58
10	MC	1	77,523	0.61	0.61
11	XI	1	77,444	0.55	0.60
12	XI	1	77,356	0.40	0.58
13	XI	1	77,286	0.38	0.68
14	MC	1	77,493	0.46	0.59
15	XI	1	77,448	0.50	0.63
16	XI	1	77,321	0.37	0.60
17	MC	1	77,495	0.23	0.51
18	XI	1	77,389	0.68	0.60
19	XI	1	77,447	0.69	0.59
20	MC	1	77,465	0.66	0.61
21	XI	1	77,388	0.46	0.62
22	XI	1	77,320	0.43	0.61
23	MC	1	77,466	0.42	0.57
24	MC	1	77,562	0.80	0.49
25	XI	1	77,411	0.45	0.64
26	XI	1	77,398	0.35	0.47
27	MC	1	77,543	0.39	0.42
28	XI	1	77,479	0.47	0.63
29	MC	1	77,528	0.68	0.51
30	XI	1	77,482	0.42	0.62
31	XI	1	77,466	0.61	0.65
32	MC	1	77,542	0.36	0.35
33	XI	1	77,460	0.42	0.62
34	XI	1	77,516	0.80	0.60
35	MC	1	77,482	0.58	0.42

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
36	XI	1	77,464	0.58	0.54
37	XI	1	77,446	0.68	0.62
38	MC	1	77,511	0.49	0.41
39	XI	1	77,395	0.37	0.52
40	XI	1	77,489	0.78	0.62
41	MC	1	77,527	0.77	0.43
42	MX	1	77,458	0.41	0.63
43	XI	1	77,443	0.50	0.52
44	MC	1	77,515	0.52	0.44
45	XI	1	77,503	0.36	0.66

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.9. Item-Level CTT Statistics, Mathematics Grade 5

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	XI	1	77,565	0.59	0.42
2	XI	1	77,534	0.55	0.59
3	MC	1	77,580	0.53	0.44
4	MX	1	77,585	0.31	0.50
5	XI	1	77,417	0.48	0.67
6	MX	1	77,578	0.39	0.36
7	XI	1	77,546	0.59	0.62
8	XI	1	77,469	0.49	0.62
9	MC	1	77,576	0.44	0.33
10	XI	1	77,503	0.46	0.66
11	XI	1	77,447	0.25	0.61
12	XI	1	77,510	0.35	0.63
13	XI	1	77,446	0.31	0.63
14	XI	1	77,447	0.47	0.59
15	XI	1	77,425	0.27	0.59
16	XI	1	77,398	0.44	0.70
17	XI	1	77,463	0.45	0.60
18	XI	1	77,469	0.50	0.54
19	MC	1	77,512	0.64	0.51
20	XI	1	77,327	0.49	0.68
21	MC	1	77,512	0.50	0.38
22	MC	1	77,502	0.31	0.35
23	MC	1	77,521	0.61	0.44
24	MC	1	77,595	0.70	0.37
25	MC	1	77,592	0.42	0.65
26	MC	1	77,590	0.61	0.37
27	XI	1	77,551	0.51	0.53
28	XI	1	77,559	0.43	0.68
29	MC	1	77,595	0.56	0.53
30	XI	1	77,533	0.58	0.39
31	XI	1	77,338	0.16	0.44

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
32	MC	1	77,563	0.34	0.45
33	XI	1	77,511	0.33	0.62
34	XI	1	77,565	0.32	0.50
35	XI	1	77,464	0.18	0.54
36	MC	1	77,564	0.50	0.48
37	MX	1	77,564	0.54	0.42
38	XI	1	77,491	0.37	0.61
39	MC	1	77,565	0.44	0.36
40	MC	1	77,569	0.58	0.53
41	XI	1	77,427	0.49	0.67
42	XI	1	77,507	0.39	0.50
43	XI	1	77,528	0.58	0.66
44	XI	1	77,504	0.58	0.60
45	XI	1	77,495	0.38	0.59

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.10. Item-Level CTT Statistics, Mathematics Grade 6

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	XI	1	77,416	0.45	0.59
2	MC	1	77,484	0.52	0.45
3	MC	1	77,494	0.64	0.26
4	XI	1	77,325	0.40	0.28
5	XI	1	77,449	0.49	0.65
6	XI	1	77,345	0.40	0.58
7	MC	1	77,461	0.27	0.63
8	XI	1	77,329	0.40	0.58
9	XI	1	77,291	0.36	0.59
10	XI	1	77,323	0.35	0.62
11	XI	1	77,378	0.30	0.64
12	XI	1	77,365	0.45	0.65
13	XI	1	77,331	0.36	0.70
14	XI	1	77,298	0.23	0.57
15	MC	1	77,424	0.38	0.27
16	XI	1	77,255	0.18	0.47
17	MC	1	77,432	0.36	0.30
18	XI	1	77,375	0.47	0.60
19	MC	1	77,379	0.41	0.42
20	XI	1	77,272	0.40	0.53
21	MC	1	77,410	0.43	0.28
22	XI	1	77,123	0.42	0.64
23	MC	1	77,415	0.36	0.46
24	XI	1	77,362	0.60	0.66
25	XI	1	77,413	0.24	0.61
26	MC	1	77,479	0.45	0.33
27	XI	1	77,364	0.52	0.66

Item Number	Item Type	Max. Points	N	P-Value	Point-Biserial
28	MC	1	77,491	0.59	0.50
29	MC	1	77,456	0.45	0.57
30	XI	1	77,388	0.50	0.28
31	XI	1	77,372	0.53	0.45
32	MC	1	77,473	0.36	0.49
33	XI	1	77,409	0.44	0.58
34	MC	1	77,449	0.34	0.63
35	XI	1	77,428	0.45	0.53
36	XI	1	77,361	0.36	0.44
37	XI	1	77,348	0.48	0.65
38	MC	1	77,448	0.37	0.51
39	XI	1	77,395	0.56	0.46
40	MC	1	77,433	0.29	0.49
41	XI	1	77,333	0.27	0.63
42	MC	1	77,433	0.42	0.47
43	XI	1	77,171	0.34	0.60
44	XI	1	77,343	0.41	0.58
45	XI	1	77,345	0.51	0.39
46	XI	1	77,290	0.42	0.63
47	MC	1	77,446	0.62	0.45

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.11. Item-Level CTT Statistics, Mathematics Grade 7

Item Number	Item Type	Max. Points	N	P-Value	Point-Biserial
1	MC	1	78,477	0.54	0.36
2	XI	1	78,214	0.29	0.56
3	XI	1	78,327	0.35	0.59
4	XI	1	78,326	0.44	0.65
5	MC	1	78,466	0.65	0.49
6	MC	1	78,463	0.22	0.25
7	MC	1	78,428	0.53	0.37
8	XI	1	78,334	0.50	0.68
9	XI	1	78,135	0.39	0.65
10	MC	1	78,432	0.36	0.34
11	XI	1	78,299	0.48	0.62
12	XI	1	78,217	0.37	0.52
13	XI	1	78,338	0.55	0.65
14	XI	1	78,231	0.18	0.45
15	XI	1	78,300	0.53	0.68
16	XI	1	78,278	0.26	0.64
17	MC	1	78,394	0.24	0.38
18	XI	1	78,390	0.39	0.58
19	MC	1	78,409	0.70	0.59
20	MC	1	78,414	0.47	0.30
21	XI	1	78,279	0.49	0.69

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
22	XI	1	78,246	0.30	0.58
23	XI	1	78,328	0.61	0.64
24	MC	1	78,484	0.81	0.43
25	XI	1	78,426	0.30	0.60
26	XI	1	78,225	0.35	0.66
27	XI	1	78,262	0.21	0.58
28	XI	1	78,220	0.34	0.50
29	XI	1	78,376	0.56	0.66
30	MC	1	78,473	0.42	0.42
31	MC	1	78,468	0.44	0.52
32	XI	1	78,383	0.26	0.65
33	XI	1	78,173	0.30	0.56
34	XI	1	78,302	0.27	0.64
35	XI	1	78,295	0.20	0.52
36	MC	1	78,440	0.34	0.47
37	XI	1	78,426	0.27	0.48
38	MC	1	78,452	0.56	0.53
39	XI	1	78,303	0.29	0.57
40	XI	1	78,145	0.20	0.63
41	XI	1	78,261	0.22	0.65
42	MC	1	78,459	0.73	0.27
43	MC	1	78,450	0.58	0.39
44	MC	1	78,442	0.53	0.40
45	XI	1	78,324	0.30	0.65
46	MC	1	78,450	0.70	0.57
47	XI	1	78,270	0.25	0.60

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.12. Item-Level CTT Statistics, Mathematics Grade 8

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
1	MC	1	82,070	0.68	0.40
2	MC	1	82,061	0.39	0.37
3	XI	1	81,794	0.38	0.62
4	MC	1	82,044	0.24	0.43
5	MC	1	82,052	0.40	0.40
6	XI	1	81,600	0.25	0.64
7	MC	1	82,044	0.53	0.50
8	MC	1	82,027	0.28	0.33
9	MC	1	82,032	0.23	0.36
10	MC	1	82,032	0.78	0.40
11	MC	1	82,041	0.27	0.29
12	MX	1	82,019	0.39	0.33
13	XI	1	81,779	0.13	0.52
14	XI	1	81,696	0.29	0.64
15	MC	1	82,009	0.37	0.36

Appendix A: Item-Level CTT Statistics

Item Number	Item Type	Max. Points	N	<i>P</i> -Value	Point-Biserial
16	MX	1	82,002	0.40	0.45
17	XI	1	81,801	0.15	0.53
18	XI	1	81,994	0.19	0.52
19	MC	1	82,012	0.48	0.41
20	XI	1	81,845	0.24	0.39
21	MC	1	82,015	0.61	0.45
22	MC	1	82,013	0.48	0.48
23	MC	1	81,996	0.40	0.52
24	MC	1	82,057	0.53	0.59
25	XI	1	82,056	0.48	0.32
26	MC	1	82,035	0.44	0.50
27	MC	1	82,045	0.43	0.32
28	MC	1	82,033	0.42	0.36
29	MC	1	82,039	0.53	0.46
30	MC	1	82,035	0.32	0.44
31	MX	1	82,031	0.39	0.50
32	MC	1	82,035	0.56	0.43
33	MC	1	82,040	0.46	0.60
34	MC	1	82,034	0.37	0.50
35	MC	1	82,022	0.35	0.39
36	XI	1	81,976	0.42	0.47
37	MC	1	82,019	0.54	0.32
38	MC	1	82,021	0.41	0.49
39	XI	1	81,977	0.38	0.49
40	XI	1	81,619	0.37	0.71
41	XI	1	81,639	0.37	0.65
42	MC	1	82,013	0.58	0.35
43	XI	1	81,703	0.32	0.62
44	MC	1	81,997	0.35	0.63
45	XI	1	81,829	0.37	0.67
46	MC	1	82,015	0.59	0.28
47	MC	1	82,005	0.47	0.50

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.13. Distractor Analysis of Multiple-Choice Items, ELA Grade 3

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
7	76.61	0.54	6.97	-0.28	9.21	-0.33	7.21	-0.23
8	44.78	0.40	21.00	-0.09	18.29	-0.21	15.94	-0.24
11	68.69	0.54	8.34	-0.26	11.03	-0.27	11.94	-0.29
12	61.61	0.59	14.92	-0.33	11.19	-0.31	12.28	-0.21
13	40.36	0.36	19.99	-0.19	18.19	-0.18	21.47	-0.07
14	67.27	0.59	11.77	-0.28	13.37	-0.32	7.60	-0.28
15	59.72	0.52	13.52	-0.22	17.45	-0.27	9.30	-0.27
18	34.45	0.26	27.97	-0.06	16.88	-0.16	20.70	-0.10
20	66.78	0.57	10.70	-0.25	10.72	-0.27	11.80	-0.33
23	47.90	0.49	16.46	-0.27	24.05	-0.16	11.59	-0.24
24	79.80	0.51	7.17	-0.27	7.20	-0.30	5.83	-0.26
25	62.15	0.49	18.25	-0.17	14.36	-0.35	5.25	-0.24
27	64.42	0.45	12.57	-0.19	17.47	-0.26	5.53	-0.24
28	38.30	0.42	13.27	-0.19	25.16	-0.16	23.28	-0.16
29	56.11	0.54	19.92	-0.21	10.75	-0.30	13.22	-0.27
32	53.19	0.36	17.52	-0.21	11.13	-0.27	18.16	-0.04
33	46.68	0.43	10.98	-0.31	12.00	-0.30	30.34	-0.04
34	46.88	0.36	12.25	-0.30	27.37	-0.03	13.50	-0.19
35	50.18	0.54	23.28	-0.32	13.11	-0.21	13.43	-0.19
36	41.82	0.38	21.12	-0.12	25.46	-0.19	11.60	-0.17
37	35.01	0.40	27.25	-0.06	18.89	-0.20	18.85	-0.22
38	39.64	0.28	27.91	-0.11	19.04	-0.10	13.41	-0.15
39	57.59	0.46	14.55	-0.27	16.49	-0.22	11.37	-0.16
40	64.21	0.52	13.06	-0.19	14.80	-0.32	7.94	-0.26
44	41.28	0.48	32.82	-0.08	15.40	-0.27	10.51	-0.33
45	55.52	0.39	12.67	-0.29	10.83	-0.26	20.99	-0.03

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.14. Distractor Analysis of Multiple-Choice Items, ELA Grade 4

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
7	76.50	0.54	8.63	-0.32	10.59	-0.34	4.28	-0.18
8	65.19	0.51	8.88	-0.21	13.54	-0.33	12.39	-0.21
10	56.67	0.42	15.94	-0.24	17.83	-0.20	9.56	-0.13
11	69.36	0.51	12.38	-0.28	9.74	-0.24	8.52	-0.25
12	74.26	0.48	4.81	-0.24	13.49	-0.27	7.44	-0.26
13	50.27	0.37	11.42	-0.21	22.02	-0.15	16.29	-0.16
14	62.84	0.41	15.11	-0.22	12.88	-0.21	9.17	-0.17
15	45.07	0.31	16.70	-0.19	16.53	-0.22	21.70	0.00
17	31.06	0.17	18.88	-0.03	36.44	-0.06	13.62	-0.11
18	36.59	0.33	18.90	-0.05	34.73	-0.15	9.79	-0.24
19	49.21	0.39	21.44	-0.23	13.90	-0.28	15.45	-0.01
20	45.83	0.47	23.93	-0.19	9.17	-0.32	21.07	-0.15
23	60.92	0.35	5.29	-0.23	29.37	-0.18	4.43	-0.18
25	45.90	0.32	12.59	-0.20	30.74	-0.09	10.76	-0.18
26	78.89	0.43	2.73	-0.17	4.78	-0.25	13.60	-0.27
27	62.36	0.42	10.04	-0.31	17.35	-0.15	10.25	-0.17
29	46.36	0.47	14.38	-0.16	15.64	-0.27	23.62	-0.19
30	47.46	0.34	12.52	-0.25	24.39	-0.10	15.63	-0.13
31	37.97	0.44	15.28	-0.16	31.24	-0.21	15.51	-0.16
32	49.47	0.51	19.03	-0.19	17.04	-0.26	14.45	-0.24
33	69.71	0.58	9.24	-0.30	14.31	-0.33	6.74	-0.25
35	59.31	0.48	19.00	-0.16	14.84	-0.31	6.85	-0.24
36	44.29	0.34	13.95	-0.15	25.88	-0.17	15.88	-0.12
37	74.50	0.51	12.02	-0.27	8.87	-0.30	4.60	-0.22
39	71.63	0.55	9.61	-0.28	11.40	-0.36	7.36	-0.20
40	38.79	0.40	28.17	-0.27	20.82	-0.10	12.23	-0.11
44	62.37	0.40	19.08	-0.13	9.41	-0.32	9.15	-0.18
45	75.06	0.51	10.20	-0.32	7.60	-0.29	7.14	-0.18

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.15. Distractor Analysis of Multiple-Choice Items, ELA Grade 5

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
7	50.66	0.36	10.99	-0.27	16.87	-0.24	21.48	-0.02
10	40.43	0.42	20.59	-0.17	28.25	-0.18	10.73	-0.18
15	57.47	0.41	13.75	-0.16	14.51	-0.31	14.27	-0.11
16	65.77	0.33	13.53	-0.09	13.40	-0.23	7.30	-0.17
19	42.33	0.36	11.89	-0.36	28.15	-0.10	17.62	-0.05
20	66.59	0.56	11.40	-0.29	17.86	-0.33	4.15	-0.24
21	56.70	0.59	15.98	-0.26	13.08	-0.31	14.23	-0.26
23	58.71	0.46	9.96	-0.30	10.79	-0.33	20.55	-0.09
24	58.94	0.42	7.47	-0.27	12.74	-0.38	20.85	-0.03
27	56.53	0.41	25.85	-0.17	11.00	-0.21	6.62	-0.25
29	52.30	0.33	17.84	-0.22	17.71	-0.15	12.15	-0.07
30	41.81	0.38	24.30	-0.11	17.72	-0.23	16.17	-0.14
32	60.19	0.38	11.13	-0.12	8.74	-0.33	19.94	-0.14
33	49.18	0.36	9.91	-0.31	32.10	-0.08	8.81	-0.18
34	46.08	0.29	15.71	-0.19	14.84	-0.29	23.36	0.07
35	61.88	0.57	15.01	-0.29	16.78	-0.32	6.33	-0.22
36	48.39	0.35	18.90	-0.03	17.33	-0.23	15.38	-0.21
38	77.30	0.54	5.61	-0.24	11.92	-0.36	5.17	-0.25
42	63.57	0.56	9.87	-0.18	16.13	-0.27	10.43	-0.38
44	55.82	0.45	12.75	-0.36	15.41	-0.24	16.02	-0.05
45	58.16	0.44	14.34	-0.23	8.76	-0.34	18.75	-0.10

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.16. Distractor Analysis of Multiple-Choice Items, ELA Grade 6

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
7	78.32	0.47	10.81	-0.28	6.40	-0.30	4.47	-0.16
8	64.95	0.46	13.78	-0.29	10.69	-0.13	10.58	-0.26
9	69.77	0.40	5.35	-0.21	8.15	-0.31	16.74	-0.14
12	66.80	0.42	11.11	-0.14	10.87	-0.26	11.22	-0.23
13	47.41	0.36	18.65	-0.18	23.36	-0.13	10.59	-0.17
15	63.98	0.49	6.87	-0.30	13.58	-0.28	15.57	-0.18
16	37.41	0.33	11.76	-0.10	18.71	-0.21	32.13	-0.09
17	73.77	0.52	8.90	-0.22	8.64	-0.31	8.69	-0.29
21	57.02	0.36	12.69	-0.30	9.16	-0.33	21.13	0.05
22	88.98	0.44	3.94	-0.25	3.39	-0.26	3.69	-0.23
25	55.35	0.37	15.71	-0.11	23.59	-0.23	5.35	-0.20
27	44.43	0.36	14.92	-0.20	34.05	-0.10	6.59	-0.25
29	49.27	0.38	14.33	-0.15	14.44	-0.28	21.96	-0.10
30	55.68	0.34	7.54	-0.24	25.85	-0.11	10.93	-0.18
31	52.95	0.47	9.19	-0.25	23.22	-0.21	14.65	-0.20
32	66.32	0.51	15.73	-0.24	8.54	-0.25	9.40	-0.29
33	46.61	0.37	18.69	-0.22	14.17	-0.15	20.53	-0.12
34	49.57	0.52	13.58	-0.31	14.75	-0.32	22.10	-0.09
35	56.07	0.44	8.18	-0.23	24.17	-0.17	11.58	-0.25
36	74.07	0.54	8.36	-0.28	8.81	-0.32	8.76	-0.25
37	48.09	0.46	17.94	-0.18	22.09	-0.29	11.87	-0.14
39	44.23	0.38	13.18	-0.34	20.66	-0.22	21.93	0.03
42	63.15	0.47	10.22	-0.28	20.74	-0.21	5.89	-0.25
44	36.39	0.29	9.16	-0.05	41.73	-0.20	12.73	-0.07
45	74.70	0.59	10.75	-0.39	6.11	-0.32	8.43	-0.22

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.17. Distractor Analysis of Multiple-Choice Items, ELA Grade 7

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
7	64.60	0.47	8.59	-0.28	7.55	-0.27	19.27	-0.19
9	47.84	0.34	26.42	-0.21	11.30	-0.10	14.44	-0.13
12	56.88	0.51	18.48	-0.30	16.96	-0.30	7.68	-0.09
13	49.76	0.29	17.40	0.04	17.20	-0.11	15.64	-0.32
14	55.59	0.42	13.25	-0.21	20.94	-0.17	10.22	-0.21
16	37.96	0.30	15.03	-0.17	33.64	-0.02	13.37	-0.22
17	60.85	0.40	10.96	-0.22	13.02	-0.28	15.18	-0.10
19	42.63	0.31	27.49	-0.05	14.82	-0.22	15.06	-0.14
20	48.21	0.45	11.68	-0.23	13.43	-0.28	26.68	-0.13
21	56.32	0.39	15.97	-0.15	17.35	-0.22	10.37	-0.17
22	36.36	0.11	18.69	-0.09	21.27	-0.10	23.67	0.06
25	77.54	0.47	8.07	-0.27	6.67	-0.30	7.73	-0.19
26	77.27	0.39	11.19	-0.19	6.78	-0.30	4.76	-0.15
27	81.58	0.47	7.24	-0.24	5.97	-0.27	5.22	-0.26
28	71.74	0.42	5.35	-0.28	4.90	-0.26	18.01	-0.19
29	61.24	0.41	10.79	-0.23	12.18	-0.23	15.79	-0.14
30	45.91	0.31	18.96	-0.20	9.45	-0.22	25.68	-0.03
31	52.02	0.38	15.63	-0.18	17.45	-0.22	14.90	-0.11
32	55.98	0.37	19.20	-0.12	11.65	-0.21	13.17	-0.20
33	59.92	0.46	12.85	-0.25	10.90	-0.24	16.32	-0.18
35	47.97	0.36	18.98	-0.13	10.55	-0.26	22.50	-0.12
37	46.60	0.45	16.80	-0.11	17.01	-0.20	19.58	-0.27
38	61.43	0.57	16.44	-0.30	17.81	-0.34	4.33	-0.19
39	54.48	0.41	14.59	-0.10	20.04	-0.27	10.89	-0.19
40	57.05	0.39	10.05	-0.17	17.60	-0.19	15.30	-0.20
41	62.76	0.54	11.37	-0.28	14.01	-0.26	11.85	-0.25
42	70.96	0.53	12.01	-0.30	9.79	-0.26	7.24	-0.25
43	51.87	0.41	20.33	-0.16	10.88	-0.23	16.92	-0.19
44	36.82	0.33	25.04	-0.16	22.98	-0.16	15.16	-0.05
45	48.59	0.36	18.47	-0.22	8.44	-0.29	24.50	-0.04

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.18. Distractor Analysis of Multiple-Choice Items, ELA Grade 8

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
7	46.77	0.43	14.87	-0.19	17.86	-0.30	20.50	-0.08
9	61.53	0.43	14.93	-0.16	9.86	-0.25	13.67	-0.21
10	50.84	0.45	21.22	-0.20	18.15	-0.22	9.79	-0.19
12	59.81	0.43	9.07	-0.25	23.12	-0.16	7.99	-0.26
13	51.61	0.34	26.37	-0.06	14.48	-0.29	7.55	-0.15
14	43.20	0.40	29.32	-0.14	13.84	-0.18	13.64	-0.21
15	35.56	0.28	18.67	-0.10	25.47	-0.12	20.29	-0.10
16	52.99	0.47	12.88	-0.21	15.15	-0.23	18.98	-0.22
17	43.37	0.38	18.87	-0.17	23.93	-0.09	13.83	-0.24
18	44.07	0.38	19.48	-0.09	21.67	-0.20	14.79	-0.21
19	40.32	0.39	10.35	-0.25	31.89	-0.15	17.43	-0.11
20	41.78	0.41	16.39	-0.27	21.01	-0.16	20.81	-0.09
22	71.37	0.44	19.65	-0.24	6.08	-0.31	2.90	-0.19
25	77.89	0.50	9.28	-0.29	6.07	-0.28	6.76	-0.22
27	46.68	0.30	9.09	-0.17	16.81	-0.20	27.41	-0.05
28	44.33	0.50	16.63	-0.26	22.03	-0.19	17.02	-0.20
29	58.47	0.46	9.33	-0.20	14.82	-0.20	17.38	-0.25
30	41.78	0.32	13.25	-0.23	21.11	-0.13	23.86	-0.06
31	34.08	0.27	11.04	-0.19	28.11	-0.09	26.76	-0.06
32	73.26	0.49	12.29	-0.25	8.85	-0.31	5.60	-0.21
33	45.93	0.45	12.62	-0.34	17.66	-0.24	23.79	-0.05
35	67.19	0.64	10.21	-0.25	10.27	-0.36	12.34	-0.35
37	55.39	0.59	9.99	-0.25	19.31	-0.33	15.30	-0.25
39	54.33	0.56	19.10	-0.30	13.25	-0.34	13.32	-0.14
42	49.00	0.31	15.54	-0.10	14.37	-0.28	21.09	-0.05
45	82.83	0.48	8.79	-0.31	5.74	-0.31	2.63	-0.15

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.19. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 3

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
1	86.41	0.33	5.44	-0.23	2.61	-0.16	5.55	-0.15
3	38.17	0.47	37.15	-0.42	23.42	-0.03	1.26	-0.14
15	42.04	0.41	28.95	-0.25	4.77	-0.21	24.23	-0.10
16	40.59	0.49	20.62	-0.27	23.74	-0.16	15.05	-0.17
17	47.77	0.40	8.06	-0.29	16.70	-0.25	27.46	-0.06
21	72.17	0.62	12.66	-0.30	9.86	-0.37	5.31	-0.29
24	76.59	0.45	13.17	-0.33	6.16	-0.18	4.08	-0.18
27	62.83	0.46	17.62	-0.24	12.55	-0.21	7.00	-0.23
29	56.31	0.58	22.52	-0.35	15.84	-0.30	5.33	-0.13
30	76.22	0.55	9.00	-0.35	10.16	-0.27	4.62	-0.24
37	71.12	0.58	10.36	-0.36	10.75	-0.36	7.76	-0.15
39	84.65	0.43	7.14	-0.29	3.87	-0.20	4.34	-0.20

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.20. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 4

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
7	41.18	0.37	27.76	-0.19	15.50	-0.18	15.57	-0.08
9	55.51	0.58	14.59	-0.15	22.20	-0.41	7.70	-0.24
14	46.42	0.59	23.66	-0.33	17.98	-0.24	11.94	-0.19
20	66.04	0.61	9.75	-0.36	11.02	-0.33	13.18	-0.23
24	79.66	0.49	16.96	-0.42	2.37	-0.19	1.01	-0.11
27	39.35	0.42	15.66	-0.15	16.36	-0.30	28.63	-0.10
29	68.49	0.51	7.06	-0.28	17.01	-0.27	7.45	-0.24
32	35.81	0.35	10.85	-0.21	10.34	-0.16	43.00	-0.11
35	58.28	0.42	11.49	-0.18	22.54	-0.23	7.68	-0.19
38	49.13	0.41	32.26	-0.12	13.19	-0.29	5.42	-0.22
41	77.16	0.43	7.52	-0.27	8.69	-0.26	6.64	-0.14
44	52.06	0.44	10.75	-0.09	31.60	-0.27	5.58	-0.30

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.21. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 5

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
3	52.77	0.44	17.58	-0.18	21.93	-0.26	7.72	-0.16
9	43.93	0.33	21.64	-0.13	19.66	-0.16	14.77	-0.14
19	64.15	0.51	14.18	-0.26	9.69	-0.30	11.99	-0.21
21	49.92	0.38	27.97	-0.37	7.64	-0.21	14.46	0.08
22	30.91	0.35	42.31	-0.16	14.85	-0.26	11.93	0.04
24	69.60	0.37	25.71	-0.28	3.25	-0.18	1.44	-0.13
26	61.02	0.37	5.15	-0.21	10.41	-0.18	23.41	-0.19
36	49.66	0.48	28.22	-0.23	12.82	-0.21	9.30	-0.22
39	44.01	0.36	36.68	-0.09	10.51	-0.29	8.80	-0.17
40	58.11	0.53	8.36	-0.23	21.58	-0.28	11.95	-0.25

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.22. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 6

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
2	52.32	0.45	30.76	-0.25	11.44	-0.27	5.47	-0.10
3	63.62	0.26	20.18	-0.16	6.02	-0.16	10.19	-0.08
15	37.53	0.27	16.43	-0.09	26.11	-0.20	19.93	-0.03
17	36.08	0.30	24.85	-0.03	22.73	-0.31	16.33	-0.01
19	41.07	0.42	23.68	-0.28	26.23	-0.14	9.01	-0.09
21	42.80	0.28	25.58	-0.10	18.52	-0.22	13.10	-0.03
23	36.12	0.46	40.12	-0.15	20.87	-0.29	2.89	-0.16
26	45.16	0.33	23.64	-0.01	18.24	-0.29	12.95	-0.15
28	59.21	0.50	28.13	-0.32	11.40	-0.28	1.25	-0.12
29	44.86	0.57	22.19	-0.30	20.37	-0.24	12.57	-0.19
38	36.58	0.51	16.88	-0.27	14.75	-0.29	31.79	-0.09
42	42.38	0.47	16.26	-0.10	30.07	-0.32	11.30	-0.15
47	61.95	0.45	19.75	-0.28	12.54	-0.22	5.77	-0.15

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.23. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 7

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
1	53.71	0.36	14.61	-0.17	19.45	-0.17	12.23	-0.16
5	65.10	0.49	17.52	-0.26	14.01	-0.30	3.37	-0.16
6	21.93	0.25	32.98	-0.14	26.74	-0.02	18.35	-0.07
7	53.21	0.37	13.63	-0.13	17.06	-0.20	16.09	-0.18
10	35.95	0.34	13.39	-0.10	32.91	-0.26	17.76	0.00
17	24.04	0.38	28.13	-0.15	30.10	-0.19	17.72	-0.02
19	70.04	0.59	14.75	-0.35	10.63	-0.33	4.58	-0.21
20	47.17	0.30	7.00	-0.13	28.88	-0.22	16.96	-0.05
24	81.08	0.43	12.35	-0.31	3.46	-0.20	3.12	-0.16
30	41.99	0.42	15.30	-0.15	19.30	-0.17	23.41	-0.21
31	44.34	0.52	12.76	-0.27	18.26	-0.16	24.64	-0.24
36	34.10	0.47	11.30	-0.20	34.95	-0.19	19.66	-0.17
38	56.10	0.53	21.46	-0.24	15.27	-0.30	7.17	-0.21
42	72.83	0.27	1.57	-0.08	5.90	-0.15	19.69	-0.19
43	57.70	0.39	13.59	-0.23	20.56	-0.17	8.14	-0.17
44	53.23	0.40	11.34	-0.22	24.17	-0.19	11.26	-0.15
46	70.42	0.57	15.23	-0.33	7.13	-0.28	7.23	-0.27

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table A.24. Distractor Analysis of Multiple-Choice Items, Mathematics Grade 8

Item Number	Correct Option		Distractor 1		Distractor 2		Distractor 3	
	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.	%	Pt. Bis.
1	67.73	0.40	10.68	-0.24	8.94	-0.20	12.65	-0.17
2	39.36	0.37	26.96	-0.12	16.01	-0.21	17.67	-0.12
4	23.73	0.43	24.29	-0.05	24.05	-0.08	27.93	-0.28
5	40.19	0.40	15.39	-0.18	14.05	-0.21	30.37	-0.13
8	27.84	0.33	39.88	-0.17	27.97	-0.14	4.31	-0.01
9	23.03	0.36	39.66	-0.07	17.92	-0.17	19.38	-0.13
10	78.41	0.40	8.33	-0.23	6.54	-0.19	6.71	-0.21
11	26.67	0.29	11.59	-0.12	44.39	-0.18	17.35	0.00
15	37.38	0.36	30.25	-0.16	20.19	-0.20	12.17	-0.06
21	60.90	0.45	17.38	-0.23	12.53	-0.26	9.19	-0.16
22	47.65	0.48	22.26	-0.17	16.67	-0.26	13.42	-0.22
23	40.30	0.52	25.33	-0.26	20.59	-0.22	13.78	-0.17
24	52.63	0.59	9.33	-0.25	18.14	-0.30	19.90	-0.26
26	44.04	0.50	7.86	-0.19	23.39	-0.26	24.71	-0.21
27	42.76	0.32	13.68	-0.19	28.99	-0.16	14.57	-0.06
28	41.89	0.36	8.02	-0.12	30.17	-0.26	19.92	-0.07
29	52.51	0.46	12.97	-0.16	15.10	-0.24	19.42	-0.23
30	31.67	0.44	19.88	-0.21	25.27	-0.12	23.18	-0.16
32	55.72	0.43	13.91	-0.25	23.93	-0.19	6.44	-0.18
34	36.60	0.50	21.66	-0.19	20.91	-0.21	20.82	-0.19
35	35.03	0.39	25.82	-0.22	30.54	-0.19	8.61	-0.01
37	53.69	0.32	16.02	-0.09	18.73	-0.19	11.56	-0.17
38	41.14	0.49	14.67	-0.14	21.85	-0.19	22.34	-0.27
42	57.76	0.35	17.92	-0.21	14.88	-0.11	9.45	-0.18
44	34.86	0.63	19.39	-0.18	23.59	-0.28	22.16	-0.26
46	59.13	0.29	11.55	-0.06	15.08	-0.23	14.24	-0.11
47	46.92	0.50	19.32	-0.21	25.30	-0.31	8.46	-0.11

Note. The item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Appendix B: ITEM-LEVEL IRT STATISTICS

This appendix includes the following item-level IRT statistics:

- Table B.1 – Table B.12 present the IRT statistics, including item type, Rasch difficulty, standard error (SE) of Rasch, and infit values.
- Table B.13 – Table B.24 present the raw-to-scale score conversion tables.
- Figure B.1 – Figure B.18 present the item-person map for each post-equated operational form.
- Figure B.19 – Figure B.54 present the test characteristic curve (TCC) and conditional standard error of measurement (CSEM) curve for each post-equated operational form.
- Figure B.55 – Figure B.72 present the scree plot from the principal component analysis (PCA) for each operational form. The scree plot shows only the first 10 components.

Table B.1. Item-Level IRT Statistics, ELA Grade 3

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	OE	2.3017	0.0087	0.74
2	OE	2.4937	0.0094	0.82
3	OE	-0.5632	0.0079	0.83
4	OE	2.3061	0.0088	0.78
5	OE	2.5353	0.0090	0.83
6	OE	-0.7491	0.0083	0.76
7	MC	-1.5835	0.0094	0.83
8	MC	0.1903	0.0083	1.09
9	MC	1.5322	0.0097	1.12
10	MX	-0.2552	0.0060	1.18
11	MC	-1.0901	0.0087	0.87
12	MC	-0.6569	0.0084	0.84
13	MC	0.2279	0.0083	1.13
14	MC	-0.8899	0.0085	0.81
15	MC	-0.5722	0.0083	0.93
16	MC	1.0918	0.0090	1.08
17	MX	1.5143	0.0096	1.11
18	MC	0.7924	0.0086	1.23
19	XI	-0.3018	0.0082	0.99
20	MC	-1.0014	0.0086	0.85
21	MX	-0.3143	0.0066	1.17
22	MX	-1.1192	0.0062	1.05
23	MC	0.0264	0.0082	0.99
24	MC	-1.8086	0.0098	0.84
25	MC	-0.7228	0.0084	0.96
26	MX	0.3883	0.0083	0.92
27	MC	-0.8480	0.0085	1.00
28	MC	0.5370	0.0084	1.06
29	MC	-0.4010	0.0083	0.92
30	MX	0.6011	0.0085	1.07

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
31	MX	0.7871	0.0086	0.99
32	MC	-0.2485	0.0082	1.15
33	MC	0.0905	0.0082	1.06
34	MC	-0.1070	0.0082	1.16
35	MC	-0.0919	0.0082	0.92
36	MC	0.4854	0.0084	1.13
37	MC	0.7947	0.0086	1.08
38	MC	0.3615	0.0083	1.23
39	MC	-0.6059	0.0083	1.02
40	MC	-0.8358	0.0085	0.93
41	MX	1.1330	0.0090	0.86
42	XI	0.5727	0.0085	0.82
43	MC	-0.1811	0.0083	0.90
44	MC	0.3763	0.0083	0.99
45	MC	-0.3706	0.0083	1.11
46	MX	-0.4582	0.0060	0.88
47	MX	0.7709	0.0068	1.35

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.2. Item-Level IRT Statistics, ELA Grade 4

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	OE	1.6473	0.0074	0.73
2	OE	1.7245	0.0075	0.79
3	OE	-0.3533	0.0079	0.73
4	OE	2.1945	0.0086	0.96
5	OE	2.3682	0.0096	0.79
6	OE	-0.1034	0.0084	0.75
7	MC	-1.0433	0.0093	0.82
8	MC	-0.3735	0.0084	0.91
9	MC	0.0219	0.0082	0.87
10	MC	0.0754	0.0082	1.04
11	MC	-0.6068	0.0086	0.90
12	MC	-0.8981	0.0090	0.89
13	MC	0.3901	0.0081	1.10
14	MC	-0.2710	0.0083	1.04
15	MC	0.8749	0.0083	1.21
16	MX	0.9695	0.0055	1.47
17	MC	1.3307	0.0086	1.27
18	MC	1.0915	0.0084	1.12
19	MC	0.2998	0.0081	1.09
20	MC	0.7958	0.0082	1.01
21	MX	0.0309	0.0058	1.09
22	MX	0.3502	0.0060	1.10
23	MC	-0.1222	0.0083	1.11

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
24	MC	0.8056	0.0083	0.97
25	MC	0.5815	0.0082	1.16
26	MC	-1.2374	0.0096	0.95
27	MC	-0.1009	0.0082	1.01
28	MX	0.5217	0.0052	1.12
29	MC	0.6853	0.0082	0.99
30	MC	0.6029	0.0082	1.14
31	MC	1.0400	0.0084	1.01
32	MC	0.4876	0.0081	0.94
33	MC	-0.6259	0.0087	0.82
34	MC	0.4611	0.0082	0.82
35	MC	-0.0609	0.0082	0.96
36	MC	0.7066	0.0082	1.14
37	MC	-0.9141	0.0091	0.88
38	MC	1.2479	0.0086	0.99
39	MC	-0.7389	0.0088	0.83
40	MC	0.7343	0.0082	1.03
41	MC	1.3287	0.0087	0.97
42	MX	1.8545	0.0094	0.89
43	MC	1.1477	0.0085	0.85
44	MC	-0.1457	0.0083	1.04
45	MC	-0.9502	0.0091	0.86
46	MX	0.2945	0.0058	1.13
47	MX	-0.7305	0.0062	1.02

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.3. Item-Level IRT Statistics, ELA Grade 5

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	OE	1.4397	0.0077	0.79
2	OE	1.4125	0.0076	0.80
3	OE	-1.1037	0.0085	0.77
4	OE	1.8855	0.0083	0.77
5	OE	2.1084	0.0088	0.73
6	OE	-1.0956	0.0090	0.80
7	MC	-0.0881	0.0082	1.15
8	MC	0.5748	0.0083	0.93
9	MX	1.5753	0.0094	1.10
10	MC	0.8288	0.0085	1.07
11	XI	0.5023	0.0083	1.03
12	MX	1.3922	0.0091	1.01
13	MC	-0.4448	0.0084	0.93
14	XI	0.4586	0.0083	0.96
15	MC	-0.1823	0.0082	1.07
16	MC	-0.7216	0.0085	1.15

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
17	MX	-0.1997	0.0082	0.82
18	MX	0.6657	0.0084	1.05
19	MC	0.5174	0.0083	1.13
20	MC	-0.7676	0.0086	0.87
21	MC	-0.2308	0.0083	0.86
22	MC	1.2005	0.0088	1.07
23	MC	-0.3365	0.0083	1.01
24	MC	-0.2379	0.0083	1.05
25	MX	-0.0142	0.0059	1.08
26	MX	-0.9021	0.0062	1.08
27	MC	-0.2224	0.0082	1.08
28	MC	0.3303	0.0082	0.97
29	MC	0.0809	0.0082	1.19
30	MC	0.3012	0.0082	1.11
31	XI	1.6204	0.0094	1.07
32	MC	-0.3456	0.0083	1.10
33	MC	0.1600	0.0082	1.14
34	MC	0.3206	0.0082	1.24
35	MC	-0.5069	0.0084	0.87
36	MC	0.2012	0.0082	1.15
37	MX	0.0152	0.0082	1.01
38	MC	-1.4313	0.0095	0.84
39	MC	-0.0669	0.0082	0.92
40	MC	-0.0183	0.0082	0.83
41	MX	1.5012	0.0092	0.92
42	MC	-0.7338	0.0086	0.91
43	MC	0.1683	0.0082	0.88
44	MC	-0.1059	0.0082	1.02
45	MC	-0.3074	0.0083	1.04
46	MX	-0.5504	0.0061	1.08
47	MX	-0.0551	0.0058	1.00

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.4. Item-Level IRT Statistics, ELA Grade 6

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	OE	1.1014	0.0071	0.75
2	OE	1.5260	0.0076	0.78
3	OE	-1.0407	0.0087	0.73
4	OE	0.8086	0.0074	0.78
5	OE	1.2670	0.0079	0.79
6	OE	-1.3560	0.0097	0.73
7	MC	-1.3381	0.0095	0.89
8	MC	-0.5108	0.0085	0.98
9	MC	-0.7668	0.0087	1.02

Appendix B: Item-Level IRT Statistics

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
10	MC	0.5508	0.0082	1.02
11	MX	0.2267	0.0082	0.98
12	MC	-0.7990	0.0087	1.07
13	MC	0.3566	0.0082	1.13
14	MC	0.8305	0.0084	1.07
15	MC	-0.4561	0.0084	0.94
16	MC	0.8816	0.0084	1.14
17	MC	-1.0592	0.0091	0.88
18	MX	1.6025	0.0092	0.98
19	MC	0.7713	0.0083	0.87
20	XI	1.8046	0.0096	1.13
21	MC	-0.2555	0.0083	1.14
22	MC	-2.4212	0.0126	0.91
23	MX	0.5899	0.0060	1.21
24	MX	0.5820	0.0059	1.11
25	MC	-0.0765	0.0082	1.12
26	MX	1.6406	0.0093	1.08
27	MC	0.4923	0.0082	1.13
28	OE	0.0161	0.0082	0.87
29	MC	0.0998	0.0082	1.11
30	MC	-0.1186	0.0082	1.17
31	MC	0.0724	0.0082	1.00
32	MC	-0.6416	0.0086	0.92
33	MC	0.4263	0.0082	1.12
34	MC	0.2460	0.0082	0.94
35	MC	-0.0889	0.0082	1.04
36	MC	-1.0088	0.0090	0.83
37	MC	0.3218	0.0082	1.00
38	MX	0.5600	0.0082	0.88
39	MC	0.5212	0.0082	1.10
40	MC	1.3498	0.0089	1.01
41	MX	0.1497	0.0082	0.76
42	MC	-0.4650	0.0084	0.98
43	MX	0.5972	0.0082	0.91
44	MC	0.9370	0.0084	1.17
45	MC	-1.1466	0.0092	0.79
46	MX	0.1710	0.0060	1.16
47	MX	-0.0494	0.0058	0.99

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.5. Item-Level IRT Statistics, ELA Grade 7

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	OE	0.9140	0.0067	0.80
2	OE	1.2060	0.0072	0.82
3	OE	-0.9982	0.0084	0.72
4	OE	0.8303	0.0071	0.79
5	OE	1.0188	0.0072	0.81
6	OE	-1.2015	0.0094	0.71
7	MC	-0.4587	0.0083	0.94
8	MC	-0.9404	0.0088	0.80
9	MC	0.3856	0.0080	1.10
10	MX	-0.0971	0.0080	1.01
11	MC	0.3505	0.0080	0.88
12	MC	-0.0625	0.0080	0.92
13	MC	0.2907	0.0080	1.17
14	MC	-0.0824	0.0080	1.03
15	MC	0.6316	0.0081	0.90
16	MC	0.8855	0.0082	1.13
17	MC	-0.2964	0.0081	1.02
18	MX	0.7006	0.0052	1.26
19	MC	0.5614	0.0080	1.12
20	MC	0.4573	0.0080	0.99
21	MC	0.2098	0.0080	1.05
22	MC	1.0876	0.0084	1.34
23	MX	0.0694	0.0060	0.99
24	MX	-0.6678	0.0060	1.03
25	MC	-1.2161	0.0093	0.89
26	MC	-1.2004	0.0092	0.95
27	MC	-1.5042	0.0099	0.87
28	MC	-0.8555	0.0087	0.96
29	MC	-0.2838	0.0081	1.02
30	MC	0.4814	0.0080	1.13
31	MC	0.1787	0.0080	1.06
32	MC	-0.0175	0.0080	1.07
33	MC	-0.2165	0.0081	0.97
34	MC	0.7146	0.0081	0.88
35	MC	0.3790	0.0080	1.08
36	MC	0.1548	0.0080	0.81
37	MC	0.1822	0.0080	0.99
38	MC	-0.2786	0.0081	0.84
39	MC	0.1344	0.0080	1.02
40	MC	0.0143	0.0080	1.04
41	MC	-0.5070	0.0083	0.92
42	MC	-0.7703	0.0086	0.85
43	MC	0.1479	0.0080	1.02
44	MC	0.8596	0.0082	1.08
45	MC	0.3485	0.0080	1.08

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
46	MX	0.4899	0.0058	1.12
47	MX	0.2932	0.0059	1.26

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis. One item for ELA Grade 7 was omitted from scoring due to an error in the stimulus.

Table B.6. Item-Level IRT Statistics, ELA Grade 8

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	OE	0.6079	0.0070	0.83
2	OE	1.3995	0.0071	0.95
3	OE	-1.4951	0.0091	0.76
4	OE	1.0561	0.0075	0.74
5	OE	1.3674	0.0072	0.84
6	OE	-1.4669	0.0093	0.74
7	MC	0.2891	0.0079	1.04
8	MX	0.7949	0.0082	0.90
9	MC	-0.3595	0.0080	1.01
10	MC	-0.0717	0.0079	1.03
11	MX	1.1532	0.0085	0.82
12	MC	-0.2582	0.0080	1.02
13	MC	0.0395	0.0079	1.15
14	MC	0.3584	0.0080	1.06
15	MC	0.8516	0.0082	1.18
16	MC	-0.0802	0.0079	0.99
17	MC	0.5562	0.0080	1.11
18	MC	0.2403	0.0079	1.09
19	MC	0.6039	0.0081	1.09
20	MC	0.5266	0.0080	1.06
21	XI	1.1928	0.0086	0.99
22	MC	-0.9556	0.0085	0.94
23	MX	-0.6283	0.0063	1.14
24	MX	-0.9164	0.0058	0.94
25	MC	-1.4505	0.0092	0.86
26	MX	1.2974	0.0087	1.01
27	MC	0.2446	0.0079	1.20
28	MC	0.4095	0.0080	0.95
29	MC	-0.3255	0.0080	1.00
30	MC	0.5422	0.0080	1.16
31	MC	0.9586	0.0083	1.21
32	MC	-1.1466	0.0088	0.91
33	MC	0.3269	0.0080	1.02
34	MX	0.0330	0.0079	0.95
35	MC	-0.7882	0.0083	0.76
36	MC	0.3195	0.0080	0.88
37	MC	-0.1583	0.0080	0.85
38	MC	0.4908	0.0080	0.98

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
39	MC	-0.1029	0.0080	0.89
40	MX	0.9240	0.0083	0.94
41	MX	1.2022	0.0086	1.02
42	MC	0.1694	0.0079	1.19
43	MC	1.4433	0.0089	0.99
44	XI	2.0303	0.0100	1.07
45	MC	-1.6580	0.0097	0.77
46	MX	-0.3994	0.0060	1.10
47	MX	-0.8097	0.0060	0.84

Note. OE = open-ended, MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.7. Item-Level IRT Statistics, Mathematics Grade 3

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	MC	-2.0209	0.0119	1.15
2	XI	0.8701	0.0088	1.00
3	MC	1.3188	0.0089	1.11
4	XI	-0.8246	0.0097	1.18
5	XI	-0.3147	0.0092	0.93
6	XI	-0.6020	0.0095	1.02
7	XI	-0.5070	0.0094	1.11
8	MX	1.5475	0.0091	1.10
9	XI	0.2634	0.0089	1.04
10	XI	0.5865	0.0088	0.87
11	XI	1.3005	0.0089	1.10
12	XI	1.6965	0.0092	0.97
13	XI	1.4490	0.0090	0.90
14	XI	1.0546	0.0088	0.91
15	MC	1.2265	0.0089	1.26
16	MC	1.0418	0.0088	1.10
17	MC	0.7365	0.0088	1.31
18	XI	0.0542	0.0090	1.13
19	XI	-0.0479	0.0090	0.90
20	XI	-0.6311	0.0095	0.79
21	MC	-0.8415	0.0098	0.90
22	XI	-0.6431	0.0095	0.79
23	XI	-0.5069	0.0094	0.83
24	MC	-1.0498	0.0100	1.12
25	MC	1.6195	0.0091	0.82
26	XI	-0.5535	0.0094	0.80
27	MC	-0.1974	0.0091	1.24
28	XI	0.1170	0.0090	0.84
29	MC	-0.0694	0.0090	1.08
30	MC	-1.0541	0.0100	0.97
31	XI	0.6796	0.0088	1.35

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
32	XI	1.7303	0.0092	0.92
33	MC	1.0490	0.0088	0.91
34	XI	0.1652	0.0089	0.87
35	XI	0.8120	0.0088	0.84
36	XI	-1.2121	0.0103	0.76
37	MC	-0.7664	0.0097	0.99
38	MC	1.0217	0.0088	1.24
39	MC	-1.9301	0.0117	1.07
40	XI	0.8420	0.0088	0.90
41	XI	0.7363	0.0088	0.94
42	XI	-0.5346	0.0094	0.87
43	XI	0.8303	0.0088	0.84
44	MC	0.0749	0.0090	0.82
45	XI	1.0196	0.0088	0.96

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.8. Item-Level IRT Statistics, Mathematics Grade 4

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	XI	0.8284	0.0090	0.83
2	XI	0.3677	0.0089	0.96
3	XI	-1.4982	0.0099	0.89
4	XI	1.6709	0.0098	0.94
5	MX	-0.5185	0.0090	1.22
6	XI	-0.5871	0.0090	1.09
7	MC	0.7896	0.0090	1.34
8	XI	-1.0555	0.0094	1.03
9	MC	-0.2253	0.0089	1.01
10	MC	-0.5781	0.0090	0.91
11	XI	-0.0581	0.0089	0.97
12	XI	0.6264	0.0089	0.98
13	XI	0.8866	0.0091	0.79
14	MC	0.3267	0.0089	0.97
15	XI	0.1355	0.0088	0.92
16	XI	0.9143	0.0091	0.93
17	MC	1.9347	0.0102	0.96
18	XI	-1.0077	0.0094	0.91
19	XI	-1.0726	0.0094	0.92
20	MC	-1.0490	0.0094	0.93
21	XI	0.3615	0.0089	0.92
22	XI	0.5094	0.0089	0.95
23	MC	0.5728	0.0089	1.02
24	MC	-1.9035	0.0105	0.99
25	XI	0.4224	0.0089	0.89
26	XI	1.0501	0.0092	1.17
27	MC	0.7618	0.0090	1.24

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
28	XI	0.3307	0.0089	0.92
29	MC	-1.0561	0.0094	1.10
30	XI	0.5770	0.0089	0.92
31	XI	-0.5636	0.0090	0.85
32	MC	1.1661	0.0093	1.37
33	XI	0.6010	0.0089	0.93
34	XI	-1.9376	0.0106	0.76
35	MC	-0.3783	0.0089	1.30
36	XI	-0.3969	0.0090	1.05
37	XI	-1.0534	0.0094	0.85
38	MC	0.1626	0.0088	1.31
39	XI	0.9581	0.0091	1.07
40	XI	-1.7648	0.0103	0.74
41	MC	-1.9371	0.0106	1.23
42	MX	0.6635	0.0090	0.88
43	XI	0.1922	0.0088	1.12
44	MC	-0.0153	0.0088	1.26
45	XI	1.0134	0.0092	0.79

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.9. Item-Level IRT Statistics, Mathematics Grade 5

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	XI	-0.7129	0.0086	1.17
2	XI	-0.4958	0.0086	0.90
3	MC	-0.1505	0.0086	1.18
4	MX	0.8991	0.0092	1.03
5	XI	-0.0308	0.0086	0.81
6	MX	0.4529	0.0088	1.27
7	XI	-0.6994	0.0086	0.84
8	XI	-0.1190	0.0086	0.87
9	MC	0.1439	0.0086	1.33
10	XI	0.0524	0.0086	0.82
11	XI	1.2433	0.0096	0.79
12	XI	0.6744	0.0089	0.86
13	XI	0.9615	0.0092	0.84
14	XI	-0.0237	0.0086	0.94
15	XI	1.0037	0.0093	0.83
16	XI	-0.0492	0.0086	0.75
17	XI	0.1014	0.0086	0.92
18	XI	-0.2098	0.0086	1.01
19	MC	-1.0212	0.0088	1.02
20	XI	-0.1130	0.0086	0.79
21	MC	-0.1986	0.0086	1.26
22	MC	0.9375	0.0092	1.24
23	MC	-0.8393	0.0087	1.12

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
24	MC	-1.3704	0.0091	1.19
25	MC	0.2690	0.0087	0.84
26	MC	-0.8369	0.0087	1.24
27	XI	-0.2638	0.0086	1.01
28	XI	0.3424	0.0087	0.81
29	MC	-0.3931	0.0086	1.00
30	XI	-0.6364	0.0086	1.22
31	XI	2.0818	0.0112	1.00
32	MC	0.7110	0.0090	1.11
33	XI	0.7354	0.0090	0.84
34	XI	0.8232	0.0091	1.04
35	XI	1.9380	0.0109	0.88
36	MC	-0.1844	0.0086	1.11
37	MX	-0.4154	0.0086	1.20
38	XI	0.5884	0.0089	0.89
39	MC	0.1392	0.0086	1.30
40	MC	-0.6667	0.0086	1.02
41	XI	-0.1443	0.0086	0.80
42	XI	0.4315	0.0088	1.06
43	XI	-0.6577	0.0086	0.78
44	XI	-0.6830	0.0086	0.87
45	XI	0.5131	0.0088	0.94

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.10. Item-Level IRT Statistics, Mathematics Grade 6

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	XI	-0.1295	0.0085	0.91
2	MC	-0.5536	0.0084	1.11
3	MC	-1.3975	0.0088	1.36
4	XI	0.1527	0.0087	1.37
5	XI	-0.3812	0.0084	0.80
6	XI	0.1290	0.0086	0.92
7	MC	0.9272	0.0094	0.82
8	XI	0.1530	0.0087	0.93
9	XI	0.3692	0.0088	0.90
10	XI	0.5220	0.0090	0.89
11	XI	0.6930	0.0091	0.81
12	XI	-0.1532	0.0085	0.82
13	XI	0.3941	0.0088	0.74
14	XI	1.2044	0.0099	0.89
15	MC	0.2708	0.0087	1.36
16	XI	1.6132	0.0107	0.99
17	MC	0.3665	0.0088	1.32
18	XI	-0.1050	0.0085	0.90
19	MC	0.0743	0.0086	1.17

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
20	XI	0.1516	0.0087	1.00
21	MC	-0.0243	0.0086	1.37
22	XI	0.0197	0.0086	0.84
23	MC	0.3646	0.0088	1.09
24	XI	-1.0345	0.0085	0.74
25	XI	1.1871	0.0098	0.83
26	MC	-0.1634	0.0085	1.29
27	XI	-0.5250	0.0084	0.77
28	MC	-0.9353	0.0085	1.00
29	MC	-0.1403	0.0085	0.93
30	XI	-0.4024	0.0084	1.35
31	XI	-0.5798	0.0084	1.10
32	MC	0.3523	0.0088	1.05
33	XI	-0.0054	0.0086	0.93
34	MC	0.4783	0.0089	0.84
35	XI	-0.1735	0.0085	1.00
36	XI	0.4728	0.0089	1.15
37	XI	-0.3027	0.0085	0.80
38	MC	0.3367	0.0088	1.02
39	XI	-0.7318	0.0084	1.04
40	MC	0.9524	0.0095	1.07
41	XI	0.9330	0.0094	0.82
42	MC	0.0472	0.0086	1.10
43	XI	0.5261	0.0090	0.90
44	XI	-0.2134	0.0085	0.91
45	XI	-0.4756	0.0084	1.18
46	XI	0.0398	0.0086	0.84
47	MC	-1.0889	0.0086	1.06

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.11. Item-Level IRT Statistics, Mathematics Grade 7

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	MC	-0.5590	0.0086	1.32
2	XI	0.9548	0.0095	0.98
3	XI	0.5075	0.0090	0.94
4	XI	-0.2400	0.0086	0.84
5	MC	-1.2689	0.0088	1.02
6	MC	1.4576	0.0102	1.36
7	MC	-0.5669	0.0086	1.32
8	XI	-0.3725	0.0086	0.78
9	XI	0.2534	0.0089	0.85
10	MC	0.4638	0.0090	1.36
11	XI	-0.1984	0.0087	0.89
12	XI	0.6692	0.0092	1.15
13	XI	-0.6491	0.0086	0.80

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
14	XI	1.6358	0.0106	1.03
15	XI	-0.5526	0.0086	0.77
16	XI	1.1109	0.0097	0.84
17	MC	1.2896	0.0100	1.20
18	XI	0.2634	0.0089	0.97
19	MC	-1.5791	0.0091	0.81
20	MC	-0.2147	0.0086	1.45
21	XI	-0.3265	0.0086	0.75
22	XI	0.8825	0.0094	0.95
23	XI	-0.7828	0.0086	0.78
24	MC	-2.3712	0.0102	0.93
25	XI	0.9471	0.0095	0.97
26	XI	0.5019	0.0090	0.84
27	XI	1.5046	0.0103	0.89
28	XI	0.5576	0.0091	1.09
29	XI	-0.7337	0.0086	0.77
30	MC	0.0916	0.0088	1.24
31	MC	-0.0491	0.0087	1.08
32	XI	1.1255	0.0097	0.82
33	XI	0.6839	0.0092	0.94
34	XI	1.2389	0.0099	0.89
35	XI	1.6322	0.0106	0.99
36	MC	0.3188	0.0089	1.12
37	XI	1.0104	0.0096	1.06
38	MC	-0.6670	0.0086	1.01
39	XI	0.8655	0.0094	0.95
40	XI	1.6575	0.0106	0.78
41	XI	1.4946	0.0103	0.78
42	MC	-1.5448	0.0090	1.25
43	MC	-0.9877	0.0087	1.26
44	MC	-0.5406	0.0086	1.26
45	XI	0.8638	0.0094	0.83
46	MC	-1.6029	0.0091	0.82
47	XI	1.2310	0.0099	0.89

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.12. Item-Level IRT Statistics, Mathematics Grade 8

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
1	MC	-1.6084	0.0082	0.98
2	MC	0.0068	0.0083	1.19
3	XI	-0.0419	0.0083	0.82
4	MC	0.8348	0.0094	1.05
5	MC	-0.1720	0.0081	1.10
6	XI	0.7564	0.0093	0.81
7	MC	-0.8414	0.0079	0.93

Appendix B: Item-Level IRT Statistics

Item Number	Item Type	Rasch Difficulty	SE	MNSQ Infit
8	MC	0.5570	0.0090	1.21
9	MC	0.8865	0.0095	1.14
10	MC	-2.4121	0.0094	1.00
11	MC	0.6845	0.0092	1.25
12	MX	-0.1094	0.0082	1.20
13	XI	1.8480	0.0120	0.91
14	XI	0.4998	0.0089	0.83
15	MC	-0.1728	0.0081	1.13
16	MX	-0.1611	0.0081	1.05
17	XI	1.6330	0.0113	0.86
18	XI	1.2029	0.0102	0.95
19	MC	-0.5810	0.0079	1.06
20	XI	0.8433	0.0095	1.09
21	MC	-1.2456	0.0079	0.96
22	MC	-0.5707	0.0079	0.97
23	MC	-0.1838	0.0081	0.95
24	MC	-0.8242	0.0079	0.83
25	XI	-0.5994	0.0079	1.18
26	MC	-0.3835	0.0080	0.96
27	MC	-0.3013	0.0080	1.20
28	MC	-0.3429	0.0080	1.14
29	MC	-0.8182	0.0079	1.00
30	MC	0.1770	0.0085	1.03
31	MX	0.1225	0.0084	1.04
32	MC	-0.9816	0.0079	1.02
33	MC	-0.5156	0.0079	0.83
34	MC	0.1202	0.0084	1.01
35	MC	0.1515	0.0084	1.13
36	XI	-0.2976	0.0080	1.01
37	MC	-0.8777	0.0079	1.15
38	MC	-0.2507	0.0081	0.99
39	XI	-0.2529	0.0081	0.97
40	XI	-0.0018	0.0083	0.71
41	XI	-0.0249	0.0083	0.80
42	MC	-1.0844	0.0079	1.10
43	XI	0.3671	0.0087	0.86
44	MC	0.1201	0.0084	0.82
45	XI	0.0565	0.0083	0.79
46	MC	-1.1546	0.0079	1.16
47	MC	-0.5740	0.0079	0.96

Note. MC = multiple-choice, MX = multi-part, XI = technology-enhanced. Item number does not indicate item location on an operational test form, as field test items were embedded on the form but not included in the analysis.

Table B.13. Raw-to-Scale Score Conversion, ELA Grade 3

Form	Raw Score	Scale Score	CSEM	Performance Level
1	2	2395	21	1
1	3	2395	21	1
1	4	2395	21	1
1	5	2406	18	1
1	6	2416	16	1
1	7	2423	15	1
1	8	2430	14	1
1	9	2436	13	1
1	10	2441	12	1
1	11	2445	12	1
1	12	2450	11	1
1	13	2454	11	1
1	14	2457	10	1
1	15	2461	10	1
1	16	2464	10	1
1	17	2467	10	1
1	18	2471	10	1
1	19	2474	10	1
1	20	2477	9	1
1	21	2480	9	1
1	22	2482	9	1
1	23	2485	9	1
1	24	2488	9	1
1	25	2491	9	1
1	26	2494	9	1
1	27	2497	9	2
1	28	2499	9	2
1	29	2502	9	2
1	30	2504	9	2
1	31	2507	9	2
1	32	2510	9	3
1	33	2513	9	3
1	34	2516	9	3
1	35	2519	9	3
1	36	2522	10	3
1	37	2525	10	3
1	38	2528	10	3
1	39	2531	10	3
1	40	2535	10	3
1	41	2538	10	3
1	42	2542	11	4
1	43	2546	11	4
1	44	2550	11	4
1	45	2554	12	4
1	46	2559	12	4
1	47	2564	13	4
1	48	2570	14	4

Appendix B: Item-Level IRT Statistics

Form	Raw Score	Scale Score	CSEM	Performance Level
1	49	2576	14	4
1	50	2584	16	4
1	51	2593	17	4
1	52	2604	19	4
1	53	2605	20	4
1	54	2605	20	4
1	55	2605	20	4
1	56	2605	20	4
2	2	2395	21	1
2	3	2395	21	1
2	4	2395	21	1
2	5	2405	18	1
2	6	2414	16	1
2	7	2422	15	1
2	8	2429	13	1
2	9	2434	13	1
2	10	2439	12	1
2	11	2444	12	1
2	12	2448	11	1
2	13	2452	11	1
2	14	2456	10	1
2	15	2459	10	1
2	16	2463	10	1
2	17	2466	10	1
2	18	2469	10	1
2	19	2472	10	1
2	20	2475	9	1
2	21	2478	9	1
2	22	2481	9	1
2	23	2484	9	1
2	24	2487	9	1
2	25	2489	9	1
2	26	2492	9	1
2	27	2495	9	1
2	28	2498	9	2
2	29	2500	9	2
2	30	2503	9	2
2	31	2506	9	2
2	32	2509	9	3
2	33	2511	9	3
2	34	2514	9	3
2	35	2517	9	3
2	36	2520	10	3
2	37	2523	10	3
2	38	2526	10	3
2	39	2530	10	3
2	40	2533	10	3
2	41	2536	10	3

Form	Raw Score	Scale Score	CSEM	Performance Level
2	42	2541	11	4
2	43	2544	11	4
2	44	2548	11	4
2	45	2552	12	4
2	46	2557	12	4
2	47	2562	13	4
2	48	2568	14	4
2	49	2575	15	4
2	50	2583	16	4
2	51	2592	18	4
2	52	2605	21	4
2	53	2605	21	4
2	54	2605	21	4
2	55	2605	21	4
2	56	2605	21	4

Table B.14. Raw-to-Scale Score Conversion, ELA Grade 4

Form	Raw Score	Scale Score	CSEM	Performance Level
1	2	2400	25	1
1	3	2400	25	1
1	4	2410	22	1
1	5	2423	18	1
1	6	2432	16	1
1	7	2440	14	1
1	8	2446	13	1
1	9	2452	13	1
1	10	2457	12	1
1	11	2461	11	1
1	12	2466	11	1
1	13	2469	11	1
1	14	2473	10	1
1	15	2476	10	1
1	16	2480	10	1
1	17	2483	10	1
1	18	2486	9	1
1	19	2489	9	1
1	20	2491	9	1
1	21	2494	9	1
1	22	2497	9	1
1	23	2499	9	1
1	24	2502	9	1
1	25	2505	9	1
1	26	2507	9	1
1	27	2510	9	2
1	28	2512	9	2
1	29	2514	9	2
1	30	2517	9	2
1	31	2519	9	2

Appendix B: Item-Level IRT Statistics

Form	Raw Score	Scale Score	CSEM	Performance Level
1	32	2522	9	2
1	33	2524	9	3
1	34	2527	9	3
1	35	2529	9	3
1	36	2532	9	3
1	37	2534	9	3
1	38	2537	9	3
1	39	2540	9	3
1	40	2543	9	3
1	41	2546	9	3
1	42	2549	10	3
1	43	2552	10	3
1	44	2555	10	3
1	45	2559	11	4
1	46	2563	11	4
1	47	2567	11	4
1	48	2571	12	4
1	49	2576	12	4
1	50	2582	13	4
1	51	2588	14	4
1	52	2595	15	4
1	53	2604	17	4
1	54	2610	18	4
1	55	2610	18	4
1	56	2610	18	4
1	57	2610	18	4
2	2	2400	25	1
2	3	2400	25	1
2	4	2410	22	1
2	5	2423	18	1
2	6	2433	16	1
2	7	2440	14	1
2	8	2447	13	1
2	9	2452	13	1
2	10	2457	12	1
2	11	2462	11	1
2	12	2466	11	1
2	13	2470	11	1
2	14	2474	10	1
2	15	2477	10	1
2	16	2480	10	1
2	17	2483	10	1
2	18	2487	9	1
2	19	2489	9	1
2	20	2492	9	1
2	21	2495	9	1
2	22	2498	9	1
2	23	2501	9	1

Form	Raw Score	Scale Score	CSEM	Performance Level
2	24	2503	9	1
2	25	2506	9	1
2	26	2508	9	1
2	27	2511	9	2
2	28	2513	9	2
2	29	2516	9	2
2	30	2518	9	2
2	31	2521	9	2
2	32	2524	9	3
2	33	2526	9	3
2	34	2529	9	3
2	35	2531	9	3
2	36	2534	9	3
2	37	2537	9	3
2	38	2540	9	3
2	39	2543	9	3
2	40	2546	10	3
2	41	2549	10	3
2	42	2552	10	3
2	43	2555	10	3
2	44	2559	11	4
2	45	2563	11	4
2	46	2567	11	4
2	47	2571	12	4
2	48	2576	12	4
2	49	2582	13	4
2	50	2587	14	4
2	51	2594	15	4
2	52	2602	16	4
2	53	2610	17	4
2	54	2610	17	4
2	55	2610	17	4
2	56	2610	17	4
2	57	2610	17	4

Table B.15. Raw-to-Scale Score Conversion, ELA Grade 5

Form	Raw Score	Scale Score	CSEM	Performance Level
1	2	2419	23	1
1	3	2419	23	1
1	4	2423	22	1
1	5	2436	18	1
1	6	2446	16	1
1	7	2454	14	1
1	8	2460	13	1
1	9	2465	12	1
1	10	2470	12	1
1	11	2475	11	1
1	12	2479	11	1

Appendix B: Item-Level IRT Statistics

Form	Raw Score	Scale Score	CSEM	Performance Level
1	13	2483	11	1
1	14	2486	10	1
1	15	2490	10	1
1	16	2493	10	1
1	17	2496	10	1
1	18	2499	9	1
1	19	2502	9	1
1	20	2505	9	1
1	21	2508	9	1
1	22	2510	9	1
1	23	2513	9	1
1	24	2516	9	1
1	25	2518	9	1
1	26	2521	9	2
1	27	2524	9	2
1	28	2526	9	2
1	29	2529	9	2
1	30	2532	9	2
1	31	2534	9	2
1	32	2537	9	2
1	33	2540	9	2
1	34	2543	9	3
1	35	2545	9	3
1	36	2548	9	3
1	37	2551	9	3
1	38	2554	10	3
1	39	2557	10	3
1	40	2560	10	3
1	41	2564	10	3
1	42	2567	10	3
1	43	2571	11	3
1	44	2575	11	3
1	45	2579	12	4
1	46	2584	12	4
1	47	2589	13	4
1	48	2595	13	4
1	49	2601	14	4
1	50	2608	15	4
1	51	2617	17	4
1	52	2628	19	4
1	53	2629	20	4
1	54	2629	20	4
1	55	2629	20	4
2	2	2419	23	1
2	3	2419	23	1
2	4	2424	22	1
2	5	2437	18	1
2	6	2447	16	1

Form	Raw Score	Scale Score	CSEM	Performance Level
2	7	2454	14	1
2	8	2461	13	1
2	9	2466	13	1
2	10	2471	12	1
2	11	2476	11	1
2	12	2480	11	1
2	13	2484	11	1
2	14	2487	10	1
2	15	2491	10	1
2	16	2494	10	1
2	17	2497	10	1
2	18	2500	10	1
2	19	2503	9	1
2	20	2506	9	1
2	21	2509	9	1
2	22	2512	9	1
2	23	2514	9	1
2	24	2517	9	1
2	25	2520	9	2
2	26	2523	9	2
2	27	2525	9	2
2	28	2528	9	2
2	29	2530	9	2
2	30	2533	9	2
2	31	2536	9	2
2	32	2539	9	2
2	33	2541	9	2
2	34	2544	9	3
2	35	2547	9	3
2	36	2550	9	3
2	37	2553	10	3
2	38	2556	10	3
2	39	2559	10	3
2	40	2562	10	3
2	41	2566	10	3
2	42	2569	11	3
2	43	2573	11	3
2	44	2578	11	4
2	45	2582	12	4
2	46	2587	12	4
2	47	2592	13	4
2	48	2598	14	4
2	49	2604	15	4
2	50	2612	16	4
2	51	2621	18	4
2	52	2629	19	4
2	53	2629	19	4
2	54	2629	19	4
2	55	2629	19	4

Table B.16. Raw-to-Scale Score Conversion, ELA Grade 6

Form	Raw Score	Scale Score	CSEM	Performance Level
1	2	2431	23	1
1	3	2431	23	1
1	4	2433	22	1
1	5	2447	18	1
1	6	2457	16	1
1	7	2464	15	1
1	8	2471	14	1
1	9	2477	13	1
1	10	2482	12	1
1	11	2487	12	1
1	12	2491	11	1
1	13	2495	11	1
1	14	2499	10	1
1	15	2502	10	1
1	16	2506	10	1
1	17	2509	10	1
1	18	2512	10	1
1	19	2515	9	1
1	20	2518	9	1
1	21	2521	9	1
1	22	2524	9	1
1	23	2527	9	1
1	24	2529	9	1
1	25	2532	9	2
1	26	2535	9	2
1	27	2537	9	2
1	28	2540	9	2
1	29	2543	9	2
1	30	2545	9	2
1	31	2548	9	2
1	32	2551	9	2
1	33	2554	9	3
1	34	2556	9	3
1	35	2559	9	3
1	36	2562	9	3
1	37	2565	9	3
1	38	2568	10	3
1	39	2571	10	3
1	40	2574	10	3
1	41	2578	10	3
1	42	2581	10	3
1	43	2585	11	3
1	44	2589	11	3
1	45	2593	11	3
1	46	2598	12	4
1	47	2603	13	4
1	48	2608	13	4

Appendix B: Item-Level IRT Statistics

Form	Raw Score	Scale Score	CSEM	Performance Level
1	49	2615	14	4
1	50	2622	15	4
1	51	2630	17	4
1	52	2641	19	4
1	53	2641	19	4
1	54	2641	19	4
1	55	2641	19	4
2	2	2431	22	1
2	3	2431	22	1
2	4	2432	22	1
2	5	2445	18	1
2	6	2455	16	1
2	7	2463	15	1
2	8	2470	14	1
2	9	2475	13	1
2	10	2481	12	1
2	11	2485	12	1
2	12	2490	11	1
2	13	2494	11	1
2	14	2497	11	1
2	15	2501	10	1
2	16	2504	10	1
2	17	2508	10	1
2	18	2511	10	1
2	19	2514	10	1
2	20	2517	9	1
2	21	2520	9	1
2	22	2523	9	1
2	23	2525	9	1
2	24	2528	9	1
2	25	2532	9	2
2	26	2534	9	2
2	27	2536	9	2
2	28	2539	9	2
2	29	2542	9	2
2	30	2544	9	2
2	31	2547	9	2
2	32	2550	9	2
2	33	2553	9	3
2	34	2555	9	3
2	35	2558	9	3
2	36	2561	9	3
2	37	2564	9	3
2	38	2567	10	3
2	39	2570	10	3
2	40	2573	10	3
2	41	2577	10	3
2	42	2580	10	3

Form	Raw Score	Scale Score	CSEM	Performance Level
2	43	2584	11	3
2	44	2588	11	3
2	45	2592	11	3
2	46	2597	12	4
2	47	2602	12	4
2	48	2607	13	4
2	49	2613	14	4
2	50	2620	15	4
2	51	2629	17	4
2	52	2639	19	4
2	53	2641	19	4
2	54	2641	19	4
2	55	2641	19	4

Table B.17. Raw-to-Scale Score Conversion, ELA Grade 7

Form	Raw Score	Scale Score	CSEM	Performance Level
1	2	2438	23	1
1	3	2438	23	1
1	4	2442	22	1
1	5	2455	18	1
1	6	2464	16	1
1	7	2472	14	1
1	8	2478	13	1
1	9	2484	12	1
1	10	2489	12	1
1	11	2493	11	1
1	12	2497	11	1
1	13	2501	11	1
1	14	2505	10	1
1	15	2508	10	1
1	16	2511	10	1
1	17	2514	10	1
1	18	2517	9	1
1	19	2520	9	1
1	20	2523	9	1
1	21	2526	9	1
1	22	2528	9	1
1	23	2531	9	1
1	24	2534	9	1
1	25	2536	9	1
1	26	2539	9	1
1	27	2541	9	1
1	28	2544	9	2
1	29	2546	9	2
1	30	2549	9	2
1	31	2551	9	2
1	32	2554	9	2
1	33	2556	9	2

Appendix B: Item-Level IRT Statistics

Form	Raw Score	Scale Score	CSEM	Performance Level
1	34	2559	9	2
1	35	2561	9	3
1	36	2564	9	3
1	37	2566	9	3
1	38	2569	9	3
1	39	2572	9	3
1	40	2575	9	3
1	41	2578	10	3
1	42	2581	10	3
1	43	2584	10	3
1	44	2588	10	3
1	45	2591	11	3
1	46	2595	11	3
1	47	2600	11	4
1	48	2604	12	4
1	49	2609	13	4
1	50	2615	14	4
1	51	2621	15	4
1	52	2629	16	4
1	53	2639	18	4
1	54	2648	21	4
1	55	2648	21	4
1	56	2648	21	4
2	2	2438	23	1
2	3	2438	23	1
2	4	2441	22	1
2	5	2454	18	1
2	6	2464	16	1
2	7	2471	14	1
2	8	2478	13	1
2	9	2483	12	1
2	10	2488	12	1
2	11	2493	11	1
2	12	2497	11	1
2	13	2501	11	1
2	14	2504	10	1
2	15	2508	10	1
2	16	2511	10	1
2	17	2514	10	1
2	18	2517	9	1
2	19	2520	9	1
2	20	2523	9	1
2	21	2525	9	1
2	22	2528	9	1
2	23	2531	9	1
2	24	2533	9	1
2	25	2536	9	1
2	26	2538	9	1

Form	Raw Score	Scale Score	CSEM	Performance Level
2	27	2541	9	1
2	28	2543	9	2
2	29	2546	9	2
2	30	2548	9	2
2	31	2551	9	2
2	32	2553	9	2
2	33	2556	9	2
2	34	2558	9	2
2	35	2561	9	3
2	36	2563	9	3
2	37	2566	9	3
2	38	2569	9	3
2	39	2571	9	3
2	40	2574	9	3
2	41	2577	10	3
2	42	2580	10	3
2	43	2584	10	3
2	44	2587	10	3
2	45	2591	11	3
2	46	2594	11	3
2	47	2600	11	4
2	48	2603	12	4
2	49	2608	13	4
2	50	2614	13	4
2	51	2620	14	4
2	52	2628	16	4
2	53	2637	18	4
2	54	2648	21	4
2	55	2648	21	4
2	56	2648	21	4

Table B.18. Raw-to-Scale Score Conversion, ELA Grade 8

Form	Raw Score	Scale Score	CSEM	Performance Level
1	2	2448	22	1
1	3	2448	22	1
1	4	2448	22	1
1	5	2461	18	1
1	6	2471	16	1
1	7	2479	15	1
1	8	2485	14	1
1	9	2491	13	1
1	10	2496	12	1
1	11	2501	12	1
1	12	2505	11	1
1	13	2509	11	1
1	14	2513	11	1
1	15	2517	10	1
1	16	2520	10	1

Appendix B: Item-Level IRT Statistics

Form	Raw Score	Scale Score	CSEM	Performance Level
1	17	2524	10	1
1	18	2527	10	1
1	19	2530	10	1
1	20	2533	9	1
1	21	2536	9	1
1	22	2539	9	1
1	23	2542	9	1
1	24	2544	9	1
1	25	2547	9	1
1	26	2551	9	2
1	27	2553	9	2
1	28	2555	9	2
1	29	2558	9	2
1	30	2561	9	2
1	31	2564	9	2
1	32	2566	9	2
1	33	2569	9	2
1	34	2572	9	3
1	35	2575	9	3
1	36	2578	9	3
1	37	2581	9	3
1	38	2584	10	3
1	39	2587	10	3
1	40	2590	10	3
1	41	2593	10	3
1	42	2597	10	3
1	43	2600	11	3
1	44	2604	11	4
1	45	2608	11	4
1	46	2613	12	4
1	47	2617	12	4
1	48	2623	13	4
1	49	2628	14	4
1	50	2635	15	4
1	51	2643	16	4
1	52	2653	18	4
1	53	2658	20	4
1	54	2658	20	4
1	55	2658	20	4
2	2	2448	22	1
2	3	2448	22	1
2	4	2448	22	1
2	5	2461	18	1
2	6	2471	16	1
2	7	2478	15	1
2	8	2485	14	1
2	9	2491	13	1
2	10	2496	12	1

Appendix B: Item-Level IRT Statistics

Form	Raw Score	Scale Score	CSEM	Performance Level
2	11	2501	12	1
2	12	2505	11	1
2	13	2509	11	1
2	14	2513	11	1
2	15	2517	10	1
2	16	2520	10	1
2	17	2524	10	1
2	18	2527	10	1
2	19	2530	10	1
2	20	2533	10	1
2	21	2536	9	1
2	22	2539	9	1
2	23	2542	9	1
2	24	2545	9	1
2	25	2548	9	1
2	26	2551	9	2
2	27	2553	9	2
2	28	2556	9	2
2	29	2559	9	2
2	30	2561	9	2
2	31	2564	9	2
2	32	2567	9	2
2	33	2570	9	2
2	34	2572	9	3
2	35	2575	9	3
2	36	2578	9	3
2	37	2581	9	3
2	38	2584	10	3
2	39	2587	10	3
2	40	2591	10	3
2	41	2594	10	3
2	42	2598	10	3
2	43	2601	11	3
2	44	2605	11	4
2	45	2609	11	4
2	46	2614	12	4
2	47	2619	12	4
2	48	2624	13	4
2	49	2630	14	4
2	50	2637	15	4
2	51	2646	17	4
2	52	2656	19	4
2	53	2658	19	4
2	54	2658	19	4
2	55	2658	19	4

Table B.19. Raw-to-Scale Score Conversion, Mathematics Grade 3

Raw Score	Scale Score	CSEM	Performance Level
0	3395	24	1
1	3395	24	1
2	3401	22	1
3	3415	19	1
4	3425	16	1
5	3433	15	1
6	3440	14	1
7	3447	13	1
8	3452	13	1
9	3457	12	1
10	3462	12	1
11	3466	11	1
12	3470	11	1
13	3474	11	1
14	3478	11	1
15	3482	10	1
16	3485	10	1
17	3489	10	1
18	3492	10	1
19	3495	10	2
20	3499	10	2
21	3502	10	2
22	3505	10	2
23	3508	10	2
24	3512	10	2
25	3515	10	2
26	3518	10	2
27	3522	10	2
28	3525	10	2
29	3528	10	2
30	3532	10	3
31	3535	10	3
32	3539	11	3
33	3543	11	3
34	3547	11	3
35	3551	11	3
36	3556	12	3
37	3561	12	3
38	3566	13	3
39	3573	14	4
40	3578	15	4
41	3586	16	4
42	3596	18	4
43	3605	21	4
44	3605	21	4
45	3605	21	4

Table B.20. Raw-to-Scale Score Conversion, Mathematics Grade 4

Raw Score	Scale Score	CSEM	Performance Level
0	3435	22	1
1	3435	22	1
2	3435	22	1
3	3448	19	1
4	3458	16	1
5	3466	15	1
6	3473	14	1
7	3479	13	1
8	3485	13	1
9	3490	12	1
10	3495	12	1
11	3499	11	1
12	3503	11	1
13	3507	11	1
14	3511	11	1
15	3515	10	1
16	3518	10	1
17	3522	10	1
18	3525	10	1
19	3528	10	1
20	3532	10	2
21	3535	10	2
22	3538	10	2
23	3541	10	2
24	3545	10	2
25	3548	10	2
26	3551	10	2
27	3554	10	2
28	3558	10	2
29	3562	10	3
30	3565	10	3
31	3568	10	3
32	3572	11	3
33	3576	11	3
34	3580	11	3
35	3584	11	3
36	3588	12	3
37	3593	12	3
38	3598	13	3
39	3606	14	4
40	3611	15	4
41	3619	16	4
42	3629	18	4
43	3642	22	4
44	3645	23	4
45	3645	23	4

Table B.21. Raw-to-Scale Score Conversion, Mathematics Grade 5

Raw Score	Scale Score	CSEM	Performance Level
0	3478	25	1
1	3478	25	1
2	3487	22	1
3	3500	18	1
4	3510	16	1
5	3517	15	1
6	3524	13	1
7	3530	13	1
8	3535	12	1
9	3539	12	1
10	3544	11	1
11	3548	11	1
12	3551	11	1
13	3555	10	1
14	3559	10	1
15	3563	10	2
16	3565	10	2
17	3568	10	2
18	3571	10	2
19	3574	10	2
20	3577	9	2
21	3580	9	2
22	3583	9	2
23	3586	9	2
24	3589	9	2
25	3592	10	2
26	3596	10	3
27	3599	10	3
28	3602	10	3
29	3605	10	3
30	3608	10	3
31	3612	10	3
32	3615	10	3
33	3619	11	3
34	3623	11	3
35	3627	11	3
36	3631	12	3
37	3636	12	4
38	3641	13	4
39	3647	14	4
40	3654	15	4
41	3662	16	4
42	3672	18	4
43	3685	22	4
44	3688	23	4
45	3688	23	4

Table B.22. Raw-to-Scale Score Conversion, Mathematics Grade 6

Raw Score	Scale Score	CSEM	Performance Level
0	3512	25	1
1	3512	25	1
2	3520	22	1
3	3534	18	1
4	3543	16	1
5	3551	14	1
6	3557	13	1
7	3563	13	1
8	3568	12	1
9	3572	11	1
10	3577	11	1
11	3581	11	1
12	3584	10	1
13	3588	10	1
14	3591	10	1
15	3594	10	1
16	3597	10	1
17	3600	9	1
18	3603	9	2
19	3606	9	2
20	3609	9	2
21	3612	9	2
22	3615	9	2
23	3618	9	2
24	3620	9	2
25	3623	9	2
26	3626	9	2
27	3629	9	3
28	3632	9	3
29	3634	9	3
30	3637	9	3
31	3640	10	3
32	3644	10	3
33	3647	10	3
34	3650	10	3
35	3654	10	3
36	3657	11	3
37	3661	11	3
38	3665	11	4
39	3670	12	4
40	3675	13	4
41	3681	13	4
42	3687	14	4
43	3695	16	4
44	3704	18	4
45	3717	22	4
46	3722	23	4
47	3722	23	4

Table B.23. Raw-to-Scale Score Conversion, Mathematics Grade 7

Raw Score	Scale Score	CSEM	Performance Level
0	3529	23	1
1	3529	23	1
2	3533	22	1
3	3547	19	1
4	3557	16	1
5	3565	15	1
6	3572	14	1
7	3578	13	1
8	3584	13	1
9	3589	12	1
10	3593	12	1
11	3598	11	1
12	3602	11	1
13	3606	11	1
14	3610	11	1
15	3613	10	1
16	3617	10	1
17	3620	10	1
18	3624	10	1
19	3627	10	1
20	3630	10	2
21	3633	10	2
22	3636	10	2
23	3640	10	2
24	3643	10	2
25	3646	10	2
26	3649	10	2
27	3652	10	3
28	3655	10	3
29	3658	10	3
30	3662	10	3
31	3665	10	3
32	3668	10	3
33	3672	10	3
34	3676	11	3
35	3680	11	4
36	3683	11	4
37	3687	11	4
38	3692	12	4
39	3697	12	4
40	3702	13	4
41	3708	14	4
42	3714	15	4
43	3722	16	4
44	3732	18	4
45	3739	20	4
46	3739	20	4
47	3739	20	4

Table B.24. Raw-to-Scale Score Conversion, Mathematics Grade 8

Raw Score	Scale Score	CSEM	Performance Level
0	3566	22	1
1	3566	22	1
2	3566	22	1
3	3578	18	1
4	3588	16	1
5	3596	15	1
6	3603	14	1
7	3608	13	1
8	3614	12	1
9	3618	12	1
10	3623	11	1
11	3627	11	1
12	3631	11	1
13	3634	10	1
14	3638	10	1
15	3641	10	1
16	3644	10	1
17	3647	10	1
18	3650	10	2
19	3653	9	2
20	3656	9	2
21	3659	9	2
22	3662	9	2
23	3665	9	2
24	3668	9	2
25	3671	9	2
26	3674	9	3
27	3677	9	3
28	3680	9	3
29	3683	10	3
30	3686	10	3
31	3689	10	3
32	3692	10	3
33	3696	10	3
34	3699	10	3
35	3703	11	3
36	3707	11	4
37	3711	11	4
38	3715	12	4
39	3720	12	4
40	3725	13	4
41	3731	14	4
42	3737	15	4
43	3745	16	4
44	3755	18	4
45	3769	22	4
46	3776	25	4
47	3776	25	4

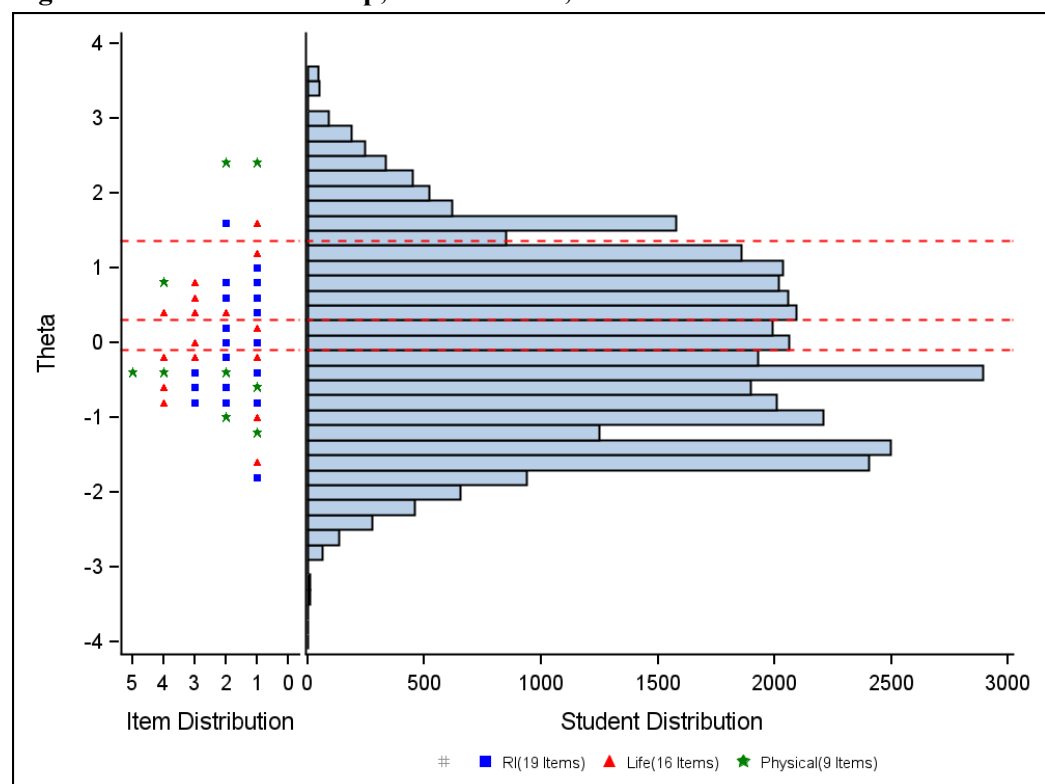
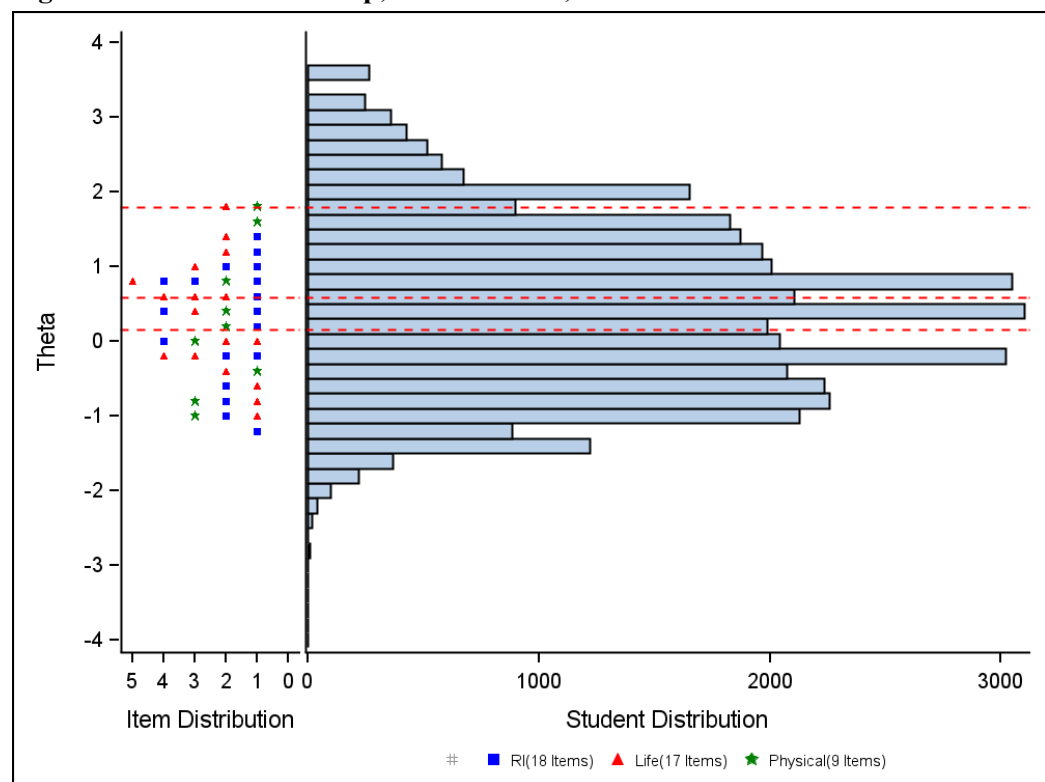
Figure B.1. Item-Person Map, ELA Grade 3, Form 1**Figure B.2. Item-Person Map, ELA Grade 4, Form 1**

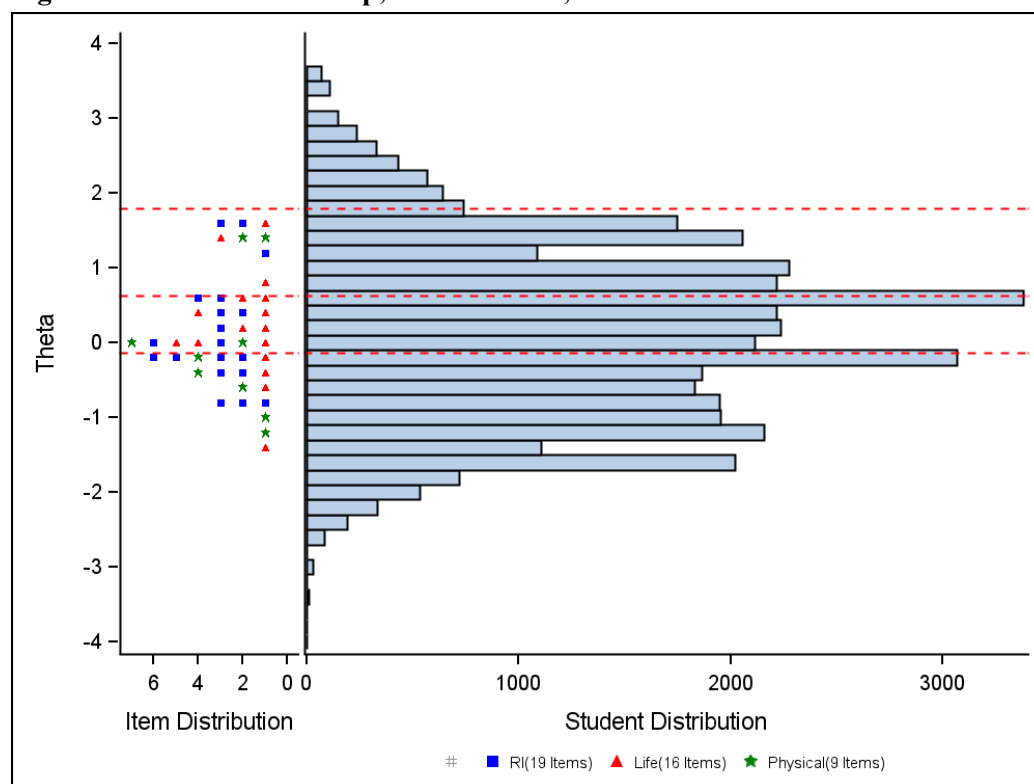
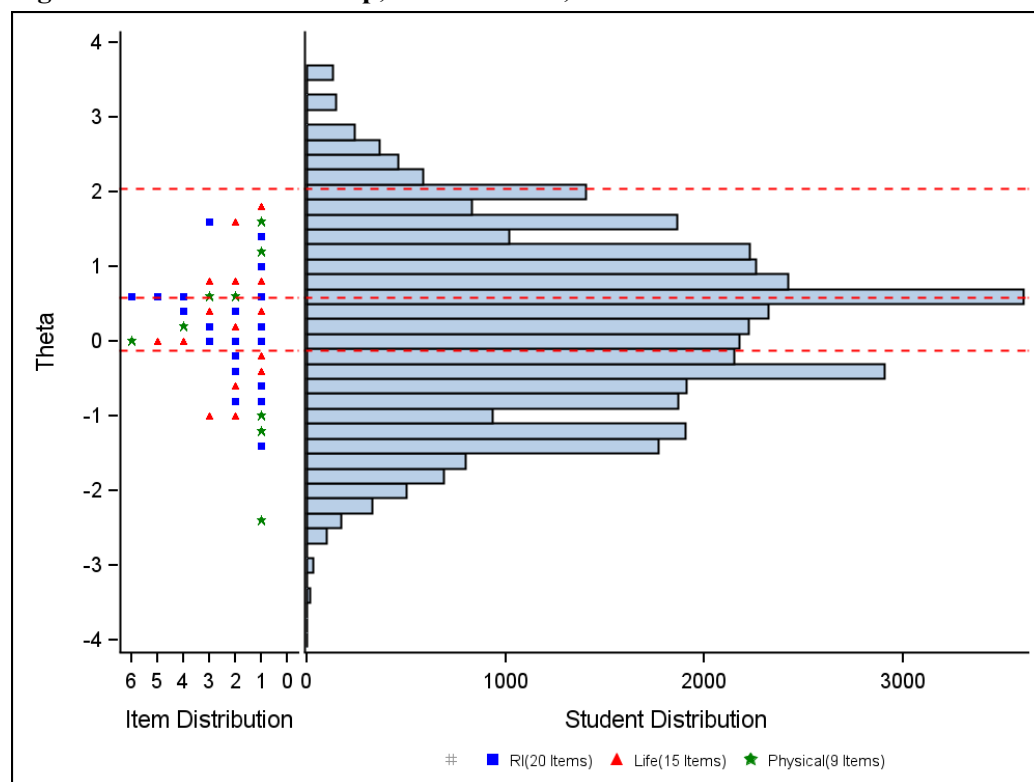
Figure B.3. Item-Person Map, ELA Grade 5, Form 1**Figure B.4. Item-Person Map, ELA Grade 6, Form 1**

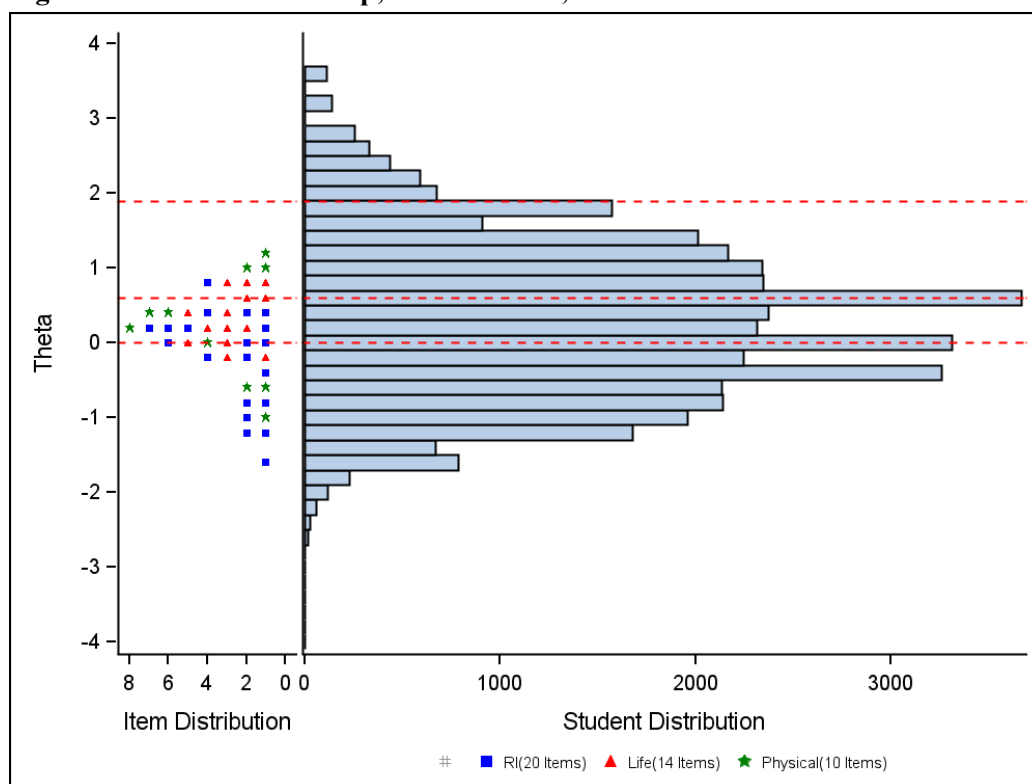
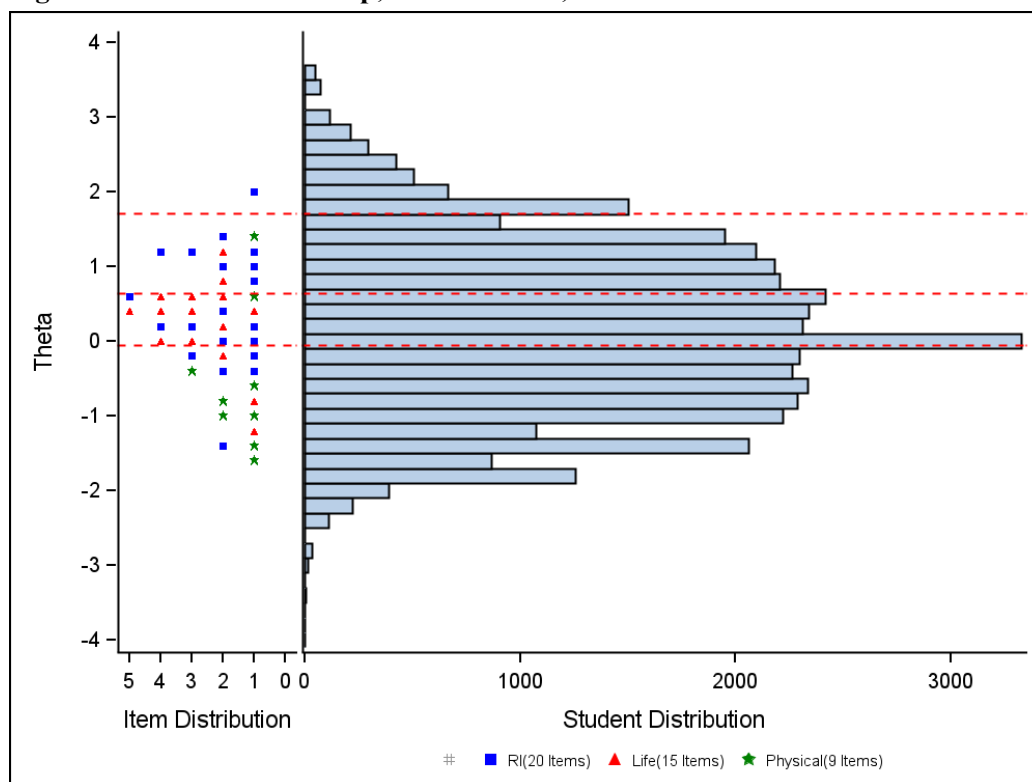
Figure B.5. Item-Person Map, ELA Grade 7, Form 1**Figure B.6. Item-Person Map, ELA Grade 8, Form 1**

Figure B.7. Item-Person Map, ELA Grade 3, Form 2

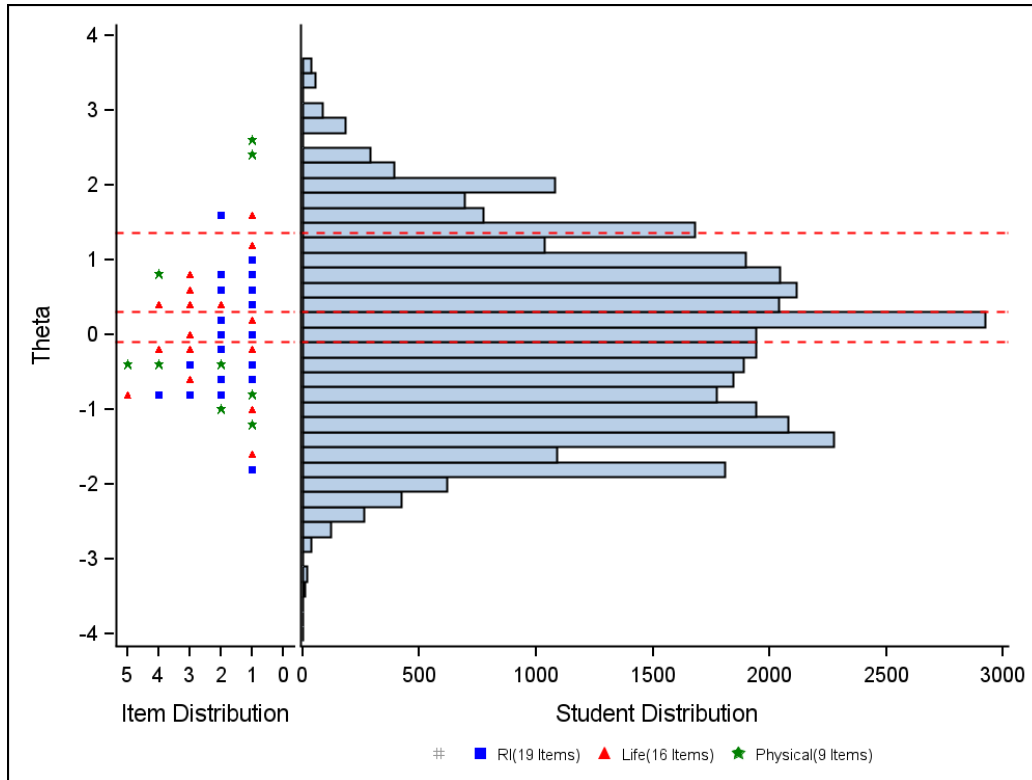


Figure B.8. Item-Person Map, ELA Grade 4, Form 2

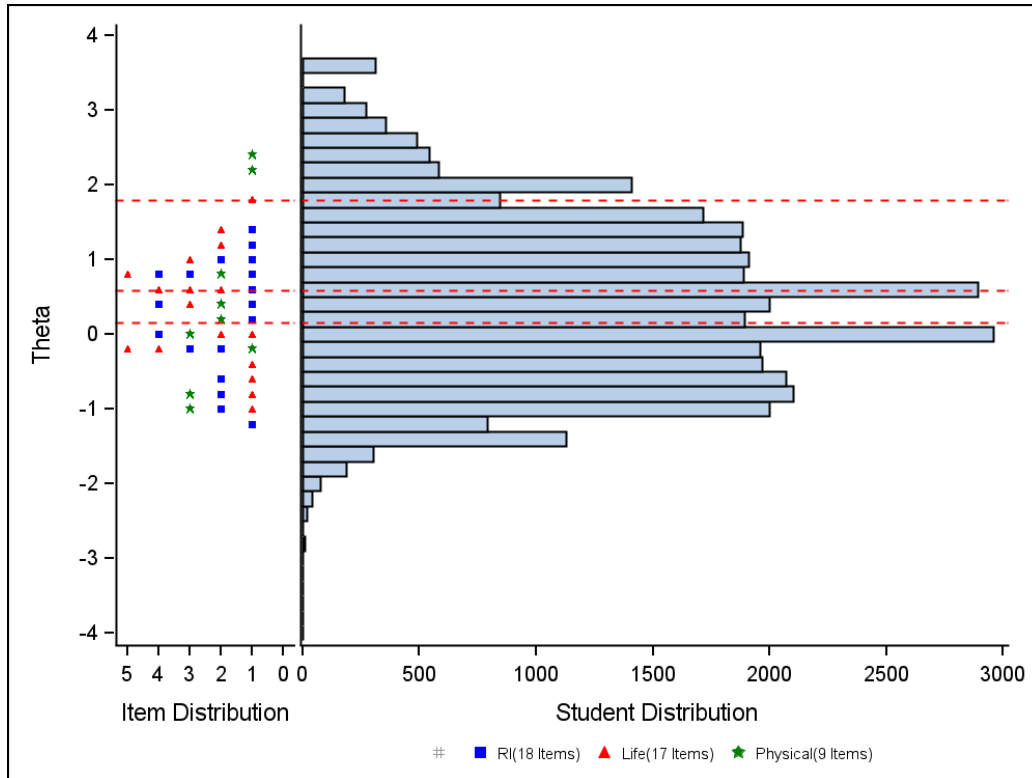


Figure B.9. Item-Person Map, ELA Grade 5, Form 2

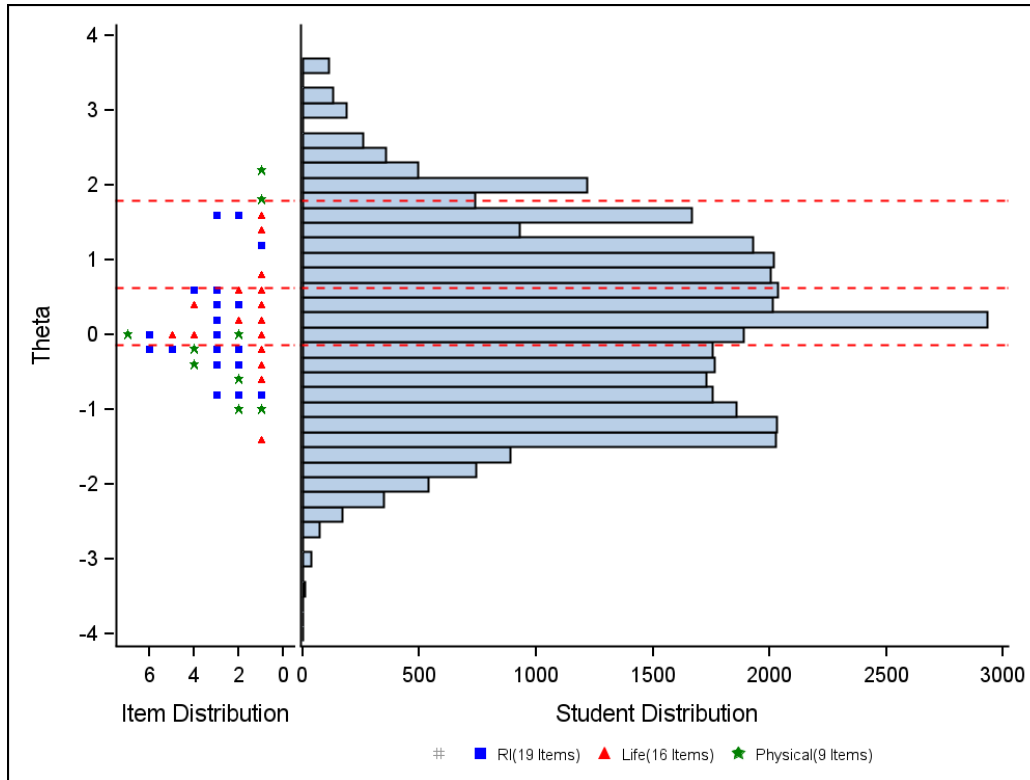


Figure B.10. Item-Person Map, ELA Grade 6, Form 2

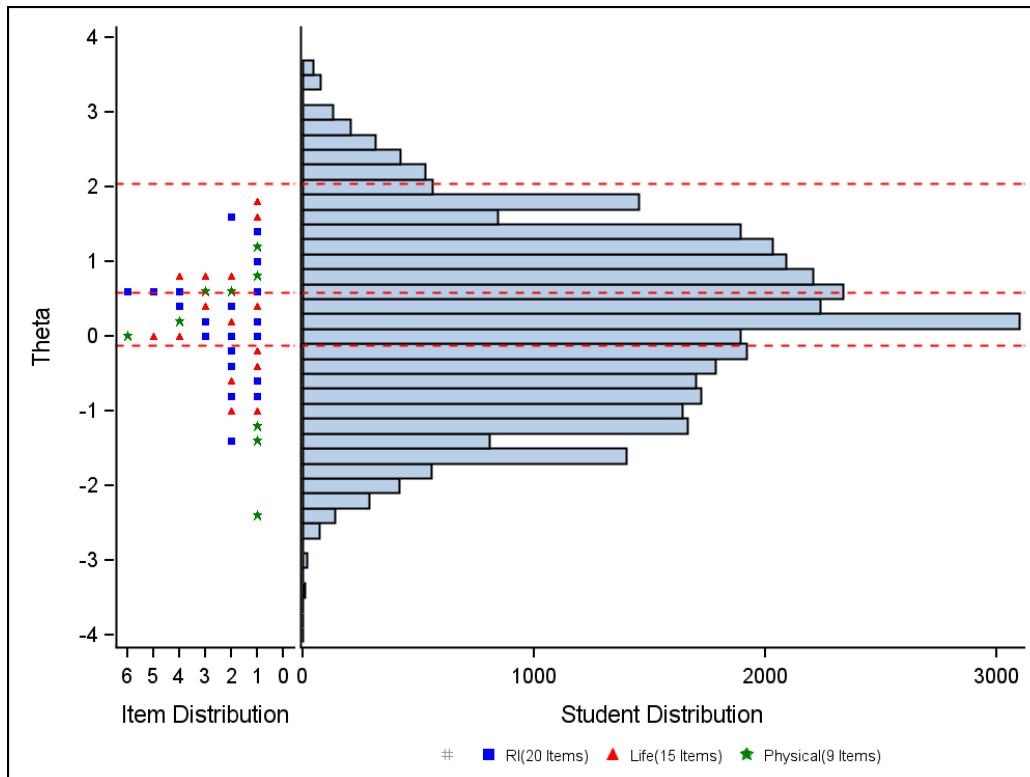


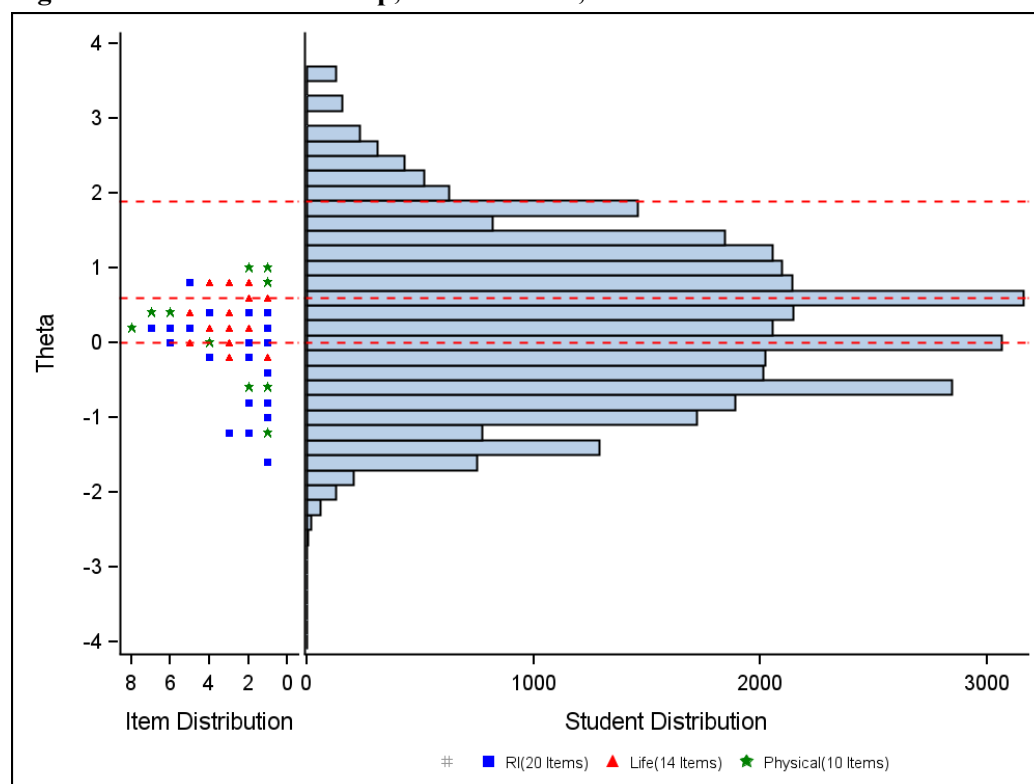
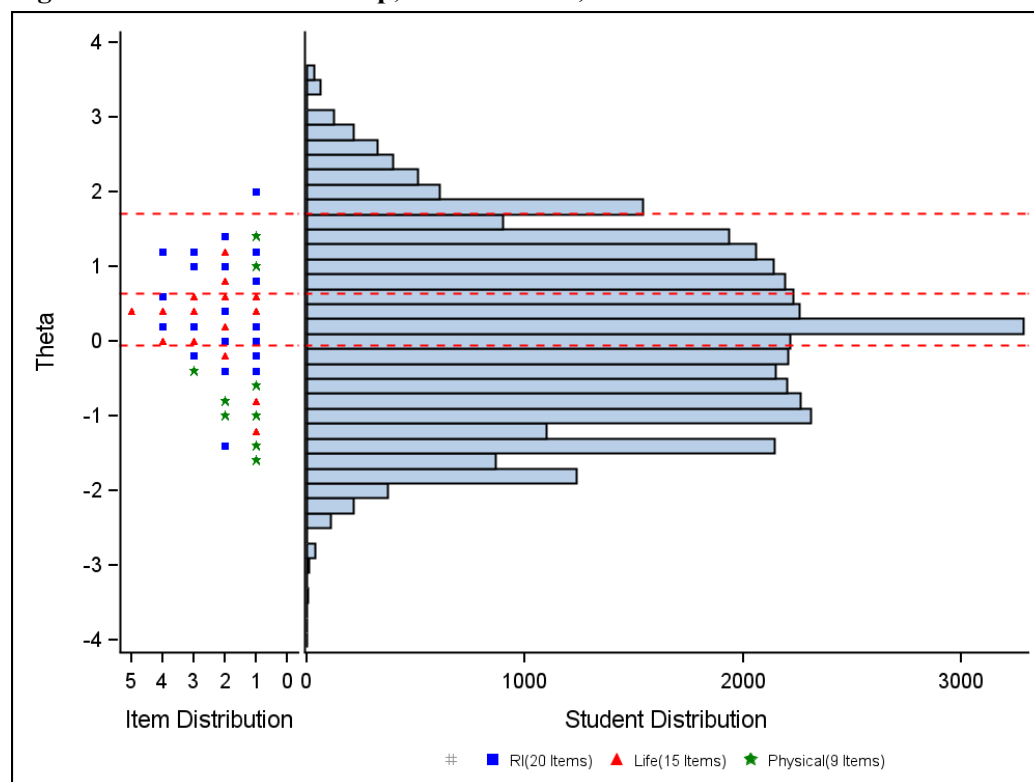
Figure B.11. Item-Person Map, ELA Grade 7, Form 2**Figure B.12. Item-Person Map, ELA Grade 8, Form 2**

Figure B.13. Item-Person Map, Mathematics Grade 3

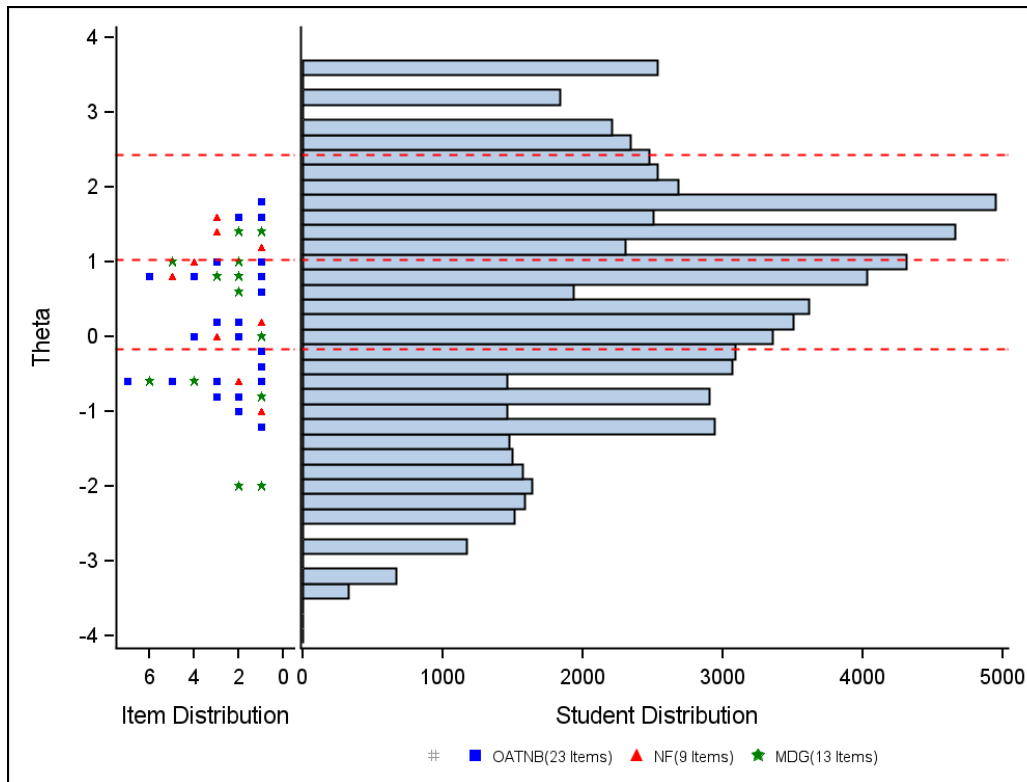


Figure B.14. Item-Person Map, Mathematics Grade 4

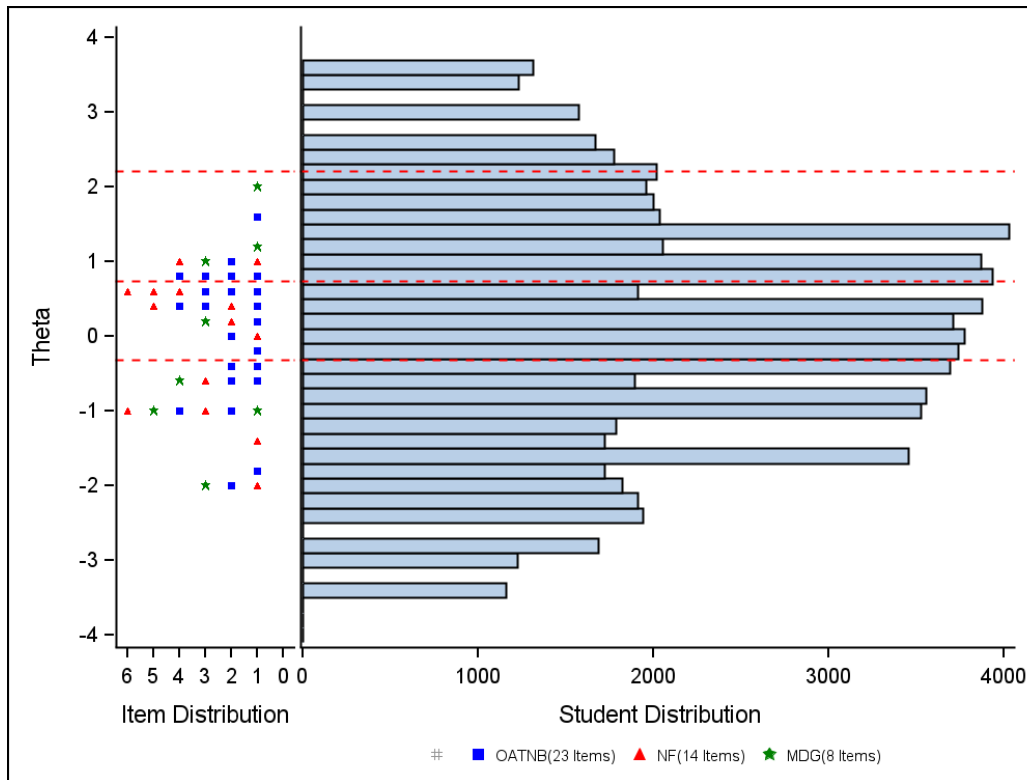


Figure B.15. Item-Person Map, Mathematics Grade 5

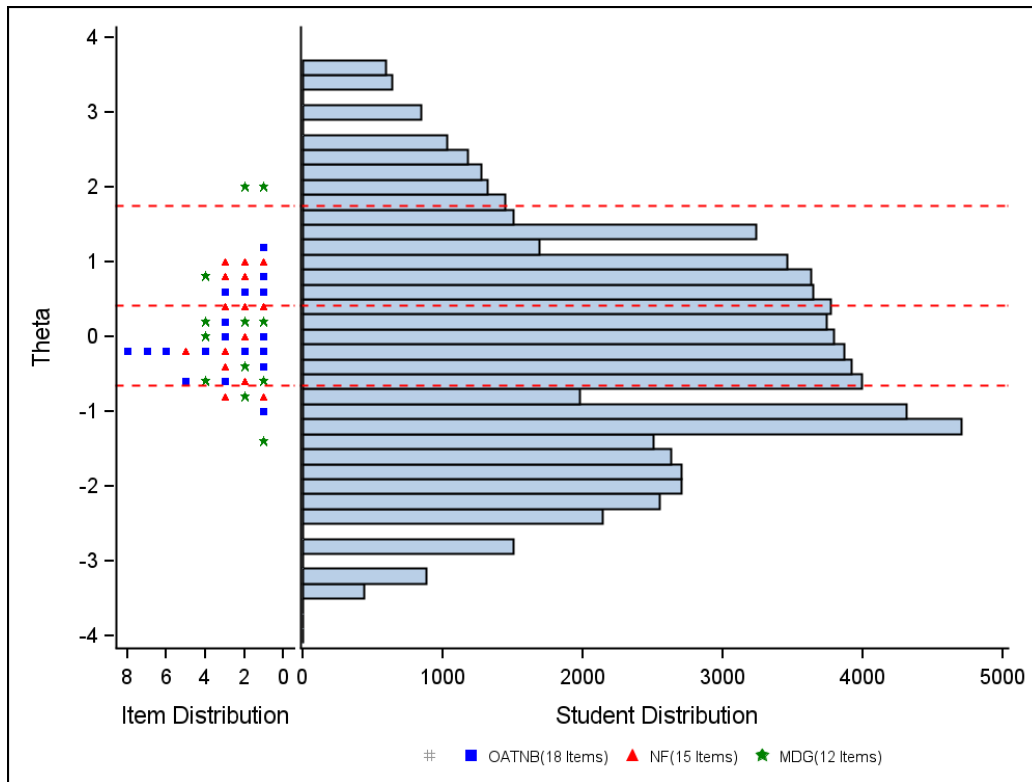


Figure B.16. Item-Person Map, Mathematics Grade 6

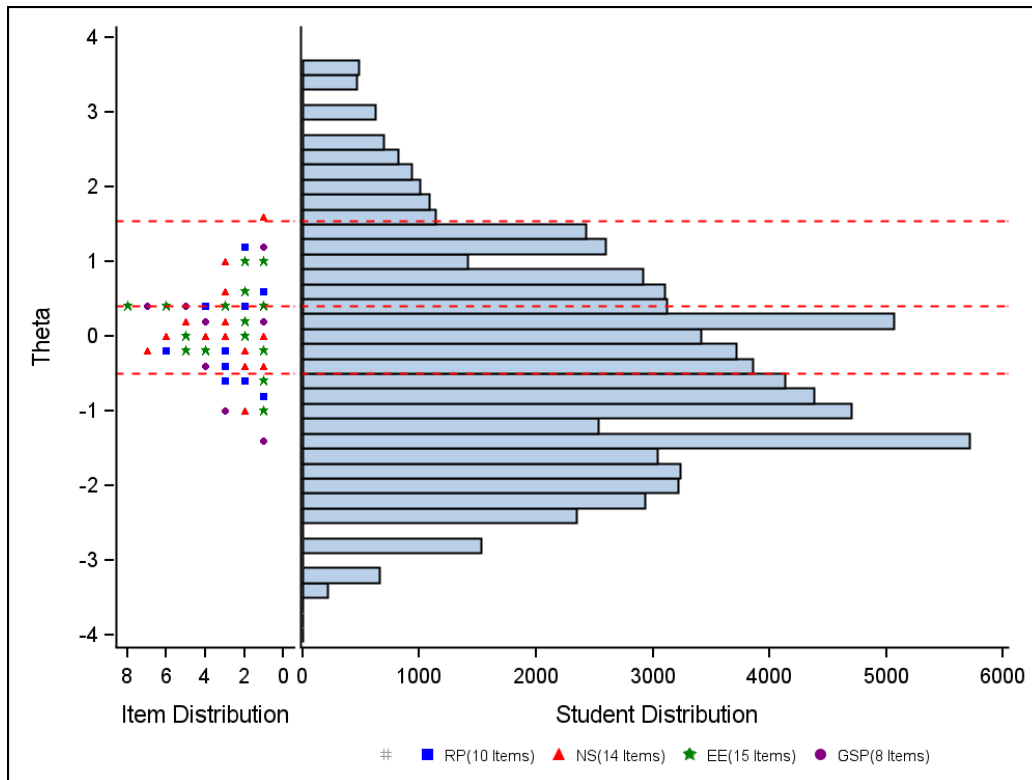


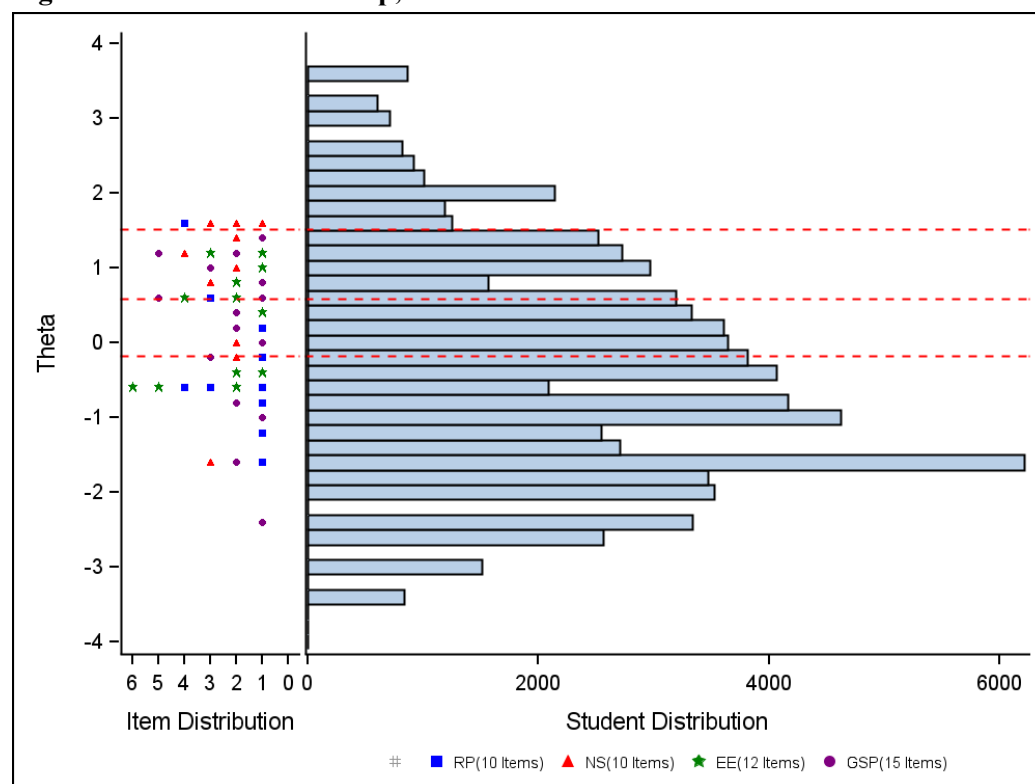
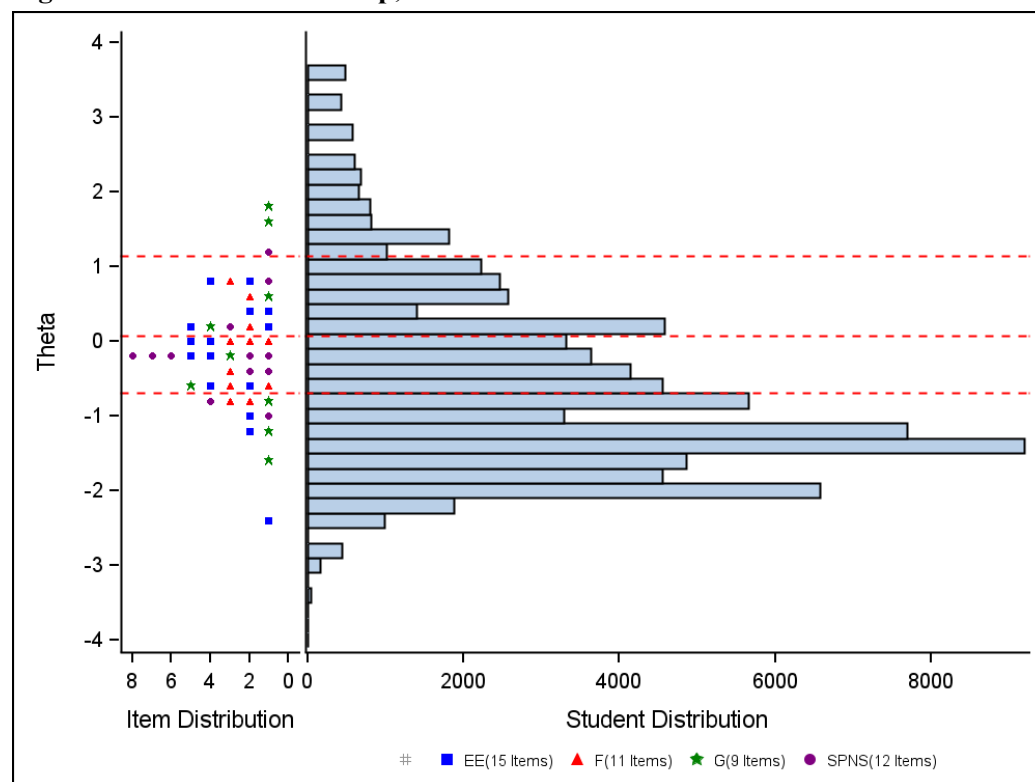
Figure B.17. Item-Person Map, Mathematics Grade 7**Figure B.18. Item-Person Map, Mathematics Grade 8**

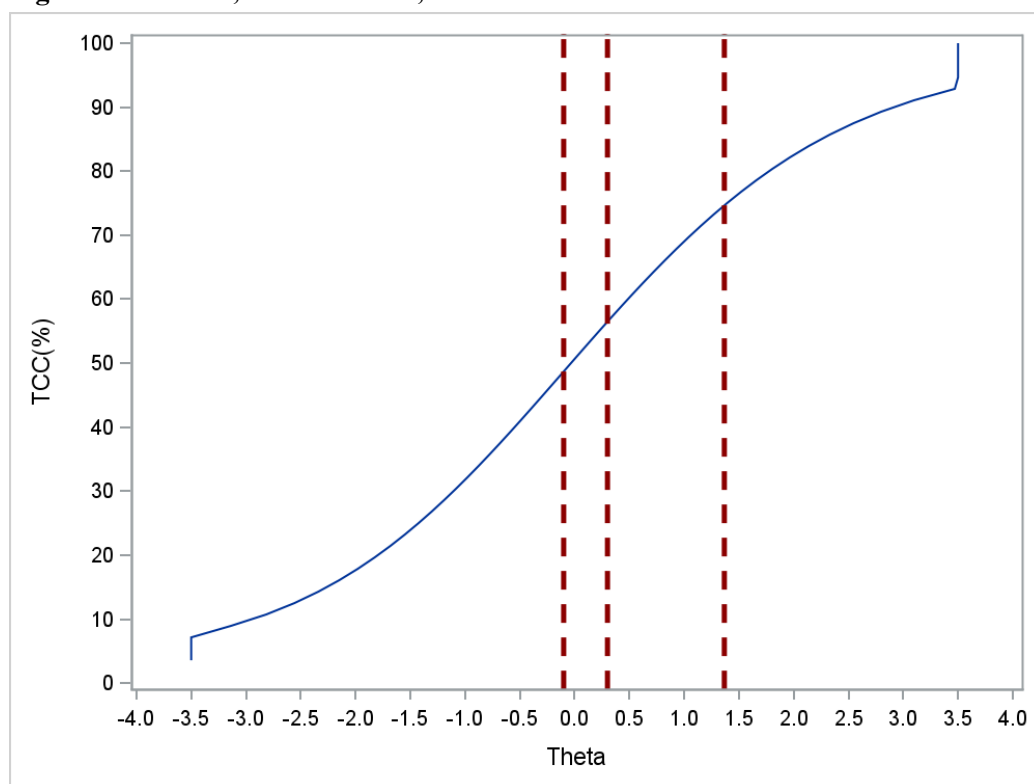
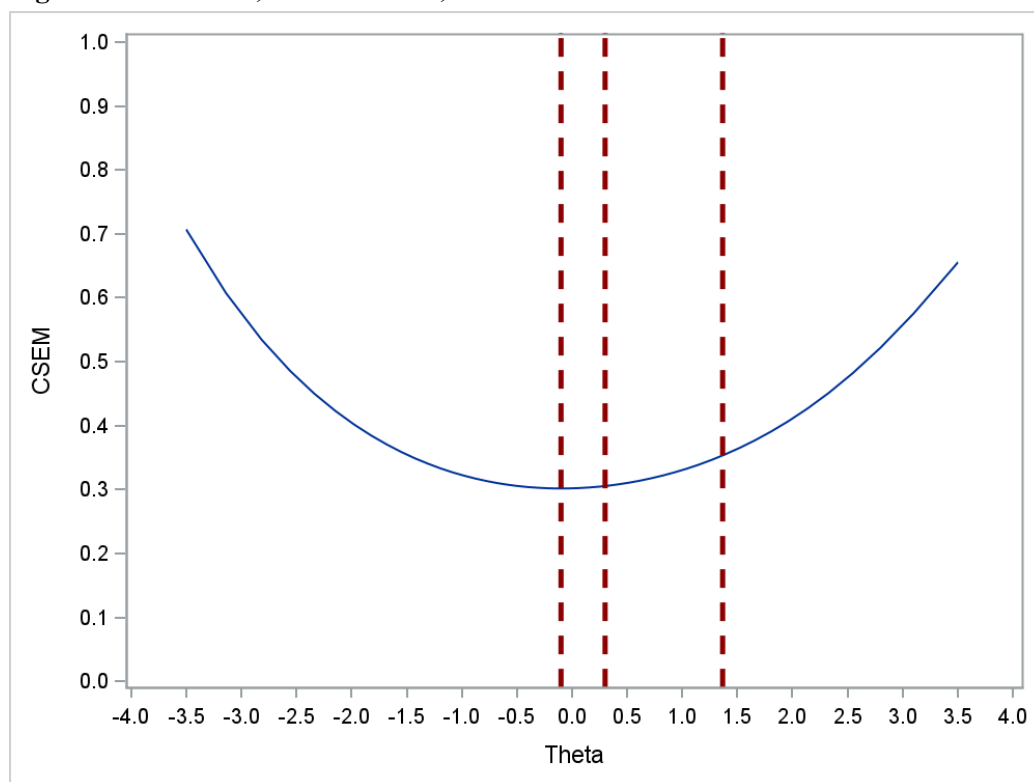
Figure B.19. TCC, ELA Grade 3, Form 1**Figure B.20. CSEM, ELA Grade 3, Form 1**

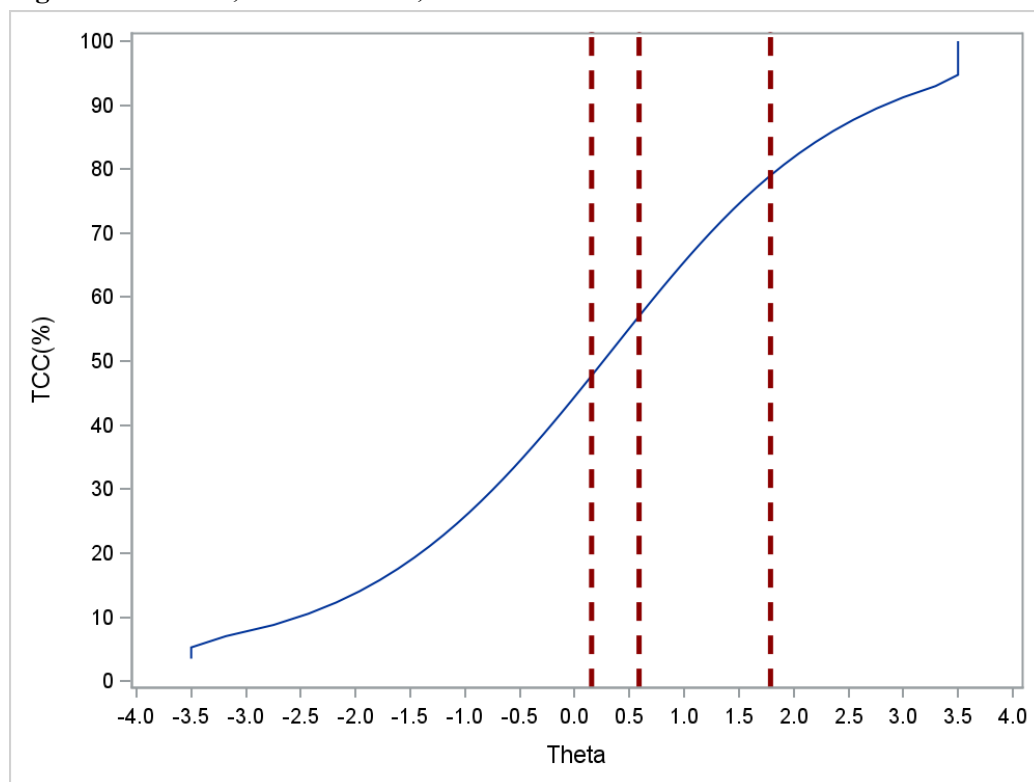
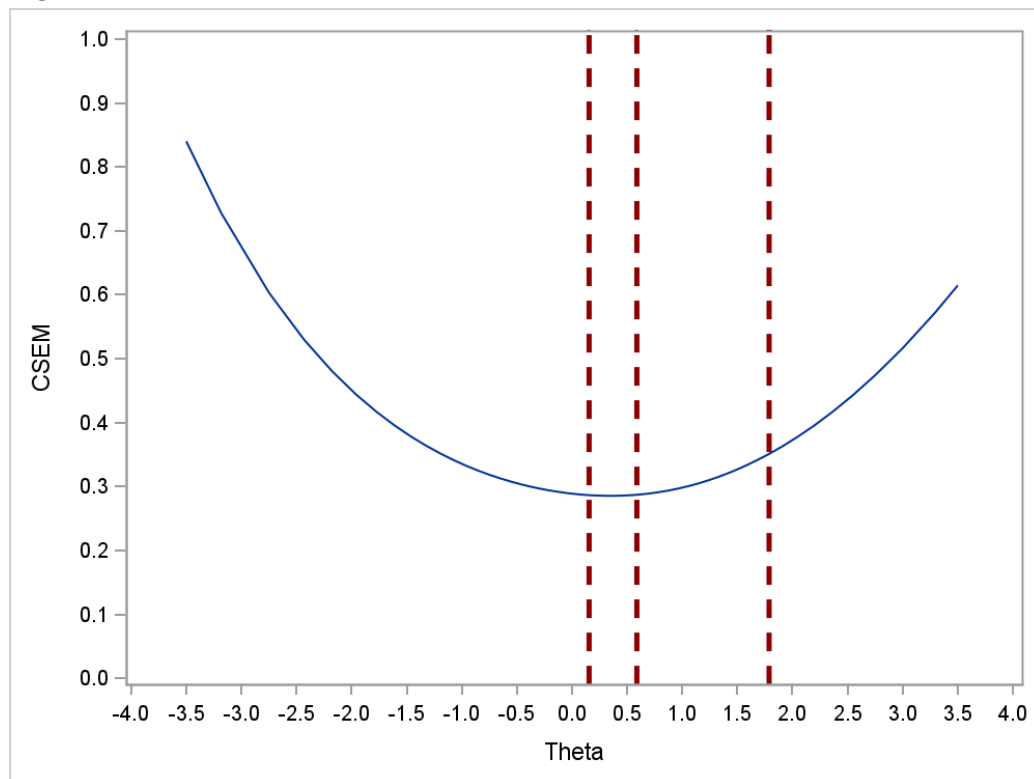
Figure B.21. TCC, ELA Grade 4, Form 1**Figure B.22. CSEM, ELA Grade 4, Form 1**

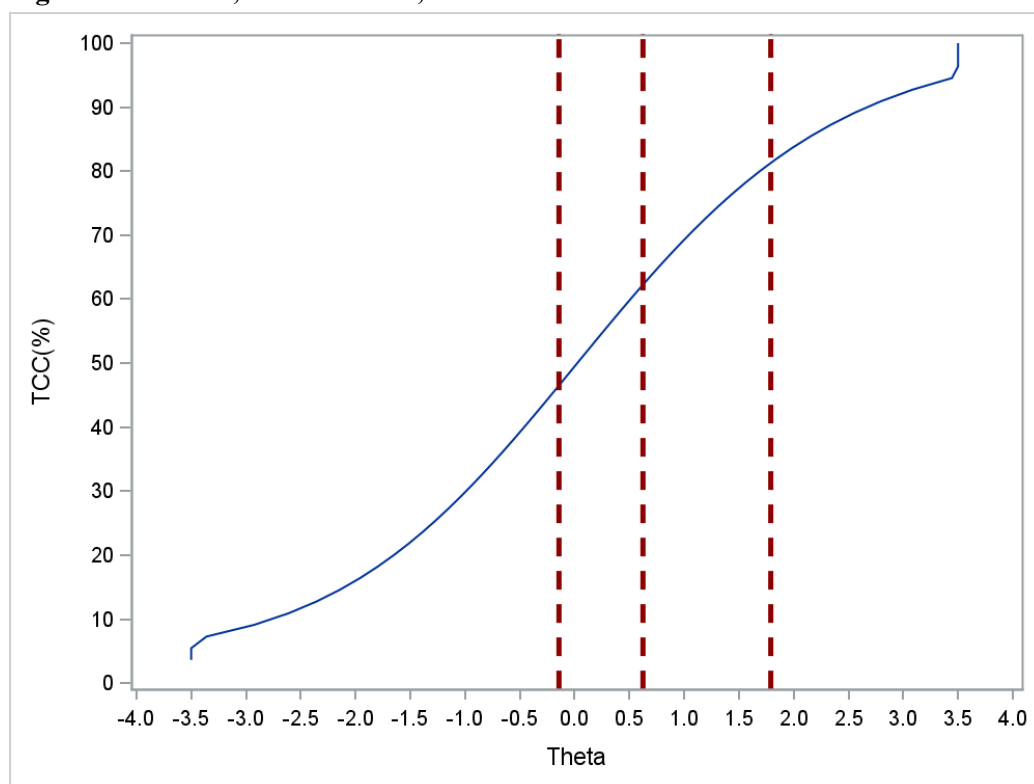
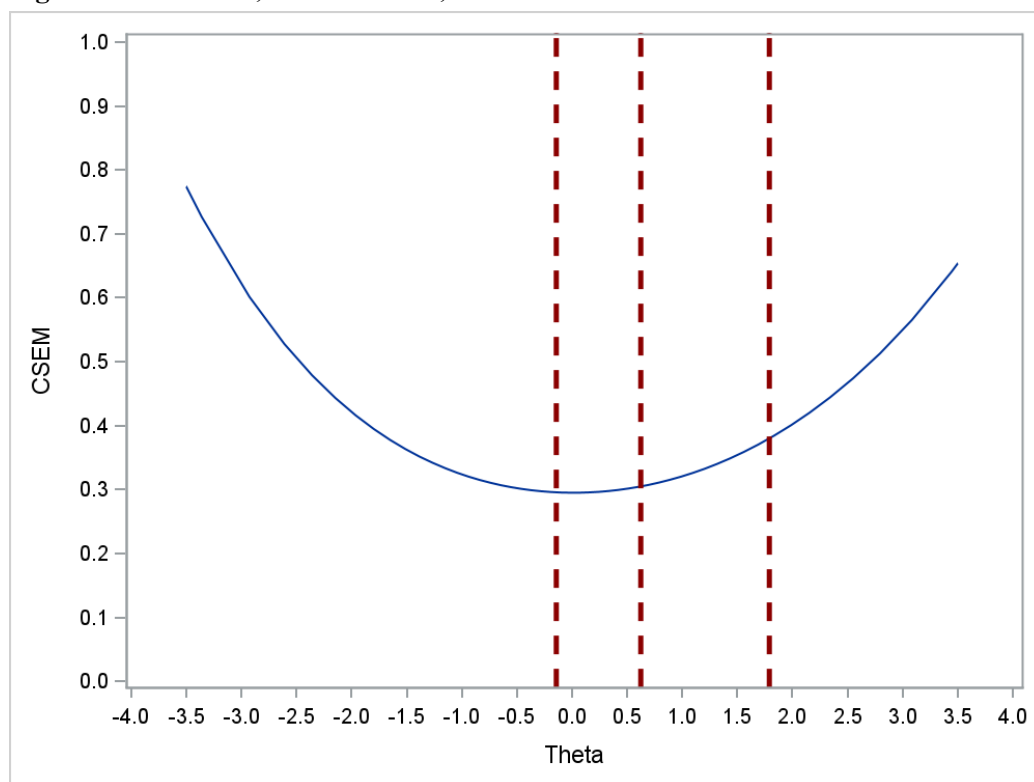
Figure B.23. TCC, ELA Grade 5, Form 1**Figure B.24. CSEM, ELA Grade 5, Form 1**

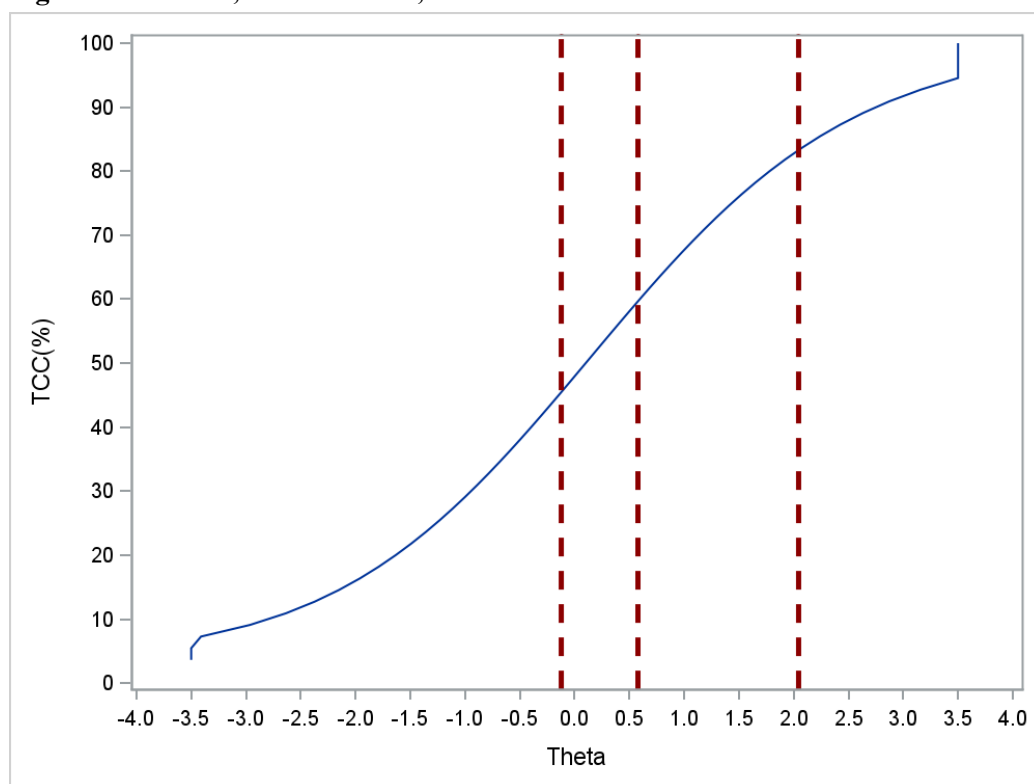
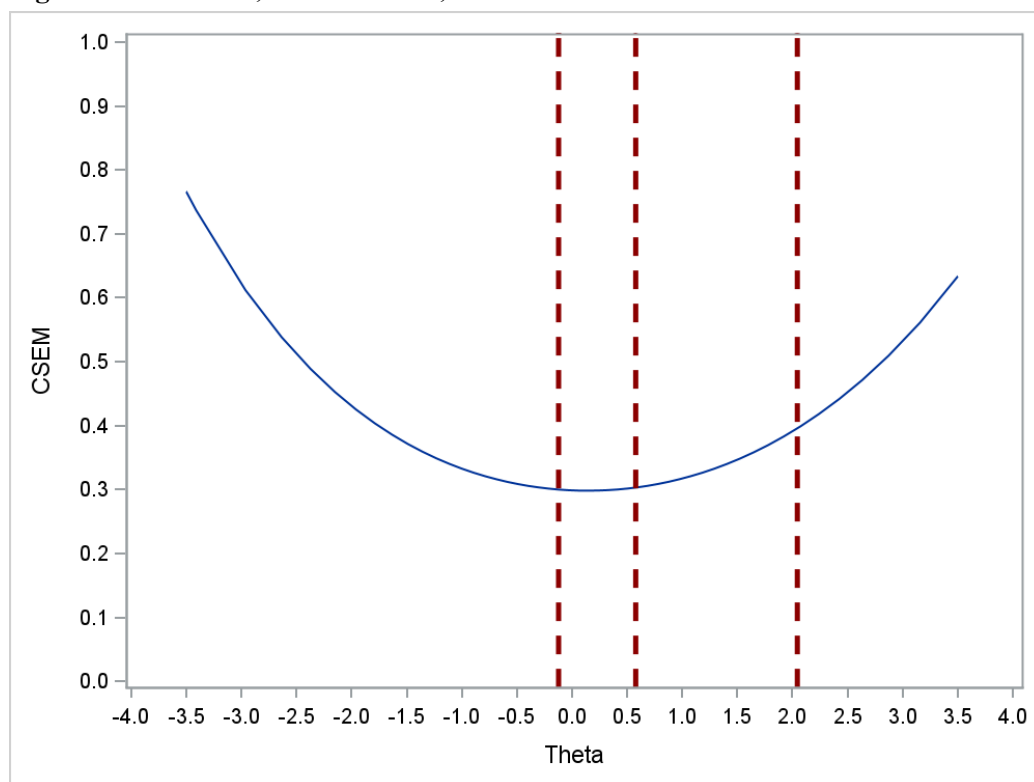
Figure B.25. TCC, ELA Grade 6, Form 1**Figure B.26. CSEM, ELA Grade 6, Form 1**

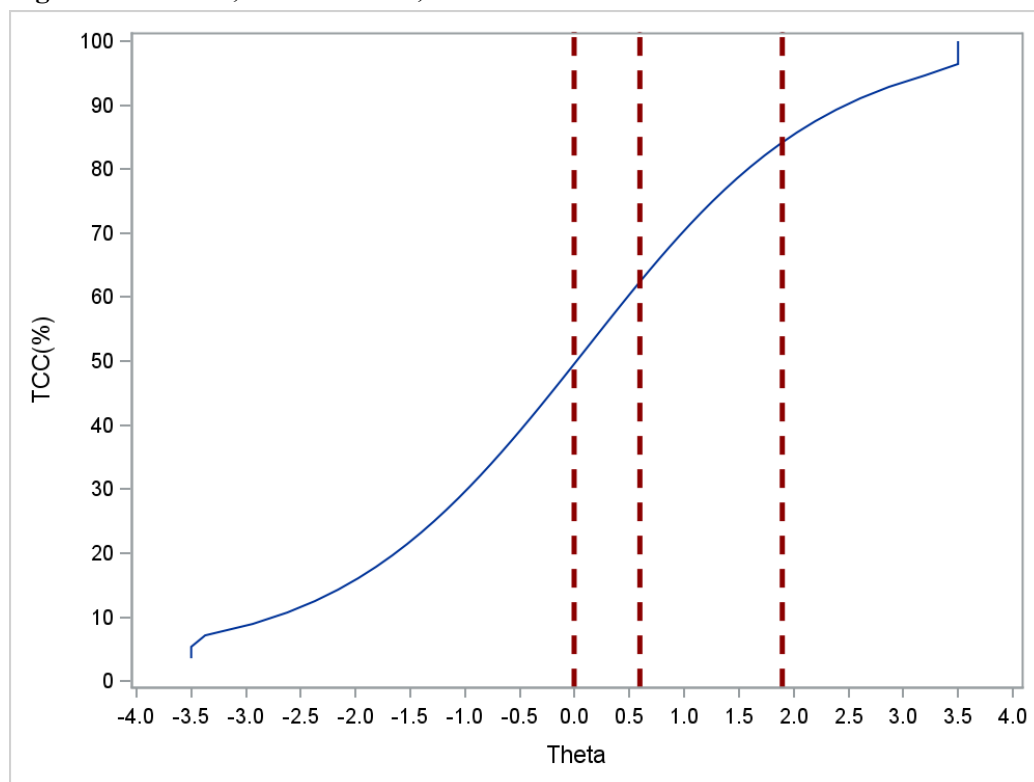
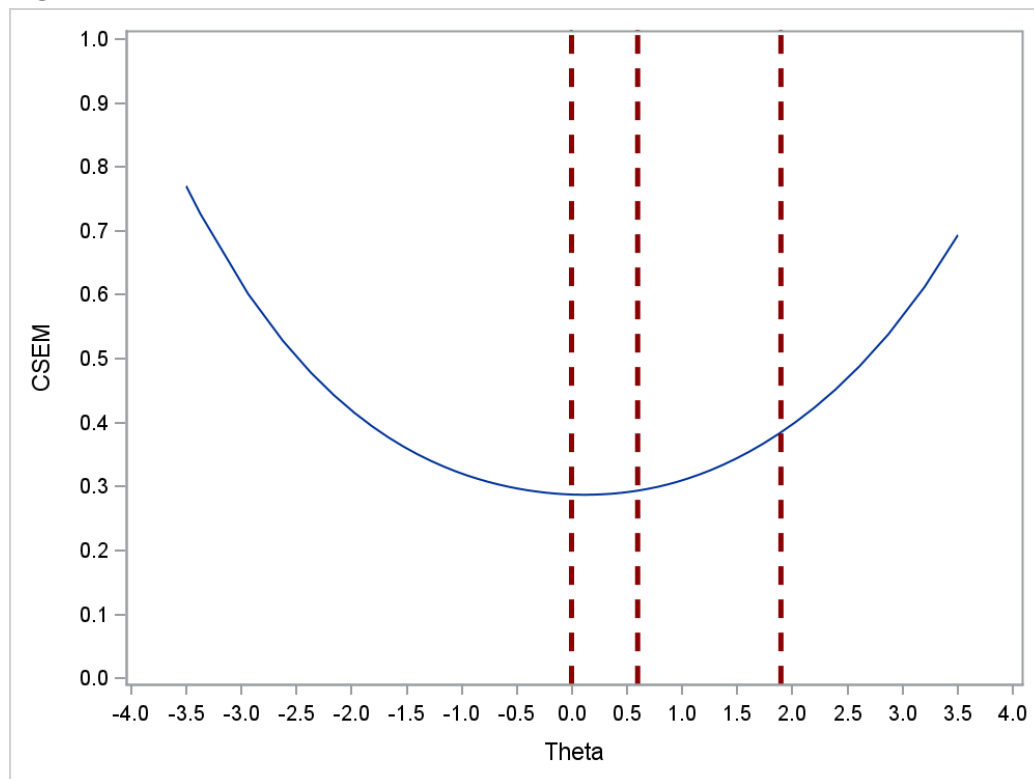
Figure B.27. TCC, ELA Grade 7, Form 1**Figure B.28. CSEM, ELA Grade 7, Form 1**

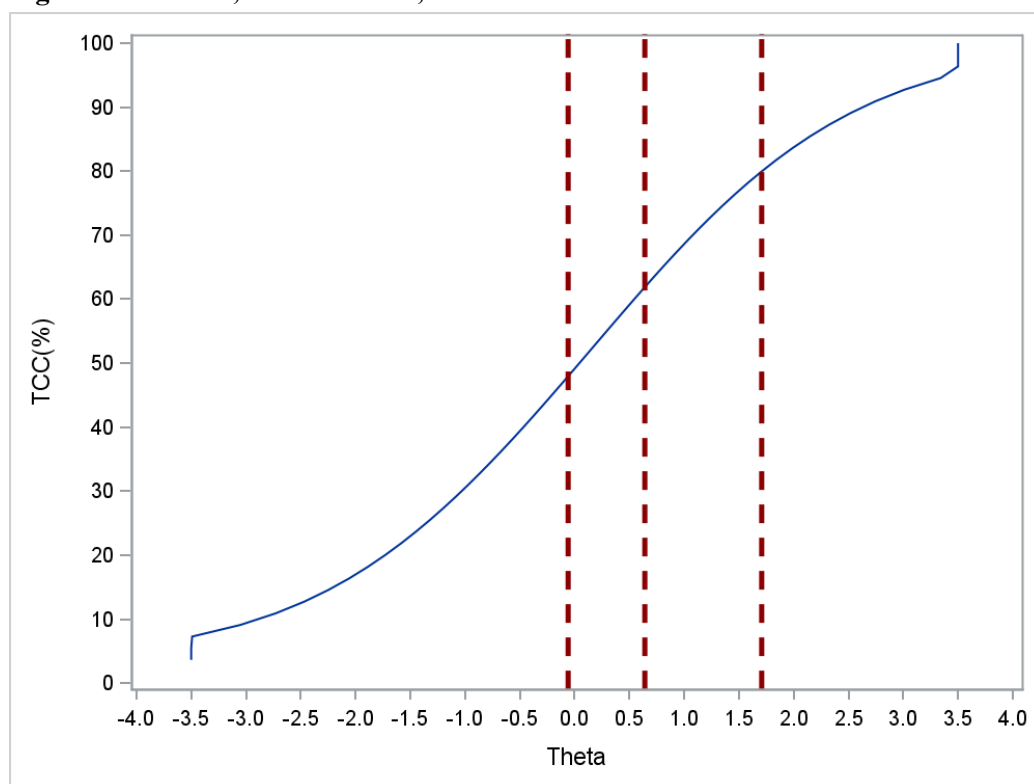
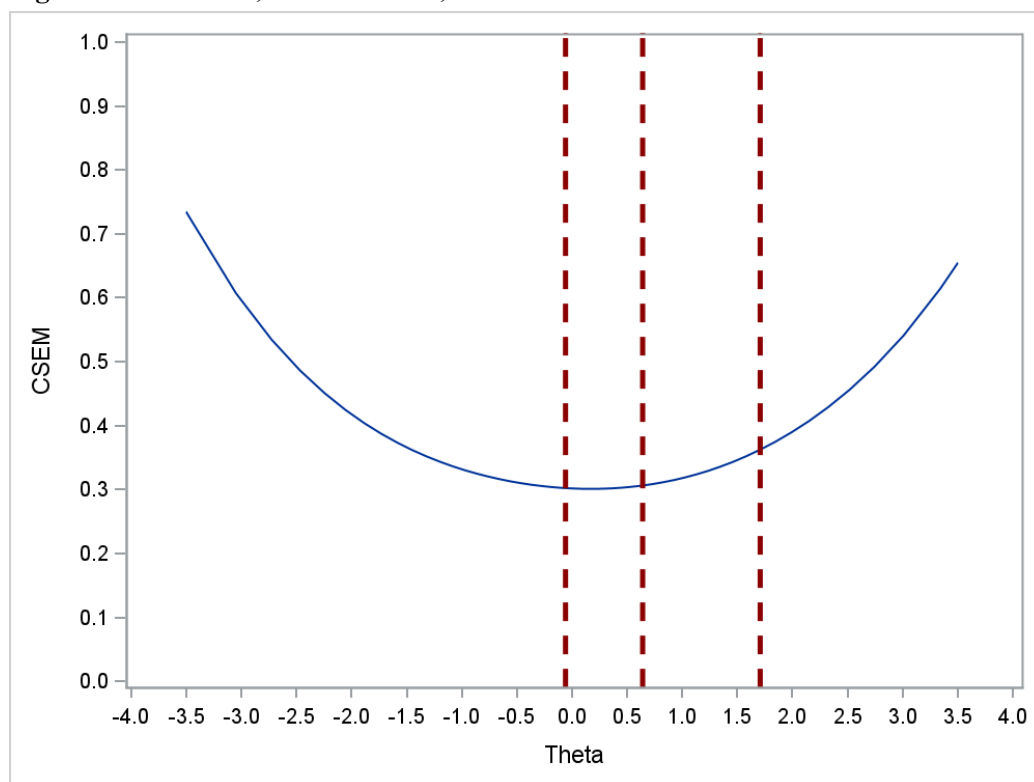
Figure B.29. TCC, ELA Grade 8, Form 1**Figure B.30. CSEM, ELA Grade 8, Form 1**

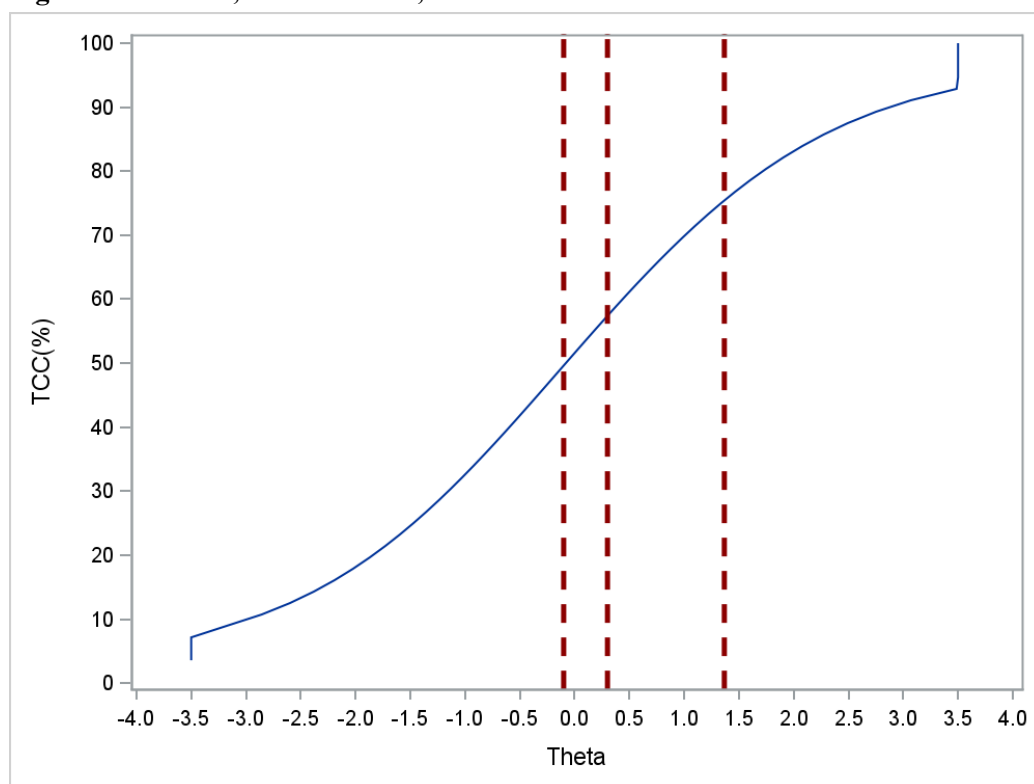
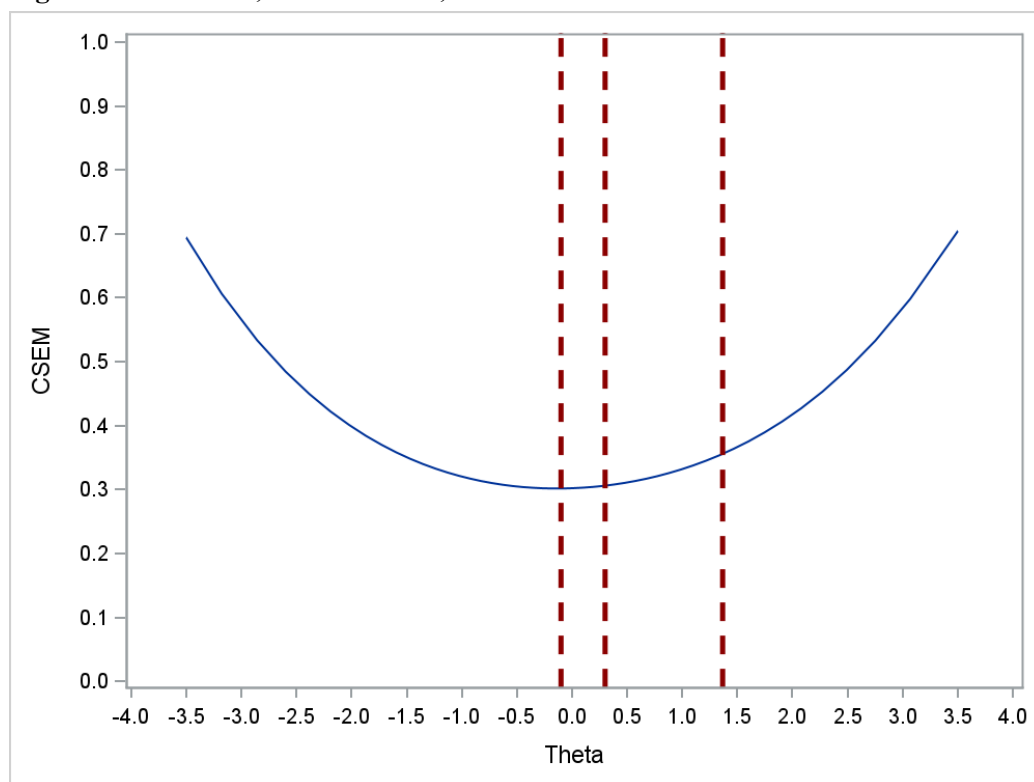
Figure B.31. TCC, ELA Grade 3, Form 2**Figure B.32. CSEM, ELA Grade 3, Form 2**

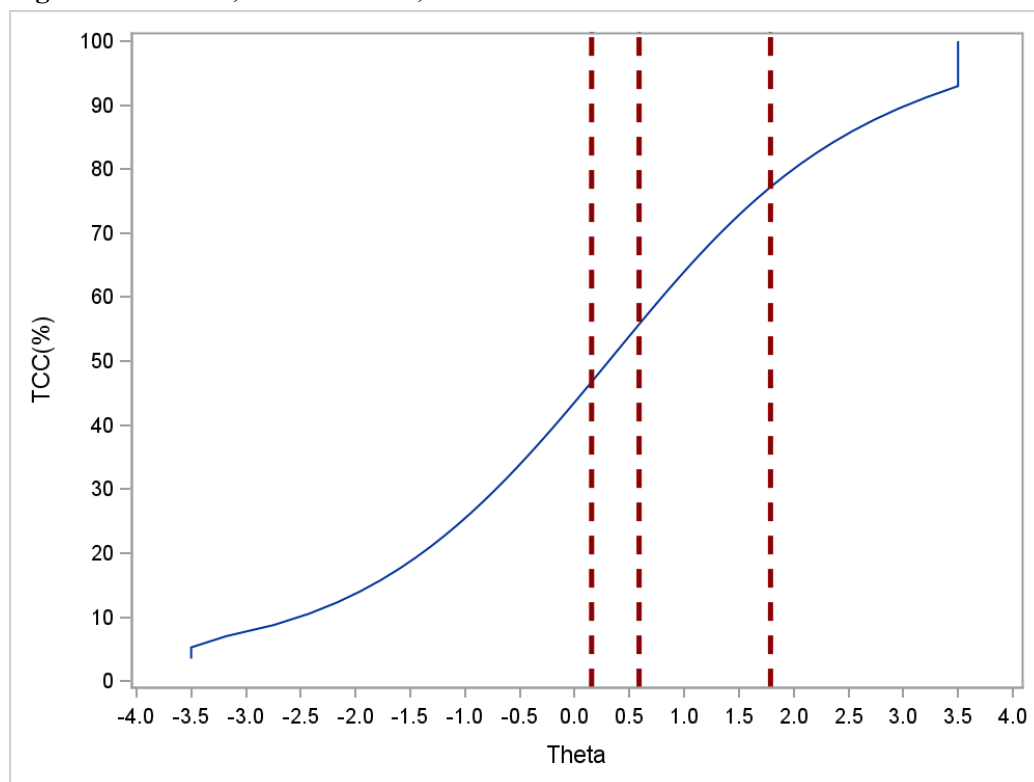
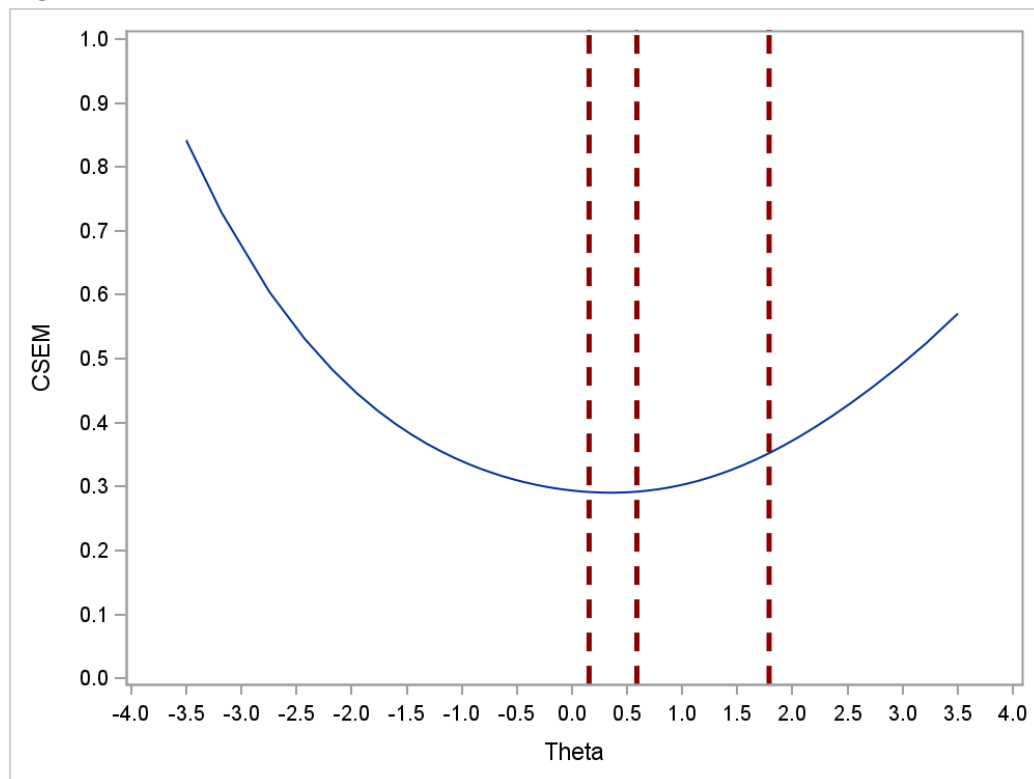
Figure B.33. TCC, ELA Grade 4, Form 2**Figure B.34. CSEM, ELA Grade 4, Form 2**

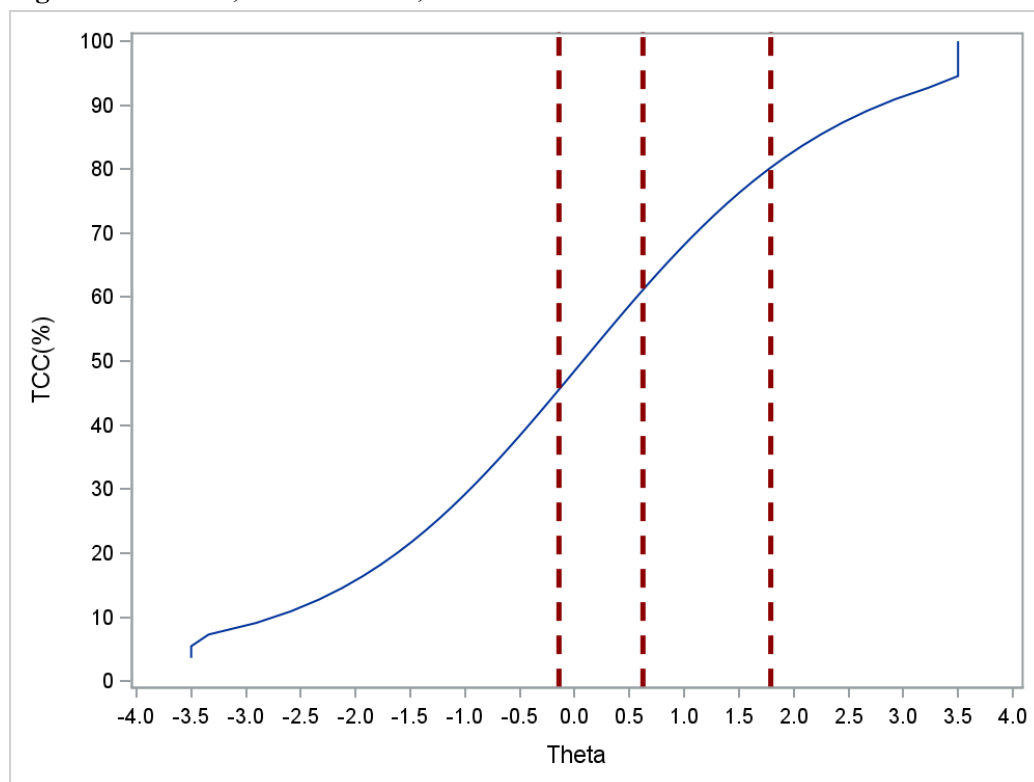
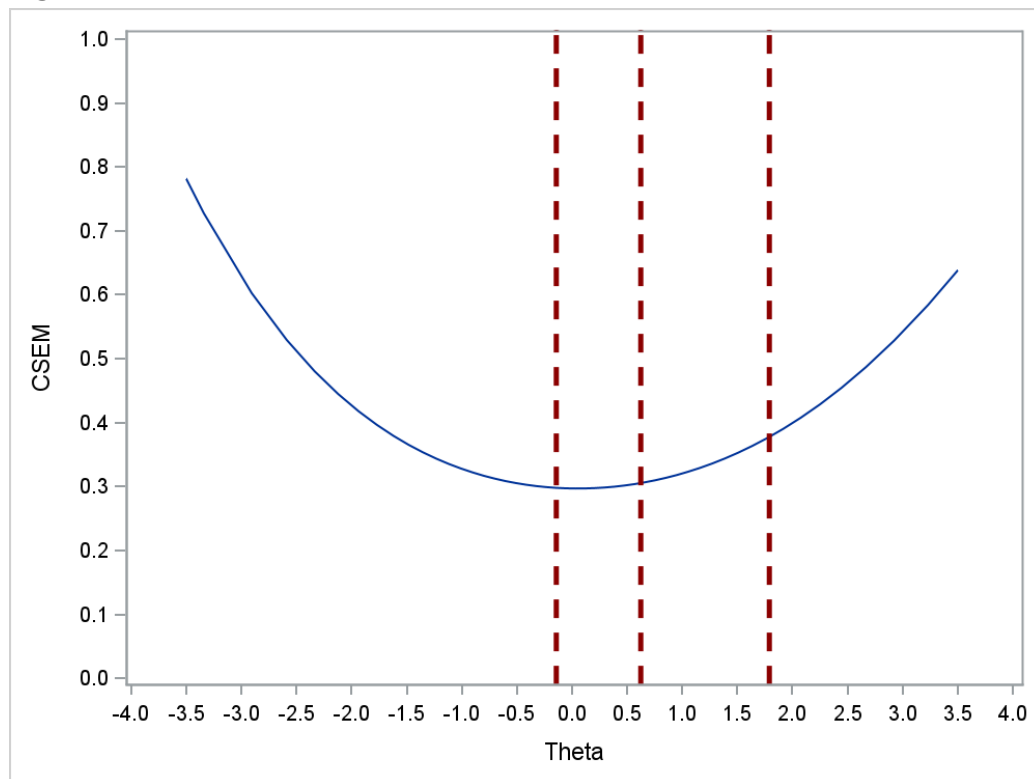
Figure B.35. TCC, ELA Grade 5, Form 2**Figure B.36. CSEM, ELA Grade 5, Form 2**

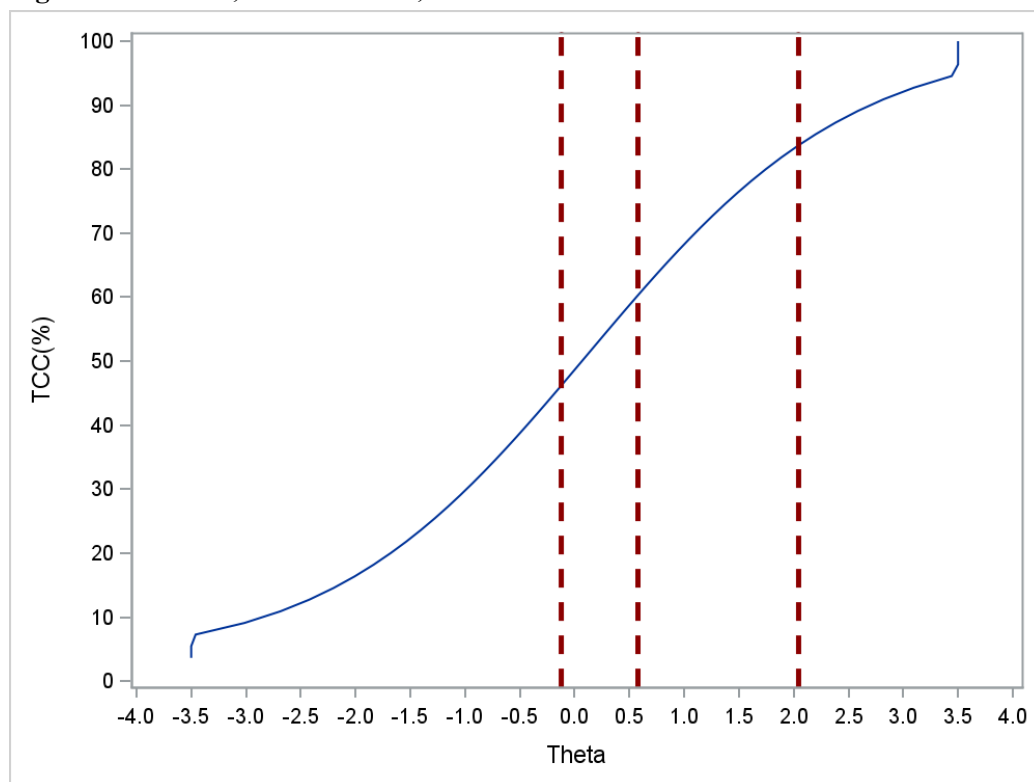
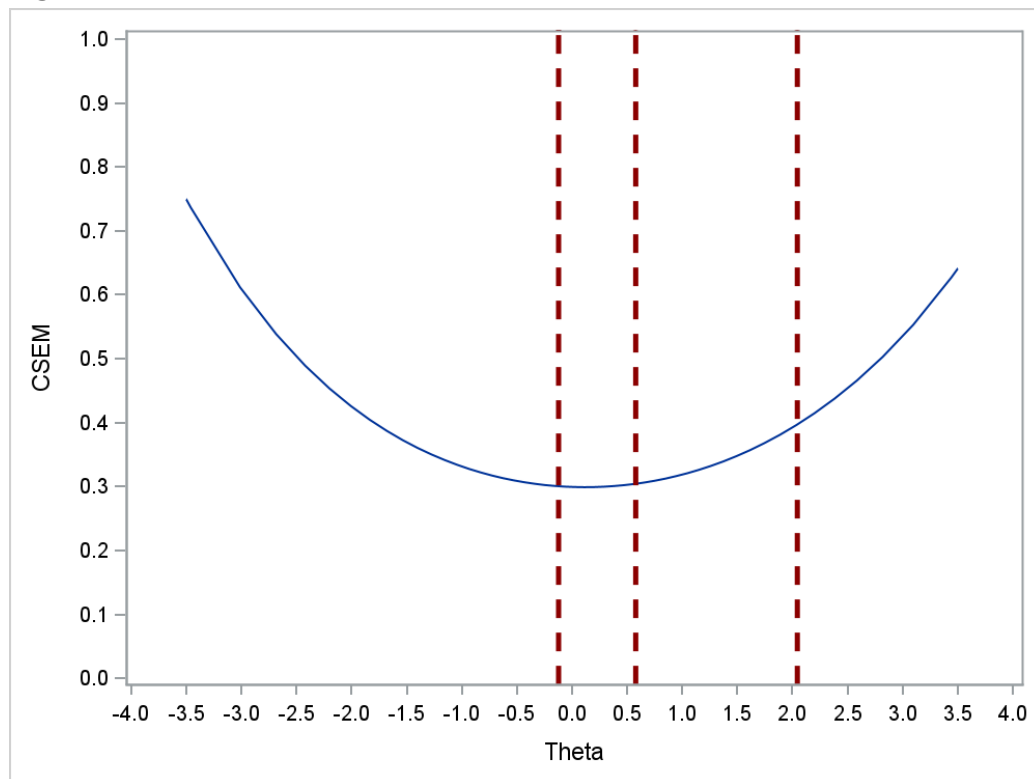
Figure B.37. TCC, ELA Grade 6, Form 2**Figure B.38. CSEM, ELA Grade 6, Form 2**

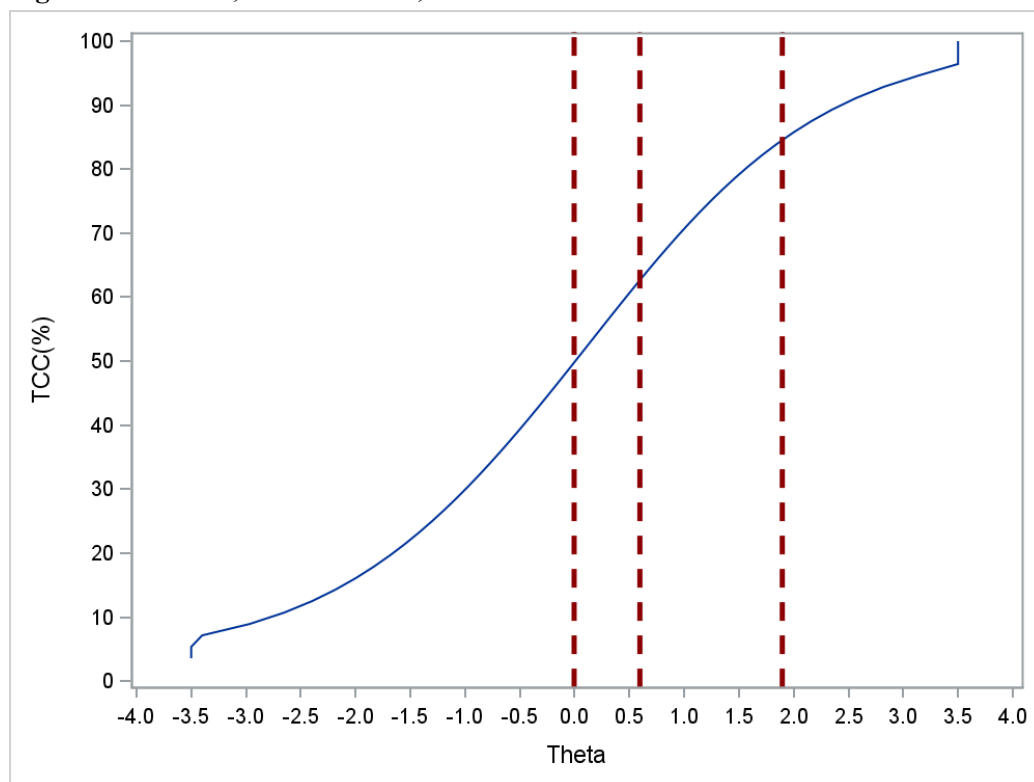
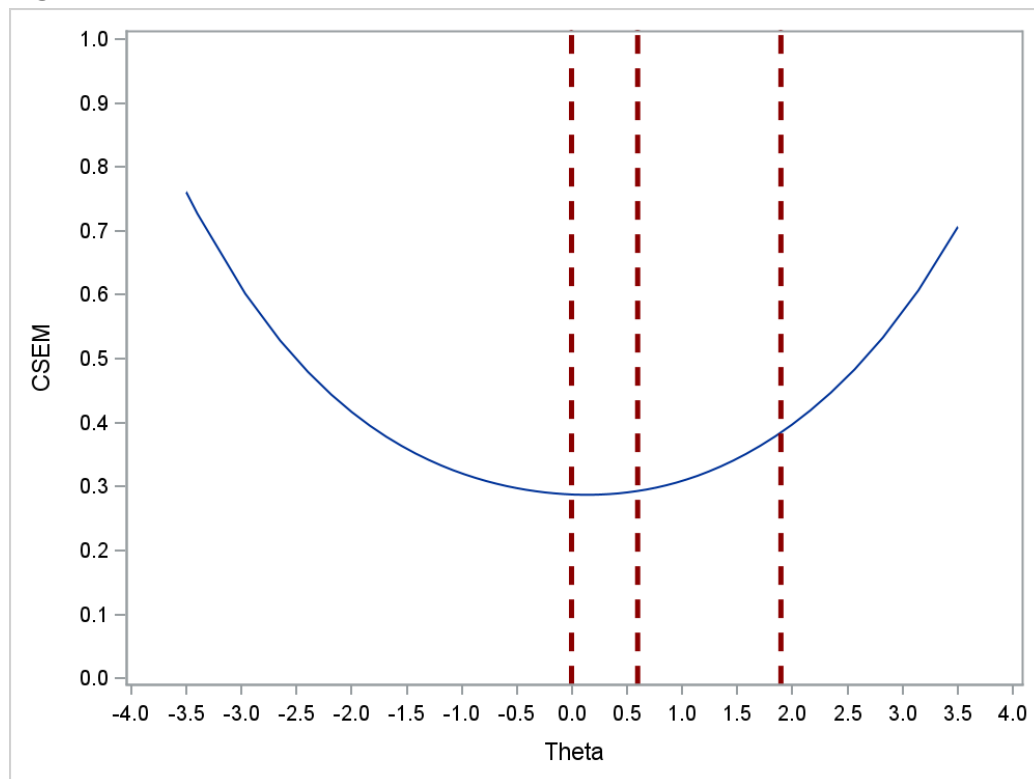
Figure B.39. TCC, ELA Grade 7, Form 2**Figure B.40. CSEM, ELA Grade 7, Form 2**

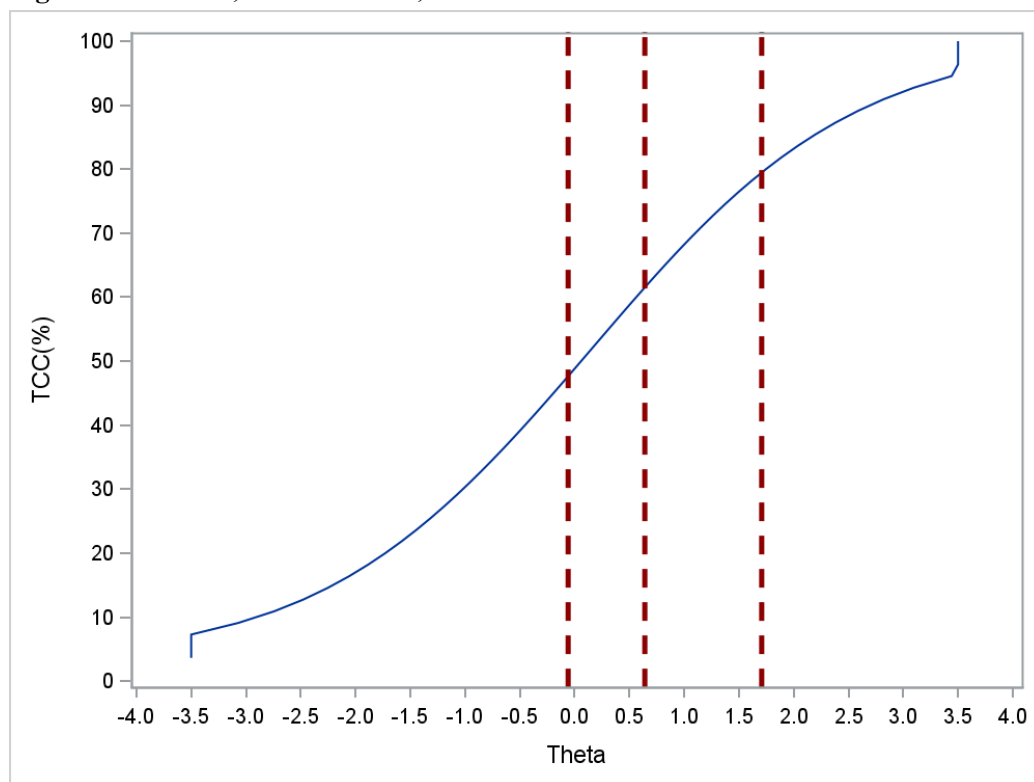
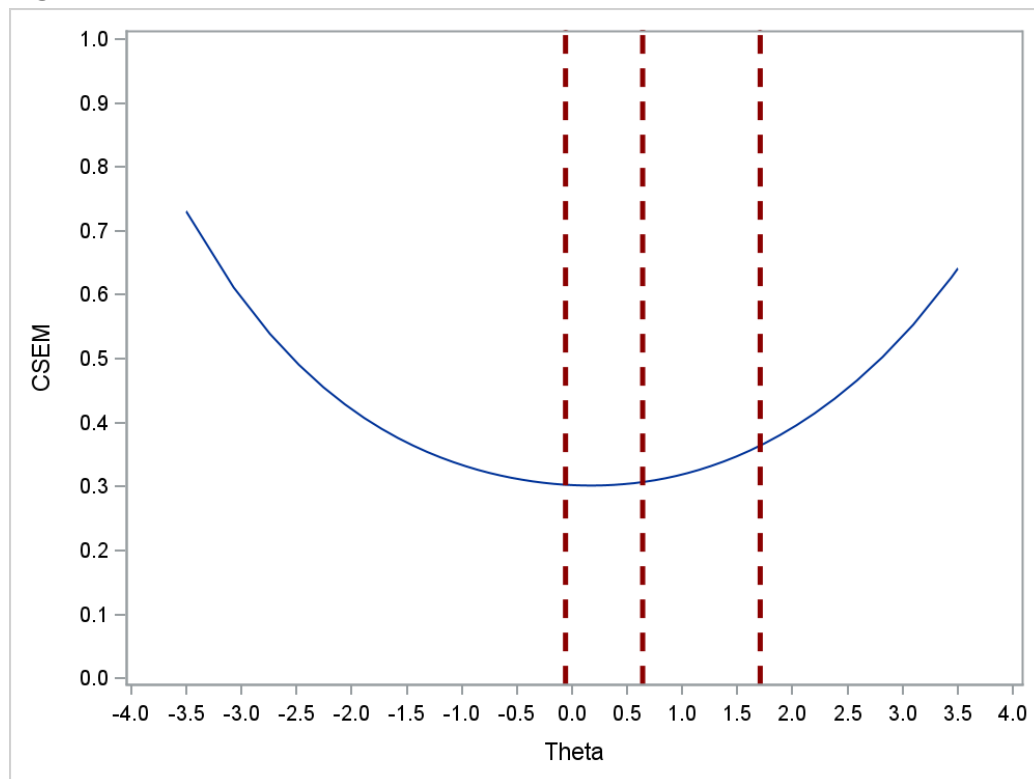
Figure B.41. TCC, ELA Grade 8, Form 2**Figure B.42. CSEM, ELA Grade 8, Form 2**

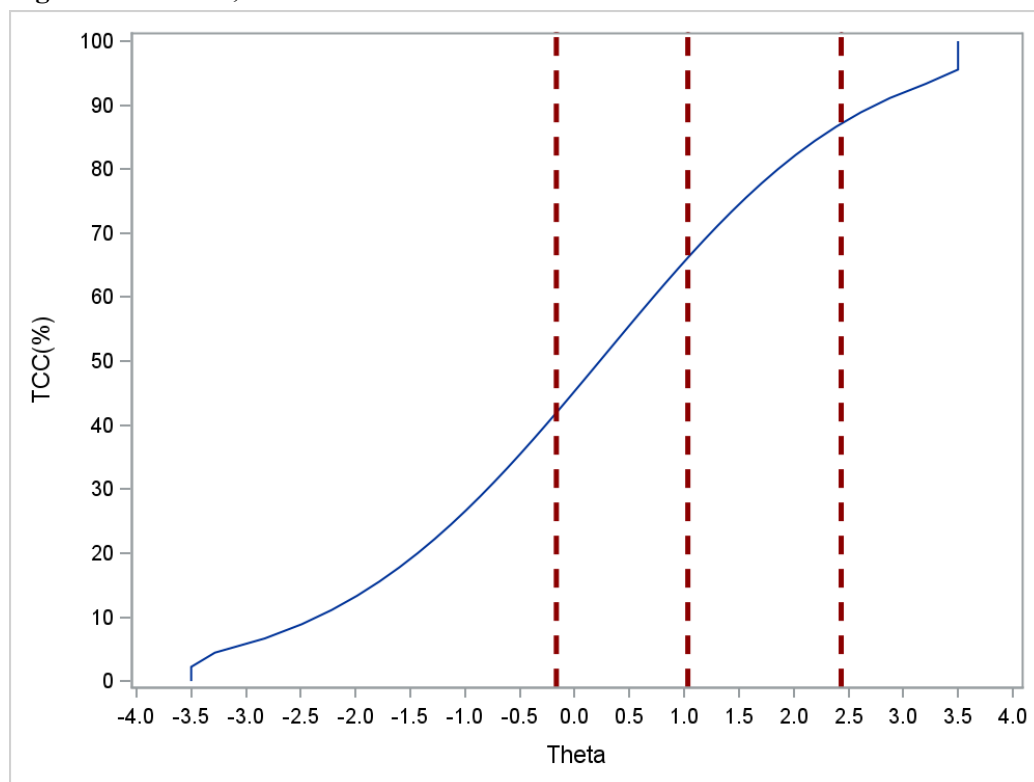
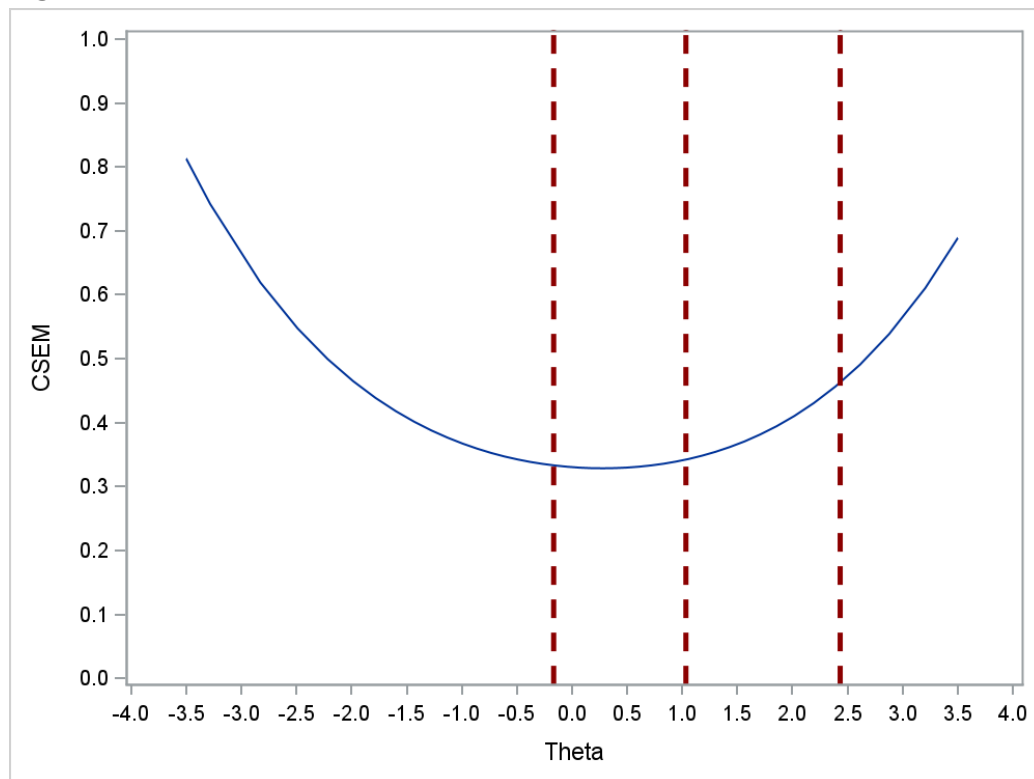
Figure B.43. TCC, Mathematics Grade 3**Figure B.44. CSEM, Mathematics Grade 3**

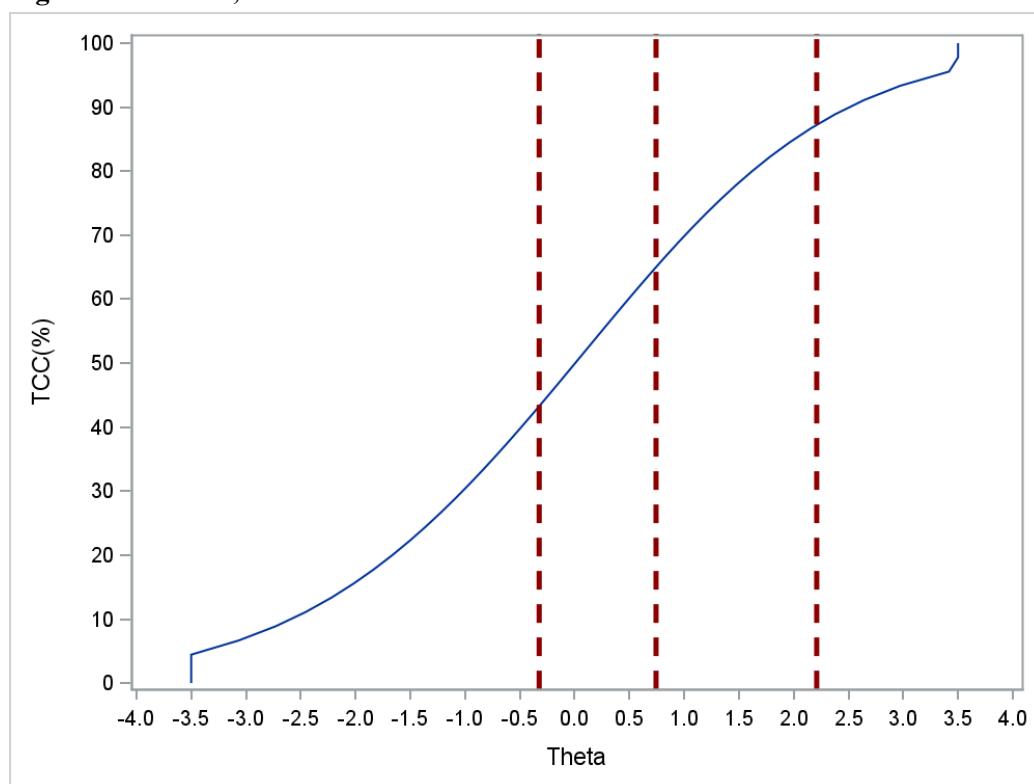
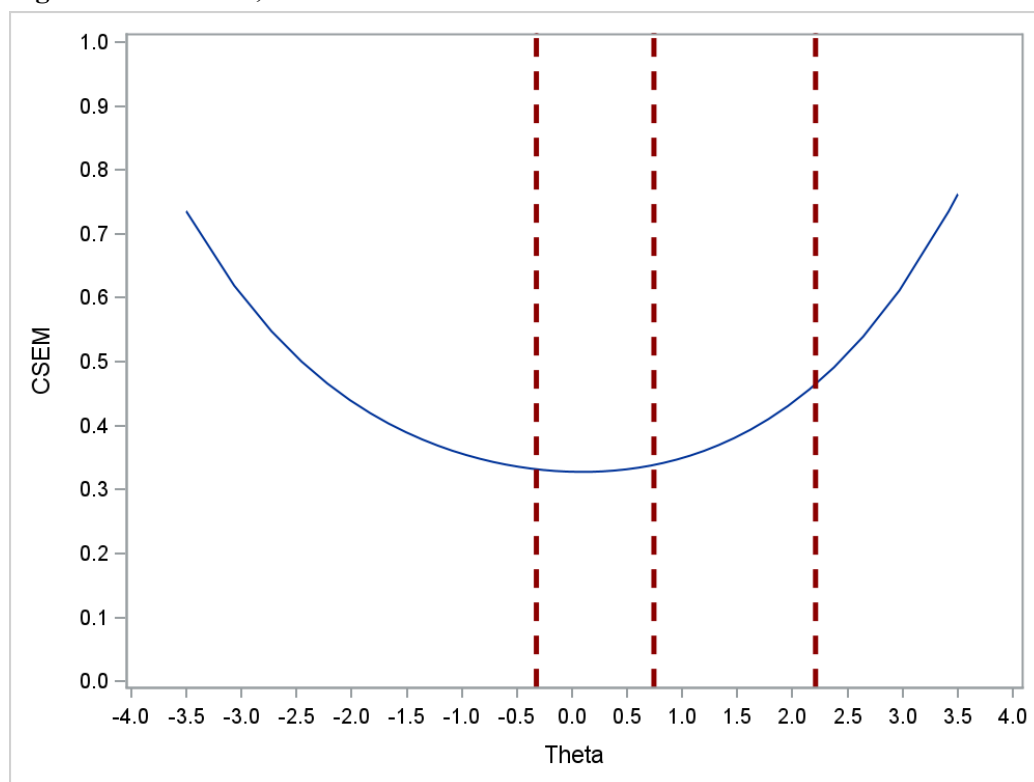
Figure B.45. TCC, Mathematics Grade 4**Figure B.46. CSEM, Mathematics Grade 4**

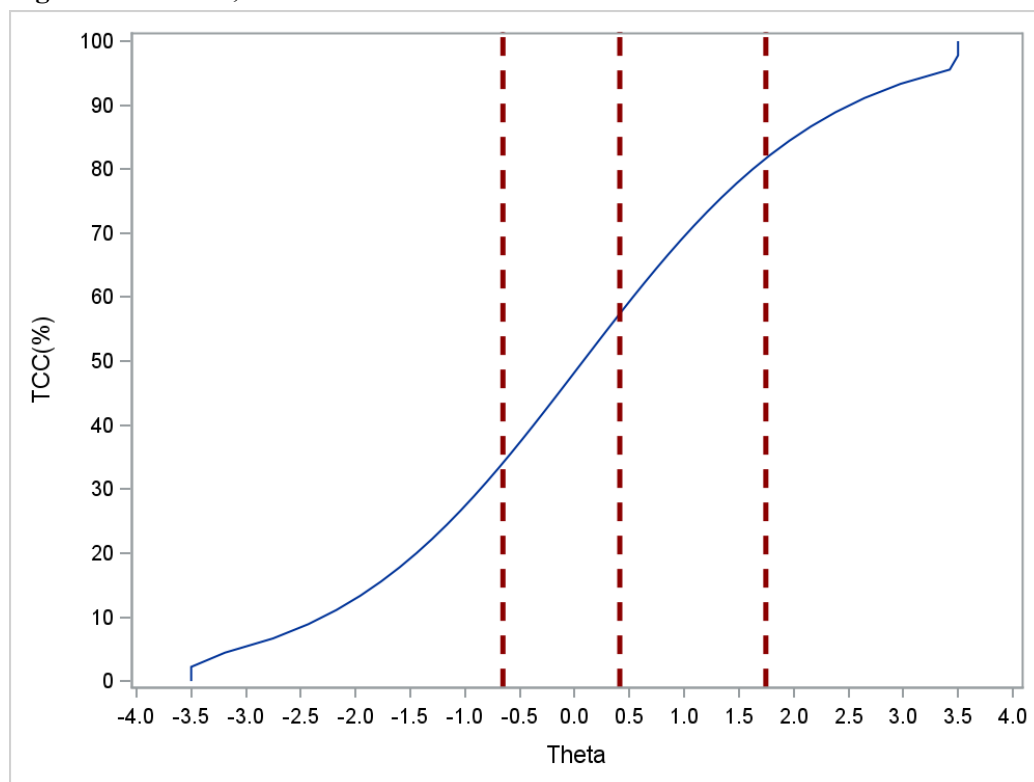
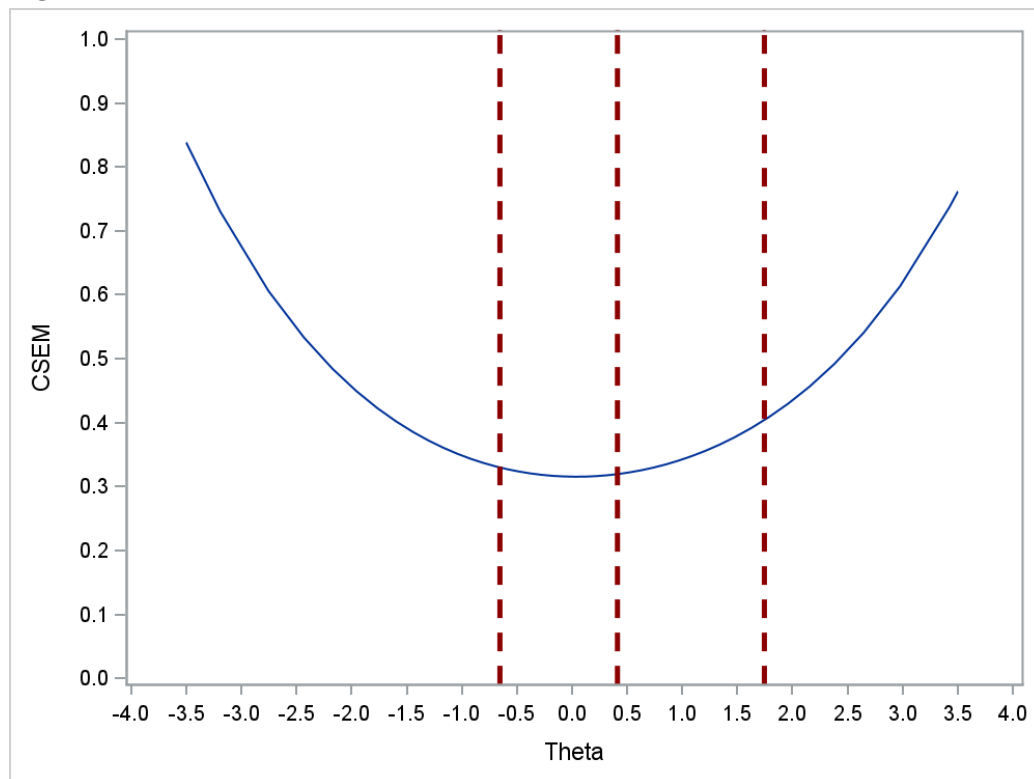
Figure B.47. TCC, Mathematics Grade 5**Figure B.48. CSEM, Mathematics Grade 5**

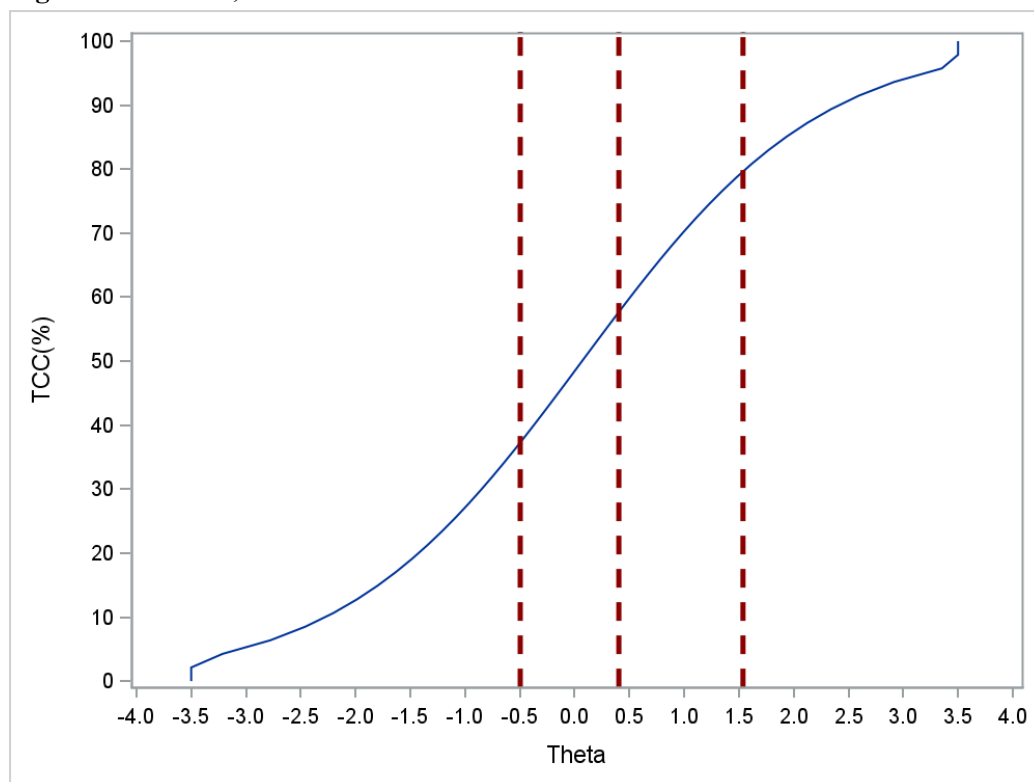
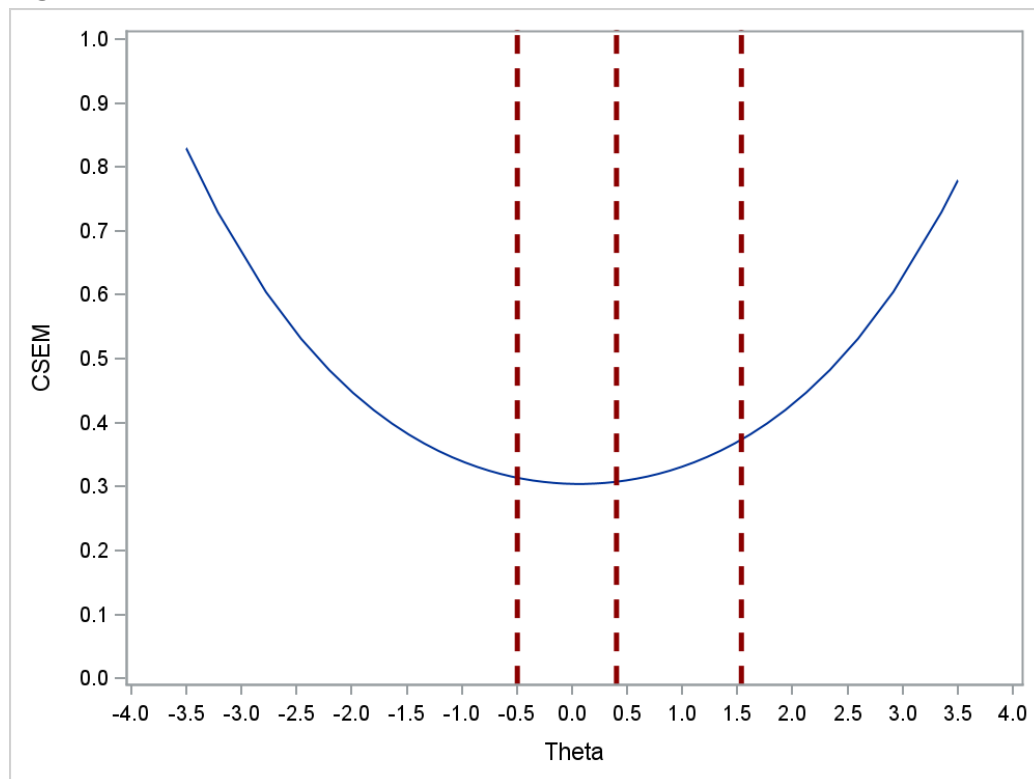
Figure B.49. TCC, Mathematics Grade 6**Figure B.50. CSEM, Mathematics Grade 6**

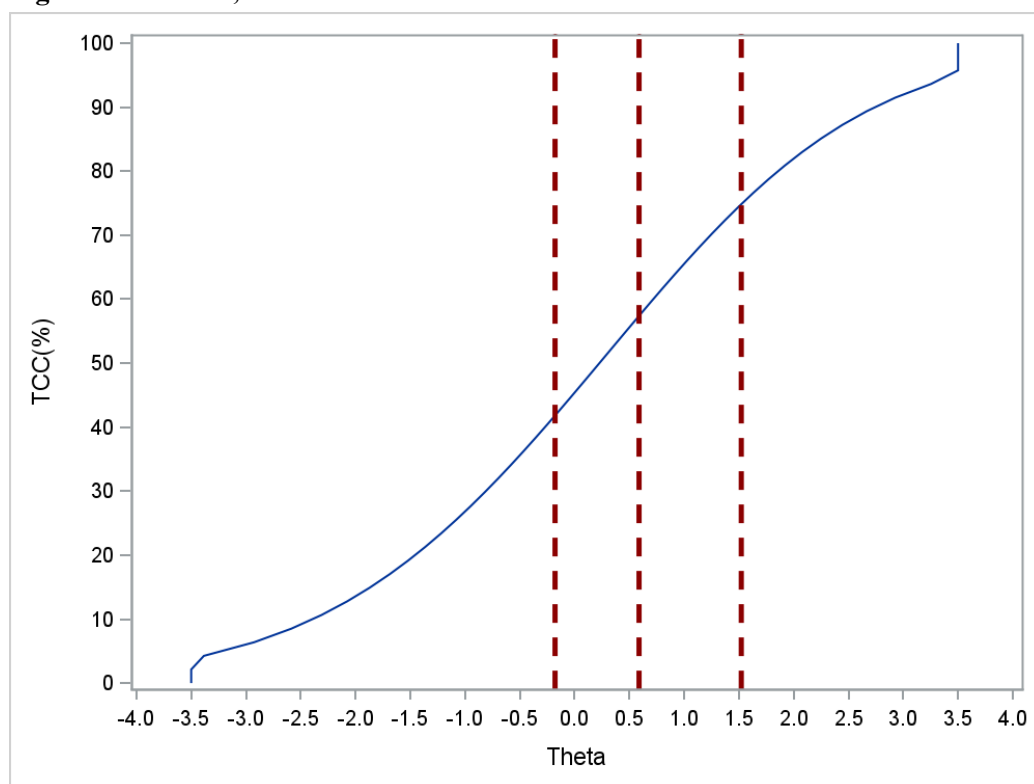
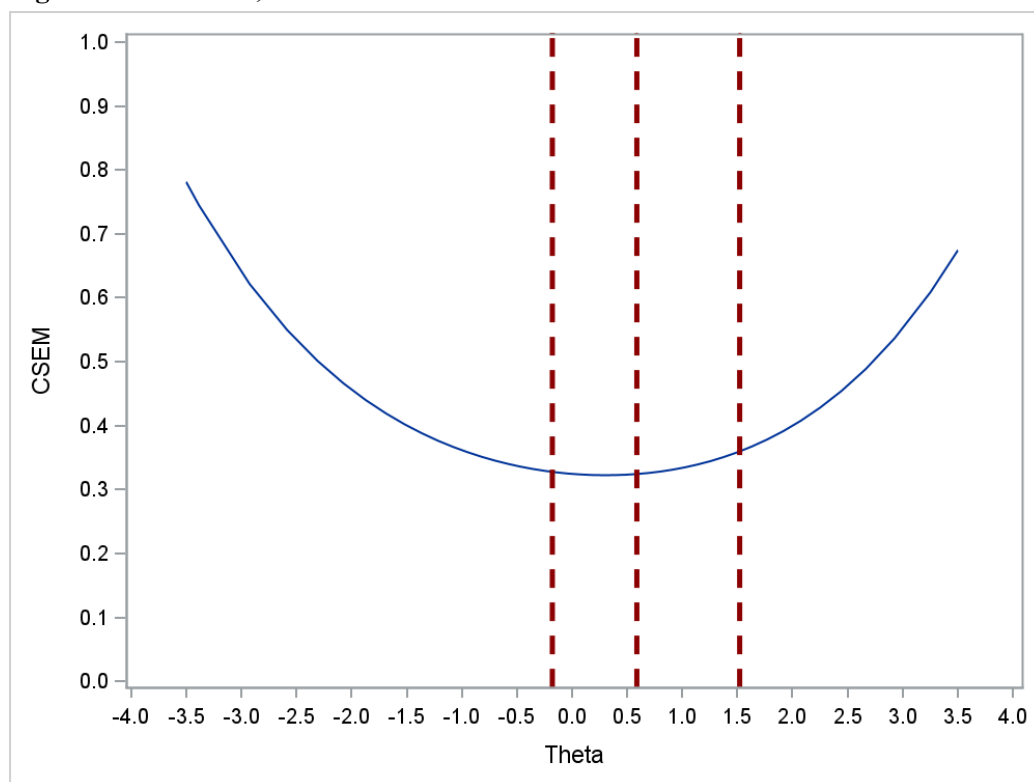
Figure B.51. TCC, Mathematics Grade 7**Figure B.52. CSEM, Mathematics Grade 7**

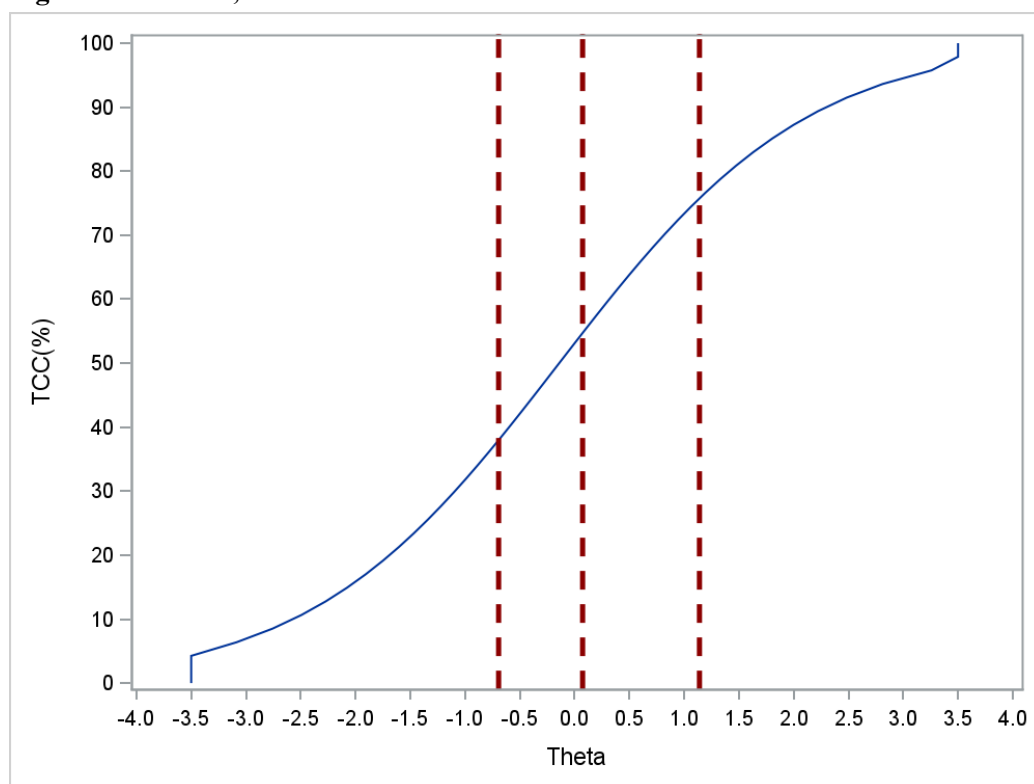
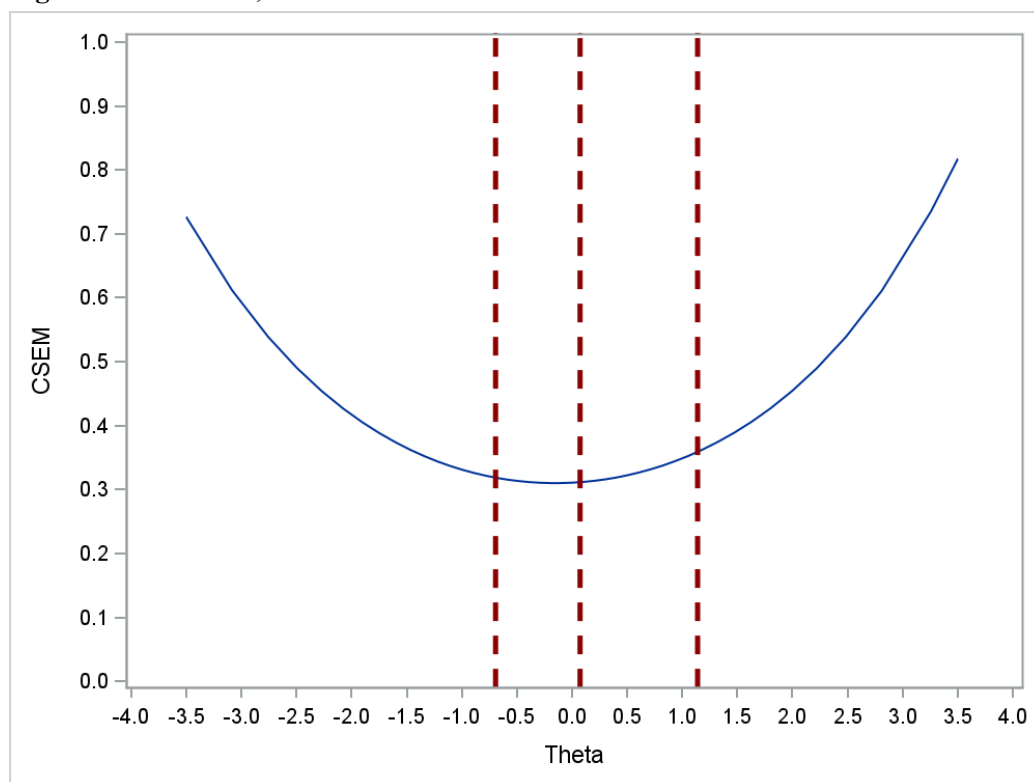
Figure B.53. TCC, Mathematics Grade 8**Figure B.54. CSEM, Mathematics Grade 8**

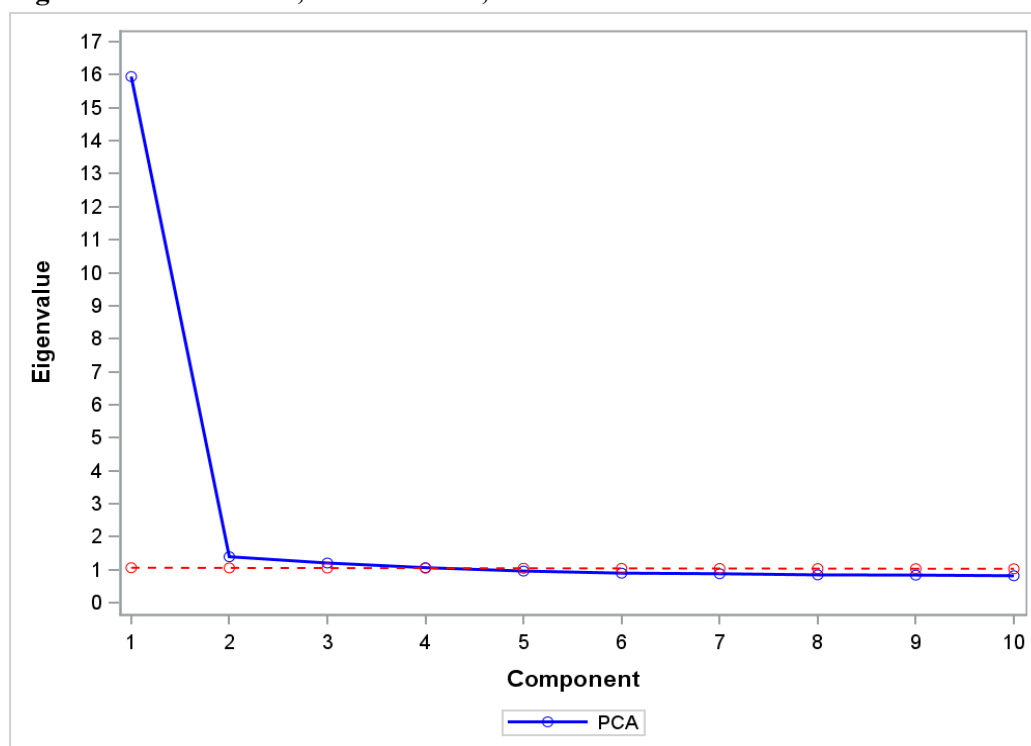
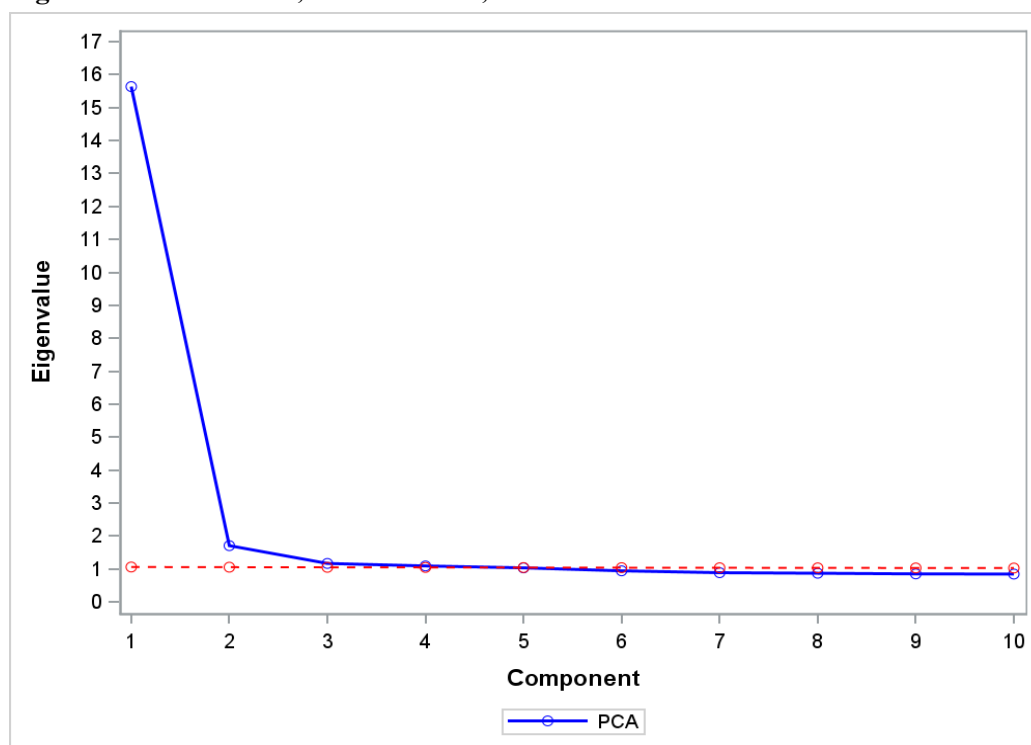
Figure B.55. Scree Plot, ELA Grade 3, Form 1**Figure B.56. Scree Plot, ELA Grade 4, Form 1**

Figure B.57. Scree Plot, ELA Grade 5, Form 1

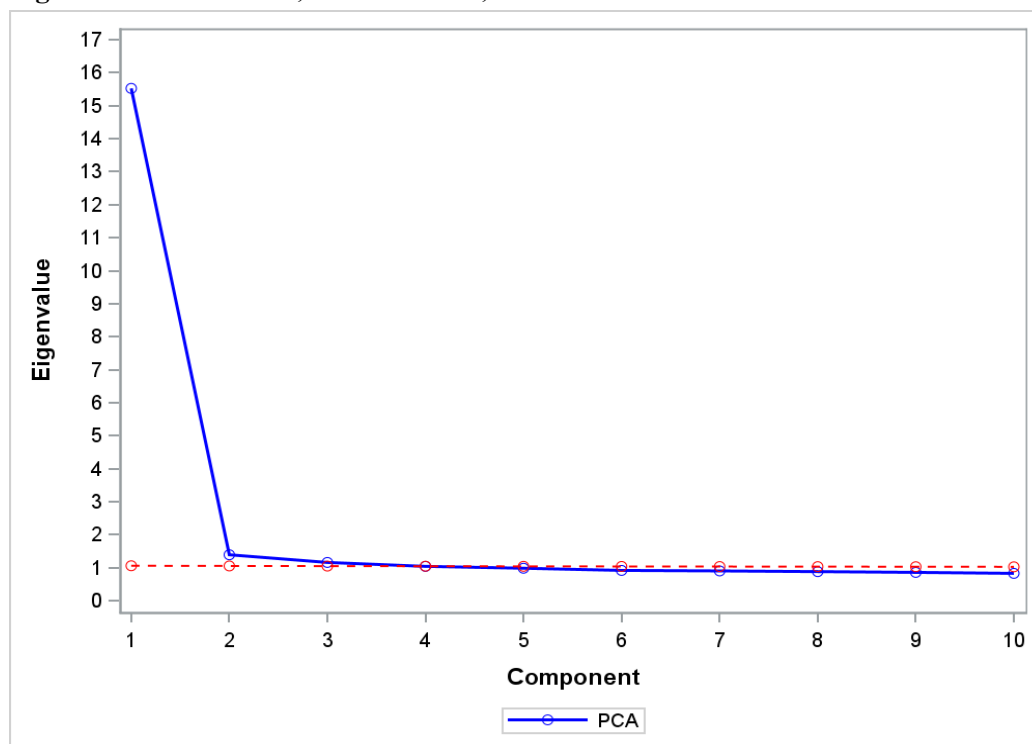


Figure B.58. Scree Plot, ELA Grade 6, Form 1

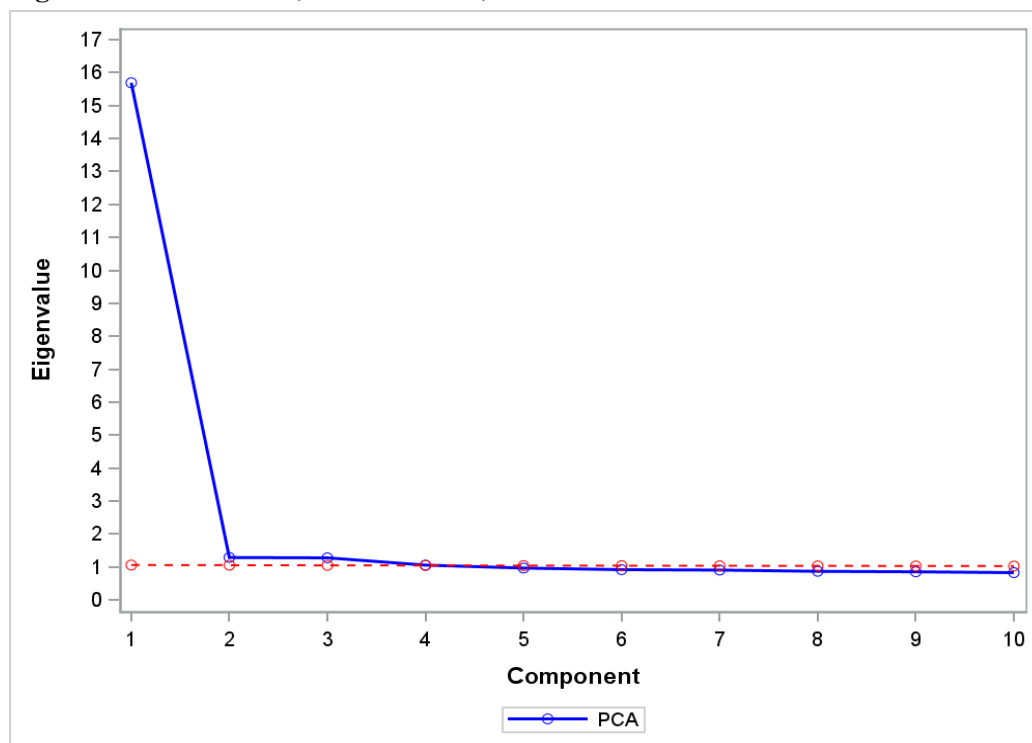


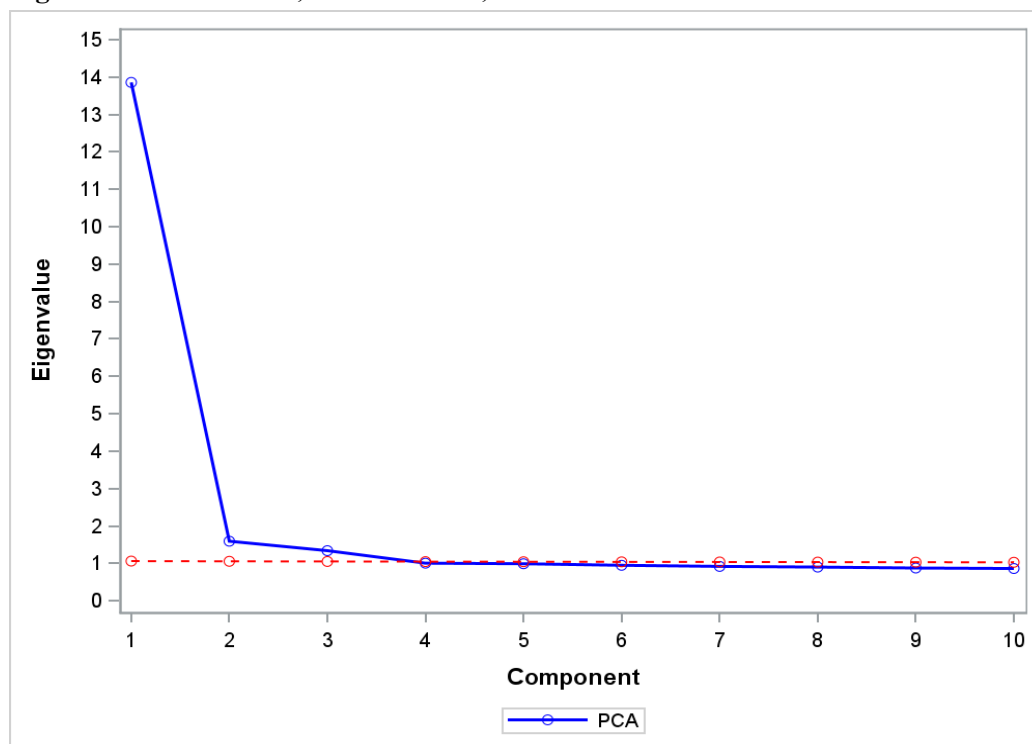
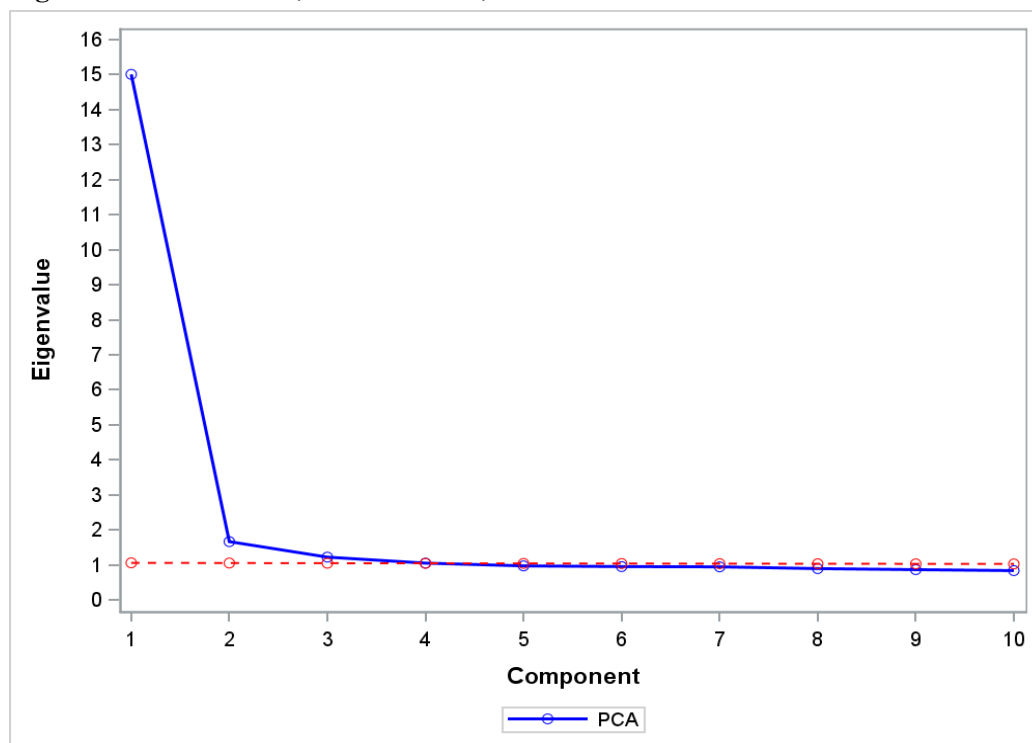
Figure B.59. Scree Plot, ELA Grade 7, Form 1**Figure B.60. Scree Plot, ELA Grade 8, Form 1**

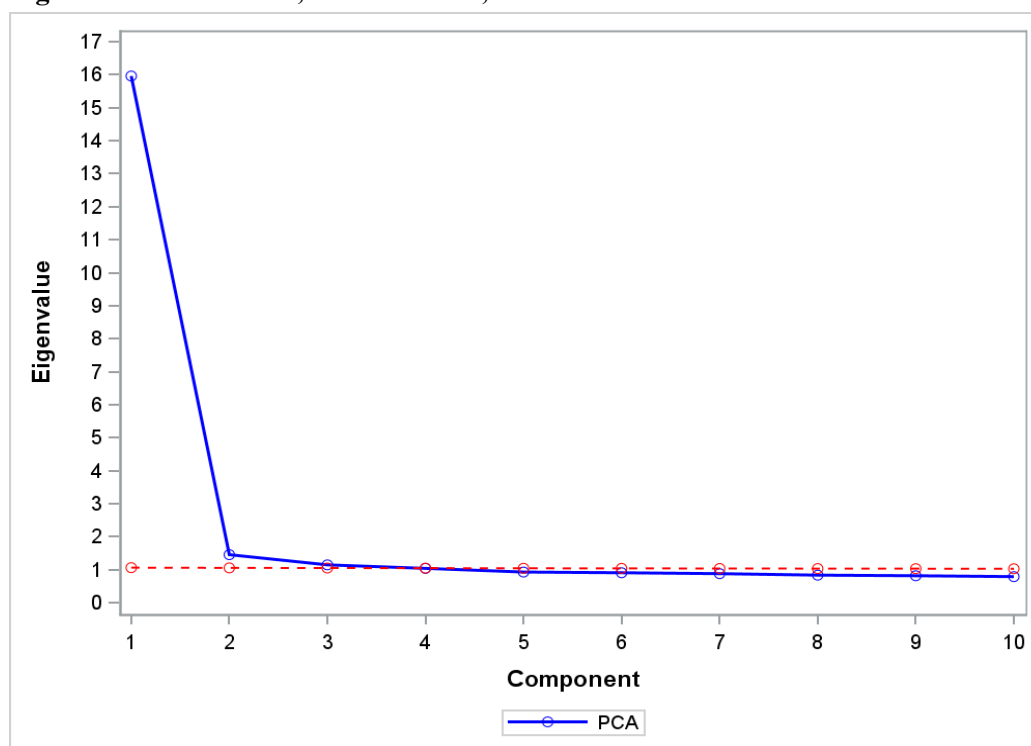
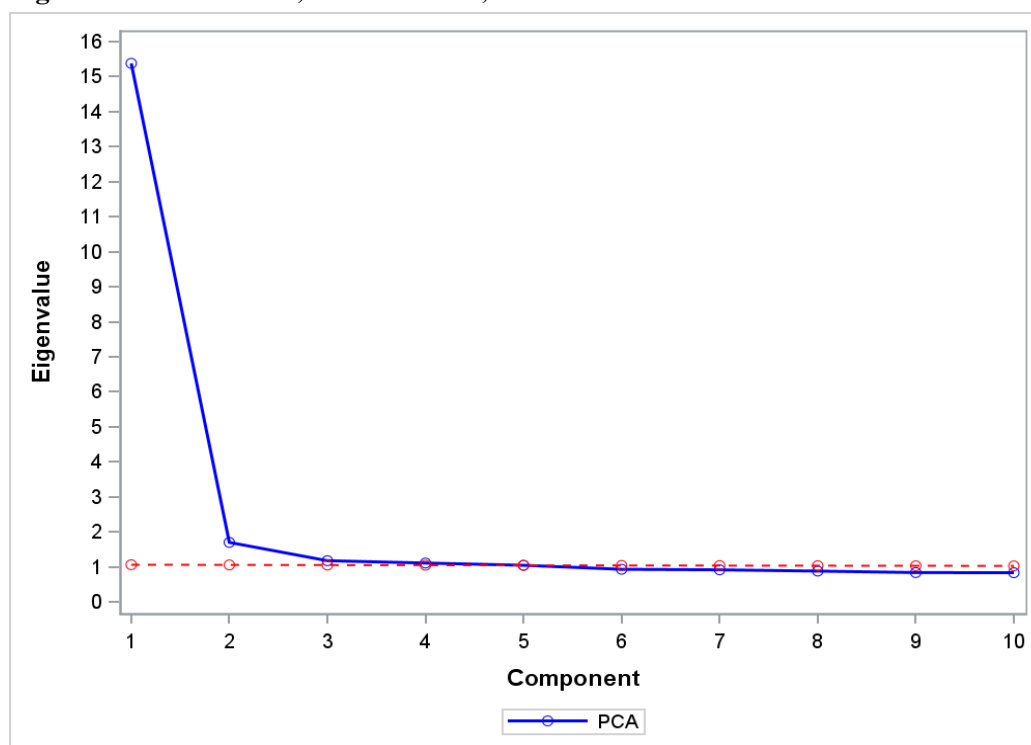
Figure B.61. Scree Plot, ELA Grade 3, Form 2**Figure B.62. Scree Plot, ELA Grade 4, Form 2**

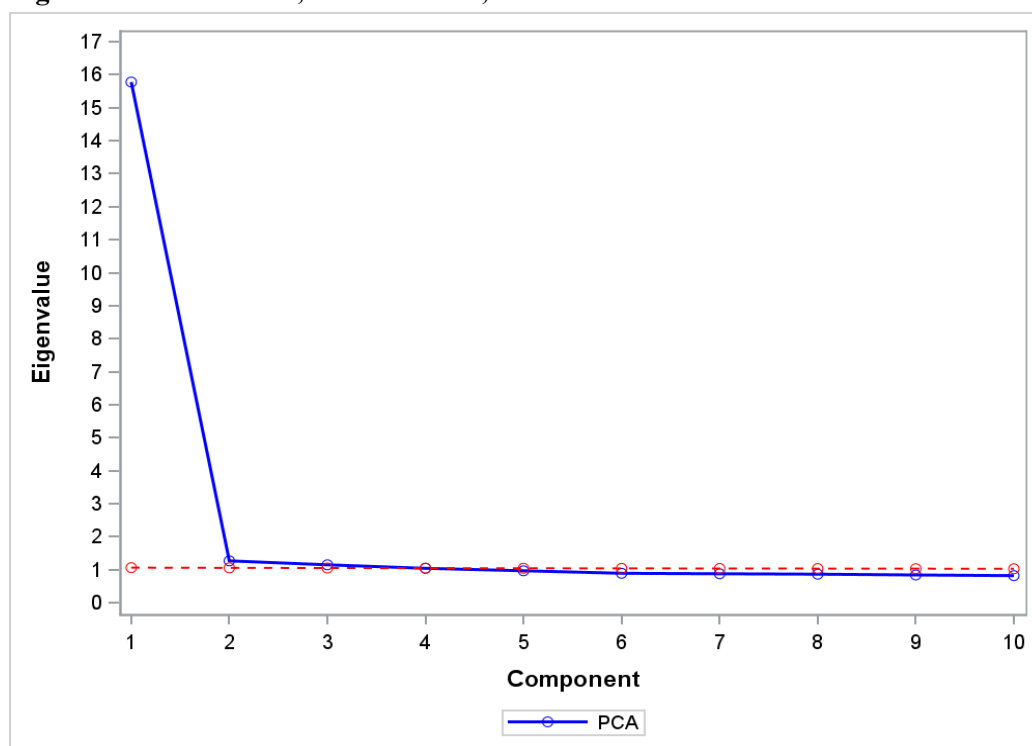
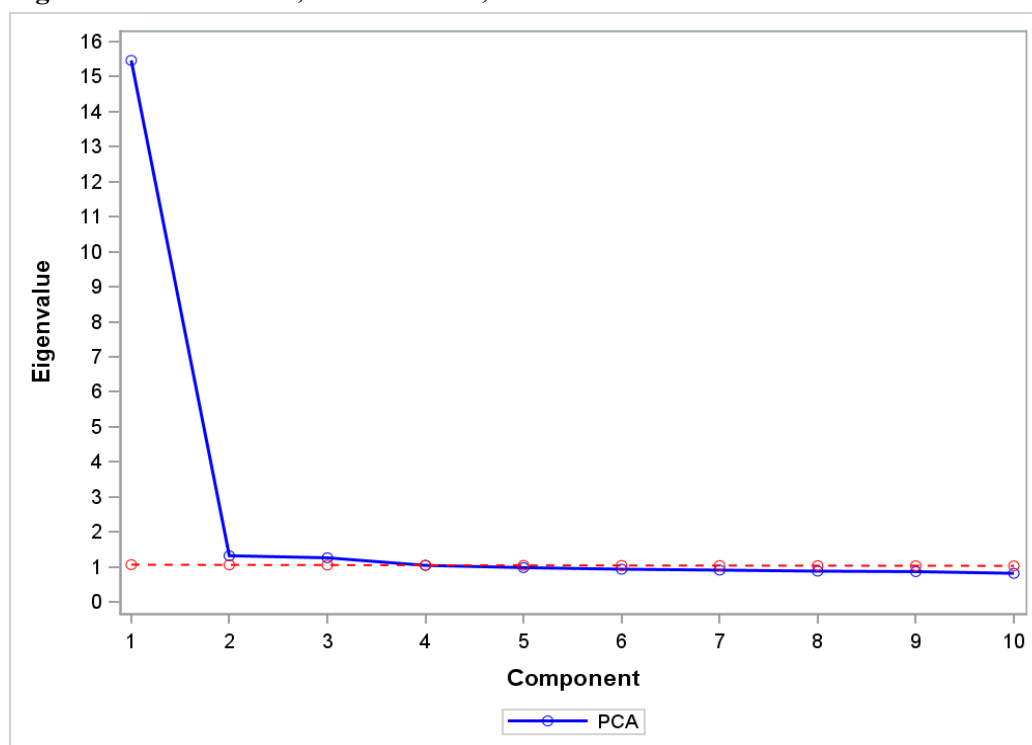
Figure B.63. Scree Plot, ELA Grade 5, Form 2**Figure B.64. Scree Plot, ELA Grade 6, Form 2**

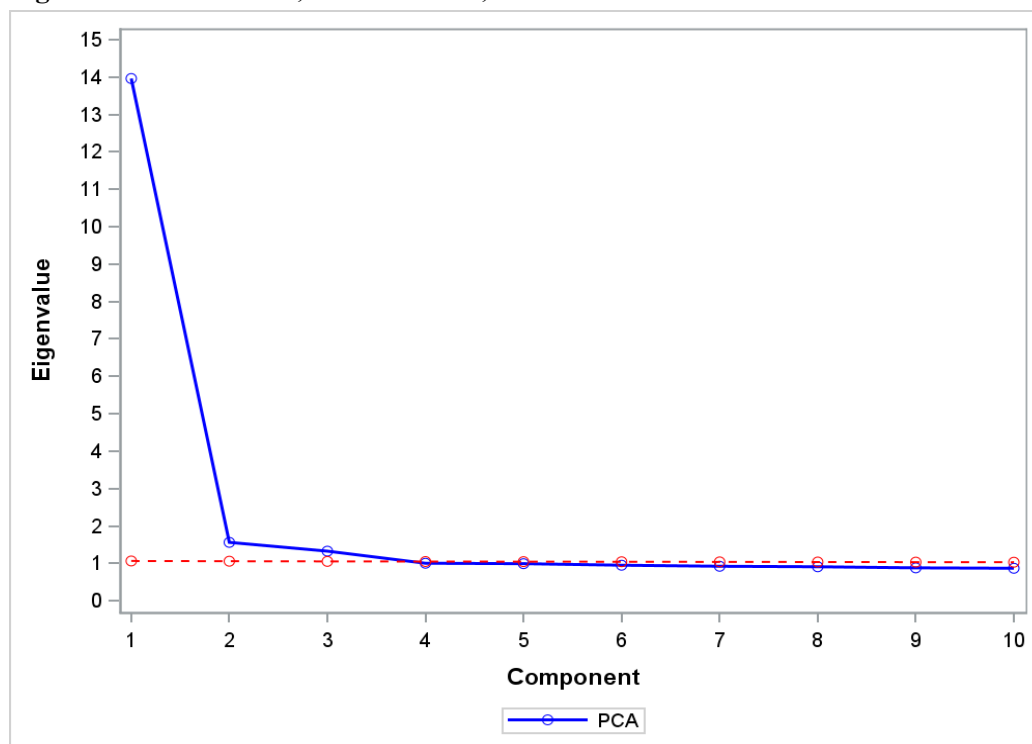
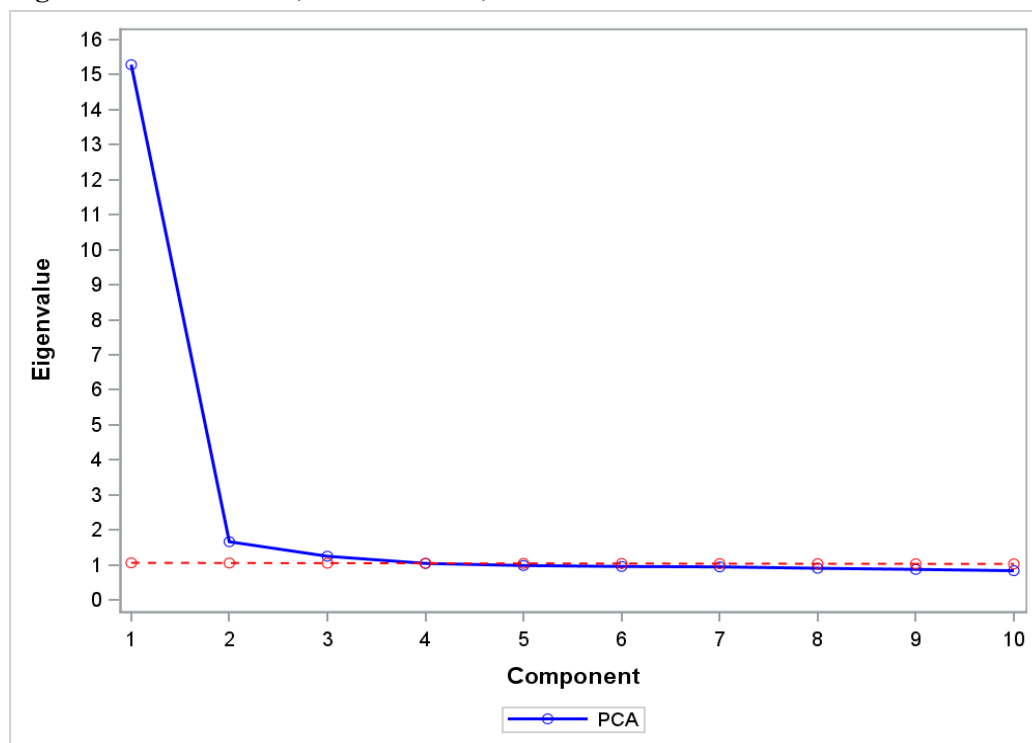
Figure B.65. Scree Plot, ELA Grade 7, Form 2**Figure B.66. Scree Plot, ELA Grade 8, Form 2**

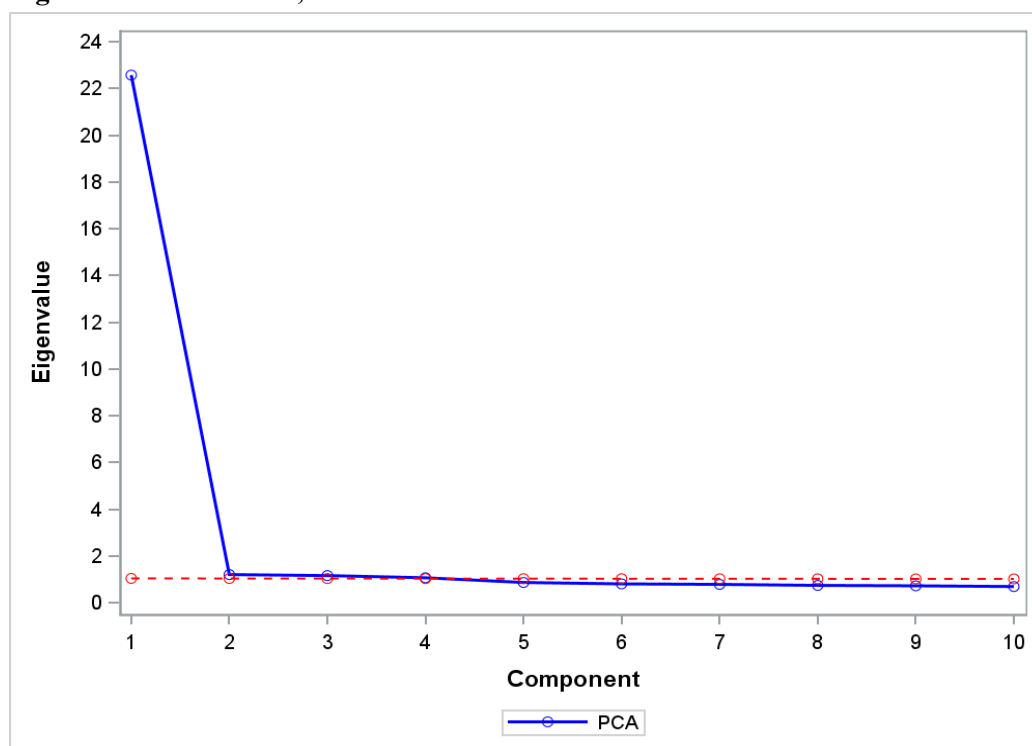
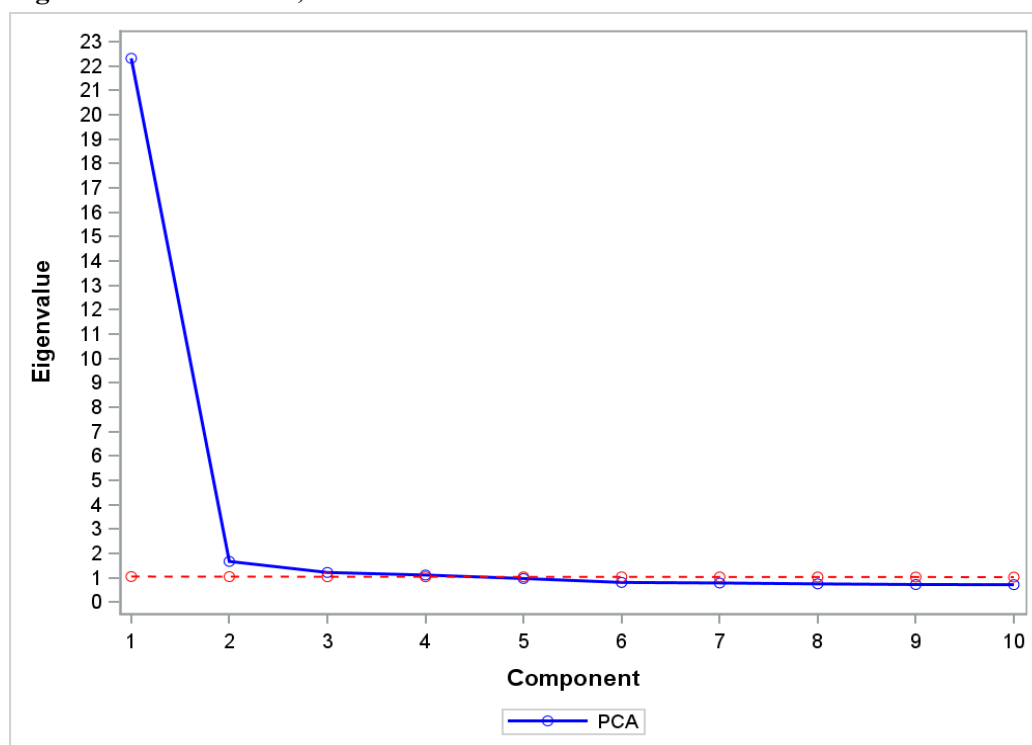
Figure B.67. Scree Plot, Mathematics Grade 3**Figure B.68. Scree Plot, Mathematics Grade 4**

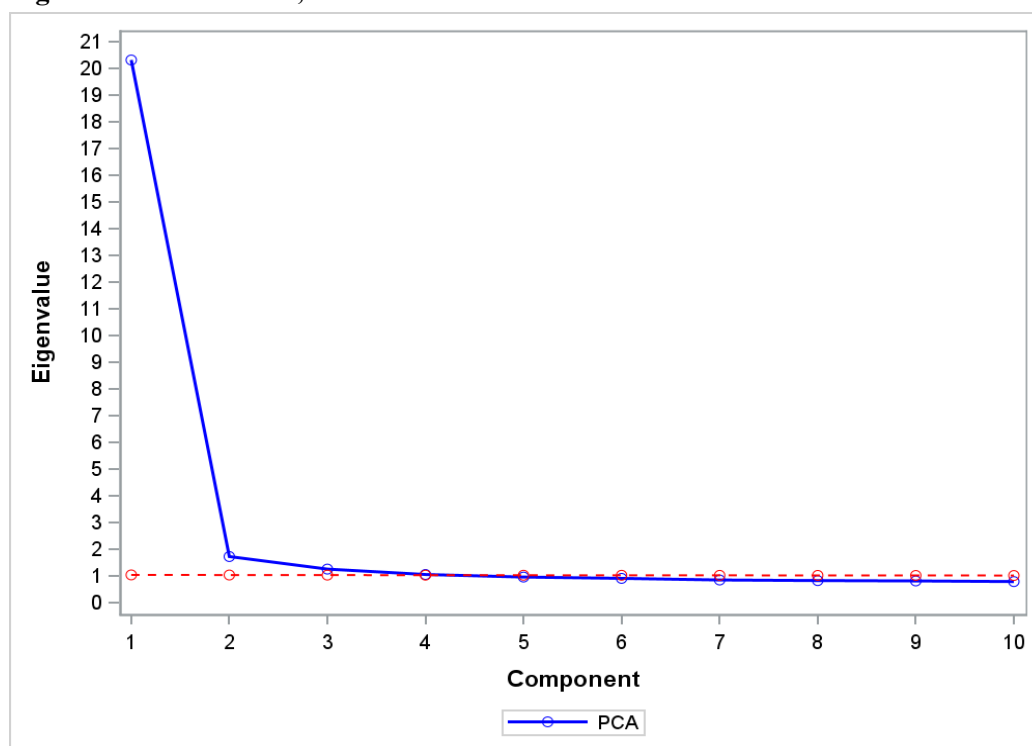
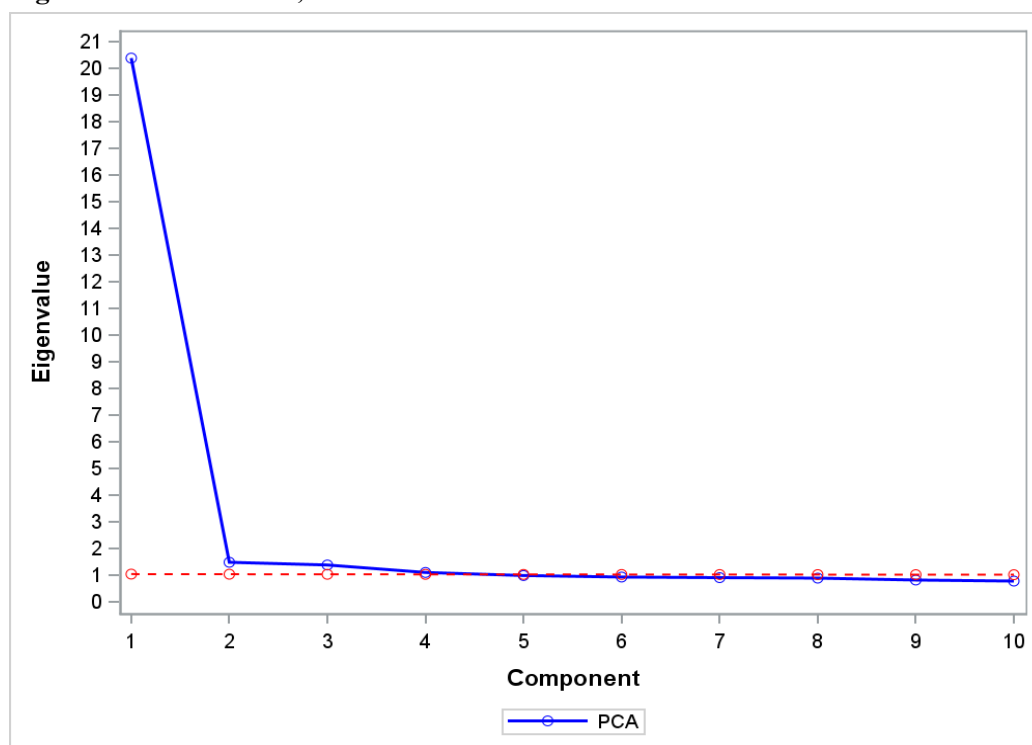
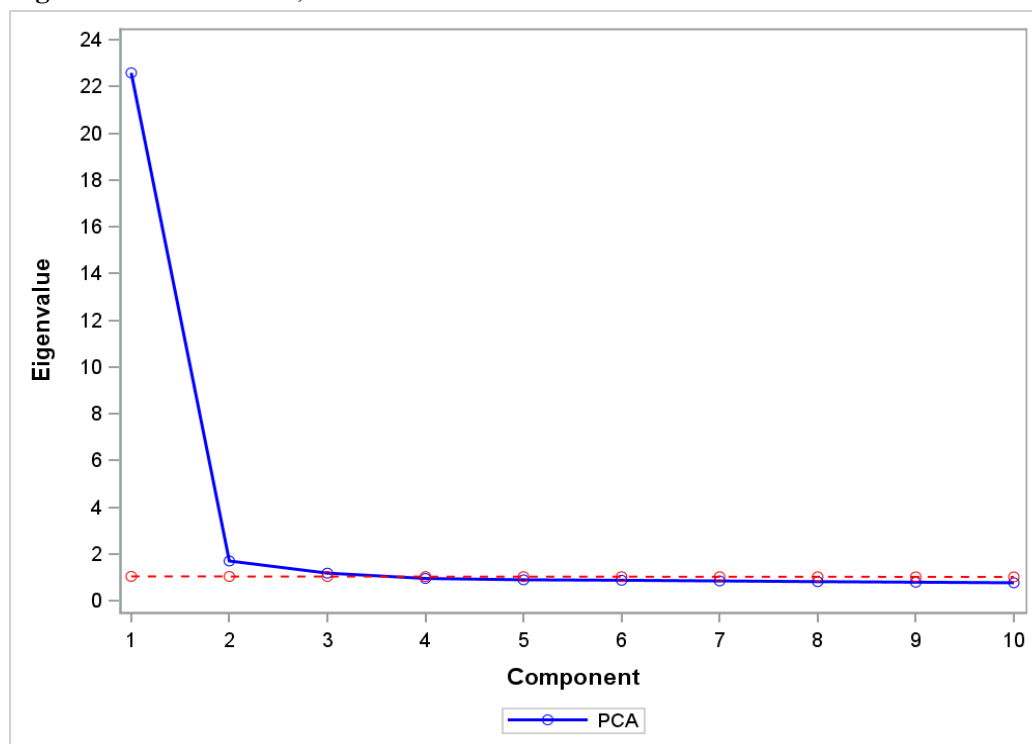
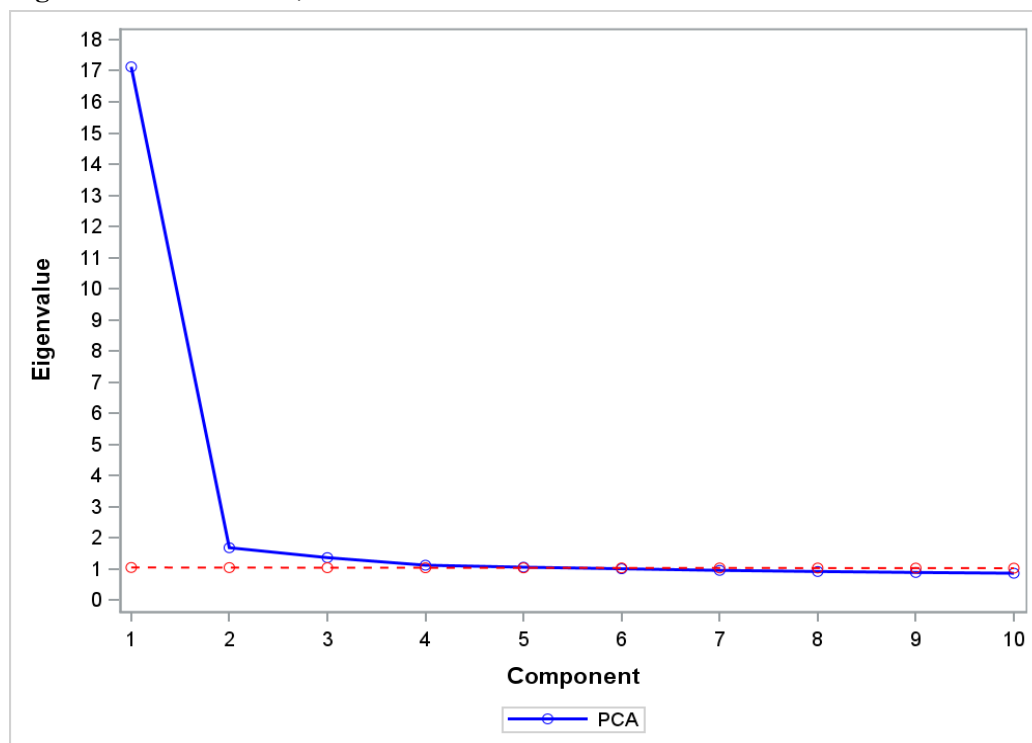
Figure B.69. Scree Plot, Mathematics Grade 5**Figure B.70. Scree Plot, Mathematics Grade 6**

Figure B.71. Scree Plot, Mathematics Grade 7**Figure B.72. Scree Plot, Mathematics Grade 8**

Appendix C: ADMINISTRATION RESULTS

This appendix presents the Spring 2023 AASA results for all students and subgroups by gender, ethnicity (Hispanic or Not-Hispanic), race, and special education, English learner (EL), and low socioeconomic status. Specifically:

- Table C.1 – Table C.12 present the overall results by subgroup, including the sample size, mean and standard deviation (SD) of the total scale score (SS), and percentage of students at each performance level overall.
- Figure C.1 – Figure C.12 present histograms of the total scale score distribution.

Table C.1. Test Results by Subgroup, ELA Grade 3

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	80,814	2498.68	35.51	47.4	11.7	27.1	13.8
Male	40,905	2496.66	35.66	49.6	11.5	26.1	12.8
Female	39,909	2500.75	35.24	45.1	11.9	28.2	14.8
Hispanic	38,838	2489.92	32.99	57.9	11.8	22.6	7.8
Non-Hispanic	41,976	2506.79	35.83	37.6	11.7	31.3	19.4
American Indian	4,346	2479.61	29.28	71.9	10.3	14.4	3.4
Asian	2,885	2521.43	34.47	21.9	9.4	37.3	31.4
Black or African American	5,814	2489.05	32.39	59.2	12.0	21.7	7.1
Multi-racial	5,226	2504.35	35.17	39.8	12.4	31.1	16.7
Native Hawaiian or Other Pacific Islander	467	2497.90	32.71	47.3	13.9	29.1	9.6
White	62,002	2499.38	35.41	46.4	11.8	27.7	14.1
Missing	74	2509.62	39.46	29.7	18.9	23.0	28.4
Special Education	12,514	2473.55	31.52	77.7	6.8	11.1	4.4
English Learner (EL)	9,055	2466.50	22.70	88.6	6.0	5.0	0.4
Low Socioeconomic Status (SES)	38,364	2487.19	32.02	61.2	11.7	20.8	6.3
Migrant	424	2476.79	29.54	75.2	9.4	12.3	3.1

Note. SS = scale score, 1 = *Minimally Proficient*, 2 = *Partially Proficient*, 3 = *Proficient*, 4 = *Highly Proficient*

Table C.2. Test Results by Subgroup, ELA Grade 4

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	80,659	2519.03	33.92	41.1	14.2	30.5	14.3
Male	40,783	2516.73	34.16	44.2	13.9	28.6	13.3
Female	39,876	2521.38	33.51	37.9	14.4	32.4	15.3
Hispanic	38,320	2510.74	31.05	50.8	14.9	26.2	8.0
Non-Hispanic	42,339	2526.53	34.66	32.3	13.5	34.4	19.9
American Indian	4,398	2502.29	28.18	62.6	14.4	18.9	4.1
Asian	2,945	2539.79	34.65	19.3	10.4	37.1	33.1
Black or African American	5,915	2509.15	30.81	53.1	14.6	25.1	7.3
Multi-racial	5,059	2524.01	33.81	34.5	14.1	34.3	17.1
Native Hawaiian or Other Pacific Islander	452	2515.71	32.30	44.2	17.7	27.0	11.1
White	61,821	2519.80	33.78	40.0	14.3	31.2	14.5
Missing	69	2513.62	39.23	52.2	5.8	23.2	18.8
Special Education	12,218	2492.81	28.71	76.2	8.7	11.4	3.7
English Learner (EL)	8,727	2489.03	21.25	83.0	9.5	7.0	0.4
Low Socioeconomic Status (SES)	38,418	2508.37	30.18	53.8	15.1	24.6	6.5
Migrant	447	2500.81	30.49	66.0	10.3	18.1	5.6

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.3. Test Results by Subgroup, ELA Grade 5

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	80,917	2528.74	35.21	40.2	22.6	28.3	8.9
Male	41,123	2525.53	35.44	43.8	21.9	26.5	7.7
Female	39,794	2532.05	34.65	36.4	23.4	30.1	10.1
Hispanic	38,120	2520.21	32.58	49.5	23.9	22.1	4.5
Non-Hispanic	42,797	2536.33	35.71	31.9	21.5	33.8	12.8
American Indian	4,423	2510.25	30.42	63.1	20.8	13.6	2.5
Asian	2,972	2550.71	33.49	17.3	17.7	44.1	20.9
Black or African American	5,797	2517.23	32.04	53.8	22.2	20.6	3.4
Multi-racial	5,034	2533.72	34.84	34.4	23.1	31.7	10.8
Native Hawaiian or Other Pacific Islander	432	2527.11	33.10	43.5	23.1	25.5	7.9
White	62,195	2529.67	35.03	38.8	23.0	29.0	9.1
Missing	64	2539.78	45.33	34.4	10.9	31.3	23.4
Special Education	11,699	2499.99	30.30	76.2	13.1	8.6	2.1
English Learner (EL)	8,041	2495.33	23.47	83.9	12.3	3.6	0.2
Low Socioeconomic Status (SES)	37,763	2517.52	31.97	53.0	23.3	20.1	3.6
Migrant	462	2511.17	33.41	60.6	18.0	19.3	2.2

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.4. Test Results by Subgroup, ELA Grade 6

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	81,369	2543.19	34.76	35.7	22.4	35.5	6.5
Male	41,180	2539.72	34.97	39.7	21.9	33.0	5.4
Female	40,189	2546.75	34.18	31.6	22.8	38.0	7.6
Hispanic	38,415	2534.71	32.62	44.9	23.7	28.3	3.1
Non-Hispanic	42,954	2550.78	34.86	27.5	21.1	41.8	9.5
American Indian	4,595	2524.26	30.68	58.3	22.7	17.6	1.3
Asian	2,941	2565.90	32.84	13.6	16.7	52.1	17.6
Black or African American	5,783	2532.81	32.58	47.0	23.4	26.8	2.7
Multi-racial	4,903	2548.54	33.55	29.4	21.8	41.1	7.8
Native Hawaiian or Other Pacific Islander	505	2542.00	33.18	37.4	22.2	35.6	4.8
White	62,575	2544.06	34.51	34.5	22.6	36.3	6.6
Missing	67	2544.85	44.87	35.8	13.4	37.3	13.4
Special Education	11,034	2512.68	30.14	74.3	14.3	10.5	1.0
English Learner (EL)	6,933	2506.96	23.99	82.7	13.0	4.3	0
Low Socioeconomic Status (SES)	37,775	2532.15	32.20	48.0	23.4	26.0	2.5
Migrant	455	2521.99	31.02	62.0	18.5	18.2	1.3

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.5. Test Results by Subgroup, ELA Grade 7

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	82,061	2552.90	31.66	38.9	19.6	32.8	8.7
Male	41,753	2549.83	31.56	42.6	19.1	31.1	7.2
Female	40,308	2556.07	31.46	35.1	20.1	34.6	10.2
Hispanic	39,074	2545.35	29.28	47.9	20.7	27.1	4.3
Non-Hispanic	42,987	2559.75	32.19	30.8	18.6	38.0	12.6
American Indian	4,591	2537.18	26.39	60.4	19.8	17.8	2.0
Asian	2,843	2576.01	32.91	15.4	13.7	44.5	26.3
Black or African American	5,881	2544.67	28.68	49.1	20.7	26.2	4.0
Multi-racial	4,741	2558.10	31.41	31.5	20.6	36.8	11.1
Native Hawaiian or Other Pacific Islander	456	2548.99	29.09	40.8	24.6	29.8	4.8
White	63,469	2553.40	31.44	38.0	19.6	33.7	8.6
Missing	80	2553.00	41.64	41.3	16.3	25.0	17.5
Special Education	10,397	2524.81	25.06	79.2	11.0	8.5	1.3
English Learner (EL)	7,013	2520.74	19.74	85.8	10.5	3.7	0
Low Socioeconomic Status (SES)	37,290	2543.06	28.53	51.2	20.7	24.6	3.5
Migrant	455	2532.55	25.90	64.6	20.0	14.3	1.1

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.6. Test Results by Subgroup, ELA Grade 8

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	85,232	2557.60	34.17	41.4	22.2	26.5	9.8
Male	43,762	2552.69	34.16	47.1	21.6	23.5	7.8
Female	41,470	2562.79	33.40	35.3	22.9	29.8	12.0
Hispanic	40,734	2549.89	32.31	50.7	22.5	21.2	5.7
Non-Hispanic	44,498	2564.66	34.30	32.9	22.0	31.4	13.6
American Indian	4,854	2542.23	30.04	61.5	19.9	16.0	2.6
Asian	2,863	2581.75	32.46	15.6	17.9	39.7	26.8
Black or African American	5,860	2549.49	32.28	51.1	22.3	21.1	5.5
Multi-racial	4,889	2562.67	33.34	34.6	23.6	30.3	11.5
Native Hawaiian or Other Pacific Islander	513	2552.51	32.83	48.0	21.6	23.0	7.4
White	66,164	2558.05	34.00	40.6	22.5	27.0	9.9
Missing	89	2568.65	42.61	33.7	13.5	25.8	27.0
Special Education	10,151	2526.89	27.15	80.7	11.8	6.3	1.2
English Learner (EL)	6,884	2522.72	21.37	88.1	9.5	2.4	0
Low Socioeconomic Status (SES)	37,802	2547.95	31.85	53.0	22.2	19.9	4.8
Migrant	501	2536.43	32.54	67.3	17.0	11.0	4.8

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.7. Test Results by Subgroup, Mathematics Grade 3

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	81,986	3517.02	48.72	31.2	25.5	28.3	14.9
Male	41,669	3519.20	49.95	29.8	24.4	28.9	16.9
Female	40,317	3514.77	47.31	32.7	26.7	27.7	13.0
Hispanic	39,416	3505.74	46.68	39.1	28.3	24.1	8.5
Non-Hispanic	42,570	3527.47	48.24	23.9	23.0	32.2	20.9
American Indian	4,470	3489.32	44.26	53.9	26.5	16.1	3.5
Asian	2,915	3551.39	43.02	10.8	15.1	33.7	40.4
Black or African American	5,984	3499.82	47.01	44.4	27.4	21.4	6.8
Multi-racial	5,309	3523.00	48.05	26.1	25.6	30.5	17.8
Native Hawaiian or Other Pacific Islander	474	3513.99	45.34	32.1	28.3	29.3	10.3
White	62,760	3518.54	48.03	29.7	25.7	29.4	15.1
Missing	74	3528.91	56.91	27.0	17.6	27.0	28.4
Special Education	12,926	3483.23	50.34	59.4	20.8	14.1	5.7
English Learner (EL)	9,305	3478.63	41.10	64.2	24.3	9.9	1.6
Low Socioeconomic Status (SES)	39,010	3502.29	46.47	42.1	28.0	22.7	7.2
Migrant	429	3496.89	46.22	45.0	29.4	20.3	5.4

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.8. Test Results by Subgroup, Mathematics Grade 4

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	81,480	3543.68	49.02	39.3	21.9	26.2	12.6
Male	41,293	3546.40	50.44	37.2	21.2	26.9	14.6
Female	40,187	3540.90	47.36	41.3	22.6	25.4	10.6
Hispanic	38,708	3531.28	45.94	49.2	23.0	21.3	6.6
Non-Hispanic	42,772	3554.90	49.02	30.3	20.9	30.6	18.1
American Indian	4,502	3517.99	43.75	62.0	20.4	14.0	3.6
Asian	2,968	3577.15	45.09	15.2	15.6	37.0	32.2
Black or African American	6,015	3524.43	45.08	55.5	22.2	17.2	5.1
Multi-racial	5,104	3549.48	48.15	33.3	23.5	28.6	14.6
Native Hawaiian or Other Pacific Islander	457	3536.80	46.47	44.6	22.1	24.5	8.8
White	62,366	3545.39	48.54	37.7	22.1	27.2	13.0
Missing	68	3532.76	60.53	47.1	17.6	20.6	14.7
Special Education	12,496	3508.55	46.76	69.6	15.2	11.2	3.9
English Learner (EL)	8,867	3504.80	39.08	74.4	16.2	8.3	1.1
Low Socioeconomic Status (SES)	38,859	3528.16	45.22	51.9	22.9	19.5	5.6
Migrant	447	3523.72	46.56	54.4	23.5	17.4	4.7

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.9. Test Results by Subgroup, Mathematics Grade 5

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	81,451	3578.26	44.04	36.5	27.3	25	11.2
Male	41,428	3579.61	45.75	36.3	25.6	25.2	12.8
Female	40,023	3576.85	42.16	36.7	28.9	24.8	9.5
Hispanic	38,364	3567.37	40.25	45.7	28.5	20.1	5.7
Non-Hispanic	43,087	3587.95	45.01	28.3	26.1	29.4	16.1
American Indian	4,501	3555.91	37.44	57.3	26.3	13.7	2.7
Asian	2,987	3615.07	42.72	10.7	18.9	35.8	34.5
Black or African American	5,860	3560.51	39.49	52.9	26.8	16.2	4.2
Multi-racial	5,067	3582.62	43.48	32.3	27.5	27.7	12.5
Native Hawaiian or Other Pacific Islander	434	3575.08	40.35	35.7	34.1	21.9	8.3
White	62,537	3579.43	43.49	35.0	27.7	26.0	11.3
Missing	65	3585.09	54.70	36.9	18.5	23.1	21.5
Special Education	11,825	3546.63	39.80	68.2	18.9	9.6	3.3
English Learner (EL)	8,130	3543.86	33.45	71.2	20.8	6.9	1.0
Low Socioeconomic Status (SES)	38,027	3564.38	39.70	48.8	28.1	18.1	4.9
Migrant	463	3565.78	41.53	47.1	26.6	21.2	5.2

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.10. Test Results by Subgroup, Mathematics Grade 6

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	82,066	3606.26	41.62	48.7	20.9	20.8	9.6
Male	41,566	3607.87	42.76	47.3	20.6	21.2	10.9
Female	40,500	3604.61	40.34	50.2	21.2	20.3	8.3
Hispanic	38,776	3595.50	37.30	59.9	20.2	15.2	4.7
Non-Hispanic	43,290	3615.90	42.90	38.7	21.5	25.7	14.0
American Indian	4,674	3583.43	33.78	73.0	16.1	8.7	2.1
Asian	2,956	3640.12	40.48	16.6	19.2	36.5	27.7
Black or African American	5,841	3589.38	36.33	66.5	17.9	12.2	3.5
Multi-racial	4,948	3611.34	41.27	42.9	22.7	23.4	11.0
Native Hawaiian or Other Pacific Islander	510	3603.13	38.40	49.2	25.1	18.2	7.5
White	63,062	3607.56	41.23	47.2	21.4	21.5	9.8
Missing	75	3599.59	43.29	52.0	13.3	30.7	4.0
Special Education	11,210	3575.42	34.60	80.6	10.7	6.5	2.1
English Learner (EL)	7,045	3573.21	29.47	84.8	10.2	4.0	0.9
Low Socioeconomic Status (SES)	38,137	3592.67	36.50	62.6	19.6	14.0	3.8
Migrant	463	3588.19	35.02	66.3	20.7	9.5	3.5

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.11. Test Results by Subgroup, Mathematics Grade 7

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	82,799	3625.99	46.01	54.1	15.8	15.3	14.8
Male	42,140	3628.32	47.05	51.6	15.8	16.1	16.5
Female	40,659	3623.58	44.78	56.7	15.8	14.5	13.0
Hispanic	39,451	3613.71	41.33	65.9	14.7	11.7	7.7
Non-Hispanic	43,348	3637.18	47.18	43.3	16.9	18.6	21.2
American Indian	4,690	3601.01	37.27	77.5	11.6	7.3	3.6
Asian	2,855	3666.80	48.59	22.4	12.4	19.6	45.6
Black or African American	5,949	3607.92	39.53	71.6	12.7	10	5.8
Multi-racial	4,780	3631.65	45.56	48.3	17.5	17.1	17.1
Native Hawaiian or Other Pacific Islander	457	3619.95	41.39	57.5	20.1	11.4	10.9
White	63,986	3627.31	45.40	52.6	16.4	16.1	14.9
Missing	82	3622.05	50.45	51.2	9.8	26.8	12.2
Special Education	10,598	3588.78	35.04	87.0	6.4	4.0	2.6
English Learner (EL)	7,134	3585.65	28.88	91.5	5.3	2.5	0.7
Low Socioeconomic Status (SES)	37,648	3610.66	40.28	68.5	14.4	10.8	6.3
Migrant	461	3601.84	39.65	76.6	11.9	7.2	4.3

Note. SS = scale score, 1 = Minimally Proficient, 2 = Partially Proficient, 3 = Proficient, 4 = Highly Proficient

Table C.12. Test Results by Subgroup, Mathematics Grade 8

Subgroup	N	SS Mean	SS SD	%Level 1	%Level 1	%Level 1	%Level 1
All	86,031	3653.61	37.08	53.9	19.1	16.6	10.4
Male	44,172	3654.32	38.77	53.5	18.5	16.4	11.6
Female	41,859	3652.87	35.20	54.4	19.7	16.9	9.1
Hispanic	41,156	3644.18	31.79	65.2	17.6	11.9	5.3
Non-Hispanic	44,875	3662.27	39.41	43.6	20.4	21.0	15.0
American Indian	4,949	3637.14	27.80	75.7	13.0	8.2	3.0
Asian	2,883	3690.67	45.71	20.8	15.9	26.0	37.2
Black or African American	5,924	3641.49	29.46	67.8	18.0	10.3	4.0
Multi-racial	4,917	3657.05	37.17	48.9	20.8	18.7	11.6
Native Hawaiian or Other Pacific Islander	517	3648.27	34.91	60.5	18.4	13.3	7.7
White	66,747	3654.08	36.61	52.8	19.7	17.3	10.2
Missing	94	3664.38	40.96	42.6	12.8	23.4	21.3
Special Education	10,313	3628.38	23.91	86.2	8.4	4.0	1.4
English Learner (EL)	7,017	3626.38	20.26	89.8	7.0	2.7	0.5
Low Socioeconomic Status (SES)	38,145	3642.43	30.65	67.2	17.4	10.8	4.6
Migrant	508	3635.52	30.02	79.7	9.1	7.3	3.9

Note. SS = scale score, 1 = *Minimally Proficient*, 2 = *Partially Proficient*, 3 = *Proficient*, 4 = *Highly Proficient*

Figure C.1. Total Scale Score Distribution, ELA Grade 3

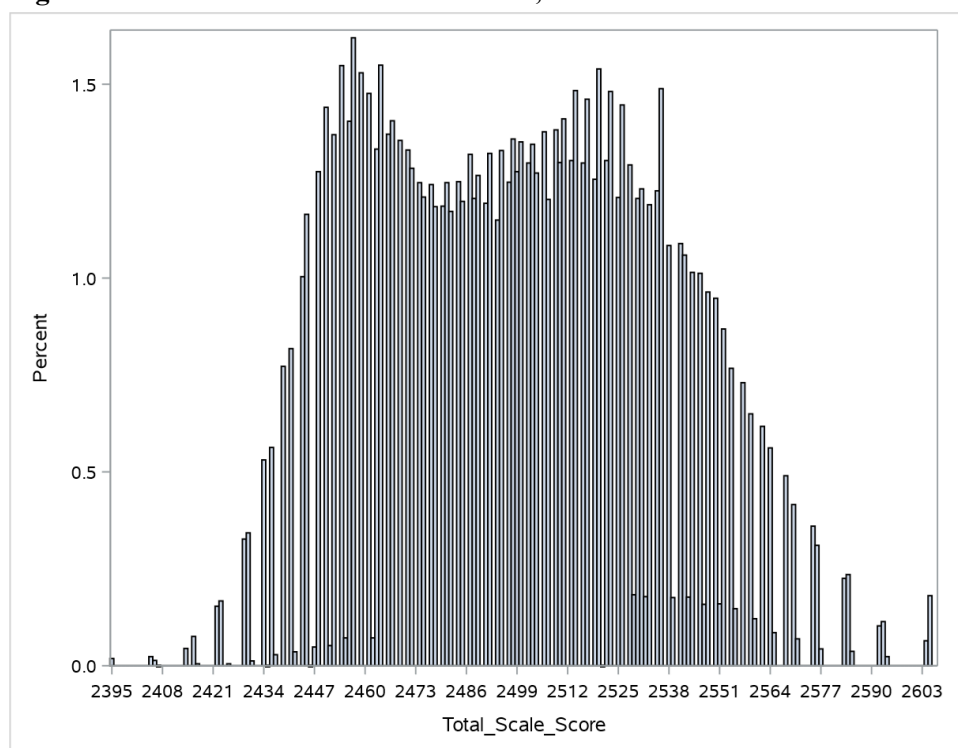


Figure C.2. Total Scale Score Distribution, ELA Grade 4

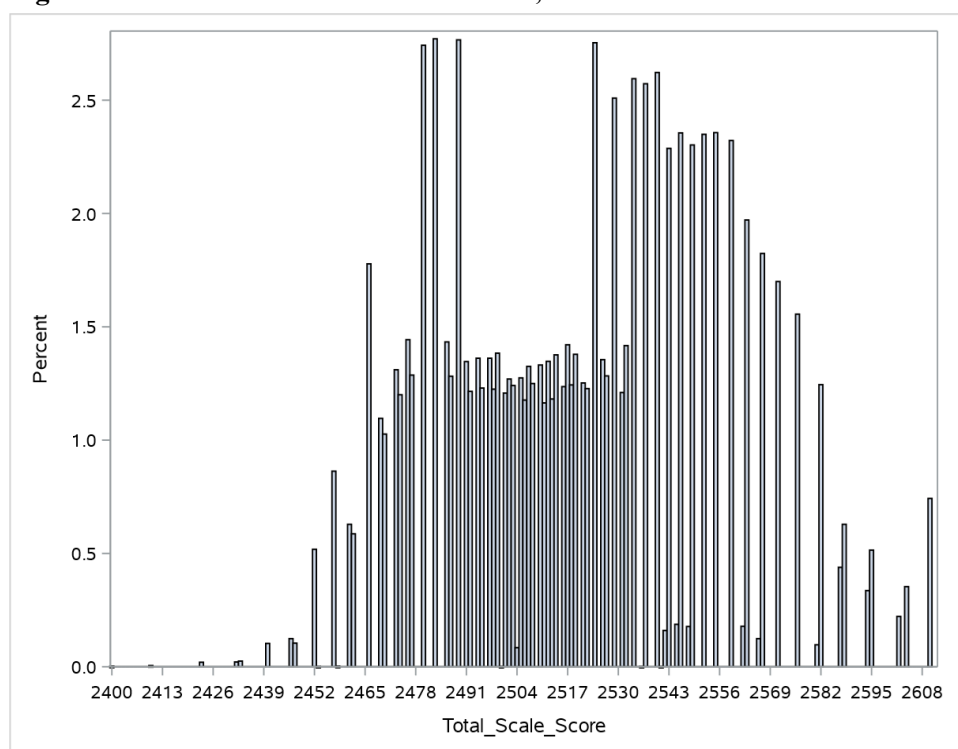


Figure C.3. Total Scale Score Distribution, ELA Grade 5

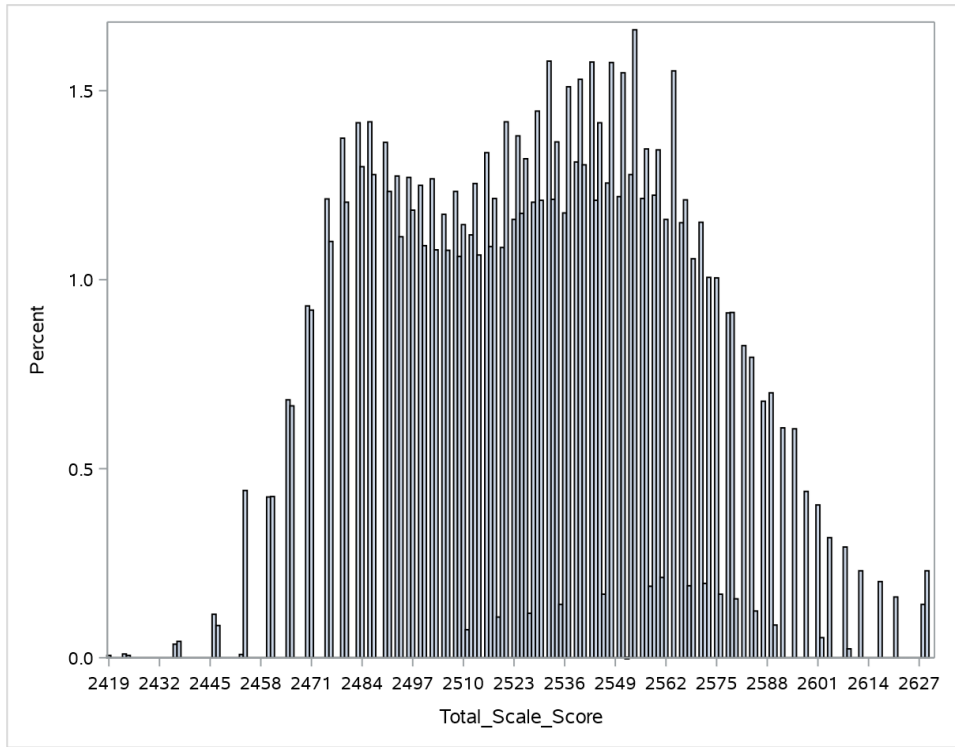


Figure C.4. Total Scale Score Distribution, ELA Grade 6

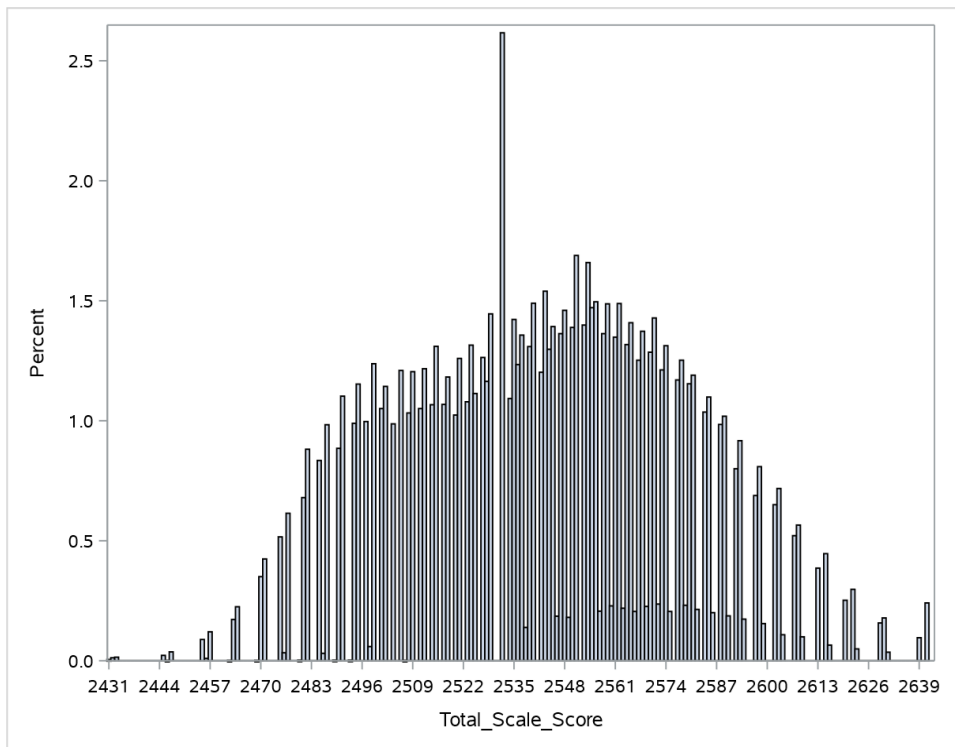


Figure C.5. Total Scale Score Distribution, ELA Grade 7

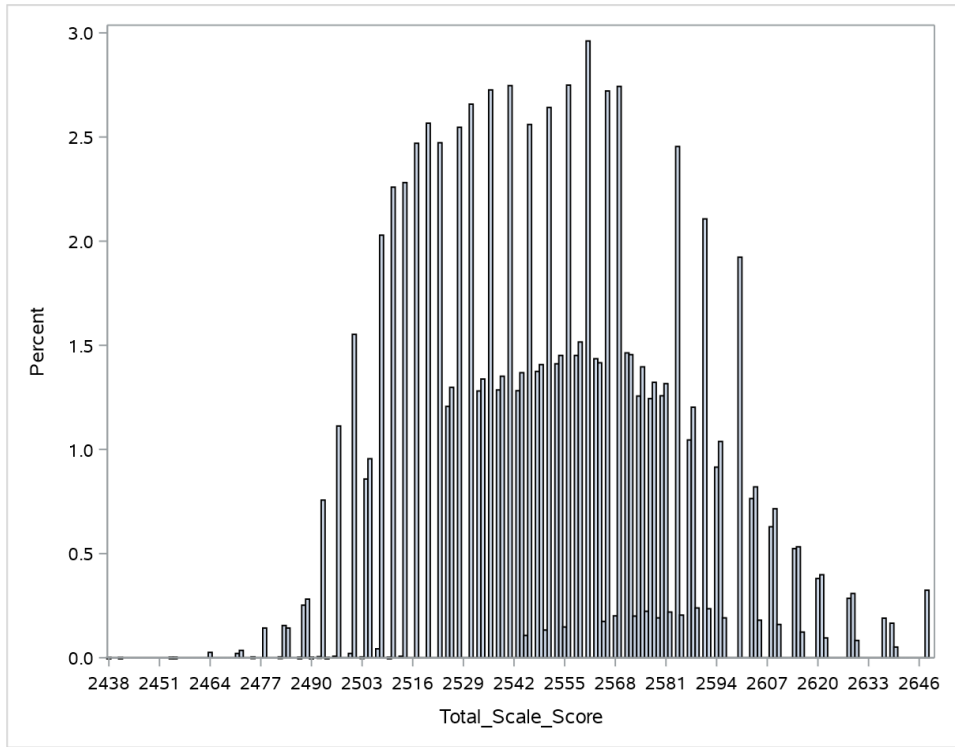


Figure C.6. Total Scale Score Distribution, ELA Grade 8

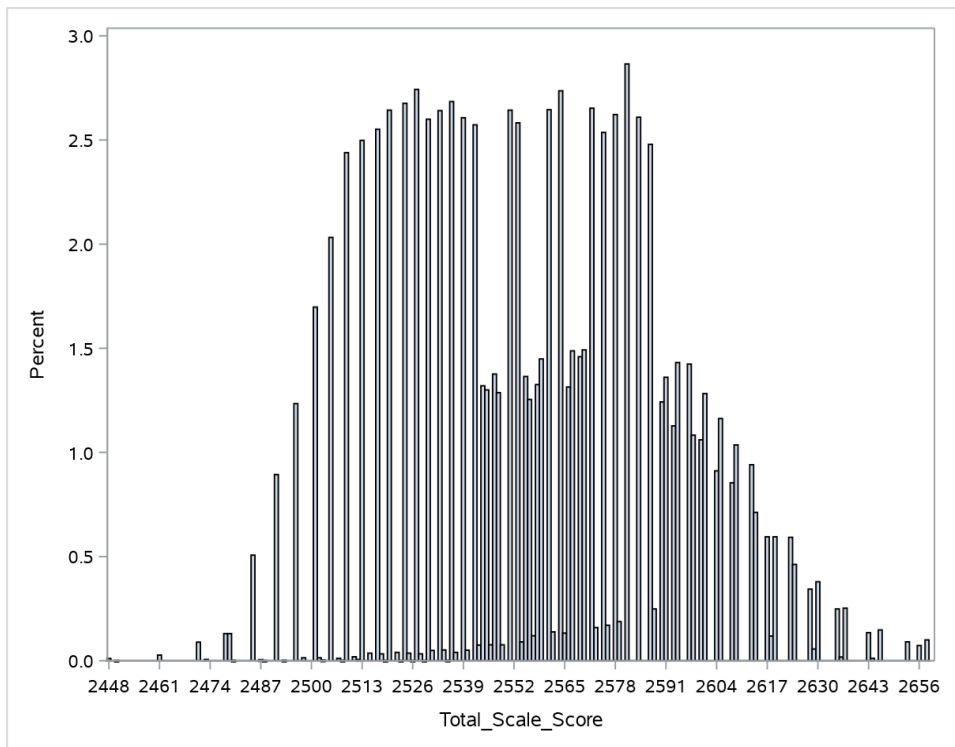


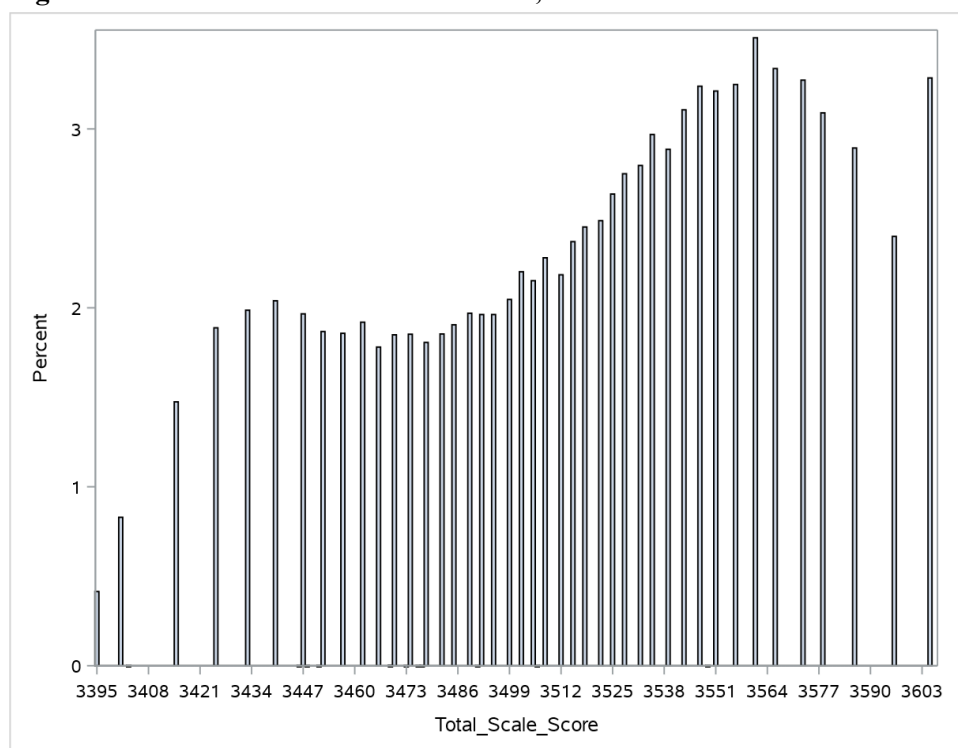
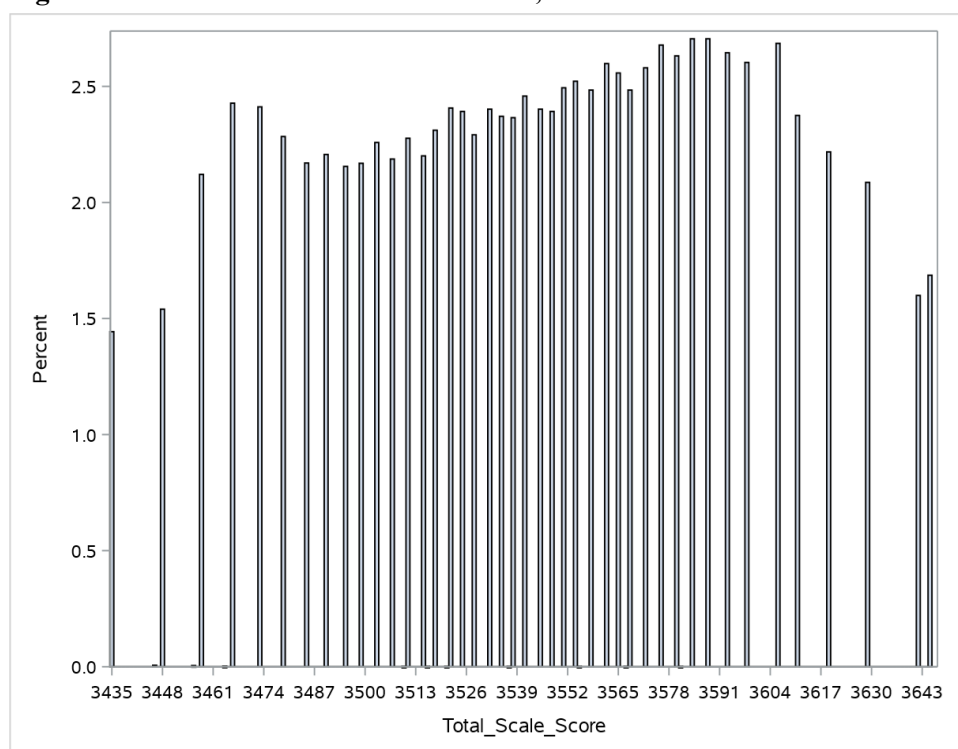
Figure C.7. Total Scale Score Distribution, Mathematics Grade 3**Figure C.8. Total Scale Score Distribution, Mathematics Grade 4**

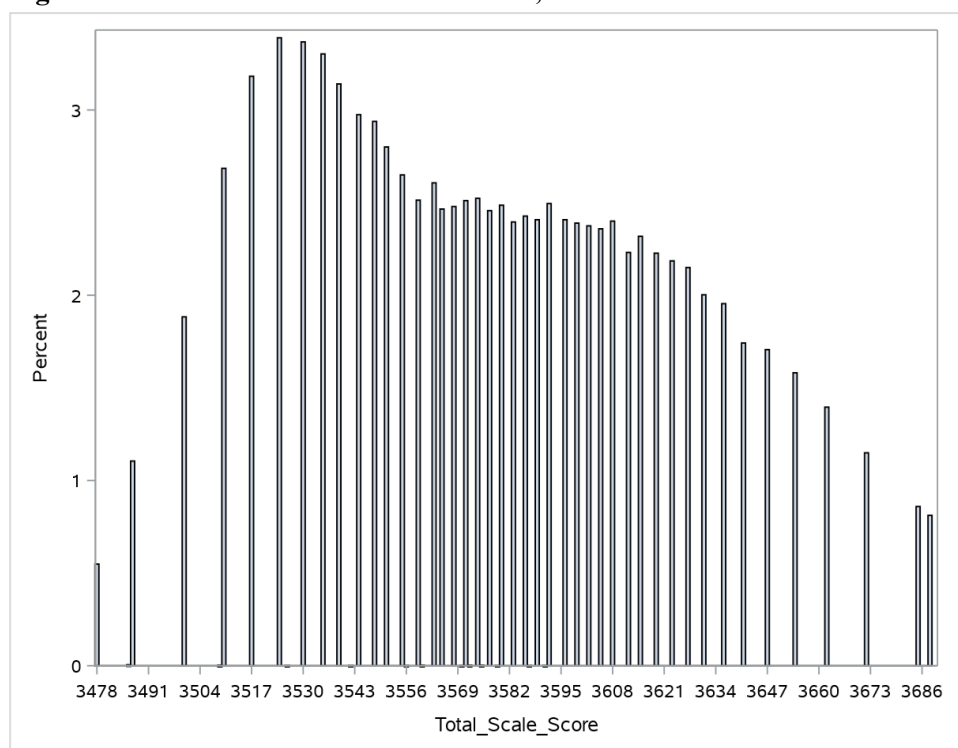
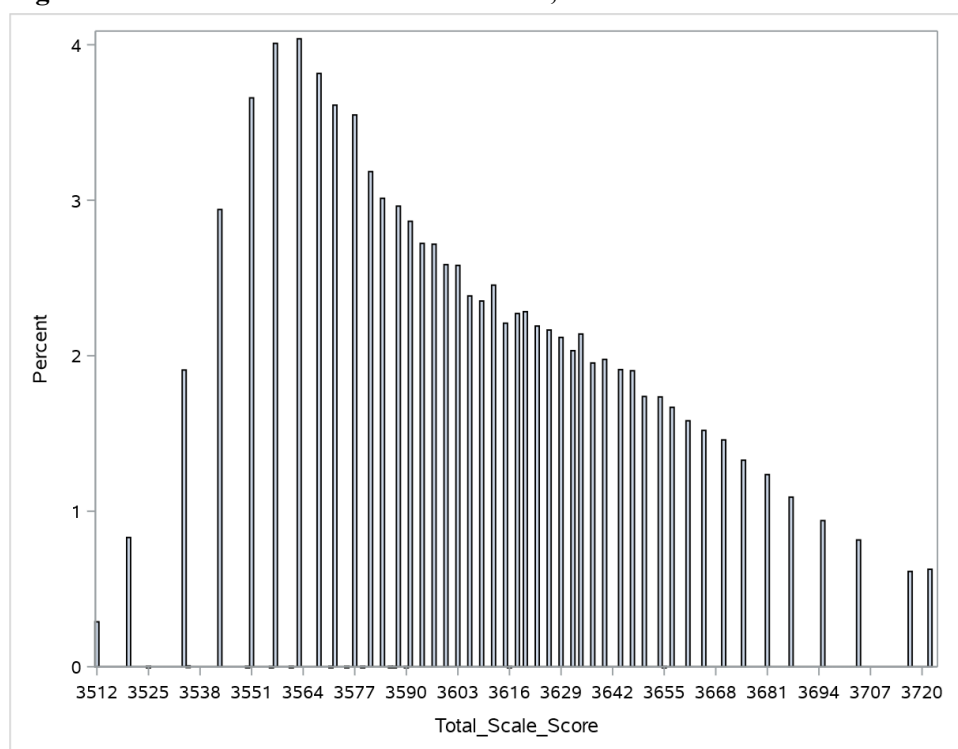
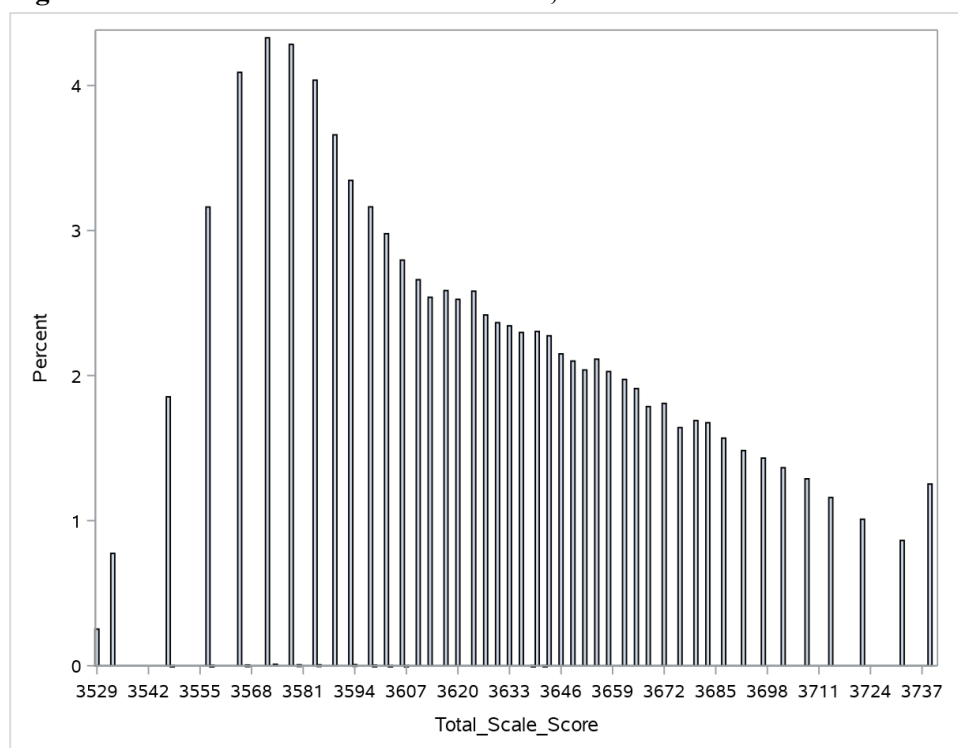
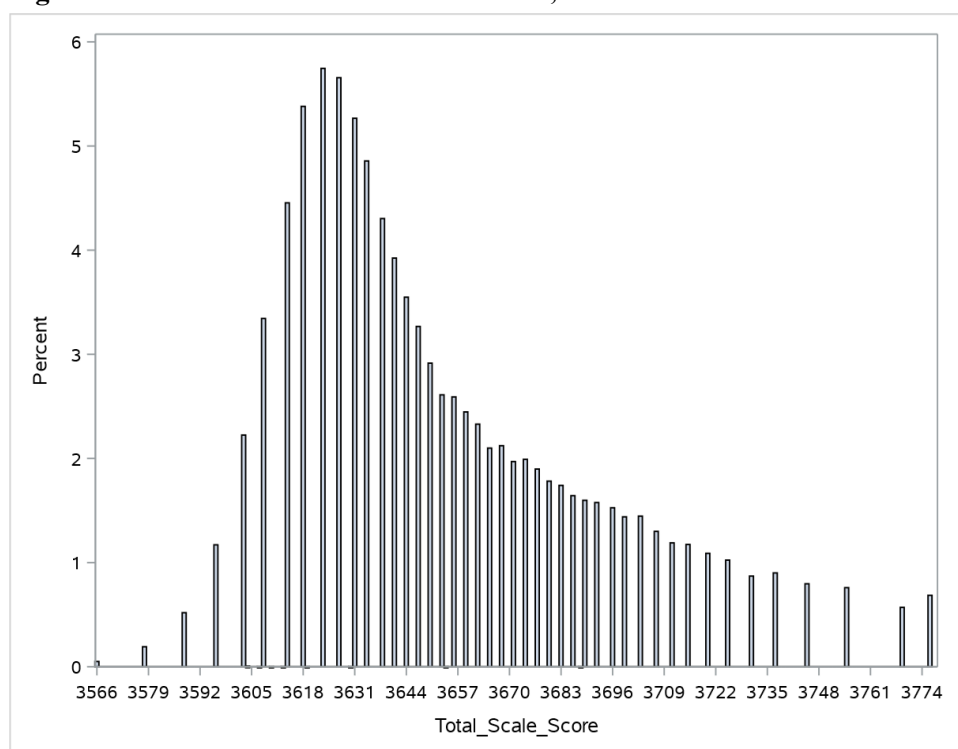
Figure C.9. Total Scale Score Distribution, Mathematics Grade 5**Figure C.10. Total Scale Score Distribution, Mathematics Grade 6**

Figure C.11. Total Scale Score Distribution, Mathematics Grade 7**Figure C.12. Total Scale Score Distribution, Mathematics Grade 8**

Appendix D: ACT GRADE 8 LINKING STUDY

This appendix presents the results of the study conducted by ACT to link AASA to the ACT Aspire scale using Spring 2022 data.

D.1. Purpose

Reporting ACT test score predictions can help Grade 8 students and stakeholders understand college readiness and plan course selection as students prepare to enter high school. ACT Aspire Mathematics and Reading items were included in the Spring 2022 administration of the Grade 8 AASA Mathematics and ELA assessments to be able to conduct a linking study and provide predictions of performance on the ACT assessment. Based on this study, students may use their AASA mathematics and ELA scores to obtain predicted ACT scores that will help students understand their predicted college readiness and plan future course work. This report summarizes the linking and prediction studies.

D.2. Background

ACT score predictions are available for ACT Aspire scale scores for Grades 3-10 for all subject areas (English, mathematics, reading, science, ELA, STEM, and Composite) and are available for the AASA ELA and Mathematics assessments. The data and methodology used for the established ACT score predictions is described in Section 14.1 of the 2020 [ACT Aspire Technical Manual](#).

Using data from the Spring 2022 AASA administration, a chained equipercentile concordance study was conducted to establish the relationship between AASA scale scores and ACT Aspire scale scores. The Aspire scores were then used to predict ACT scores. Concordance can be used to link two tests measuring similar constructs and intended for similar populations (Dorans, 2004; Holland & Dorans, 2006; Kolen & Brennan, 2014), as was the case of the AASA and Aspire assessment. The following steps were used for the concordance:

- **Phase I (before test administration):** ACT Aspire items were reviewed for content alignment to the Arizona standards and statistical characteristics for consideration as candidate items to be used in creating the link between the AASA and Aspire assessments. Statistics for the candidate items were based on the national administration of the ACT Aspire assessment. Approved items were embedded in the AASA test forms for the Spring 2022 administration.
- **Phase II (after test administration):** Item statistics (i.e., p -value, point-biserial) based on the Arizona student data were compared to the statistics of the same items based on the ACT Aspire national sample. ACT then developed the concordances between (a) AASA scale scores and the raw score for the common item set and (b) ACT Aspire scale scores and the raw score for the common item set. The AASA to ACT Aspire concordance were then generated by combining the two concordances. Based on the resulting concordance and existing ACT Aspire to ACT test score prediction table, each AASA scale score was assigned a predicted ACT score range. The resulting predictions will be used for three years until longitudinal data are available for Arizona-specific predicted ACT scores. Arizona's longitudinal data will be used to establish a direct prediction of ACT scores. For example, the Grade 8 test results from 2022 will be matched to the Grade 11 ACT test results from 2025 to predict ACT scores directly.

D.3. Results

Table D.1 – Table D.5 present the summary statistics of the ACT Aspire scores to create the short-to-long concordance and the AASA scores to create the long-to-short concordance. The lowest obtainable scale score (LOSS) and highest obtainable scale score (HOSS) for each assessment are indicated as a note in each table, along with the maximum possible raw score for both the total test and common item set.

Table D.1. Summary Statistics of ACT Aspire Scores to Create Short-to-Long Concordance—Mathematics

Score	N	Mean	SD	Min.	Max.	Pearson's Correlation	
Scale Score (A)	6,688	425.79	7.85	401	456	$r(A, B) = 0.99$	$r(A, C) = 0.91$
Total Raw Score (B)	6,688	21.38	8.72	1	51	–	$r(B, C) = 0.93$
Common-Item Raw Score (C)	6,688	8.62	4.04	0	20	–	–

Note. The LOSS is 400, and the HOSS is 456. The maximum possible raw score is 51 for the total test and 20 for the common item set.

Table D.2. Summary Statistics of ACT Aspire Scores to Create Short-to-Long Concordance—Reading

Score	N	Mean	SD	Min.	Max.	Pearson's Correlation	
Scale Score (A)	6,869	422.84	7.47	400	440	$r(A, B) = 0.99$	$r(A, C) = 0.91$
Total Raw Score (B)	6,869	15.09	6.59	0	30	–	$r(B, C) = 0.91$
Common-Item Raw Score (C)	6,869	6.90	2.97	0	12	–	–

Note. The LOSS is 400, and the HOSS is 440. The maximum possible raw score is 30 for the total test and 12 for the common item set.

Table D.3. Summary Statistics of ACT Aspire Scores to Create Short-to-Long Concordance—ELA

N	Mean	SD	Min.	Max.	Pearson's Correlation between ELA and Reading Scale Scores
6,818	425.73	6.62	404	444	0.89

Note. The LOSS is 403, and the HOSS is 447.

Table D.4. Summary Statistics of AASA Scores to Create Long-to-Short Concordance—Mathematics

Score	N	Mean	SD	Min.	Max.	Pearson's Correlation	
Scale Score (A)	82,328	3651.72	35.59	3566	3776	$r(A, B) = 0.99$	$r(A, C) = 0.91$
Total Raw Score (B)	82,328	16.85	10.07	0	47	–	$r(B, C) = 0.92$
Common-Item Raw Score (C)	82,328	8.07	4.52	0	20	–	–

Note. The LOSS is 3566, and the HOSS is 3776. The maximum possible raw score is 47 for the total test and 20 for the common item set.

Table D.5. Summary Statistics of AASA Scores to Create Long-to-Short Concordance—ELA

Score	N	Mean	SD	Min.	Max.	Pearson's Correlation	
Scale Score (A)	81,428	2557.77	32.83	2448	2658	$r(A, B) = 1.00$	$r(A, C) = 0.83$
Total Raw Score (B)	81,428	27.75	10.25	3	55	–	$r(B, C) = 0.84$
Common-Item Raw Score (C)	81,428	6.73	3.04	0	12	–	–

Note. The LOSS is 2448, and the HOSS is 2658. The maximum possible raw score is 55 for the total test and 12 for the common item set.

Table D.6 and Table D.7 present the item statistics for the common items for the ACT Aspire Mathematics and Reading and the AASA Mathematics and ELA assessments, and Table D.8 and Table D.9 present the final concordance tables.

Table D.6. Item Statistics for the Common Items—Mathematics

Item	<i>P</i> -value		Point-Biserial (Total Raw Score)		Point-Biserial (Total Common Item Score)	
	AASA	ACT Aspire	AASA	ACT Aspire	AASA	ACT Aspire
M_MC1	0.64	0.76	0.48	0.42	0.53	0.44
M_MC2	0.60	0.62	0.43	0.38	0.46	0.43
M_MC3	0.31	0.38	0.59	0.59	0.60	0.60
M_MC4	0.48	0.55	0.47	0.47	0.51	0.49
M_MC5	0.55	0.69	0.47	0.34	0.50	0.37
M_MC6	0.30	0.30	0.33	0.29	0.36	0.34
M_MC7	0.38	0.39	0.38	0.35	0.44	0.40
M_MC8	0.44	0.48	0.39	0.35	0.46	0.40
M_MC9	0.27	0.32	0.45	0.41	0.45	0.45
M_MC10	0.39	0.42	0.26	0.23	0.32	0.29
M_MC11	0.45	0.38	0.47	0.43	0.50	0.46
M_MC12	0.24	0.25	0.46	0.47	0.47	0.51
M_MC13	0.40	0.40	0.47	0.46	0.50	0.49
M_MC14	0.29	0.26	0.25	0.25	0.29	0.29
M_MC15	0.25	0.19	0.52	0.43	0.52	0.46
M_MC16	0.24	0.20	0.49	0.44	0.53	0.46
M_MC17	0.27	0.20	0.46	0.42	0.49	0.44
M_MC18	0.57	0.65	0.44	0.42	0.51	0.45
M_TE1	0.55	0.64	0.47	0.43	0.53	0.46
M_TE2	0.47	0.53	0.54	0.49	0.57	0.51

Table D.7. Item Statistics for the Common Items—Reading/ELA

Item	<i>P</i> -value		Point-Biserial (Total Raw Score)		Point-Biserial (Total Common Item Score)	
	AASA	ACT Aspire	AASA	ACT Aspire	AASA	ACT Aspire
R_MC1	0.57	0.66	0.55	0.53	0.64	0.44
R_MC2	0.52	0.58	0.31	0.37	0.43	0.43
R_MC3	0.59	0.55	0.54	0.54	0.61	0.60
R_MC4	0.65	0.70	0.32	0.45	0.42	0.49
R_MC5	0.56	0.65	0.46	0.50	0.56	0.37
R_MC6	0.72	0.64	0.56	0.60	0.62	0.34
R_MC7	0.19	0.16	0.19	0.19	0.23	0.40
R_MC8	0.67	0.57	0.55	0.59	0.61	0.40
R_MC9	0.63	0.67	0.43	0.46	0.53	0.45
R_MC10	0.54	0.63	0.49	0.51	0.59	0.29
R_MC11	0.61	0.67	0.36	0.41	0.46	0.46
R_TE1	0.49	0.40	0.54	0.53	0.60	0.51

Table D.8. Final Concordance Table—Mathematics

AASA	Common	ACT Aspire	AASA	Common	ACT Aspire
3566	0	402	3608	2	412
3567	0	402	3609	2	412
3568	0	402	3610	3	415
3569	0	402	3611	3	415
3570	0	402	3612	3	415
3571	0	402	3613	3	415
3572	0	402	3614	3	415
3573	0	402	3615	3	415
3574	0	402	3616	3	415
3575	0	402	3617	3	415
3576	0	402	3618	3	415
3577	0	402	3619	4	416
3578	0	402	3620	4	416
3579	0	402	3621	4	416
3580	0	402	3622	4	416
3581	0	402	3623	4	416
3582	0	402	3624	4	416
3583	0	402	3625	4	416
3584	1	408	3626	5	418
3585	1	408	3627	5	418
3586	1	408	3628	5	418
3587	1	408	3629	5	418
3588	1	408	3630	5	418
3589	1	408	3631	5	418
3590	1	408	3632	5	418
3591	1	408	3633	5	418
3592	1	408	3634	6	421
3593	1	408	3635	6	421
3594	1	408	3636	6	421
3595	1	408	3637	6	421
3596	1	408	3638	6	421
3597	1	408	3639	6	421
3598	1	408	3640	6	421
3599	2	412	3641	6	421
3600	2	412	3642	7	423
3601	2	412	3643	7	423
3602	2	412	3644	7	423
3603	2	412	3645	7	423
3604	2	412	3646	7	423
3605	2	412	3647	7	423
3606	2	412	3648	8	425
3607	2	412	3649	8	425

Appendix D: ACT Grade 8 Linking Study

AASA	Common	ACT Aspire	AASA	Common	ACT Aspire
3650	8	425	3693	14	436
3651	8	425	3694	14	436
3652	8	425	3695	14	436
3653	8	425	3696	14	436
3654	8	425	3697	14	436
3655	9	427	3698	15	437
3656	9	427	3699	15	437
3657	9	427	3700	15	437
3658	9	427	3701	15	437
3659	9	427	3702	15	437
3660	9	427	3703	15	437
3661	9	427	3704	15	437
3662	10	429	3705	15	437
3663	10	429	3706	16	438
3664	10	429	3707	16	438
3665	10	429	3708	16	438
3666	10	429	3709	16	438
3667	10	429	3710	16	438
3668	10	429	3711	16	438
3669	11	432	3712	16	438
3670	11	432	3713	16	438
3671	11	432	3714	17	439
3672	11	432	3715	17	439
3673	11	432	3716	17	439
3674	11	432	3717	17	439
3675	11	432	3718	17	439
3676	12	433	3719	17	439
3677	12	433	3720	17	439
3678	12	433	3721	17	439
3679	12	433	3722	17	439
3680	12	433	3723	17	439
3681	12	433	3724	17	439
3682	12	433	3725	17	439
3683	12	433	3726	18	440
3684	13	435	3727	18	440
3685	13	435	3728	18	440
3686	13	435	3729	18	440
3687	13	435	3730	18	440
3688	13	435	3731	18	440
3689	13	435	3732	18	440
3690	13	435	3733	18	440
3691	14	436	3734	18	440
3692	14	436	3735	18	440

AASA	Common	ACT Aspire	AASA	Common	ACT Aspire
3736	18	440	3757	19	443
3737	18	440	3758	19	443
3738	18	440	3759	19	443
3739	18	440	3760	19	443
3740	18	440	3761	20	451
3741	19	443	3762	20	451
3742	19	443	3763	20	451
3743	19	443	3764	20	451
3744	19	443	3765	20	451
3745	19	443	3766	20	451
3746	19	443	3767	20	451
3747	19	443	3768	20	451
3748	19	443	3769	20	451
3749	19	443	3770	20	451
3750	19	443	3771	20	451
3751	19	443	3772	20	451
3752	19	443	3773	20	451
3753	19	443	3774	20	451
3754	19	443	3775	20	451
3755	19	443	3776	20	451
3756	19	443			

Table D.9. Final Concordance Table—ELA

AASA	Common	ACT Aspire	AASA	Common	ACT Aspire
2448	0	406	2468	0	406
2449	0	406	2469	0	406
2450	0	406	2470	0	406
2451	0	406	2471	0	406
2452	0	406	2472	0	406
2453	0	406	2473	0	406
2454	0	406	2474	0	406
2455	0	406	2475	0	406
2456	0	406	2476	0	406
2457	0	406	2477	0	406
2458	0	406	2478	0	406
2459	0	406	2479	0	406
2460	0	406	2480	0	406
2461	0	406	2481	0	406
2462	0	406	2482	0	406
2463	0	406	2483	0	406
2464	0	406	2484	0	406
2465	0	406	2485	0	406
2466	0	406	2486	0	406
2467	0	406	2487	0	406
			2488	1	411

Appendix D: ACT Grade 8 Linking Study

AASA	Common	ACT Aspire	AASA	Common	ACT Aspire
2489	1	411	2537	5	423
2490	1	411	2538	5	423
2491	1	411	2539	5	423
2492	1	411	2540	5	423
2493	1	411	2541	5	423
2494	1	411	2542	5	423
2495	1	411	2543	5	423
2496	1	411	2544	5	423
2497	1	411	2545	6	424
2498	1	411	2546	6	424
2499	1	411	2547	6	424
2500	1	411	2548	6	424
2501	1	411	2549	6	424
2502	1	411	2550	6	424
2503	1	411	2551	6	424
2504	1	411	2552	6	424
2505	2	415	2553	7	426
2506	2	415	2554	7	426
2507	2	415	2555	7	426
2508	2	415	2556	7	426
2509	2	415	2557	7	426
2510	2	415	2558	7	426
2511	2	415	2559	7	426
2512	2	415	2560	7	426
2513	2	415	2561	7	426
2514	2	415	2562	8	428
2515	2	415	2563	8	428
2516	3	418	2564	8	428
2517	3	418	2565	8	428
2518	3	418	2566	8	428
2519	3	418	2567	8	428
2520	3	418	2568	8	428
2521	3	418	2569	8	428
2522	3	418	2570	8	428
2523	3	418	2571	9	430
2524	3	418	2572	9	430
2525	3	418	2573	9	430
2526	3	418	2574	9	430
2527	4	420	2575	9	430
2528	4	420	2576	9	430
2529	4	420	2577	9	430
2530	4	420	2578	9	430
2531	4	420	2579	9	430
2532	4	420	2580	9	430
2533	4	420	2581	9	430
2534	4	420	2582	9	430
2535	4	420	2583	10	432
2536	5	423	2584	10	432

Appendix D: ACT Grade 8 Linking Study

AASA	Common	ACT Aspire	AASA	Common	ACT Aspire
2585	10	432	2622	12	439
2586	10	432	2623	12	439
2587	10	432	2624	12	439
2588	10	432	2625	12	439
2589	10	432	2626	12	439
2590	10	432	2627	12	439
2591	10	432	2628	12	439
2592	10	432	2629	12	439
2593	10	432	2630	12	439
2594	10	432	2631	12	439
2595	10	432	2632	12	439
2596	10	432	2633	12	439
2597	10	432	2634	12	439
2598	10	432	2635	12	439
2599	11	434	2636	12	439
2600	11	434	2637	12	439
2601	11	434	2638	12	439
2602	11	434	2639	12	439
2603	11	434	2640	12	439
2604	11	434	2641	12	439
2605	11	434	2642	12	439
2606	11	434	2643	12	439
2607	11	434	2644	12	439
2608	11	434	2645	13	439
2609	11	434	2646	13	439
2610	11	434	2647	13	439
2611	11	434	2648	13	439
2612	11	434	2649	13	439
2613	11	434	2650	13	439
2614	11	434	2651	13	439
2615	11	434	2652	13	439
2616	11	434	2653	13	439
2617	11	434	2654	13	439
2618	11	434	2655	13	439
2619	11	434	2656	13	439
2620	12	439	2657	13	439
2621	12	439	2658	13	439

D.4. References

- Dorans, N. J. (2004). Equating, concordance, and expectation. *Applied Psychological Measurement*, 28(4), 227–246. <https://doi.org/10.1177/0146621604265031>
- Holland, P. W., & Dorans, N. J. (2006). Linking and equating. In R. L. Brennan (Ed.), *Educational measurement* (4th ed., pp. 187–220). Westport.
- Kolen, M. J., & Brennan, R. L. (2014). Test equating, scaling, and linking. *Methods and practices* (3rd ed.). Springer.