



New Meridian
Technical Report 2021–2022
Alternate Blueprint
January 13, 2023

Table of Contents

Executive Summary.....	1
Section 1: Introduction.....	4
1.1 Background.....	4
1.2 Purpose of the Operational Tests.....	5
1.3 Composition of Operational Tests.....	5
1.4 Intended Population.....	6
1.5 Groups and Organizations Involved with the Summative Assessments.....	6
1.6 Overview of the Technical Report.....	6
1.7 Glossary of Abbreviations.....	9
Section 2: Test Development.....	11
2.1 Overview of the Summative Assessments, Claims, and Design.....	11
2.1.1 English Language Arts/Literacy Assessments—Claims and Subclaims.....	11
2.1.2 Mathematics Assessments—Claims and Subclaims.....	12
2.2 Test Development Activities.....	13
2.2.1 Item Development Process.....	13
2.2.2 Item and Text Review Committees.....	14
2.2.3 Operational Test Construction.....	15
2.2.4 Linking Design of the Operational Test.....	18
2.2.5 Field Test Data Collection Overview.....	18
Section 3: Test Administration.....	20
3.1 Test Security and Administration Policies.....	20
3.1.1 Secure versus Nonsecure Materials.....	20
3.1.2 Scorable versus Nonscorable Materials.....	20
3.2 Accessibility Features and Accommodations.....	21
3.2.1 Participation Guidelines for Assessments.....	21
3.2.2 Accessibility System.....	22
3.2.3 What Are Accessibility Features?.....	22
3.2.4 Accommodations for Students with Disabilities and English Learners.....	22
3.2.5 Unique Accommodations.....	23
3.2.6 Emergency Accommodations.....	24
3.2.7 Student Refusal Form.....	24
3.3 Testing Irregularities and Security Breaches.....	24
3.4 Data Forensics Analyses.....	26
3.4.1 Response Change Analysis.....	26
3.4.2 Aberrant Response Analysis.....	26
3.4.3 Plagiarism Analysis.....	27
3.4.4 Longitudinal Performance Monitoring.....	27
3.4.5 Internet and Social Media Monitoring.....	27

3.4.6 Off-Hours Testing Monitoring.....	27
Section 4: Item Scoring.....	29
4.1 Machine-Scored Items.....	29
4.1.1 Key-Based Items.....	29
4.1.2 Rule-Based Items	29
4.2 Human or Handscored Items.....	30
4.2.1 Scorer Training	31
4.2.2 Scorer Qualification	33
4.2.3 Managing Scoring.....	34
4.2.4 Monitoring Scoring	34
4.3 Automated Scoring for Prose Constructed Responses.....	36
4.3.1 Concepts Related to Automated Scoring.....	37
4.3.2 Sampling Responses Used for Training IEA.....	38
4.3.3 Primary Criteria for Evaluating IEA Performance.....	39
4.3.4 Contingent Primary Criteria for Evaluating IEA Performance.....	39
4.3.5 Applying Smart Routing	39
4.3.6 Evaluation of Secondary Criteria for Evaluating IEA Performance	40
4.3.7 Inter-Rater Agreement for Prose Constructed Response	42
Section 5: Classical Item Analysis	43
5.1 Overview.....	43
5.2 Data Screening Criteria.....	43
5.3 Description of Classical Item Analysis Statistics	43
5.4 Summary of Classical Item Analysis Flagging Criteria.....	45
5.5 Classical Item Analysis Results.....	46
Section 6: Differential Item Functioning.....	49
6.1 Overview.....	49
6.2 DIF Procedures	49
6.3 Operational Analysis DIF Comparison Groups	51
6.4 Operational Differential Item Functioning Results.....	52
Section 7: IRT Model and Parameters.....	54
7.1 Overview.....	54
7.2 Two-Parameter Logistic/Generalized Partial Credit Model	54
7.3 Summary Statistics and Distributions from IRT Analyses.....	54
7.3.1 IRT Summary Statistics for English Language Arts/Literacy.....	54
7.3.2 IRT Summary Statistics for Mathematics	55
Section 8: Performance Level Setting.....	57
8.1 Performance Standards.....	57
8.2 Performance Levels and Policy Definitions	57
8.3 Performance Level Setting Process for the Assessment System.....	59

8.3.1 Research Studies	60
8.3.2 Pre-Policy Meeting.....	60
8.3.3 Performance Level Setting Meetings.....	60
8.3.4 Post-Policy Reasonableness Review	61
Section 9: Quality Control Procedures.....	63
9.1 Quality Control of the Item Bank.....	63
9.2 Quality Control of Test Form Development	63
9.3 Quality Control of Test Materials	64
9.4 Quality Control of Scanning.....	65
9.5 Quality Control of Image Editing	66
9.6 Quality Control of Answer Document Processing and Scoring	66
9.7 Quality Control of Psychometric Processes.....	67
Section 10: Operational Test Forms.....	69
Section 11: Student Characteristics	70
11.1 Overview of Test Taking Population.....	70
11.2 Rules for Inclusion of Students in Analyses.....	70
11.3 Students by Grade/Course, Mode, and Gender.....	70
11.4 Demographics.....	72
Section 12: Scale Scores.....	73
12.1 Operational Test Content (Claims and Subclaims)	73
12.1.1 English Language Arts/Literacy	73
12.1.2 Mathematics	75
12.2 Establishing the Reporting Scales.....	75
12.2.1 Summative Score Scale and Performance Levels	75
12.2.2 ELA/L Reading and Writing Claim Scale	77
12.2.3 Subclaims Scale	77
12.3 Creating Conversion Tables	78
12.4 Score Distributions	80
12.4.1 Score Distributions for English Language Arts/Literacy.....	80
12.4.2 Scale Score Cumulative Frequencies for ELA/L.....	85
12.4.3 Summary Scale Score Statistics for ELA/L Groups	85
12.4.4 Score Distributions for Mathematics.....	89
12.4.5 Scale Score Cumulative Frequencies for Mathematics.....	89
12.4.6 Summary Scale Score Statistics for Mathematics Groups	91
12.5 Interpreting Claim Scores and Subclaim Scores	93
12.5.1 Interpreting Claim Scores	93
12.5.2 Interpreting Subclaim Scores.....	93
Section 13: Reliability.....	94
13.1 Overview.....	94

13.2 Reliability and SEM Estimation.....	95
13.2.1 Raw Score Reliability Estimation.....	95
13.2.2 Scale Score Reliability Estimation	96
13.3 Reliability Results for Total Group.....	97
13.3.1 Raw Score Reliability Results	97
13.3.2 Scale Score Reliability Results.....	98
13.4 Reliability Results for Subgroups of Interest	99
13.4.1 Reliability Results for Gender	99
13.4.2 Reliability Results for Ethnicity	99
13.4.3 Reliability Results for Special Education Needs	100
13.4.4 Reliability Results for Students Taking Accommodated Forms	100
13.4.5 Reliability Results of Students Taking Translated Forms	100
13.5 Reliability Results for ELA/L Claims and Subclaims	103
13.6 Reliability Results for Mathematics Subclaims.....	105
13.7 Reliability of Classification.....	107
13.7.1 English Language Arts/Literacy.....	107
13.7.2 Mathematics	108
13.8 Inter-Rater Agreement	109
Section 14: Validity.....	111
14.1 Overview.....	111
14.2 Evidence Based on Test Content	111
14.3 Evidence Based on Internal Structure	113
14.3.1 Intercorrelations	113
14.3.2 Reliability.....	122
14.3.3 Local Item Dependence	122
14.4 Evidence Based on Relationships to Other Variables.....	126
14.5 Evidence from Special Studies.....	129
14.5.1 Content Alignment Studies	129
14.5.2 Benchmarking Study	131
14.5.3 Longitudinal Study of External Validity of Performance Levels (Phase 1)	132
14.5.4 Mode and Device Comparability Studies.....	133
14.5.5 Quality Testing Standards	134
14.6 Evidence Based on Response Processes.....	143
14.7 Interpretations of Test Scores	144
14.8 Evidence Based on the Consequences to Testing	144
14.9 Summary.....	145
Section 15: Student Growth Measures.....	147
15.1 Norm Groups	147
15.2 Student Growth Percentile Estimation.....	150

15.3 Student Growth Percentile Results/Model Fit for Total Group	150
15.4 Student Growth Percentile Results for Subgroups of Interest	152
15.4.1 SGP Results for Gender	152
15.4.2 SGP Results for Ethnicity	152
15.4.3 SGP Results for Special Instructional Needs	153
15.4.4 SGP Results for Students Taking Spanish Forms	154
References	155
Appendices	158
Appendix 6: Summary of Differential Item Function (DIF) Results	158
Appendix 7.1: Pre-Equated IRT Results for Spring 2022 English Language Arts/Literacy (ELA/L)	175
Appendix 7.2: Pre-Equated IRT Results for Spring 2022 Mathematics	176
Appendix 11: Students by Grade/Subject and Mode, for Each State	179
Appendix 12.1: Form Composition	205
Appendix 12.2: Threshold Scores and Scaling Constants	211
Appendix 12.3: IRT Test Characteristic Curves, CSEM Curves, and Information Curves	215
Appendix 12.4: Scale Score Cumulative Frequencies	232
Appendix 12.5: Subgroup Scale Score Performance	249
Appendix 13.1: Reliability by Content and Grade/Subject	274
Appendix 13.2: Reliability of Classification by Content and Grade/Subject	291
Appendix 14: Quality Testing Standards	297
Appendix 15: Growth	323

List of Tables

Table 1.1 Glossary of Abbreviations and Acronyms	9
Table 4.1 Training Materials Used During Scoring	32
Table 4.2 Mathematics Qualification Requirements	34
Table 4.3 Scoring Hierarchy Rules	34
Table 4.4 Scoring Validity Agreement Requirements	35
Table 4.5 Inter-Rater Agreement Expectations and Results	36
Table 4.6 Comparison Groups	41
Table 4.7 Prose Constructed Response Average Agreement Indices by Test	42
Table 5.1 Summary of Pre-Administration p-Values for ELA/L Operational Items by Grade and Mode	46
Table 5.2 Summary of p-Values for Mathematics Operational Items by Grade and Mode	47
Table 5.3 Summary of Pre-Administration Item-Total Polyserial Correlations for ELA/L Operational Items by Grade and Mode	47
Table 5.4 Summary of Item-Total Correlations for Mathematics Operational Items by Grade and Mode	48
Table 6.1 DIF Categories for Dichotomous Selected-Response and Constructed-Response Items	51
Table 6.2 DIF Categories for Polytomous Constructed-Response Items	51
Table 6.3 Traditional DIF Comparison Groups	51
Table 6.4 Pre-Administration Differential Item Functioning for ELA/L Grade 3	52

Table 6.5 Pre-Administration Differential Item Functioning for Mathematics Grade 3	53
Table 7.1 Pre-Equated IRT Parameter Estimates Summary for All Items for ELA/L by Grade	55
Table 7.2 Pre-Equated IRT Parameter Distribution by Year for All Items for ELA/L by Grade	55
Table 7.3 Pre-Equated IRT Parameter Estimates Summary for All Items for Mathematics by Grade/Course	56
Table 7.4 Pre-Equated IRT Parameter Distribution by Year for All Items for Mathematics by Grade/Course	56
Table 8.1 Performance Level Setting Committee Meetings and Dates	62
Table 10.1 Number of Core Operational Forms per Grade/Subject and Mode for ELA/L and Mathematics	69
Table 11.1 ELA/L Students by Grade and Mode: All States Combined	71
Table 11.2 Mathematics Students by Grade/Course and Mode: All States Combined	71
Table 11.3 Spanish-Language Mathematics Students by Grade/Course and Mode: All States Combined	71
Table 12.1 Form Composition for ELA/L Grade 3	74
Table 12.2 Contribution of Prose Constructed-Response Items to ELA/L	74
Table 12.3 Mathematics Form Composition for Grade 3	75
Table 12.4 Calculating Scaling Constants for Reading and Writing Claim Scores	77
Table 12.5 Subgroup Performance for ELA/L Scale Scores: Grade 3	86
Table 12.6 Subgroup Performance for ELA/L Scale Scores: Grade 10	87
Table 12.7 Subgroup Performance for Mathematics Scale Scores: Grade 3	91
Table 12.8 Subgroup Performance for Mathematics Scale Scores: Algebra I	91
Table 13.1 Summary of ELA/L Test Reliability Estimates for Total Group	97
Table 13.2 Summary of Mathematics Test Reliability Estimates for Total Group	97
Table 13.3 Summary of ELA/L Test Pre-Equated Scale Score Reliability Estimates for Total Group	98
Table 13.4 Summary of Mathematics Test Scale Score Reliability Estimates for Total Group	98
Table 13.5 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 3	101
Table 13.6 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 3	102
Table 13.7 Descriptions of ELA/L Claims and Subclaims	103
Table 13.8 Average ELA/L Reliability Estimates for Subscores	104
Table 13.9 Average Mathematics Reliability Estimates for Subscores	106
Table 13.10 Reliability of Classification: Summary for ELA/L	108
Table 13.11 Reliability of Classification: Grade 3 ELA/L	108
Table 13.12 Reliability of Classification: Summary for Mathematics	109
Table 13.13 Inter-Rater Agreement Expectations and Results	110
Table 14.1 Average Intercorrelations and Reliability between Grade 3 ELA/L Subclaims	115
Table 14.2 Average Intercorrelations and Reliability between Grade 4 ELA/L Subclaims	115
Table 14.3 Average Intercorrelations and Reliability between Grade 5 ELA/L Subclaims	116
Table 14.4 Average Intercorrelations and Reliability between Grade 6 ELA/L Subclaims	116
Table 14.5 Average Intercorrelations and Reliability between Grade 7 ELA/L Subclaims	117
Table 14.6 Average Intercorrelations and Reliability between Grade 8 ELA/L Subclaims	117
Table 14.7 Average Intercorrelations and Reliability between Grade 9 ELA/L Subclaims	118
Table 14.8 Average Intercorrelations and Reliability between Grade 10 ELA/L Subclaims	118
Table 14.9 Average Intercorrelations and Reliability between Grade 3 Mathematics Subclaims	119
Table 14.10 Average Intercorrelations and Reliability between Grade 4 Mathematics Subclaims	119
Table 14.11 Average Intercorrelations and Reliability between Grade 5 Mathematics Subclaims	119
Table 14.12 Average Intercorrelations and Reliability between Grade 6 Mathematics Subclaims	120
Table 14.13 Average Intercorrelations and Reliability between Grade 7 Mathematics Subclaims	120

Table 14.14 Average Intercorrelations and Reliability between Grade 8 Mathematics Subclaims	120
Table 14.15 Average Intercorrelations and Reliability between Algebra I Subclaims.....	121
Table 14.16 Average Intercorrelations and Reliability between Geometry Subclaims	121
Table 14.17 Average Intercorrelations and Reliability between Algebra II Subclaims	121
Table 14.18 Conditions Used in LID Investigation and Results.....	125
Table 14.19 Summary of Q3 Values for ELA/L Grade 4 and Integrated Mathematics II (Spring 2015)	125
Table 14.20 Correlations between ELA/L and Mathematics for Grade 3	127
Table 14.21 Correlations between ELA/L and Mathematics for Grade 4	127
Table 14.22 Correlations between ELA/L and Mathematics for Grade 5	127
Table 14.23 Correlations between ELA/L and Mathematics for Grade 6	127
Table 14.24 Correlations between ELA/L and Mathematics for Grade 7	127
Table 14.25 Correlations between ELA/L and Mathematics for Grade 8	128
Table 14.26 Correlations between ELA/L and Mathematics for High School.....	128
Table 14.27 Correlations between ELA/L Reading and Mathematics for High School.....	128
Table 14.28 Correlations between ELA/L Writing and Mathematics for High School.....	128
Table 14.29 Prior Grades Used in ELA/L Matching	135
Table 14.30 Prior Grades/Courses Used in Mathematics Matching.....	135
Table 14.31 ELA/L Matching Sample Size Results	136
Table 14.32 Mathematics Matching Sample Size Results.....	137
Table 15.1 ELA/L Grade-Level Progressions for One- and Two-Year Prior Test Scores	148
Table 15.2 Mathematics Grade-Level Progressions for One- and Two-Year Prior Test Scores.....	148
Table 15.3 Algebra I Grade/Content Area Progressions for One- and Two-Year Prior Test Scores.....	148
Table 15.4 Geometry Grade/Content Area Progressions for One- and Two-Year Prior Test Scores	148
Table 15.5 Algebra II Grade/Content Area Progressions for One- and Two-Year Prior Test Scores	149
Table 15.6 Integrated Mathematics I Grade/Content Area Progressions for One- and Two-Year Prior Test Scores.....	149
Table 15.7 Integrated Mathematics II Grade/Content Area Progressions for One- and Two-Year Prior Test Scores.....	149
Table 15.8 Integrated Mathematics III Grade/Content Area Progressions for One- and Two-Year Prior Test Scores.....	149
Table 15.9 State-Specific SGP Progressions	149
Table 15.10 Summary of ELA/L SGP Estimates for Total Group	151
Table 15.11 Summary of Mathematics SGP Estimates for Total Group	151
Table 15.12 Summary of SGP Estimates for Subgroups: Grade 4 ELA/L	153
Table 15.13 Summary of SGP Estimates for Subgroups: Grade 4 Mathematics.....	154
Table A.6.1 Pre-Administration Differential Item Functioning for ELA/L Grade 3	158
Table A.6.2 Pre-Administration Differential Item Functioning for ELA/L Grade 4	159
Table A.6.3 Pre-Administration Differential Item Functioning for ELA/L Grade 5	160
Table A.6.4 Pre-Administration Differential Item Functioning for ELA/L Grade 6	161
Table A.6.5 Pre-Administration Differential Item Functioning for ELA/L Grade 7	162
Table A.6.6 Pre-Administration Differential Item Functioning for ELA/L Grade 8	163
Table A.6.7 Pre-Administration Differential Item Functioning for ELA/L Grade 9	164
Table A.6.8 Pre-Administration Differential Item Functioning for ELA/L Grade 10	165
Table A.6.9 Pre-Administration Differential Item Functioning for Mathematics Grade 3.....	166

Table A.6.10 Pre-Administration Differential Item Functioning for Mathematics Grade 4.....	167
Table A.6.11 Pre-Administration Differential Item Functioning for Mathematics Grade 5.....	168
Table A.6.12 Pre-Administration Differential Item Functioning for Mathematics Grade 6.....	169
Table A.6.13 Pre-Administration Differential Item Functioning for Mathematics Grade 7.....	170
Table A.6.14 Pre-Administration Differential Item Functioning for Mathematics Grade 8.....	171
Table A.6.15 Pre-Administration Differential Item Functioning for Algebra I	172
Table A.6.16 Pre-Administration Differential Item Functioning for Geometry.....	173
Table A.6.17 Pre-Administration Differential Item Functioning for Algebra II.....	174
Table A.7.1 Pre-Equated IRT Summary Parameter Estimates for All Items for ELA/L by Grade	175
Table A.7.2 Pre-Equated IRT Summary Parameter Estimates for All Items for Mathematics by Grade/Subject	176
Table A.11.1 All ELA/L Test Takers, by State and Grade.....	179
Table A.11.2 All Mathematics Test Takers, by State and Grade	181
Table A.11.3 All Spanish-Language Mathematics Test Takers, by State and Grade	183
Table A.11.4 All States Combined: ELA/L Test Takers, by Grade, Mode, and Gender	185
Table A.11.5 All States Combined: Mathematics Test Takers, by Grade, Mode, and Gender.....	186
Table A.11.6 All States Combined: Spanish-Language Mathematics Test Takers, by Grade, Mode, and Gender	187
Table A.11.7 Demographic Information: Grade 3 ELA/L, Overall and by State	188
Table A.11.8 Demographic Information: Grade 4 ELA/L, Overall and by State	189
Table A.11.9 Demographic Information: Grade 5 ELA/L, Overall and by State	190
Table A.11.10 Demographic Information: Grade 6 ELA/L, Overall and by State.....	191
Table A.11.11 Demographic Information: Grade 7 ELA/L, Overall and by State.....	192
Table A.11.12 Demographic Information: Grade 8 ELA/L, Overall and by State.....	193
Table A.11.13 Demographic Information: Grade 9 ELA/L, Overall and by State.....	194
Table A.11.14 Demographic Information: Grade 10 ELA/L, Overall and by State	195
Table A.11.15 Demographic Information: Grade 3 Mathematics, Overall and by State	196
Table A.11.16 Demographic Information: Grade 4 Mathematics, Overall and by State	197
Table A.11.17 Demographic Information: Grade 5 Mathematics, Overall and by State	198
Table A.11.18 Demographic Information: Grade 6 Mathematics, Overall and by State	199
Table A.11.19 Demographic Information: Grade 7 Mathematics, Overall and by State	200
Table A.11.20 Demographic Information: Grade 8 Mathematics, Overall and by State	201
Table A.11.21 Demographic Information: Algebra I, Overall and by State.....	202
Table A.11.22 Demographic Information: Geometry, Overall and by State	203
Table A.11.23 Demographic Information: Algebra II, Overall and by State	204
Table A.12.1 Form Composition for ELA/L Grade 3	205
Table A.12.2 Form Composition for ELA/L Grade 4.....	205
Table A.12.3 Form Composition for ELA/L Grade 5	206
Table A.12.4 Form Composition for ELA/L Grade 6.....	206
Table A.12.5 Form Composition for ELA/L Grade 7	206
Table A.12.6 Form Composition for ELA/L Grade 8.....	207
Table A.12.7 Form Composition for ELA/L Grade 9.....	207
Table A.12.8 Form Composition for ELA/L Grade 10	207
Table A.12.9 Form Composition for Mathematics Grade 3.....	208
Table A.12.10 Form Composition for Mathematics Grade 4.....	208

Table A.12.11 Form Composition for Mathematics Grade 5.....	208
Table A.12.12 Form Composition for Mathematics Grade 6.....	208
Table A.12.13 Form Composition for Mathematics Grade 7.....	209
Table A.12.14 Form Composition for Mathematics Grade 8.....	209
Table A.12.15 Form Composition for Algebra I.....	209
Table A.12.16 Form Composition for Geometry.....	209
Table A.12.17 Form Composition for Algebra II.....	210
Table A.12.18 Threshold Scores and Scaling Constants for ELA/L Grades 3 to 8.....	211
Table A.12.19 Threshold Scores and Scaling Constants for Mathematics Grades 3 to 8.....	212
Table A.12.20 Threshold Scores and Scaling Constants for High School ELA/L.....	213
Table A.12.21 Threshold Scores and Scaling Constants for High School Mathematics.....	213
Table A.12.22 Scaling Constants for Reading and Writing Grades 3 to 11.....	214
Table A.12.23 Scale Score Cumulative Frequencies: ELA/L Grade 3.....	232
Table A.12.24 Scale Score Cumulative Frequencies: ELA/L Grade 4.....	233
Table A.12.25 Scale Score Cumulative Frequencies: ELA/L Grade 5.....	234
Table A.12.26 Scale Score Cumulative Frequencies: ELA/L Grade 6.....	235
Table A.12.27 Scale Score Cumulative Frequencies: ELA/L Grade 7.....	236
Table A.12.28 Scale Score Cumulative Frequencies: ELA/L Grade 8.....	237
Table A.12.29 Scale Score Cumulative Frequencies: ELA/L Grade 9.....	238
Table A.12.30 Scale Score Cumulative Frequencies: ELA/L Grade 10.....	239
Table A.12.31 Scale Score Cumulative Frequencies: Mathematics Grade 3.....	240
Table A.12.32 Scale Score Cumulative Frequencies: Mathematics Grade 4.....	241
Table A.12.33 Scale Score Cumulative Frequencies: Mathematics Grade 5.....	242
Table A.12.34 Scale Score Cumulative Frequencies: Mathematics Grade 6.....	243
Table A.12.35 Scale Score Cumulative Frequencies: Mathematics Grade 7.....	244
Table A.12.36 Scale Score Cumulative Frequencies: Mathematics Grade 8.....	245
Table A.12.37 Scale Score Cumulative Frequencies: Algebra I.....	246
Table A.12.38 Scale Score Cumulative Frequencies: Geometry.....	247
Table A.12.39 Scale Score Cumulative Frequencies: Algebra II.....	248
Table A.12.40 Subgroup Performance for ELA/L Scale Scores: Grade 3.....	249
Table A.12.41 Subgroup Performance for ELA/L Scale Scores: Grade 4.....	251
Table A.12.42 Subgroup Performance for ELA/L Scale Scores: Grade 5.....	253
Table A.12.43 Subgroup Performance for ELA/L Scale Scores: Grade 6.....	255
Table A.12.44 Subgroup Performance for ELA/L Scale Scores: Grade 7.....	257
Table A.12.45 Subgroup Performance for ELA/L Scale Scores: Grade 8.....	259
Table A.12.46 Subgroup Performance for ELA/L Scale Scores: Grade 9.....	261
Table A.12.47 Subgroup Performance for ELA/L Scale Scores: Grade 10.....	263
Table A.12.48 Subgroup Performance for Mathematics Scale Scores: Grade 3.....	265
Table A.12.49 Subgroup Performance for Mathematics Scale Scores: Grade 4.....	266
Table A.12.50 Subgroup Performance for Mathematics Scale Scores: Grade 5.....	267
Table A.12.51 Subgroup Performance for Mathematics Scale Scores: Grade 6.....	268
Table A.12.52 Subgroup Performance for Mathematics Scale Scores: Grade 7.....	269
Table A.12.53 Subgroup Performance for Mathematics Scale Scores: Grade 8.....	270
Table A.12.54 Subgroup Performance for Mathematics Scale Scores: Algebra I.....	271

Table A.12.55 Subgroup Performance for Mathematics Scale Scores: Geometry	272
Table A.12.56 Subgroup Performance for Mathematics Scale Scores: Algebra II	273
Table A.13.1 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 3.....	274
Table A.13.2 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 4.....	275
Table A.13.3 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 5.....	276
Table A.13.4 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 6.....	277
Table A.13.5 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 7.....	277
Table A.13.6 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 8.....	279
Table A.13.7 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 9.....	280
Table A.13.8 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 10.....	281
Table A.13.9 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 3	282
Table A.13.10 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 4.....	283
Table A.13.11 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 5.....	284
Table A.13.12 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 6.....	285
Table A.13.13 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 7.....	286
Table A.13.14 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 8.....	287
Table A.13.15 Summary of Test Reliability Estimates for Subgroups: Algebra I.....	288
Table A.13.16 Summary of Test Reliability Estimates for Subgroups: Geometry	289
Table A.13.17 Summary of Test Reliability Estimates for Subgroups: Algebra II	290
Table A.13.18 Reliability of Classification: Grade 3 ELA/L.....	291
Table A.13.19 Reliability of Classification: Grade 4 ELA/L.....	291
Table A.13.20 Reliability of Classification: Grade 5 ELA/L.....	291
Table A.13.21 Reliability of Classification: Grade 6 ELA/L.....	292
Table A.13.22 Reliability of Classification: Grade 7 ELA/L.....	292
Table A.13.23 Reliability of Classification: Grade 8 ELA/L.....	292
Table A.13.24 Reliability of Classification: Grade 9 ELA/L.....	293
Table A.13.25 Reliability of Classification: Grade 10 ELA/L.....	293
Table A.13.26 Reliability of Classification: Grade 3 Mathematics	293
Table A.13.27 Reliability of Classification: Grade 4 Mathematics	294
Table A.13.28 Reliability of Classification: Grade 5 Mathematics	294
Table A.13.29 Reliability of Classification: Grade 6 Mathematics	294
Table A.13.30 Reliability of Classification: Grade 7 Mathematics	295
Table A.13.31 Reliability of Classification: Grade 8 Mathematics	295
Table A.13.32 Reliability of Classification: Algebra I.....	295
Table A.13.33 Reliability of Classification: Geometry	296
Table A.13.34 Reliability of Classification: Algebra II	296
Table A.14.1 ELA/L Grade 6 Form 1 Matching Results	297
Table A.14.2 Mathematics Grade 6 Form 1 Matching Results.....	298
Table A.14.3 ELA/L Grade 10 Form 1 Matching Results	299
Table A.14.4 Distributions of P-Value Differences* for ELA/L	303
Table A.14.5 Distributions of P-Value Differences* for Mathematics	303
Table A.14.6 Distributions of Polyserial Differences* for ELA/L.....	307
Table A.14.7 Distributions of Polyserial Differences* for Mathematics	307
Table A.14.8 DIF Category Cross Tabulations for ELA/L	307

Table A.14.9 DIF Category Crosstabulations for Mathematics Grades 3–8 and Algebra I.....	307
Table A.14.10 DIF Category Crosstabulations for Algebra II and Geometry	308
Table A.14.11 ELA/L Reliability	308
Table A.14.12 ELA/L Raw Score Standard Error of Measurement.....	308
Table A.14.13 ELA/L Scale Score Standard Error of Measurement.....	308
Table A.14.14 Mathematics Reliability	309
Table A.14.15 Mathematics Raw Score Standard Error of Measurement.....	309
Table A.14.16 Mathematics Scale Score Standard Error of Measurement	309
Table A.14.17 ELA/L Scale Score Descriptive Statistics	310
Table A.14.18 Mathematics Scale Score Descriptive Statistics	310
Table A.14.19 ELA/L Writing Claim Score Descriptive Statistics	310
Table A.14.20 Reading Claim Score Descriptive Statistics.....	311
Table A.14.21 ELA/L Subclaim Distributions.....	311
Table A.14.22 Mathematics Subclaim Distributions.....	311
Table A.14.23 ELA/L Subclaim Distribution Comparison: Effect Size.....	312
Table A.14.24 Mathematics Subclaim Distribution Comparison: Effect Size	312
Table A.14.25 ELA/L Longitudinal Scale Score Comparison: Original to Current	312
Table A.14.26 ELA/L Longitudinal Scale Score Comparison: Original to Original.....	313
Table A.14.27 Mathematics Longitudinal Scale Score Comparison: Original to Current.....	313
Table A.14.28 Mathematics Longitudinal Scale Score Comparison: Original to Original	314
Table A.14.29 ELA/L Longitudinal Regression.....	314
Table A.14.30 Mathematics Longitudinal Regression.....	314
Table A.14.31 ELA/L Grade 3 Performance Level Comparison.....	315
Table A.14.32 Mathematics Grade 3 Performance Level Comparison	315
Table A.14.33 Performance Level Comparison Summary: Effect Sizes	315
Table A.14.34 College and Career Readiness Comparison Summary: Effect Sizes	316
Table A.14.35 ELA/L Classification Accuracy.....	316
Table A.14.36 ELA/L Classification Consistency	316
Table A.14.37 Mathematics Classification Accuracy	317
Table A.14.38 Mathematics Classification Consistency	317
Table A.14.39 ELA/L Grade 6 Performance Level Comparison.....	317
Table A.14.40 Mathematics Grade 6 Performance Level Comparison	318
Table A.14.41 Performance Level Comparison Summary: Effect Sizes	318
Table A.14.42 ELA/L Reading Claim Reliability	318
Table A.14.43 ELA/L Writing Claim Reliability	319
Table A.14.44 ELA/L Reading Information (RI) Subclaim Reliability	319
Table A.14.45 ELA/L Reading Literature (RL) Subclaim Reliability	319
Table A.14.46 ELA/L Reading Vocabulary (RV) Subclaim Reliability	320
Table A.14.47 ELA/L Writing Knowledge and Conventions (WKL) Subclaim Reliability	320
Table A.14.48 ELA/L Written Expression (WE) Subclaim Reliability	320
Table A.14.49 Mathematics Subclaim A Reliability.....	321
Table A.14.50 Mathematics Subclaim B Reliability.....	321
Table A.14.51 Mathematics Subclaim C Reliability.....	321
Table A.14.52 Mathematics Subclaim D Reliability.....	322

Table A.15.1 Summary of Student Growth Percentile Estimates for Subgroups: Grade 4 ELA/L.....	323
Table A.15.2 Summary of Student Growth Percentile Estimates for Subgroups: Grade 5 ELA/L.....	324
Table A.15.3 Summary of Student Growth Percentile Estimates for Subgroups: Grade 6 ELA/L.....	325
Table A.15.4 Summary of Student Growth Percentile Estimates for Subgroups: Grade 7 ELA/L.....	326
Table A.15.5 Summary of Student Growth Percentile Estimates for Subgroups: Grade 8 ELA/L.....	327
Table A.15.6 Summary of Student Growth Percentile Estimates for Subgroups: Grade 4 Mathematics	328
Table A.15.7 Summary of Student Growth Percentile Estimates for Subgroups: Grade 5 Mathematics	329
Table A.15.8 Summary of Student Growth Percentile Estimates for Subgroups: Grade 6 Mathematics	330
Table A.15.9 Summary of Student Growth Percentile Estimates for Subgroups: Grade 7 Mathematics	331
Table A.15.10 Summary of Student Growth Percentile Estimates for Subgroups: Grade 8 Mathematics.....	332
Table A.15.11 Summary of Student Growth Percentile Estimates for Subgroups: Algebra I.....	333
Table A.15.12 Summary of Student Growth Percentile Estimates for Subgroups: Geometry	334
Table A.15.13 Summary of Student Growth Percentile Estimates for Subgroups: Algebra II	335

List of Figures

Figure 12.1 Test Characteristic Curves, Conditional Standard Error of Measurement Curves, and Information Curves for ELA/L Grade 3.....	80
Figure 12.2 Distributions of ELA/L Scale Scores: Grades 3–10.....	82
Figure 12.3 Distributions of Reading Scale Scores: Grades 3–10	84
Figure 12.4 Distributions of Writing Scale Scores: Grades 3–10	85
Figure 12.5 Distributions of Mathematics Scale Scores: Grades 3–8	90
Figure 12.6 Distributions of Mathematics Scale Scores: High School	90
Figure 14.1 Comparison of Internal Consistency by Item and Cluster (Testlet)	124
Figure 14.2 Distribution of Q3 Values for Grade 4 ELA/L (Spring 2015)	125
Figure 14.3 Distribution of Q3 Values for Integrated Mathematics II (Spring 2015).....	126
Figure 14.4 ELA/L Grades 3–6 P-Values.....	138
Figure 14.5 Mathematics Grades 3–6 P-Values	139
Figure A.12.1 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 3.....	215
Figure A.12.2 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 4.....	216
Figure A.12.3 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 5.....	217
Figure A.12.4 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 6.....	218
Figure A.12.5 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 7.....	219
Figure A.12.6 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 8.....	220
Figure A.12.7 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 9.....	221
Figure A.12.8 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 10.....	222
Figure A.12.9 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 3.....	223
Figure A.12.10 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 4.....	224
Figure A.12.11 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 5.....	225
Figure A.12.12 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 6.....	226
Figure A.12.13 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 7.....	227
Figure A.12.14 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 8.....	228
Figure A.12.15 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Algebra I.....	229
Figure A.12.16 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Geometry	230
Figure A.12.17 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Algebra II	231
Figure A.14.1 ELA/L Grades 3–6 P-Values.....	300
Figure A.14.2 ELA/L Grades 7–8 P-Values.....	300
Figure A.14.3 ELA/L Grade 10 P-Values	301
Figure A.14.4 Mathematics Grades 3–6 P-Values	301
Figure A.14.5 Mathematics Grades 7–8 and Algebra I P-Values	302
Figure A.14.6 Algebra II and Geometry P-Values.....	302
Figure A.14.7 Polyserial Correlations ELA/L Grades 3–6	304
Figure A.14.8 Polyserial Correlations ELA/L Grades 7–8	304

Figure A.14.9 Polyserial Correlations ELA/L Grade 10.....	305
Figure A.14.10 Polyserial Correlations Mathematics Grades 3–6.....	305
Figure A.14.11 Polyserial Correlations Mathematics Grades 7–8 and Algebra I	306
Figure A.14.12 Polyserial Correlations Algebra II and Geometry	306

Executive Summary

The purpose of this report is to describe the technical qualities of the 2021–2022 operational administration of the English language arts/literacy (ELA/L) and mathematics assessments in grades 3 through 8 and high school. Due to the outbreak of the global COVID-19 pandemic, the spring 2020 administration was suspended in March 2020 and ultimately canceled for all participating states. At that time, testing only occurred for a small number of students in grades 3 through 8 in Illinois, although other states had planned to administer tests in grades 3 through 8 as well as in high school. For spring 2021, two participating states canceled their administration, while Illinois, the Department of Defense Education Activity, and the Bureau of Indian Education administered the assessments. Illinois provided the option to test either in spring 2021 or in fall 2021. The forms administered were the same in spring and fall. For spring 2022, all participating states and agencies (Illinois, the Department of Defense Education Activity, the District of Columbia, and New Jersey) administered the assessments.

Committees of educators, state education agency staff, and national experts led the work in the development of the summative assessments that are aligned to the Common Core State Standards (CCSS). These summative assessments are intended to measure more complex skills like critical thinking, persuasive writing, and problem-solving. New Meridian assumes the responsibility for management of the summative assessments, as well as item development and forms construction. New Meridian, in coordination with multiple states and vendors, developed an alternate form of the summative assessment to meet the provision for shorter testing times requested by several states. Through extensive research and guidance from the Technical Advisory Committee, the alternate blueprint was made available in spring 2019.

The ELA/L assessments focus on reading and comprehending a range of sufficiently complex texts independently and writing effectively when analyzing text. The ELA/L assessments contain literary and informational texts; each passage set has four to eight brief comprehension and vocabulary questions. ELA/L constructed-response items include three types of tasks: literary analysis, narrative writing, and research simulation. For each task, students are instructed to read one or more texts, answer several brief questions, and then write an essay based on the material they read.

The mathematics assessments contain tasks that measure a combination of conceptual understanding, applications, skills, and procedures. Mathematics constructed-response items consist of tasks designed to assess a student's ability to use mathematics to solve real-life problems. Some of the tasks require students to describe how they solved a problem, while other tasks measure conceptual understanding and the ability to apply concepts by means of selected-response or technology-enhanced items. In addition, students are required to demonstrate their skills and knowledge by answering innovative selected-response and short-answer questions that measure concepts and skills.

In both content areas, students also demonstrate their acquired skills and knowledge by answering selected-response items and fill-in-the-blank questions. Each assessment consists of multiple units, and additionally, one of the mathematics units is split into two sections: a noncalculator section and a calculator section.

The summative assessments are designed to achieve several purposes. First, the tests are intended to provide evidence to determine whether students are on track for college- and career-readiness. Second, the tests are structured to access the full range of CCSS and measure the total breadth of student performance. Finally, the tests are designed to provide data to help inform classroom instruction, student interventions, and professional development.

This technical report includes the following topics:

- background and purpose of the assessments
- test development of items and forms
- test administration, security, and scoring
- student characteristics
- classical item analyses and differential item functioning
- reliability and validity of scores
- item response theory calibration and scaling
- performance level setting
- development of the score reporting scales and student performance
- student growth measures
- quality control procedures

The information provided in this technical report 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 *Standards for Educational and Psychological Testing* (American Educational Research Association [AERA], American Psychological Association [APA], and National Council on Measurement in Education [NCME], 2014).

<This page intentionally left blank>

Section 1: Introduction

1.1 Background

States associated with the Partnership for Assessment of Readiness for College and Careers (PARCC) came together in early 2010 with a shared vision of ensuring that all students—regardless of income, family background, or geography—have equal access to a world-class education that will prepare them for success after high school in college and/or careers. The goal was to develop new assessments that tie into more rigorous academic expectations and help prepare students for success in college and the workforce, as well as to provide information back to teachers and parents about where students are on their path to success. Calling on the expertise of thousands of teachers, higher education faculty, and other educators in multiple states, the resulting assessment system is a high-quality set of summative assessments, diagnostic assessments, formative tasks, and other support materials for teachers including professional development and communications tools.

The partnership develops and administers next-generation assessments that, compared to traditional K–12 assessments, more accurately measure student progress toward college- and career-readiness. The assessments are aligned to the Common Core State Standards (CCSS) and include both English language arts/literacy (ELA/L) assessments (grades 3 through 11) and mathematics assessments (grades 3 through 8 and high school). Compared to traditional standardized tests, these assessments are intended to measure more complex skills like critical thinking, persuasive writing, and problem-solving.

In 2013, the PARCC Governing Board launched Parcc Inc., a nonprofit organization designed to support the successful delivery of the tests in 2014–2017 and the long-term success of the multistate partnership. States continued to govern decisions about the assessment system; the nonprofit organization was their “agent” for overseeing the many vendors involved in the assessment system, coordinating the multiple work groups and committees (including Governing Board meetings), managing the intellectual property, overseeing the research agenda and the Technical Advisory Committee, and developing and launching the multiple nonsummative tools.

Summative assessments for the first operational administration were constructed in 2014. Ten states and the District of Columbia participated in the first administration of the summative assessments during the 2014–2015 school year. Six states, the Bureau of Indian Education, and District of Columbia participated in the second administration in school year 2015–2016. Five states, the Bureau of Indian Education, the Department of Defense Education Activity, and the District of Columbia participated in the third administration in school year 2016–2017. Four states, the Bureau of Indian Education, the Department of Defense Education Activity, and the District of Columbia participated in the fourth administration in school year 2017–2018.

Following the Parcc Inc. contract ending in June 2017, participating states and agencies released the intellectual property of the contract to the Council of Chief State School Officers, and also contracted with New Meridian to manage the intellectual property and provide item development, forms construction, and governance. Starting in August 2017, New Meridian oversaw item development, data review for field-test items, and test construction activities.

New Meridian, in coordination with multiple states and vendors, developed an alternate form of the summative assessment to meet the provision for shorter testing times requested by several states. Through extensive research and guidance from the Technical Advisory Committee, the alternate blueprint was available in spring 2019 in addition to the original blueprint. New Meridian’s state-centric solution to educational assessment

allowed states the flexibility of selecting the assessment solution that best fit their specific needs. For the academic year 2018–2019, participating states and agencies included the Bureau of Indian Education, the District of Columbia, Illinois, New Jersey, and New Mexico. For the academic years 2019–2020 and 2020–2021, participating states and agencies included the Bureau of Indian Education, the District of Columbia, the Department of Defense Education Activity, Illinois, and New Jersey. Most testing in spring 2020 was canceled due to COVID-19, with the exception of a small number of students in Illinois who tested prior to the closure of schools, and some states canceled administration in spring 2021. For the academic year 2021–2022, participating states and agencies included the District of Columbia, the Department of Defense Education Activity, Illinois, and New Jersey. All participating states and agencies administered assessments in spring 2022.

The purpose of this technical report is to describe the operational administration of the summative assessments in the 2021–2022 academic year, including test form construction, test administration, item scoring, student characteristics, classical item analysis results, reliability results, evidence of validity, item response theory (IRT) calibrations and scaling, performance level setting procedure, growth measures, and quality control procedures.

1.2 Purpose of the Operational Tests

The summative assessments are designed to achieve several purposes. First, the assessments are intended to provide evidence to determine whether students are on track for college- and career-readiness. Second, the assessments are structured to access the full range of CCSS and measure the total breadth of student performance. Finally, the assessments are designed to provide data to help inform classroom instruction, student interventions, and professional development.

1.3 Composition of Operational Tests

Each operational test form is constructed to reflect the test blueprint in terms of content, standards measured, and item types. Sets of common items, included to provide data to support horizontal linking across test forms within a grade and content area, are proportionally representative of the operational test blueprint. The summative assessment is a mixed-format test. The current summative assessments are administered in either computer-based test (CBT) or paper-based test (PBT) format.

The ELA/L assessments focus on reading and comprehending a range of sufficiently complex texts independently and writing effectively when analyzing text. The ELA/L assessments contain literary and informational texts; each passage set has four to eight brief comprehension and vocabulary questions. ELA/L constructed-response items include three types of tasks: literary analysis, narrative writing, and research simulation. For each task, students are instructed to read one or more texts, answer several brief questions, and then write an essay based on the material they read.

The mathematics assessments contain tasks that measure a combination of conceptual understanding, applications, skills, and procedures. Mathematics constructed-response items consist of tasks designed to assess a student's ability to use mathematics to solve real-life problems. Some of the tasks require students to describe how they solved a problem, while other tasks measure conceptual understanding and the ability to apply concepts by means of selected-response or technology-enhanced items. In addition, students are required to demonstrate their skills and knowledge by answering innovative selected-response and short-answer questions that measure concepts and skills.

In both content areas, students also demonstrate their acquired skills and knowledge by answering selected-response items and fill-in-the-blank questions. Each assessment consists of multiple units; additionally, one of the mathematics units is split into two sections: a noncalculator section and a calculator section.

1.4 Intended Population

The tests are intended for students taking ELA/L in grades 3 through 11, and/or mathematics in grades 3 through 8, as well as students taking high school mathematics (i.e., Algebra I, Geometry, Algebra II, and Integrated Mathematics I–II). For these students, the tests measure whether students are meeting state academic standards and mastering the knowledge and skills needed to progress in their K–12 education and beyond. In academic year 2021–2022, assessments of ELA/L 11 and Integrated Mathematics I–III were not administered.

1.5 Groups and Organizations Involved with the Summative Assessments

New Meridian is a nonprofit organization that assumed the responsibility for the management of the assessments in 2017, as well as the responsibility for item development and forms construction of the assessments.

Committees of educators, state education agency staff, and national experts lead the work of the assessments. These committees include:

- the Governing Board, which makes major policy and operational decisions;
- the Technical Advisory Committee, which helps ensure all assessments will provide reliable results to inform valid instructional and accountability decisions;
- the State Lead Council, which coordinates all aspects of development of the summative assessment system and serves as the conduit to the Technical Advisory Committee and the Governing Board; and
- ELA/L, mathematics, and accessibility and accommodation features operational working groups.

Pearson serves as the primary contractor for the operational administration and is responsible for producing all testing materials, packaging and distribution, receiving and scanning of materials, and scoring, as well as program management and customer service. In addition, test and item development activities are conducted by Pearson under the guidance and oversight of New Meridian.

Pearson Psychometrics is responsible for all psychometric analyses of the operational test data. This includes classical item analyses, differential item functioning analyses, item calibrations based on IRT, scaling, and development of all conversion tables.

1.6 Overview of the Technical Report

This report begins by providing explanations of the test form construction process, test administration, and scoring of the test items. Subsequent sections of the report present descriptions of student characteristics, results of classical item analyses, IRT calibrations and scaling, performance level setting procedure, quality control procedures, results of students' scale score analyses, results of reliability analyses, evidence of validity, and measures of student growth.

The technical report contains the following sections:

Section 2 – Test Development

This section describes the test design and the procedures followed during the development of operational test forms.

Section 3 – Test Administration

This section presents the operational administration schedule, information regarding test security and confidentiality, accessibility features and accommodations, and testing irregularities and security breaches.

Section 4 – Item Scoring

The key-based and rule-based processes for machine-scored items, as well as the training and monitoring processes for human-scored items, are provided in this section.

Section 5 – Classical Item Analysis

The classical item-level statistics calculated for the operational test data, the flagging criteria used to identify items that performed differently than expected, and the results of these analyses are presented in this section.

Section 6 – Differential Item Functioning

In this section, the methods for conducting differential item functioning analyses as well as corresponding flagging criteria are described. This is followed by definitions of the comparison groups and subsequent results for the comparison groups.

Section 7 – Item Response Theory Model and Parameters

This section presents the information related to the IRT models used and the descriptive statistics of the item parameters. Note that all tests delivered in 2020 employed a pre-equated model, in which previously estimated item parameters are used to generate scoring tables.

Section 8 – Performance Level Setting

Performance levels and policy definitions, as well as the processes followed to establish performance level thresholds, are described in this section.

Section 9 – Quality Control Procedures

All aspects of quality control are presented in this section. These activities range from quality assurance of item banking, test form construction, and all testing materials to quality control of scanning, image editing, and scoring. This is followed by a detailed description of the steps taken to ensure that all psychometric analyses were of the highest quality.

Section 10 – Operational Test Forms

This section describes the operational test forms, including high-level blueprints for the assessments.

Section 11 – Student Characteristics

This section describes the composition of test forms; rules for inclusion of students in analyses; distributions of students by grade, mode, and gender; and distributions of demographic variables of interest.

Section 12 – Scale Scores

This section provides an overview of the claims and subclaims, describes the development of the reporting scales and conversion tables, and presents scale score distributions. Finally, information regarding the interpretation of claim scores and subclaim scores is presented.

Section 13 – Reliability

The results of scale score reliability and internal consistency reliability analyses and corresponding standard errors of measurement, for each grade, content area, and mode (CBT or PBT) for all students, and for subgroups of interest, is provided in this section. This is followed by reliability results for subscores and reliability of classification (i.e., decision accuracy and decision consistency). Finally, expectations and results for inter-rater agreement for handscored items are summarized.

Section 14 – Validity

Validity evidence based on analyses of the internal structure of the tests is provided in this section. Correlations between subscores are reported by grade, content area, and mode (CBT or PBT) for all students.

Section 15 – Student Growth Measures

This section provides details on student growth percentiles (SGP). Information about the model, model fit, and SGP averages at the overall level for all students, and for subgroups of interest, are provided in this section.

References

Appendices

To facilitate utility, tables in the appendices are numbered sequentially according to the section represented by the tables. For example, the first appendix table for Section 6 is numbered A.6.1, the second appendix table for Section 6 is numbered A.6.2, and so on.

1.7 Glossary of Abbreviations

Table 1.1 Glossary of Abbreviations and Acronyms

Abbreviation/Acronym	Definition
A1	Algebra I
A2	Algebra II
AAF	Accessibility, Accommodations, and Fairness
ABBI	Assessment Banking for Building and Interoperability
AERA	American Educational Research Association
AIQ	Assessment and Information
AIS	Average Item Score
AmerIndian	American Indian/Alaska Native
APA	American Psychological Association
ASC	Additional and Supporting Content (Mathematics)
ASL	American Sign Language
CBT	Computer-Based Test
CCSS	Common Core State Standards
COVID-19	Coronavirus Disease 2019
CR	Constructed Response
CSEM	Conditional Standard Error of Measurement
DIF	Differential Item Functioning
EconDis	Economically disadvantaged
EBSS	Evidence-Based Standard Setting
ELA/L	English Language Arts/Literacy
EL	English Learner
ELL	English Language Learner
ELN	Not an English learner
ELY	English Learner
ePEN2	Electronic Performance Evaluation Network second generation
ESEA	Elementary and Secondary Education Act
FRL	Free or Reduced-Price Lunch
GO	Geometry
HOSS	Highest Obtainable Scale Score
HumRRO	Human Resources Research Organization
IA	Item Analysis
IDEA	Individuals with Disabilities Education Act
IEA	Intelligent Essay Assessor
IEP	Individualized Education Program
INF	Information Curve
IRT	Item Response Theory
K-12	Kindergarten to Grade 12
LEA	Local Education Agency
LID	Local Item Dependence
LOSS	Lowest Obtainable Scale Score
LSA	Latent Semantic Analysis
M1	Integrated Mathematics I
M2	Integrated Mathematics II
M3	Integrated Mathematics III
MC	Major Content (Mathematics)

Abbreviation/Acronym	Definition
MH	Mantel-Haenszel
MP	Modeling Practice (Mathematics)
MR	Mathematical Reasoning
Multiracial	Multiple Races Selected
n/a	Not Applicable
NAEP	National Assessment of Educational Progress
NCLB	No Child Left Behind
NCME	National Council on Measurement in Education
NoEconDis	Not Economically Disadvantaged
n/r	Not Reported
OWG	Operational Working Group
Pacific Islander	Native Hawaiian or Pacific Islander
PARCC	Partnership for Assessment of Readiness for College and Careers
PBT	Paper-Based Test
PCR	Prose Constructed Response (ELA/L)
PEJ	Postsecondary Educators' Judgment
PLD	Performance Level Descriptor
PLS	Performance Level Setting
RD	Reading (ELA/L)
RI	Reading Information (ELA/L)
RL	Reading Literature (ELA/L)
RV	Reading Vocabulary (ELA/L)
SD	Standard Deviation
SEM	Standard Error of Measurement
SGP	Student Growth Percentile
SSMC	Single Select Multiple Choice
SWD	Students with Disabilities
SWDN	Not student with disability
SWDY	Students with Disabilities
TCC	Test Characteristic Curve
TTS	Text-to-Speech
WE	Writing Written Expression (ELA/L)
WKL	Writing Knowledge Language and Conventions (ELA/L)
WLS	Weighted Least Squares
WR	Writing (ELA/L)

Section 2: Test Development

2.1 Overview of the Summative Assessments, Claims, and Design

Aligned to the Common Core State Standards (CCSS) as articulated in the Model Content Frameworks, the summative assessments are designed to determine whether students are college- and career-ready or on track, assess the full range of the CCSS, measure the full range of student performance, and provide data to help inform instruction, interventions, and professional development. Test development is an ongoing process involving educators, researchers, psychometricians, subject matter professionals, and assessment experts who participate in the development of the test design and its underlying foundational documents; develop and review passages and items used to build the summative assessments; monitor the program for quality, accessibility, and fairness for all students; and construct, review, and score the assessments.

The summative assessments include both English language arts/literacy (ELA/L) and mathematics assessments in grades 3 through 8 and high school. The high school mathematics tests include traditional mathematics and integrated mathematics course pathways. Assessments contain selected response, brief and extended constructed response, technology-enabled and technology-enhanced items, as well as performance tasks. Technology-enabled items are single-response or constructed-response items that involve some type of digital stimulus or open-ended response box with which the students engage in answering questions. Technology-enhanced items involve specialized student interactions for collecting performance data. In other words, the act of performing the task is the way in which data is collected. Students may be asked, among other interactions, to categorize information, organize or classify data, order a series of events, plot data, generate equations, highlight text, or fill in a blank. One example of a technology-enhanced item is an interaction in which students are asked to drag response options onto a Venn diagram to show the relationship among ideas.

The summative assessments offer a wide range of accessibility features for all students and accommodations for students with disabilities (SWDs). These accommodations may include screen readers, assistive technology, braille, large print (LP), text-to-speech (TTS), and American Sign Language (ASL) video versions of the test, as well as response accommodations that allow students to respond to test items using different formats. For English learners (ELs) who are native Spanish speakers, participating states and agencies offer the mathematics assessments in Spanish, and both LP and TTS versions of the test in Spanish (refer to the Accessibility Features and Accommodations Manual for in-depth information).

2.1.1 English Language Arts/Literacy Assessments—Claims and Subclaims

The ELA/L summative assessment at each grade level consists of three task types: literary analysis, research simulation, and narrative writing. For each performance-based task, students are asked to read or view one or more texts, answer comprehension and vocabulary questions, and write an extended response that requires them to draw evidence from the text(s). The summative assessment also contains literary and informational reading passages with comprehension and vocabulary questions.

The claim structure, grounded in the CCSS, undergirds the design and development of the ELA/L summative assessments.

Master Claim: The master claim is the overall performance goal for the ELA/L Summative Assessment System—students must demonstrate that they are college- and career-ready or on track to readiness as demonstrated through reading and comprehending of grade-level texts of appropriate complexity and writing effectively when using and/or analyzing sources.

Major Claims: (1) reading and comprehending a range of sufficiently complex texts independently, and (2) writing effectively when using and/or analyzing sources.

Subclaims: The subclaims further explicate what is measured on the summative assessments and include claims about student performance on the standards and evidence outlined in the evidence tables for reading and writing (refer to the test specifications documents). The claims and evidence are grouped into the following categories:

1. Vocabulary Interpretation and Use
2. Reading Literature
3. Reading Informational Text
4. Written Expression
5. Knowledge of Language and Conventions

2.1.2 Mathematics Assessments—Claims and Subclaims

The summative mathematics assessment at each grade level includes both short- and extended-response questions focused on applying skills and concepts to solve problems that require demonstration of the mathematical practices from the CCSS with a focus on modeling and reasoning with precision. The assessments also include performance-based short-answer questions focused on conceptual understanding, procedural skills, and application.

The claim structure, grounded in the CCSS, undergirds the design and development of the summative assessments.

Master Claim: The degree to which a student is college- or career-ready or on track to being ready in mathematics. The student solves grade-level/course-level problems aligned to the Standards for Mathematical Content with connections to the Standards for Mathematical Practice.

Subclaims: The subclaims further explicate what is measured on the summative assessments and include claims about student performance on the standards and evidence outlined in the evidence statement tables for mathematics (refer to the test specifications documents). The claims and evidence are grouped into the following categories:

- **Subclaim A:** Major Content with Connections to Practices
- **Subclaim B:** Additional and Supporting Content with Connections to Practices
- **Subclaim C:** Highlighted Practices with Connections to Content: Expressing mathematical reasoning by constructing viable arguments, critiquing the reasoning of others, and/or attending to precision when making mathematical statements
- **Subclaim D:** Highlighted Practice with Connections to Content: Modeling/Application by solving real-world problems by applying knowledge and skills articulated in the standards

2.2 Test Development Activities

Test development activities began with the standards and model content frameworks. From these, more than 2,000 educators, researchers, and psychometricians have developed the test specifications documents that guide the development of test items and the composition of the tests. These documents include the College- and Career-Ready Determinations and Performance Level Descriptors, Claim Structure, Evidence Statement Tables, Blueprints, Informational Guides, Passage Selection Guidelines, Mathematics Sequencing Guidelines, Task Generation Models, Fairness and Sensitivity Guidelines, Text Selection Guidelines, and the Style Guide. Refer to the [website](#) for further information about these documents.

2.2.1 Item Development Process

Test and item development activities were conducted by Pearson under the guidance and oversight of the K–12 state leads, the Higher Education Leadership Team, the Technical Advisory Committee, the Operational Working Group (OWG) members from each of the member states, the Text and Content Item Review Committees, and staff members from New Meridian, the project manager.

Developing high-quality assessment content with authentic stimuli for computer-based tests and paper-based tests measuring rigorous standards is a complex process involving the services of many experts including assessment designers, psychometricians, managers, trainers, content providers, content experts, editors, artists, programmers, technicians, human scorers, advisors, and members of the OWGs.

Bank Analysis and Item Development Plan

The summative item bank houses passages and items at each assessed grade level and subject. The bank supports the administration of the assessments, along with item release and practice tests. Items are developed and field-tested annually. Prior to the annual item development cycle, the item development teams, in conjunction with members of the OWGs for ELA/L and mathematics, evaluated the strengths of the bank and considered the needs for future tests to establish an item development plan.

Text Selection for ELA/L

Using the Passage Selection Guidelines, English language arts subject matter experts were trained to search for appropriate passages to support an annual pool of passages for consideration. Guided by the test specifications documents, Pearson recruited, trained, and managed the contracted subject matter experts to deliver the number of texts specified in the annual asset development plan. The Passage Selection Guidelines provided a text complexity framework and guidance on selecting a variety of text types and passages that allow for a range of standards/evidence to be demonstrated to meet the assessment claims. ELA/L tests are based on authentic texts, including multimedia stimuli. Authentic texts are grade-appropriate texts that are not developed for the purposes of the assessment or to achieve a particular readability metric but reflect the original language of the authors. Pearson content experts reviewed the passages for adherence to the Passage Selection Guidelines to meet the annual asset development plan described above in the number and distribution of genres and topics prior to review and consideration by the Text Review Committee. ELA/L item development was not conducted until after texts were approved by the Text Review Committee.

Item Development

Guided by foundational documents, Pearson recruited and trained the item writers and managed the item writing to develop the number of items specified in the annual asset development plan. Prior to further committee reviews, the assessment teams at Pearson reviewed the items for content accuracy, alignment to the standards, range of difficulty, adherence to universal design principles (which maximize the participation of the widest possible range of students), bias and sensitivity, and copy editing to enable the accurate measurement of the standards.

2.2.2 Item and Text Review Committees

Members of the OWGs for ELA/L and mathematics, state-level experts, local educators, postsecondary faculty, and community members conducted rigorous reviews of every item and passage being developed for the summative assessment system to ensure all test items are of the highest quality, aligned to the standards, and fair for all student populations. All reviewers were nominated by their state education agency. The purpose of the educator reviews was to provide feedback to Pearson and participating states and agencies on the quality, accuracy, alignment, and appropriateness of the test passages and items developed annually for the summative assessments. The meetings were conducted either in person or virtually and included large group training on the expectations and processes of each meeting, followed by breakout meetings of grade/subject working committees where additional training was provided.

Text Review

The Text Review is a review and approval by the Text Review Committee of the texts eligible for item development. Participants reviewed and provided feedback to Pearson and participating states and agencies about the grade-level appropriateness, content, and potential bias concerns, and reached consensus about which texts would move forward for development. The Text Review Committee was made up of members of both Content Item Review and Bias and Sensitivity Review Committees.

Content Item Review

During Content Item Review, committees reviewed and edited test items for adherence to the foundational documents, basic universal design principles, Accessibility Guidelines, associated item metadata, and the Style Guide. Committees accessed the item content within the Pearson Assessment Banking for Building and Interoperability (ABBI) system, which previews how the passages and items will be displayed in an operational online environment. Committees also verified that the appropriate scoring rule had been applied to each item. The Content Item Review Committees were made up of OWG members and educators nominated by participating states.

Bias and Sensitivity Review

Educators and community members make up the committee that reviews items and tasks to confirm that there are no bias or sensitivity issues that would interfere with a student's ability to achieve his or her best performance. The committee reviewed items and tasks to evaluate adherence to the Fairness and Sensitivity Guidelines, and to ensure that items and tasks do not unfairly advantage or disadvantage one student or group of students over another. Bias and Sensitivity Committee members made edits and modifications to items and passages to eliminate sources of bias and improve accessibility for all students.

Editorial Review

The Editorial Review Committee consists of editors who reviewed up to 10 percent of the items and tasks. The committee reviewed the items for grammar, punctuation, clarity, and adherence to the Style Guide.

Data Review

Following the field test, educator and bias committee members met to evaluate test items and associated performance data with regard to appropriateness, level of difficulty, and potential gender, ethnic, or other bias, then recommended acceptance or rejection of each field-test item for inclusion on an operational assessment. The Data Review Committee also made recommendations that items be revised and re-field-tested. Items that were approved by the committee are eligible for use on operational summative assessments.

2.2.3 Operational Test Construction

Under the guidance in the operational test form creation specifications, Pearson constructed the operational forms to adhere to the test blueprints and the assessment goals outlined in the form creation specifications. These goals were:

- test forms designed to measure well across the full range of student ability;
- scores that are comparable among forms and across test administrations;
- scales that support classification of students into performance levels;
- maximization of the number of parallel forms;
- minimization of overexposure of items; and
- adherence to standards for validity, reliability, and fairness (*Standards for Educational and Psychological Testing*, AERA, APA, & NCME, 2014).

Each content-area and grade-level assessment was based on a specific test blueprint that guided how each test was built. Test blueprints determined the range and distribution of content, and the distribution of points across the subclaims and task types.

Multiple core forms were constructed for a given assessment to enhance test security and to support opportunity for item release. Core forms were the operational test forms consisting of only those items that counted toward a student's score. These forms were designed to facilitate psychometric equating through a common item linking strategy and to be constructed as "parallel" as possible from a content and test-taking experience. Evaluation criteria for parallelism included adherence to blueprint; sequencing of content across the forms; statistical averages and distributions for difficulty (e.g., p-value) and discrimination (e.g., polyserial correlation); item type and cognitive complexity; and passage characteristics for ELA/L including genre, topics, word count, and text complexity.

Additionally, appropriate forms were identified as accessibility and accommodated forms. The forms are accommodated to support braille, LP, human reader/human signers, assistive technology, TTS, closed-captioning, and Spanish. Human reader/human signers and Spanish are provided for mathematics assessments only. Closed captioning is provided for ELA/L assessments only.

Test Construction Activities

After the data review meetings and prior to the test construction meetings, Pearson assessment specialists constructed initial versions of all the core forms. Content specialists constructed the initial core forms based on the support documents and specific processes to achieve fair parallel forms. The following steps were used to construct the operational core forms taken to the Test Construction Committee for review:

1. constructed the online forms to match the blueprint and test construction specifications;
2. constructed the paper forms to match the blueprint and test construction specifications; and
3. constructed accommodated and accessibility forms to match the blueprint, test construction specifications, and Accessibility, Accommodations, and Fairness (AAF) constraints.

The test construction process included iterative steps between content specialists and psychometricians. Custom test construction reports generated by the Pearson psychometric team provided information on adherence to blueprint and statistical averages/distributions of item difficulty and discrimination describing the forms and allowing comparison of the forms. These reports facilitated content changes to better achieve the test construction goals. Equating across operational forms within an administration was accomplished by repeating core items across forms. Linking across administrations for operational forms was accomplished by including prior operational items on the current operational test forms.

Pearson assessment specialists identified forms for each grade/subject suitable for use as the accommodated forms. Pearson psychometrics reviewed the psychometric properties of each of the accommodated forms with respect to the required criteria. The content of these forms was also reviewed by Pearson accessibility specialists allowing for content changes prior to the Test Construction Committee meetings.

These test construction activities provided significant inputs to commence the meetings including:

- the proposed items for the initial operational core forms and the accommodated forms described above;
- reports describing each form and comparing parallel forms; and
- recommended accommodated forms.

Test Form Verification Meeting to Review Test Construction Inputs

Members of the Content Item Review Committees and the AAF OWG participated in the building of operational core forms that met the summative assessment requirements. In that process, they met in an in-person meeting to review and make recommendations for changes so that test forms conformed to both the content and psychometric requirements of the assessment.

Accommodated Form Review Process

In addition to participating in many of the development activities, including the Text Review and the Bias and Sensitivity Review meetings, the AAF OWG reviewed the proposed accommodated forms at the Test Construction Committee meeting for accessibility to make sure that the content can be accommodated for SWDs and ELs without changing the underlying measured construct.

Forms were identified to support the following accommodations:

Accommodated Base 1

- Spanish paper (also serves Spanish LP, Spanish human reader paper)
- Spanish human reader/human signer online
- base accommodated paper (serves braille, LP, human reader paper)
- human reader/human signer online
- assistive technology screen reader
- assistive technology non-screen reader
- ASL

Accommodated Base 2

- closed-captioning
- TTS first form
- Spanish online
- Spanish TTS

Accommodated Base 3 (mathematics only)

- TTS second form

Spanish is mathematics only. Closed captioning is ELA/L only.

At the conclusion of the meetings, all test forms were constructed to meet test blueprints and requirements, and if necessary, reflect the operational linking design. Each test form reflected the test blueprint in terms of content, item types, and test length, as well as *expected* difficulty and performance along the ability continuum. Linking sets were proportionally representative of the operational test blueprint. The operational core forms, linking set forms, and field-test forms were reviewed by the Forms Review Committees and approved prior to the test administration.

Spanish-Language Assessments for Mathematics

For English learners, the mathematics assessments are offered in Spanish, as well as in Spanish-language LP and TTS versions. Once the operational form was approved, the form was sent to Pearson's subcontractor, Teneo, for transadaptation of the items. Transadaptation differs from translation in that it takes into consideration the grade-level appropriateness of the words, as well as the linguistic and cultural differences that exist between speakers of two different languages. Accounting for these differences allows the item to measure the achievement of Spanish language speakers in the same way that the original version of the item does for native speakers of English. The Spanish Glossary provided guidance to the translator conducting the transadaptation in grade-level and culturally appropriate ways of transadapting the items. For the Spanish language TTS form, the alternate text (used for description and/or text in art and graphics) was transadapted from the alternate text for the English language version of the TTS form. Phonetic markup, which guides how the TTS reader pronounces content-specific words and phrases, was also applied in this process.

In addition to the expert review of potential content for all accommodated forms conducted by the AAF OWG with assistance from content experts at the test construction meetings, the transadapted forms underwent additional quality checks: a Pearson Spanish copy edit services review and approval, and an AAF OWG review and approval.

2.2.4 Linking Design of the Operational Test

To support the goal of score comparability within and across administrations and years, a hybrid approach was implemented that incorporated the strengths of common item linking and randomly equivalent groups. The use of repeated operational core items was leveraged for common item linking. In addition, all forms were available throughout the operational administration, with spiraling at the student level, leveraged to support linking through randomly equivalent groups.

The operational test forms involved various types of linking; horizontal linking and across-administration linking. Horizontal linking consisted of linking items, or common items, included in both forms in a single administration, which was the case for mathematics forms and some ELA/L forms. Across-administration linking, or year-to-year linking, consisted of common items included in two different administrations, which was used for all forms due to the pre-equated model. The placement of linking items across forms or administrations supports the development of comparable scores.

Linking item sets can be internal or external linking sets. Internal linking sets consist of common items in operational positions such that the items contribute to the students' scores. External linking sets consist of common items in positions resulting in the items not contributing to students' scores. The current linking designs included internal linking sets.

2.2.5 Field Test Data Collection Overview

Field-test items were embedded in the spring operational forms to collect data for psychometric analysis necessary to support the assessment system for future administrations. Field-test administration entailed paper and computer administration modes, with computer administration as the dominant mode. The ELA/L unit of field-test items was administered to a sample of students.

Field-test sets were constructed to balance the expected cognitive load and difficulty across forms, reflected in the number of points, distribution of task types, and balance of passages for ELA/L. Forms for each content area were spiraled at the student level. The data collection design entailed three conditions. Condition 1, which comprised the mathematics assessment, was an embedded census field-test model in which all students taking the summative assessment participated in the field test.

Under Condition 2, which comprised the ELA/L assessment, approximately one-third of the schools were sampled across some of the participating states. Students in the sampled schools or districts took forms containing ELA/L embedded field-test tasks. Schools or districts were selected so that the sample for each ELA/L assessment was representative of the general testing populations in terms of achievement (i.e., average scale score and percentage of students at Level 4 and Level 5 in the previous year) and demographics (i.e., ethnicity composition, percentage of economically disadvantaged, ELs, and SWDs). The sampling plan was created such that if a given school was part of the ELA/L field test one year (e.g., spring 2017), it would not be required to participate in the field test for the subsequent two years (e.g., spring 2018 and spring 2019).

For Condition 3, states or agencies may select to field-test two ELA/L grade levels rather than all grade levels. The grade levels selected participate in a census field test in which all students are administered the embedded field-test items. The remaining grade levels do not participate in field testing. The selected grade levels are rotated across years. In spring 2022, the District of Columbia selected three ELA/L grades (ELA/L grades 3, 6, and 10) for District-specific field testing, and all students in these three grades were included. In Illinois and New Jersey, around 50 percent of students were sampled in each grade.

Section 3: Test Administration

3.1 Test Security and Administration Policies

The administration of the summative assessment is a secure testing event. Maintaining the security of test materials before, during, and after the test administration is crucial to obtaining valid and reliable results. School Test Coordinators are responsible for ensuring that all personnel with authorized access to secure materials are trained in and subsequently act in accordance with all security requirements.

School Test Coordinators must implement chain-of-custody requirements for specified materials. School Test Coordinators are responsible for distributing materials to Test Administrators, collecting materials from Test Administrators, returning secure test materials, and securely destroying certain specified materials after testing.

The administration of the summative assessment includes both secure and nonsecure materials, and these materials are further delineated by whether they are “scorable” or “nonscorable,” depending on whether the assessments were administered via paper/pencil (i.e., paper-based assessments) or online (i.e., computer-based assessments). For the paper-based administration, students used paper-based answer documents (except in grade 3, where students responded directly in test booklets). Nearly all of the summative assessments administered during the 2021–2022 administration were online assessments (see Tables 11.1 through 11.3).

3.1.1 Secure versus Nonsecure Materials

Participating states and agencies define secure materials as those that must be closely monitored and tracked to prevent unauthorized access to or prohibited use or distribution of secure content such as test items, reading passages, student work, and so on. For paper-based tests (PBTs), secure materials include both used and unused test booklets and used scratch paper, while for computer-based tests (CBTs), secure materials include student testing tickets, secure administration scripts (e.g., mathematics read-aloud), and used scratch paper. Nonsecure materials are defined as any authorized testing materials that do not include secure content (e.g., test items or student work). These include test administration manuals, unused scratch paper, and mathematics reference sheets that have not been written upon, and so on.

3.1.2 Scorable versus Nonscorable Materials

Paper-based assessments have both scorable and nonscorable materials while computer-based assessments have only nonscorable materials. Scorable materials for paper-based assessments consist of used (including student work) test booklets (grade 3) and answer documents (grades 4 and above) only. Scorable materials must be returned to the vendor to be scored. All other materials for PBTs, such as blank (i.e., unused) test booklets, test administration manuals, scratch paper, mathematics reference sheets, and so on, are deemed nonscorable. For CBTs, there are no scorable materials as student work is submitted electronically for scoring. Thus, there are limited physical materials to return (e.g., secure administration scripts for certain accommodations).

Students taking the CBT may not have access to secure test materials before testing, including printed student testing tickets. Printed mathematics reference sheets (if applicable) and scratch paper must be new and unmarked.

Students taking the PBT may not have access to scorable or nonscorable secure test content before or after testing. Scorable secure materials that are to be provided by Test Administrators to students include test booklets (grade 3) or answer documents (grades 4 through high school). Nonscorable secure materials that are distributed by Test Administrators to paper-based testing students include large print test booklets, braille test booklets, scratch paper (paper used by students to take notes and work through items), and printed mathematics reference sheets (grades 5 through 8 and high school).

School Test Coordinators are required to maintain a tracking log to account for collection and destruction of test materials, including mathematics reference sheets and scratch paper written on by students. As part of the test administration policy, schools are required to maintain the Chain-of-Custody Form or tracking log of secure materials for at least three years unless otherwise directed by state policy. Copies of the Chain-of-Custody Form for paper-based testing are included in each local education agency (LEA) or school's test materials shipment.

Test administrators are not to have extended access to test materials before or after administration (except for certain accessibility or accommodations purposes). Test administrators must document the receipt and return of all secure test materials (used and unused) to the school test coordinator immediately after testing.

All test security and administration policies are found in the *Test Coordinator Manual* and the *Test Administrator Manuals*. State-specific policies are included in Appendix C of the *Test Coordinator Manual*.

3.2 Accessibility Features and Accommodations

3.2.1 Participation Guidelines for Assessments

All students, including students with disabilities (SWDs) and English learners (ELs), are required to participate in statewide assessments and have their assessment results be part of the state's accountability systems, with narrow exceptions for ELs in their first year in a U.S. school, and certain SWDs who have been identified by the Individualized Education Program (IEP) team to take their state's alternate assessment. Federal laws governing student participation in statewide assessments include the No Child Left Behind Act of 2001 (NCLB), the Individuals with Disabilities Education Act of 2004 (IDEA), Section 504 of the Rehabilitation Act of 1973 (reauthorized in 2008), and the Elementary and Secondary Education Act (ESEA) of 1965, as amended. All students can receive accessibility features on the summative assessments.

Four distinct groups of students may receive accommodations on the summative assessments:

1. SWDs who have an IEP;
2. students with a Section 504 plan who have a physical or mental impairment that substantially limits one or more major life activities, have a record of such an impairment, or are regarded as having such an impairment, but who do not qualify for special education services;
3. students who are ELs; and
4. students who are ELs with disabilities who have an IEP or 504 plan.

These students are eligible for accommodations intended for both SWDs and ELs. Testing accommodations for SWDs or students who are ELs must be documented according to the guidelines and requirements outlined in the *Accessibility Features and Accommodations Manual*.

3.2.2 Accessibility System

Through a combination of universal design principles and accessibility features, participating states and agencies designed an inclusive assessment system by considering accessibility from initial design through item development, field testing, and implementation of the assessments for all students, including SWDs, ELs, and ELs with disabilities. Accommodations may still be needed for some SWDs and ELs to assist in demonstrating what they know and can do. However, the accessibility features available to students should minimize the need for accommodations during testing and ensure the inclusive, accessible, and fair testing of the diverse students being assessed.

3.2.3 What Are Accessibility Features?

On computer-based assessments, accessibility features are tools or preferences that are either built into the assessment system or provided externally by Test Administrators and may be used by any student taking the summative assessments (i.e., students with and without disabilities, gifted students, ELs, and ELs with disabilities). Since accessibility features are intended for all students, they are not classified as accommodations. Students should have the opportunity to select and practice using them prior to testing to determine which are appropriate for use on the assessment. Consideration should be given to the supports a student finds helpful and consistently uses during instruction. Practice tests that include accessibility features are available for teacher and student use throughout the year.

3.2.4 Accommodations for Students with Disabilities and English Learners

It is important to ensure that performance in the classroom and on assessments is influenced minimally, if at all, by a student's disability or linguistic/cultural characteristics that may be unrelated to the content being assessed. For the summative assessments, accommodations are considered to be adjustments to the testing conditions, test format, or test administration that provide equitable access during assessments for SWDs and students who are ELs. In general, the administration of the assessment should not be the first occasion on which an accommodation is introduced to the student. To the extent possible, accommodations should

- provide equitable access during instruction and assessments;
- mitigate the effects of a student's disability;
- not reduce learning or performance expectations;
- not change the construct being assessed; and
- not compromise the integrity or validity of the assessment.

Accommodations are intended to reduce and/or eliminate the effects of a student's disability and/or English language proficiency level; however, **accommodations should never reduce learning expectations by reducing the scope, complexity, or rigor of an assessment**. Moreover, accommodations provided to a student on the summative assessments must be generally consistent with those provided for classroom instruction and classroom assessments. There are some accommodations that may be used for instruction and for formative assessments that are not allowed for the summative assessment because they impact the validity

of the assessment results—for example, allowing a student to use a thesaurus or access the internet during an assessment. There may be consequences (e.g., excluding a student’s test score) for the use of nonallowable accommodations during assessments. It is important for educators to become familiar with the participating state and agencies’ policies regarding accommodations used for assessments.

To the extent possible, accommodations should adhere to the following principles:

- Accommodations enable students to participate more fully and fairly in instruction and assessments and to demonstrate their knowledge and skills.
- Accommodations should be based upon an individual student’s needs rather than on the category of a student’s disability, level of English language proficiency alone, level of or access to grade-level instruction, amount of time spent in a general classroom, current program setting, or availability of staff.
- Accommodations should be based on a documented need in the instruction/assessment setting and should not be provided for the purpose of giving the student an enhancement that could be viewed as an unfair advantage.
- Accommodations for SWDs must be described and documented in the student’s appropriate plan (i.e., either a 504 plan or an approved IEP), and must be provided if they are listed.
- Accommodations for ELs should be described and documented.
- Students who are ELs with disabilities are eligible to receive accommodations for both SWDs and ELs.
- Accommodations should become part of the student’s program of daily instruction as soon as possible after completion and approval of the appropriate plan.
- Accommodations should not be introduced for the first time during the testing of a student.
- Accommodations should be monitored for effectiveness.
- Accommodations used for instruction should also be used, if allowable, on local district assessments and state assessments.

In the following scenarios, the school must follow each state’s policies and procedures for notifying the state assessment office:

- a student **was provided a test accommodation that was not listed** in his or her IEP/504 plan/documentation for an English learner, or
- a student **was not provided a test accommodation that was listed** in his or her IEP/504 plan/documentation for an English learner.

3.2.5 Unique Accommodations

A comprehensive list of accessibility features and accommodations was provided in the *Accessibility Features and Accommodations Manual* that are designed to increase access to the summative assessments and that will result in valid, comparable assessment scores. However, SWDs or ELs may require additional accommodations that are not already listed. Participating states and agencies individually review requests for unique accommodations in their respective states and provide a determination as to whether the accommodation would result in a valid score for the student, and if so, would approve the request.

3.2.6 Emergency Accommodations

An emergency accommodation may be appropriate for a student who incurs a temporary disabling condition that interferes with test performance shortly before or during the assessment window. A student, whether or not they already have an IEP or 504 plan, may require an accommodation as a result of a recently occurring accident or illness. Cases include a student who has a recently fractured limb (e.g., arm, wrist, or shoulder); a student whose only pair of eyeglasses has broken; or a student returning to school after a serious or prolonged illness or injury. An emergency accommodation should be given only if the accommodation will result in a valid score for the student (i.e., does not change the construct being measured by the test[s]). If the principal (or designee) determines that a student requires an emergency accommodation on the summative assessment, an Emergency Accommodation Form must be completed and maintained in the student's assessment file. If required by a state, the school may need to consult with the state or district assessment office for approval. **The parent must be notified that an emergency accommodation was provided.** If appropriate, the Emergency Accommodation Form may also be submitted to the District Assessment Coordinator to be retained in the student's central office file. Requests for emergency accommodations will be approved after it is determined that use of the accommodation would result in a valid score for the student.

3.2.7 Student Refusal Form

If a student refuses an accommodation listed in his or her IEP, 504 plan, or (if required by the member state) an EL plan, the school should document in writing that the student refused the accommodation, and the accommodation must be offered and remain available to the student during testing. This form must be completed and placed in the student's file and a copy must be sent to the parent on the day of refusal. Principals (or designee) should work with Test Administrators to determine who, if any others, should be informed when a student refuses an accommodation documented in an IEP, 504, or (if required by the member state) EL plan.

3.3 Testing Irregularities and Security Breaches

Any action that compromises test security or score validity is prohibited. These may be classified as testing irregularities or security breaches. Below are examples of activities that compromise test security or score validity (note that these lists are not exhaustive). It is highly recommended that School Test Coordinators discuss other possible testing irregularities and security breaches with Test Administrators during training.

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

Electronic Devices

- using a cell phone or other prohibited handheld electronic device (e.g., smartphone, iPod, smart watch, personal scanner) while secure test materials are still distributed, while students are testing, after a student turns in his or her test materials, or during a break

(Exception: Test Coordinators, Technology Coordinators, Test Administrators, and Proctors are permitted to use cell phones in the testing environment only in cases of emergencies or when timely administration assistance is needed; LEAs may set additional restrictions on allowable devices as needed).

Test Supervision

- coaching students during testing, including giving students verbal or nonverbal cues, hints, suggestions, or paraphrasing or defining any part of the test;

- engaging in activities (e.g., grading papers, reading a book, newspaper, or magazine) that prevent proper student supervision at all times while secure test materials are still distributed or while students are testing;
- leaving students unattended for any period of time while secure test materials are still distributed or while students are testing;
- deviating from testing time procedures;
- allowing cheating of any kind;
- providing unauthorized persons with access to secure materials;
- unlocking a test in PearsonAccess^{next} during nontesting times;
- failing to provide a student with a documented accommodation or providing a student with an accommodation that is not documented and therefore is not appropriate; and
- allowing students to test before or after the state's test administration window.

Test Materials

- losing a student test booklet or answer document;
- losing a student testing ticket;
- leaving test materials unattended or failing to keep test materials secure at all times;
- reading or viewing the passages or test items before, during, or after testing
(*Exception: Administration of a human reader/signer accessibility feature for mathematics or accommodation for English language arts/literacy, which requires a Test Administrator to access passages or test items*);
- copying or reproducing (e.g., taking a picture of) any part of the passages or test items or any secure test materials or online test forms;
- revealing or discussing passages or test items with anyone, including students and school staff, through verbal exchange, email, social media, or any other form of communication; and
- removing secure test materials from the school's campus or removing them from locked storage for any purpose other than administering the test.

Testing Environment

- allowing unauthorized visitors in the testing environment;
- failing to follow administration directions exactly as specified in the *Test Administrator Manual*; and
- displaying testing aids in the testing environment (e.g., a bulletin board containing relevant instructional materials) during testing.

All instances of security breaches and testing irregularities must be reported to the School Test Coordinator immediately. The Form to Report a Testing Irregularity or Security Breach must be completed within two school days of the incident.

If any situation occurred that could cause any part of the test administration to be compromised, schools should refer to the *Test Coordinator Manual* for each state's policy and immediately follow those steps. Instructions for the School Test Coordinator or LEA Test Coordinator to report a testing irregularity or security breach are available in the *Test Coordinator Manual*.

3.4 Data Forensics Analyses

Maintaining the validity of test scores is essential in any high-stakes assessment program, and misconduct represents a serious threat to test score validity. When used appropriately, data forensic analyses can serve as an integral component of a wider test security protocol. The results of these data forensic analyses may be instrumental in identifying potential cases of misconduct for further follow-up and investigation.

The following data forensics analyses were conducted on the operational assessments:

- Response Change Analysis
- Aberrant Response Analysis
- Plagiarism Analysis
- Longitudinal Performance Modeling
- Internet and Social Media Monitoring
- Off-Hours Testing Monitoring

An overview of each data forensics analysis method is provided next.

3.4.1 Response Change Analysis

Response change analysis looks at how often student answers are changed, focusing specifically on an excessive number of wrong answers changed to right answers. In traditional paper-based, multiple-choice testing programs, this is sometimes referred to as “erasure analysis.”¹ The rationale for erasure analysis is that a teacher or administrator who is intent on improving classroom performance might be motivated to change student responses after the answer sheets are collected. A clustered number of student answer documents from the same school or classroom with unusually high numbers of answers changed from wrong to right might provide evidence to support follow-up investigation. The response change analysis extended the traditional erasure method to account for issues specific to computer-based testing as well as the variety of item types on the summative assessments, such as partial-credit, multi-part, and multiple-select items.

3.4.2 Aberrant Response Analysis

Aberrant response pattern detection analysis looks at the unusualness of student responses compared with what would be expected. Most simply, this can be thought of as quantifying the extent to which higher-scoring students miss easy questions and lower-scoring students answer difficult questions correctly. While it would be difficult to draw a definitive inference about a single student flagged as having an aberrant response pattern, a cluster of students with aberrant response patterns within a classroom or school might warrant further investigation.

¹ The term “erasure analysis” is sometimes objected to because it is inferential rather than descriptive. A more descriptive term is “mark discrimination analysis,” which recognizes that the scanning approach makes discriminations among the darkness of selected answer choices when multiple responses to a multiple-choice item are detected during answer sheet processing.

3.4.3 Plagiarism Analysis

Plagiarism analysis compares the responses given for a group of written composition items, looking for high degrees of similarity. For the summative assessments, the primary item type of interest was the prose constructed-response tasks in the English language arts/literacy content area. This analysis was conducted for prose constructed-response tasks administered online using some of the same artificial intelligence techniques that are applied in automated essay scoring. Specifically, this method was based on latent semantic analysis (LSA) technology to detect possible plagiarism. Using LSA, the content of each constructed response was compared against the content of every other constructed response, and a measure that indicated the degree of similarity was generated for each pair of response comparison. Because LSA provided a semantic representation of language, rather than a syntactic or word-based representation, it allowed the detection of potential copying behaviors, even when students or administrators substituted synonymous words or phrases.

3.4.4 Longitudinal Performance Monitoring

Longitudinal performance modeling evaluates the performance on the summative assessments across test administrations and identifies unusual performance gains in the unit of interest (e.g., school or district). A weighted least squares (WLS) regression methodology was evaluated and recommended by the Technical Advisory Committee for implementation starting spring 2017. The WLS identified unusual changes in test performance across two consecutive administrations of the assessment. In the WLS regression approach, mean current year scale scores are regressed on mean prior year scale scores, weighting by unit sample size. Standardized residuals are calculated by dividing raw residuals by their respective standard deviations. Units with a standardized residual exceeding 3.0 are flagged for unexpected performance.

3.4.5 Internet and Social Media Monitoring

Internet and social media monitoring were conducted by Caveon, LLC. Caveon's team monitored English-language websites and searchable forums that were publicly available for suspected proxy testing solicitations and website postings that contain, or appear to contain, infringements of protected operational test content. The internet and social media outlets monitored included popular websites (such as Facebook and Twitter), blogs, discussion forums, video archives, document archives, brain dumps, auction sites, media outlets, peer-to-peer servers, and so on. Caveon's process generated regular updates that categorize identified threats by level of actual or potential risk based upon the representations made on the websites, or actual analysis of the proffered content. For example, categorizations typically ranged from "cleared" (lowest risk but bookmarked for continued monitoring) to "severe" (highest risk). Note that this process only considered potential breaches of secure item content, not violations of testing administration policies. Potential breaches were reported directly to the state(s) implicated for further action. Summary reports describing the threats were provided through notification emails.

3.4.6 Off-Hours Testing Monitoring

Off-hours testing monitoring checks for suspicious testing activities at test administration locations occurring outside of the set windows for computer-based testing sessions. Participating states and agencies established set start and end times for administering computer-based assessments. Based on these hours, authorized users (that is, users with the State Role) were allowed to override the start and end times for a test session. The off-hours testing monitoring process tracked such occurrences and logged them in an operational report, which

listed the sessions within an organization that selected to test outside the set window. States could use this report to follow up with the organizations identified in the report.

Section 4: Item Scoring

4.1 Machine-Scored Items

4.1.1 Key-Based Items

Pearson performed a key review prior to the test administration to verify that the scoring (answer) keys were correct for each item. Once the forms were constructed and approved for publication, an independent key review was performed by an experienced third-party vendor. The vendor reviewed each item and confirmed that the key was correct. If discrepancies were identified, a Pearson senior content specialist or content manager reviewed the flagged item(s) and worked with the item developers to resolve the issue.

4.1.2 Rule-Based Items

Rule-based scoring refers to item types that use various scoring models. Participating states and agencies use Question and Test Interoperability item type implementation based on scoring model rules. Examples of these item types include “choice interaction,” which presents a set of choices where one or more choices can be selected; text entry, where the response is entered in a text box; hot spot or text interaction, where an area in a graph or text in a paragraph (for example) can be highlighted; or match interaction, where an association can be made between pairs of choices in a set. These items include the scoring rules and correct responses as part of their item XML (markup language) coding.

During the initial stages of item development, Pearson staff worked closely with participating states and agencies to first delineate the rules for the scoring rubrics and then to adjust those rules based on student responses. During item studies in spring 2015, Pearson content staff received input from the staff of participating states and agencies to develop a thorough rule-based scoring process that met their needs.

Pearson worked with the item developers to review initial scoring rules created during the item development. Once the rule-based scoring process was approved, and prior to test construction, Pearson content staff worked closely with the item developers to finalize scoring rubrics for items to be scored via the rule-based scoring method. The proposed scoring rubrics were sent for review, and if any additional changes were needed or new rules added, Pearson documented and applied the requested edits.

During test construction, Pearson monitored and evaluated the scoring and updated the scoring keys/scoring rules in the item bank. After the tryout items were scored, Pearson prepared a frequency distribution of student responses for each item or task scored using a rule-based approach and compared this to the expected response based on correct answers to ensure that scoring keys and rules were appropriately applied. The content team analyzed the student response data to determine if scoring was acceptable using the item metadata and the student response file in conjunction with any potential item issues as flagged by psychometrics. These frequency distributions included an indication of right/wrong and other identifying information defined by participating states and agencies, and those items that showed a statistical anomaly, whereby the frequency distribution was outside of the expected range, were sent to content experts to verify that the items were coded with the correct key.

Following the Rule-Based Scoring Educator Committee’s review, which occurred prior to year one test construction, Pearson analyzed the feedback from the committees and made recommendations about adjustments to the scoring rubrics based on the results of the reviews. Upon submission of the results, Pearson worked with the staff of participating states and agencies to discuss these findings and determine next steps prior to the completion of scoring. In subsequent years as scoring inquiries arise throughout the process of test construction, forms creation, testing, scoring, and psychometric analysis, items with scoring discrepancies are brought before the Priority Alert Task Force for resolution. This committee consists of representatives from each state as well as the content specialists at participating states and agencies and Pearson.

Following the initial development of the rule-based scoring rubrics, Pearson has continued to monitor and evaluate new item development to ensure the scoring rules established are maintained within all item types as approved.

Pearson continues to use several avenues to monitor scoring each year. Prior to testing, a third-party key review checks operational and field test items for correct keys. Any disputed items go to a second review with Pearson content experts and anything still in question is taken before the task force for review and possible key change. During testing, Pearson creates early testing files for frequency distribution analysis whereby items for which an incorrect key receives a high distribution of responses are further evaluated for accuracy. After testing, all responses are again evaluated for the distribution of responses and potential scoring abnormalities during psychometric analysis. Any change in scoring that may be requested as a result of the psychometric analysis is also taken before the Priority Alert Task Force for decisions. These processes are the same for both paper and online modes of testing.

4.2 Human or Handscored Items

Constructed-response items were scored by human scorers in a process referred to as handscoring. Online training units were used to train all scorers. The online training units included prompts (items), passages, rubrics, training sets, and qualification sets. Scorers who successfully completed the training and qualified, demonstrating they could correctly score student responses based on the guidelines in the online training units, were permitted to score student responses using the ePEN2 (Electronic Performance Evaluation Network, second generation) scoring platform. All online and paper responses were scored within the ePEN2 system. Pearson monitored quality throughout scoring.

Pearson staff roles and responsibilities were as follows:

- Scorers applied scores to student responses.
- Scoring supervisors monitored the work of a team of scorers through review of scorer statistics and backreading, which is a review of responses scored by each scorer. When backreading, a supervisor sees the scores applied by scorers, which helps the supervisor provide additional coaching or instruction to the scorer being backread.
- Scoring directors managed the scoring quality of a subset of items and monitored the work of supervisors and scorers for their assigned items. Directors backread responses scored by supervisors and scorers as part of their quality-monitoring duties.
- English language arts/literacy (ELA/L) and mathematics content specialists managed the scoring quality and monitored the work of the scoring directors.

- The project manager documented the procedures, identified risks, and managed day-to-day administrative matters.
- A portfolio manager provided oversight for the entire scoring process.

All Pearson employees involved in the scoring or the supervision of scoring possessed at least a four-year college degree.

4.2.1 Scorer Training

Key steps in the development of scorer training materials were rangefinding and rangefinder review meetings where educators and administrators from states met to interpret the scoring rubrics and determine consensus scores for student responses. Rangefinding meetings were held prior to scoring field-test items, and rangefinder review meetings were held prior to scoring operational items.

At rangefinding meetings, educators and administrators from states reviewed student responses and used scoring rubrics to determine consensus scores. Those responses scored in rangefinding were used to create field test scorer training sets. After items were selected for operational testing, educators and administrators attended rangefinder review meetings to review and approve proposed operational scorer training sets.

When developing scorer training materials, Pearson scoring directors carefully reviewed detailed notes and records from rangefinding and rangefinder review committee meetings. Training sets were developed using the responses scored by the committees and additional suitable student response samples (as needed). All scorer training sets were reviewed and approved prior to scorer training.

During training, scorers reviewed training sets of scored student responses with annotations that explained the rationale for the score assigned. The anchor set was the primary reference for scorers as they internalized the rubric during training. Each anchor set consisted of responses that were clear examples of student performance at each score point. The responses selected were representative of typical approaches to the task and arranged to reflect a continuum of performance. All scorers had access to the anchor set when they were training and scoring and were directed to refer to it regularly during scoring.

Practice sets were used in training to help trainees practice applying the scoring guidelines. Scorers reviewed the anchor sets, scored the practice sets, and then were able to compare their assigned scores for the practice sets to the actual assigned scores to help them learn.

Qualification sets were used to confirm that scorers understood how to score student responses accurately. Qualification sets were composed of responses that were clear examples of score points. Scorers were required to meet specified agreement percentages on qualification sets in order to score student responses.

Pearson has developed two types of training sets to train scorers: prototype and abbreviated sets. Prototype training sets were complete training sets consisting of anchor, practice, and qualification sets (refer to 4.2.2 for information on the qualification process). In ELA/L, there was one prototype training set per task type (Research Simulation Task, Literary Analysis Task, and Narrative Writing Task) at each of the nine grade levels (grades 3 through 11). In mathematics, a prototype training set was built for a grouping of similar items for a total of approximately three to four prototype sets per grade level or course.

The prototype training approach promoted consistency in scoring, as each subsequent abbreviated training set for the ELA/L task type or mathematics item grouping was based on the prototype. Once a prototype was

chosen, full training materials were developed for that item, and at each grade level, scorers were trained to score a particular item type using the prototype training materials for that type.

Abbreviated training sets were prepared for all items not selected for prototype training sets. The abbreviated training sets included an anchor set and two practice sets so scorers could internalize the scoring standards for these new items, which were similar to prototype items they had previously scored.

Anchor and practice sets for both prototype and abbreviated items included annotations for each response. Annotations are formal written explanations of the score for each student response.

Table 4.1 details the composition of the anchor sets, practice sets, and qualification sets.

Table 4.1 Training Materials Used During Scoring

Training Set Development	
Description	Specification
Anchor Set	
The anchor set is the primary reference for scorers as they internalize the rubric during training. All scorers have access to the anchor set when they are training and scoring, and are directed to refer to it regularly.	<p>The anchor set for mathematics prototype items comprises three annotated responses per score point.</p> <p>The anchor set for subsequent abbreviated items for mathematics comprise one to three annotated responses per score point.</p>
The anchor set comprises clear examples of student performance at each score point. The responses selected may be representative of typical approaches to the task or arranged to reflect a continuum of performance.	The anchor sets for ELA/L prototype items comprise three annotated responses per score point. Anchor sets for prototype items include separate complete anchor sets for each applicable scoring trait (Reading Comprehension and Written Expression and Conventions for Research Simulation and Literary Analysis Tasks, Written Expression for Narrative Writing Tasks, and Knowledge of Language and Conventions for all task types).
Practice Sets	
Practice sets are used to help trainees develop experience in independently applying the scoring guide (the rubric) to student responses. Some of these responses clearly reinforce the scoring guidelines presented in the anchor set. Other responses are selected because they are more difficult to evaluate, fall near the boundary between two score categories, or represent unusual approaches to the task.	<p>The practice sets for mathematics prototype and abbreviated items include two to three sets of ten annotated responses.</p> <p>ELA/L practice sets for prototype items include two sets of five annotated responses and two sets of 10 annotated responses.</p>
The practice sets provide guidance and practice for trainees in defining the line between score categories, as well as applying the scoring criteria to a wider range of types of responses.	The subsequent ELA/L practice sets for abbreviated items include two sets of ten annotated responses.

Training Set Development	
Description	Specification
Qualification Sets	
Qualification sets are used to confirm that scorer trainees understand the scoring criteria and are able to assign scores to student responses accurately. The responses in these sets are selected to reinforce the application of the scoring criteria illustrated in the anchor set.	<p>The qualification sets for mathematics prototype items include three sets of 10 responses each (not annotated).</p> <p>The subsequent mathematics abbreviated items for mathematics do not include qualification sets.</p>
Scorer trainees must demonstrate acceptable performance on these sets by meeting a predetermined standard for accuracy in order to qualify to score. Pearson scoring staff defined and documented qualifying standards in conjunction with participating states and agencies prior to scoring.	<p>The qualification sets for ELA/L prototype items include three sets of 10 responses each (not annotated).</p> <p>The subsequent ELA/L abbreviated items do not include qualification sets.</p>

4.2.2 Scorer Qualification

In order to score items, scorers were required to show that they were able to apply scoring methodology accurately through a qualification process. Scorers were asked to apply scores to three qualification sets consisting of 10 responses each. ELA/L scorers applied a score for each trait on each response in the qualification sets. Literary Analysis and Research Simulation Tasks each had two traits: the Reading Comprehension and Written Expression trait and the Conventions trait. The Narrative Writing Task had two traits: Written Expression and Conventions. Mathematics scorers applied a score for each part of an item that was a constructed response. The number of constructed-response parts for each mathematics item ranged from one to four. Scorers were required to match the approved score at a percentage agreed to by participating states and agencies in order to qualify.

For ELA/L qualification, scorers were required to meet the following three conditions:

1. On at least one of the three qualifying sets, at least 70 percent of the ratings on each of the two scoring traits (considered separately) must agree exactly with the approved scores.
2. On at least two of the three qualifying sets, at least 70 percent of the ratings (combined across the three scoring traits) must agree exactly with the approved scores.
3. Combining over the three qualifying sets and across the two scoring traits, at least 96 percent of the ratings must be within one point of the approved scores.

For mathematics qualification, the requirements were based on the item types and score point ranges. Because mathematics items can have one or more scoring traits, a scorer needed to achieve the requirements as set forth in Table 4.2 separately for each scoring trait (when applicable to the item).

Table 4.2 Mathematics Qualification Requirements

Category	Score Point Range	Perfect Agreement	Within One Point
2	0–1	90%	100%
3	0–2	80%	96%
4	0–3	70%	96%
5	0–4	70%	95%
6	0–5	70%	95%
7	0–6	70%	95%

On at least two of the three qualifying sets, a scorer was required to meet the “perfect agreement” percentage indicated in the table above for each category. “Perfect agreement” was achieved when the scores applied exactly matched the approved scores. Over the three qualifying sets, a scorer was required to meet the “within one point” percentage indicated in the table above for each category. The average is exclusive to each trait, so an item with multiple scoring traits would have multiple trait rating averages within one point of the approved score.

4.2.3 Managing Scoring

Pearson created a handscoring specifications document that detailed the handscoring schedule, customer requirements, rangefinding plans, quality management plans, item information, and staffing plans for each scoring administration.

4.2.4 Monitoring Scoring

Second Scoring

During scoring, Pearson’s ePEN2 scoring system automatically and randomly distributed a minimum of 10 percent of student responses for second scoring; scorers had no indication whether a response had been scored previously. Humans applied the second score for all mathematics items. Second scoring for ELA/L was performed either by human scorers or by Pearson’s Intelligent Essay Assessor. If the first and second scores applied were nonadjacent, a third and occasionally a fourth score were assigned to resolve scorer disagreements. When a resolution score (i.e., third score) was nonadjacent to one or both of the first and second scores, the content specialist or scoring director would apply an adjudication score (fourth score). If a response was scored more than once, the rules in Table 4.3 were applied to determine the final score.

Table 4.3 Scoring Hierarchy Rules

Score Type	Rank	Final Score Calculation
Adjudication	1	If an adjudication score is assigned, this is the final score.
Resolution	2	If no adjudication score is assigned, this is the final score.
Backread	3	If no adjudication or resolution score is assigned, the latest backreading score is the final score.
Human First Score	4	If no adjudication, resolution, or backreading score is assigned, this is the final score.
Human Second Score	5	If no adjudication, resolution, backreading, or human first score is assigned, this is the final score.
Intelligent Essay Assessor Score	6	If no human score is assigned, this is the final score.

Backreading

Backreading was one of the major responsibilities of Pearson Scoring Supervisors and a primary tool for proactively guarding against scorer drift, where scorers score responses in comparison to one another instead of in comparison to the training responses. Scoring supervisory staff used the ePEN2 backreading tool to review scores assigned to individual student responses by any given scorer in order to confirm that the scores were correctly assigned and to give feedback and remediation to individual scorers. Pearson backread approximately 5 percent of the handscored responses. Backreading scores did not override the original score but were used to monitor scorer performance.

Validity

Validity responses are prescored responses strategically interspersed in the pool of live responses. These responses were not distinguishable from any other responses so that scorers were not aware they were scoring validity responses rather than live responses. The use of validity responses provided an objective measure that helped ensure that scorers were applying the same standards throughout the project. In addition, validity was at times shared with scorers in a process known as “validity as review.” Validity as review provided scorers automated, immediate feedback: a chance to review responses they mis-scored, with reference to the correct score and a brief explanation of that score. One validity response was sent to scorers for every 25 “live” responses scored.

Validity agreement requirements for scorers are listed in Table 4.4. Scorers had to meet the required validity agreement percentages to continue working on the project. Scorers who did not maintain expected agreement statistics were given a series of interventions culminating in a targeted calibration set: a test of scorer knowledge. Scorers who did not pass targeted calibration were removed from scoring the item, and all the scores they assigned were deleted.

Table 4.4 Scoring Validity Agreement Requirements

Subject	Score Point Range	Perfect Agreement	Within One Point^{1*}
Mathematics	0–1	90%	96%
Mathematics	0–2	80%	96%
Mathematics	0–3	70%	96%
Mathematics	0–4	65%	95%
Mathematics	0–5	65%	95%
Mathematics	0–6	65%	95%
ELA/L	Multi-trait	65%	96%

*A numerical score compared to a blank or condition code score will have a disagreement greater than 1 point.

Calibration Sets

Calibration sets are special sets created during scoring to help train scorers on particular areas of concern or focus. Scoring directors used calibration sets to reinforce rangefinding standards, introduce scoring decisions, or address scoring issues and trends. Calibration was used either to correct a scoring issue or trend, or to continue scorer training by introducing a scoring decision. Calibration was administered regularly throughout scoring.

Inter-Rater Agreement

Inter-rater agreement is the agreement between the first and second scores assigned to student responses and is the measure of how often scorers agree with each other. Pearson scoring staff used inter-rater agreement statistics as one factor in determining the needs for continuing training and intervention on both individual and group levels. Inter-rater agreement expectations are shown in Table 4.5.

Table 4.5 Inter-Rater Agreement Expectations and Results

Subject	Score Point Range	Perfect Agreement Expectation	Perfect Agreement Result	Within One Point Expectation*	Within One Point Result
Mathematics	0–1	90%	98%	96%	100%
Mathematics	0–2	80%	97%	96%	100%
Mathematics	0–3	70%	96%	96%	99%
Mathematics	0–4	65%	94%	95%	99%
Mathematics	0–5	65%	91%	95%	98%
Mathematics	0–6	65%	95%	95%	98%
ELA/L	Multi-trait	65%	83%	96%	100%

*A numerical score compared to a blank or condition code score will have a disagreement greater than 1 point.

Pearson’s ePEN2 scoring system included comprehensive inter-rater agreement reports that allowed supervisory personnel to monitor both individual and group performance. Based on reviews of these reports, scoring experts targeted individuals for increased backreading and feedback, and if necessary, retraining.

The perfect agreement rate for mathematics responses scored by two scorers ranged from 76 to 100 percent and the within-one-point rate ranged from 96 to 100 percent. For all ELA/L responses scored by two scorers, the perfect agreement rate ranged from 69 percent to 100 percent and the within-one-point rate ranged from 97 percent to 100 percent.

The results by grade level for ELA/L are provided in Section 4.3.7, “Inter-rater Agreement for Prose Constructed Response.”

4.3 Automated Scoring for Prose Constructed Responses

Automated scoring performed by Pearson’s Intelligent Essay Assessor (IEA) was the default option for scoring the summative assessment’s online prose constructed-response (PCR) tasks. Under the default option, it was assumed that operational scores for approximately 90 percent of the online PCR responses would be assigned by IEA for the spring administration. The operational scores for the remaining online responses were assigned by human scorers. Human scoring was applied to responses that were scored while IEA was being trained as well as to additional responses routed to human scoring when there was uncertainty about the automated scores.

For 10 percent of responses, a second “reliability” score was assigned. The purpose of the reliability score was to provide data for evaluating the consistency of scoring, which is done by evaluating scoring agreement. When IEA provided the first score of record, the second reliability score was a human score.

4.3.1 Concepts Related to Automated Scoring

The discussion below describes concepts related to automated scoring.

Continuous Flow

Continuous flow scoring results in an integrated connection between human scoring and automated scoring. It refers to a system of scoring where either an automated score, a human score, or both can be assigned based on a predetermined asynchronous operational flow.

Training of IEA Using Operational Data

Continuous flow scoring facilitates the training of IEA using human scores assigned to operational online data collected early in the administration. Once IEA obtains sufficient data to train, it can be “turned on” and becomes the primary source of scoring (although human scoring continues for the 10 percent reliability sample and other responses that may be routed accordingly).

Smart Routing

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. Smart routing can be applied prompt by prompt to the extent needed to meet scoring quality criteria for automated scoring.

Quality Criteria for Evaluating Automated Scoring

The state leads approved specific quality criteria for evaluating automated scoring. The primary evaluation criteria for IEA was based on responses to validity papers with “known” scores assigned by experts. For each prompt scored, a set of validity papers is used to monitor the human-scoring process over time. Validity papers are seeded into human scoring throughout the administration. The expectation is that IEA can score validity papers at least as accurately as humans can.

Additional measures of inter-rater agreement for evaluating automated scoring were proposed based on the research literature (Williamson et al., 2012). These measures were previously utilized in Pearson’s automated scoring research and include Pearson correlation, kappa, quadratic-weighted kappa, exact agreement, and standardized mean difference. These measures are computed between pairs of human scores, as well as between IEA and humans, to evaluate how performance was the same or different. Criteria for evaluating the training of IEA given these measures include the following:

- Pearson correlation between IEA-human should be within 0.1 of human-human.
- Kappa between IEA-human should be within 0.1 of human-human.
- Quadratic-weighted kappa between IEA-human should be within 0.1 of human-human.
- Exact agreement between IEA-human should be within 5.25 percent of human-human.
- Standardized mean difference between IEA-human should be less than 0.15.

The specific criteria for evaluating IEA included both primary and secondary criteria and are noted below:

- Primary Criteria—Based on responses to validity papers: With smart routing applied as needed, IEA agreement is as good as or better than human agreement for each trait score.

- Contingent Primary Criteria—Based on the training responses if validity responses are not available: With smart routing applied as needed, IEA-human exact agreement is within 5.25 percent of human-human exact agreement for each trait score.
- Secondary Criteria—Based on the training responses: With smart routing applied as needed, IEA-human differences on statistical measures for each trait score are within the Williamson et al. (2012) tolerances for subgroups with at least 50 responses.

Hierarchy of Assigned Scores for Reporting

When multiple scores are assigned for a given response, the following hierarchy determines which score was reported operationally:

- The IEA score is reported if it is the only score assigned.
- If an IEA score and a human score are assigned, the human score is reported.
- If a first human score and a second human score are assigned, the first human score is reported.
- If a backread score and human and/or IEA scores are assigned, the backread score is reported if there is no resolution or adjudication score assigned.
- If a resolution score is assigned and an adjudicated score is not assigned, the resolution score is reported (note that if nonadjacent scores are encountered, responses are automatically routed to resolution).
- If an adjudicated score is assigned, it is reported (note that if a resolution score is nonadjacent to the other scores assigned, responses are automatically routed to adjudication).

4.3.2 Sampling Responses Used for Training IEA

For prompts trained using 2022 operational data, the early performance of human scoring was closely monitored to verify that an appropriate set of data would be available for training IEA. In particular, several characteristics of the human scoring data were monitored, including:

- exact agreement between human scorers (the goal was for this to be at least 65 percent for each trait);
- exact agreement between human scores conditioned on score point (the goal was for this to be at least 50 percent for each trait);
- the number of responses at each score point (the goal was to have at least 40 responses at the highest score points in the training samples used by IEA); and
- the number of responses with two human scores assigned (note that IEA “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. For other prompts, special sampling approaches were used to increase the numbers of responses that received top scores. In addition, a healthy percentage of responses were backread during the sampling period and these scores as well as double human scores were all part of the data used to train IEA.

4.3.3 Primary Criteria for Evaluating IEA Performance

The primary criteria for evaluating IEA performance are based on evaluating validity papers and is stated as follows: With smart routing applied as needed, IEA agreement is as good as or better than human agreement for each trait score.

To operationalize the primary criteria for a given prompt, the following general steps are undertaken:

1. Determine agreement of the human scores with the validity papers for each trait.
2. Calculate agreement of the IEA scores with the validity papers for each trait.
3. Compare the IEA validity agreement with the human agreement.
4. If the IEA validity agreement is greater than or equal to the human agreement for each trait, IEA can be deployed operationally.

In addition to looking at overall validity agreement, conditional agreement was also examined. In general, it was desirable for IEA to exceed 65 percent agreement at every score point as well as be close to or exceed the human validity agreement at each score point.

4.3.4 Contingent Primary Criteria for Evaluating IEA Performance

For many of the prompts trained in 2022, it was not possible to utilize human-scored validity responses in evaluating IEA performance. In these cases, IEA was evaluated based on IEA-human exact agreement for each trait score and compared to agreement based on responses that were double-scored by humans. A portion of the data was held out for evaluating IEA-human exact agreement according to the following steps:

1. Determine exact agreement of the two human scores with each other for each trait.
2. Calculate agreement of the IEA scores with the human scores for each trait.
3. Compare the IEA-human agreement with the human-human agreement.
4. If the IEA-human agreement is within 5.25 percent of the human-human agreement, IEA can be deployed operationally.

In addition to the overall comparison, the following performance thresholds were targeted in the test data set: (1) at least 65 percent overall IEA-human agreement; and (2) 50 percent IEA-human agreement by score point (i.e., conditioned on the human score). These targets went beyond the contingent primary criteria approved by the state leads.

4.3.5 Applying Smart Routing

With smart routing, the quality of automated scoring can be increased by routing responses that are more likely to disagree with a human score to receive an additional human score.

When human scorers read a paper, they typically apply integer scores based on a scoring rubric. When there is strong agreement between two independent human readers, the readers might both assign a score of 3 such that the average score over both raters is also a 3 (i.e., $(3+3)/2 = 3$). IEA simulates this behavior, but because its scores come from an artificial intelligence algorithm, it generates continuous (i.e., decimalized) scores. In this case, the IEA score might be a 2.9 or 3.1. When human readers disagree on the score for a paper, say one reader gives the paper a score of 3 and another reader gives the paper a score of 4, the average of the two scores

would be 3.5 (i.e., $3+4=7/2=3.5$). For this paper, IEA would likely provide a score between 3 and 4, say 3.4 or 3.6. Because this continuous score needs to be rounded to an integer score for reporting, it might be reported as a 3 or a 4, depending on the rounding rules. Smart routing involves routing those responses with “in between” IEA scores to additional human scoring because the nature of the responses suggests there may be less confidence in the IEA score. Since these “in between” IEA scores are based on modeling human scores, it follows that human scores may be less certain as well, and thus such responses tend to be the ones that it makes sense to have double-scored and possibly to resolve if the IEA and human scores are nonadjacent.

Smart routing was utilized as needed to help IEA achieve targeted quality metrics (e.g., validity agreement or agreement with human scorers). Smart routing involved the application of the following four steps:

1. The continuous IEA score for each of the two trait scores was rounded to the nearest score interval of 0.2, starting from zero. For example, IEA scores between 0 and 0.1 were rounded to an interval score of 0, scores between 0.1 and 0.3 were rounded to an interval score of 0.2, scores between 0.3 and 0.5 were rounded to an interval score of 0.4, and so on.
2. Within each of these intervals, the percentage of exact agreement between IEA integer scores and the human scores was calculated for each trait.
3. For each prompt, agreement rates were evaluated by rounding interval. Those intervals for which the agreement rates were below a designated threshold for either trait were identified.
4. Once IEA scoring was implemented, responses within intervals for which IEA-human agreement was below the designated threshold were routed for additional human scoring.

In training IEA, the scoring models without smart routing were evaluated first by applying either the primary validity criteria or the contingent criteria as described in Section 4.3.4. For those prompts that did not meet these criteria, increasing smart routing thresholds were applied in an iterative fashion to filter scores and evaluate the remaining scores against the criteria. That is, in any one iteration a particular smart routing threshold was applied such that only scores falling in intervals for which exact agreement exceeded the threshold were included in evaluating the criteria. If the primary or contingent criteria were not met with this level of smart routing, an increased smart routing threshold was applied iteratively until the primary or contingent criteria were met, or the maximum threshold reached. If the criteria were still not met after a maximum threshold was applied, different models were investigated and/or additional human scoring data utilized until an IEA scoring model was found that met the criteria.

4.3.6 Evaluation of Secondary Criteria for Evaluating IEA Performance

The secondary criteria for evaluating IEA performance involved comparing agreement indices for IEA-human scoring for various demographic subgroups. Because of the importance of protecting personally identifiable information, student demographic data is stored and managed separately from the performance scoring data. For this reason, it was not possible to evaluate subgroup performance in real time as IEA was being trained.

For those prompts trained on early operational data, attempts were made to prioritize the data being returned from the field to include data from states or districts where more diverse populations of students were anticipated. In addition, requests for additional human scores were made to increase the likelihood that there would be sufficient numbers of responses with two human scores for most of the demographic subgroups of interest.

Once IEA was trained and deployed, scoring sets used in training were matched to demographic information so that agreement between IEA and human scorers could be evaluated across subgroups. The analysis was conducted for the ten comparison groups outlined in Table 4.6.

Table 4.6 Comparison Groups

Group Type	Comparison Groups
Sex	Female Male
Ethnicity	American Indian/Alaska Native Asian Black/African American Hispanic/Latino Native Hawaiian or Other Pacific Islander Two or More Races White
Special Instructional Needs	English Language Learners (ELL) Students with Disabilities (SWD) Economically Disadvantaged

IEA-human agreement indices were calculated for all cases with an IEA score and at least one human score. Human-human agreement was calculated for all cases with two human scores.

To evaluate the training of IEA for subgroups, the following criteria approved by the state leads for subgroups with at least 50 IEA-human scores and at least 50 human-human scores were applied:

- Pearson correlation between IEA-human should be within 0.1 of human-human.
- Kappa between IEA-human should be within 0.1 of human-human.
- Quadratic-weighted kappa between IEA-human should be within 0.1 of human-human.
- Exact agreement between IEA-human should be within 5.25 percent of human-human.
- Standardized mean difference between IEA-human should be less than ± 0.15 (this criterion was applied to subgroups with at least 50 IEA-human scores).

Although it was not expected that these criteria would be met for all subgroups for all prompts, if results of the evaluation between IEA and human scoring for subgroups for any prompt indicated that IEA performance persistently failed on the criteria listed above, consideration would be given to resetting the responses scored by IEA and reverting to human scoring until such time that an alternate IEA model could be established with improved subgroup performance.

In addition to the secondary criteria approved by the State Leads, the performance of IEA was compared to the following targets on the various measures for subgroups with at least 50 responses:

- Pearson correlation between IEA-human should be 0.70 or above.
- Kappa between IEA-human should be 0.40 or above.
- Quadratic-weighted kappa between IEA-human should be 0.70 or above.
- Exact agreement between IEA-human should be 65 percent or above.

These targets were not intended to be directly applied in decisions about whether to deploy IEA operationally or not. Such targets may or may not be met by human scoring for any particular prompt and/or subgroup, and

if they are not met by human scoring, they are unlikely to be met by IEA scoring. Nevertheless, comparisons to these targets provided additional information about IEA performance (and human scoring) in an absolute sense.

4.3.7 Inter-Rater Agreement for Prose Constructed Response

This section presents the inter-rater agreement for operational results for the online PCR tasks by trait and grade level. PCR items are scored on two traits: (1) Reading Comprehension and Written Expression and (2) Knowledge of Language and Conventions for Research Simulation for Literary Analysis Tasks and (1) Written Expression and (2) Knowledge of Language and Conventions for the Narrative Task.

For 10 percent of responses, a second “reliability” score was assigned. The purpose of the reliability score is to provide data for evaluating the consistency of scoring, which is done by evaluating scoring agreement. Inter-rater agreement is the agreement between the first and second scores assigned to student responses and is the measure of how often scorers agree with each other. 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. Inter-rater agreement expectations are provided in Table 4.5 in Section 4.2.4. For ELA/L PCR traits, the expectation for agreement is an inter-rater agreement of 65 percent or higher between two scorers. When IEA provided the first score of record, the second reliability score was a human score. For a subset of responses, the first and second score were both human scores.

Table 4.7 presents the average agreement across the PCRs for each grade level by trait. The number of prompts included in the analyses is listed for each grade level. The agreement indices (exact agreement, kappa, quadratic-weighted kappa, and Pearson correlation) were calculated separately by PCR for each trait (Reading Comprehension and Written Expression or Written Expression and Conventions). For each grade level, the agreement indices were averaged across the PCRs. Table 4.7 presents the average count and the average for the agreement indices.

The exact agreement for the PCR traits is above the criteria of a 65 percent agreement rate for all PCRs. The strength of agreement between raters is moderate to substantial agreement as defined by Landis and Koch (1977) for all PCRs. The quadratic-weighted kappa (QW Kappa) distinguishes between differences in ratings that are close to each other versus larger differences. The weighted kappa is substantial to almost perfect agreement for all grades. The Pearson correlations (r) ranged from 0.75 to 0.95.

Table 4.7 Prose Constructed Response Average Agreement Indices by Test

Test	Number of PCRs	Count	Written Expression				Conventions			
			Exact	Kappa	QW Kappa	r	Exact	Kappa	QW Kappa	r
ELA03	4	26875	74.43	0.57	0.74	0.75	75.1	0.6	0.78	0.78
ELA04	4	26078	73.63	0.61	0.81	0.81	73.08	0.6	0.8	0.8
ELA05	4	28513	73.45	0.6	0.81	0.81	73.55	0.6	0.81	0.82
ELA06	3	51516	75.93	0.65	0.86	0.86	76.5	0.66	0.85	0.85
ELA07	3	74197	74.17	0.64	0.87	0.87	73.83	0.64	0.86	0.86
ELA08	4	29737	75.3	0.66	0.9	0.9	75.03	0.66	0.88	0.88
ELA09	4	4954	81.95	0.75	0.91	0.91	80.95	0.74	0.9	0.9
ELA10	3	1155	89.5	0.85	0.95	0.95	88.83	0.84	0.94	0.94

Section 5: Classical Item Analysis

5.1 Overview

This section describes the results of the classical item analysis conducted for data obtained from the operational test items. All English language arts/literacy (ELA/L) and mathematics assessments were pre-equated. The item statistics provided in this section were from prior operational administrations and reflect the statistics that were used in test construction and for score reporting for some states and agencies. Item analysis serves two purposes: to inform item exclusion decisions for item response theory analysis and to provide item statistics for the item bank.

Item analysis included data from the following types of items: key-based selected-response items, rule-based machine-scored items, and handscored constructed-response items. For each item, the analysis produced item difficulty, item discrimination, and item response frequencies.

5.2 Data Screening Criteria

Item analyses were conducted by test form based on administration mode. In preparation for item analysis, student response files were processed to verify that the data were free of errors. Pearson Customer Data Quality staff ran predefined checks on all data files and verified that all fields and data needed to perform the statistical analyses were present and within expected ranges.

Before beginning item analysis, Pearson performed the following data screening operations:

1. All records with an invalid form number were excluded.
2. All records that were flagged as “void” were excluded.
3. All records where the student attempted fewer than 25 percent of items were excluded.
4. For students with more than one valid record, the record with the higher raw score was chosen.
5. Records for students with administration issues or anomalies were excluded.

5.3 Description of Classical Item Analysis Statistics

A set of classical item statistics were computed for each operational item by form and by administration mode. Each statistic was designed to evaluate the performance of each item.

The following statistics and associated flagging rules were used to identify items that were not performing as expected:

Classical Item Difficulty Indices (p-value and average item score)

When constructing tests, a wide range of item difficulties is desired (i.e., from easy to hard items) so that students of all ability levels can be assessed with precision. At the operational stage, item difficulty statistics are used by test developers to build forms that meet desired test difficulty targets.

For dichotomously scored items, item difficulty is indicated by its p-value, which is the proportion of students who answered that item correctly. The range for p-values is from .00 to 1.00. Items with high p-values are easy items and those with low p-values are difficult items. Dichotomously scored items were flagged for review if the p-value was above .95 (i.e., too easy) or below .25 (i.e., too difficult).

For polytomously scored items, difficulty is indicated by the average item score (AIS). The AIS can range from .00 to the maximum total possible points for an item. To facilitate interpretation, the AIS values for polytomously scored items are often expressed as percentages of the maximum possible score, which are equivalent to the p-values of dichotomously scored items. Polytomously scored items were flagged for review if the p-value was above .95 or below .25.

Percentage of Students Choosing Each Response Option

Selected-response items on the summative assessments refer primarily to single-select multiple-choice scored items. These items require that the student select a response from a number of answer options. These statistics for single-select multiple-choice items indicate the percentage of students who select each of the answer options and the percentage that omit the item. The percentages are also computed for the high-performing subgroup of students who scored at the top 20 percent on the assessment. An item was flagged for review if more high-performing students chose an incorrect option than the correct response. Such a result could indicate that the item has multiple correct answers or is miskeyed.

Item-Total Correlation

This statistic describes the relationship between students' performance on a specific item and their performance on the total test. The item-total correlation is usually referred to as the item discrimination index. For operational item analysis, the total score on the assessment was used as the total test score. The item-total correlation was calculated for both selected-response items and constructed-response items as an estimate of the correlation between an observed continuous variable and an unobserved continuous variable hypothesized to underlie the variable with ordered categories (Olsson et al., 1982). Item-total correlations can range from -1.00 to 1.00. Desired values are positive and larger than .15. Negative item-total correlations indicate that low-ability students perform better on an item than high-ability students, an indication that the item may be potentially flawed. Item-total correlations below .15 were flagged for review.

Distractor-Total Correlation

For selected-response items, this estimate describes the relationship between selecting an incorrect response (i.e., a distractor) for a specific item and performance on the total test. The item-total correlation is calculated for the distractors. Items with distractor-total correlations above .00 were flagged for review as these items may have multiple correct answers, be miskeyed, or have other content issues.

Percentage of Students Omitting or Not Reaching Each Item

For both selected-response and constructed-response items, this statistic is useful for identifying problems with test features such as testing time and item/test layout. Typically, if students have an adequate amount of testing time, approximately 95 percent of students should attempt to answer each question on the test. A distinction is made between "omit" and "not reached" for items without responses.

- An item is considered "omit" if the student responded to subsequent items.
- An item is considered "not reached" if the student did not respond to any subsequent items.

Patterns of high omit or not-reached rates for items located near the end of a test section may indicate that students did not have adequate time. Items with high omit rates were flagged. Omit rates for constructed-response items tend to be higher than for selected-response items. Therefore, the omit rate for flagging individual items was 5 percent for selected-response items and 15 percent for constructed-response items. If a student omitted an item, then the student received a score of 0 for that item and was included in the n-count for that item. However, if an item was near the end of the test and classified as not reached, the student did not receive a score and was not included in the n-count for that item.

Distribution of Item Scores

For constructed-response items, examination of the distribution of scores is helpful to identify how well the item is functioning. If no students' responses are assigned the highest possible score point, this may indicate that the item is not functioning as expected (e.g., the item could be confusing, poorly worded, or just unexpectedly difficult), the scoring rubric is flawed, and/or students did not have an opportunity to learn the content. In addition, if all or most students score at the extreme ends of the distribution (e.g., 0 and 2 for a 3-category item), this may indicate that there are problems with the item or the rubric so that students can receive either full credit or no credit at all, but not partial credit.

The raw score frequency distributions for constructed-response items were computed to identify items with few or no observations at any score points. Items with no observations or a low percentage (i.e., less than 3 percent) of students obtaining any score point were flagged. In addition, constructed-response items were flagged if they had U-shaped distributions, with high frequencies for extreme scores and very low frequencies for middle score categories.

5.4 Summary of Classical Item Analysis Flagging Criteria

In summary, items are flagged for review if the item analysis yielded any of the following results:

1. p-value above .95 for dichotomous items or polytomous items
2. p-value below .25 for dichotomous items or polytomous items
3. item-total correlation below .15
4. any distractor-total correlation above .00
5. greater number of high-performing students (top 20 percent) choosing a distractor rather than the keyed response
6. high percentage of omits: above 5 percent for selected-response items and above 15 percent for constructed-response items
7. high percentage that did not reach the item: above 5 percent for selected-response items and above 15 percent for constructed-response items
8. constructed-response items with a score value obtained by less than 3 percent of responses

The procedure was for Pearson's psychometric staff to review any flagged items and submit them to the Priority Alert Task Force to decide if the items were problematic and should be excluded from scoring.

5.5 Classical Item Analysis Results

This section presents tables summarizing the analyses for items on the spring operational forms. All assessments were pre-equated, meaning that the scoring was based on item parameters estimated using data from earlier administrations. Item analysis results in this section are the item statistics from prior administrations that were used to make decisions during the test construction process and for scoring.

- Table 5.1 presents pre-administration p-value information by grade for the ELA/L operational items.
- Table 5.2 presents pre-administration p-value information by grade/course for the mathematics operational items.
- Table 5.3 presents pre-administration item-total correlations by grade for the ELA/L operational items.
- Table 5.4 presents pre-administration item-total correlations by grade/course for the mathematics operational items.

An operational item may appear on multiple test forms. The tables list unique item counts for an assessment and the reported item statistics may be based on student responses across multiple occurrences of an item.

Spoiled or “do not score” items were excluded from the total test score in item analysis. These items were removed from scoring because of item performance, technical scoring issues, content concerns, or multiple/no correct answers.

Some forms in the spring 2022 administration were based on previous administrations, with many of them being reused from the spring 2021 administration; therefore, the item analyses for these forms were reported in the associated technical reports.

Table 5.1 Summary of Pre-Administration p-Values for ELA/L Operational Items by Grade and Mode

Grade	N of Unique Items	Mean p-Value	SD p-Value	Min. p-Value	Max. p-Value	Median p-Value
3	44	0.47	0.18	0.20	0.78	0.44
4	66	0.46	0.15	0.18	0.82	0.46
5	62	0.45	0.15	0.17	0.84	0.42
6	50	0.48	0.14	0.18	0.75	0.47
7	55	0.45	0.14	0.23	0.82	0.41
8	64	0.47	0.12	0.22	0.80	0.45
9	66	0.46	0.14	0.14	0.83	0.43
10	52	0.44	0.13	0.19	0.81	0.44

Note. SD = standard deviation.

Table 5.2 Summary of p-Values for Mathematics Operational Items by Grade and Mode

Grade	N of Unique Items	Mean p-Value	SD p-Value	Min. p-Value	Max. p-Value	Median p-Value
3	90	0.59	0.21	0.18	0.95	0.63
4	88	0.53	0.19	0.23	0.94	0.51
5	91	0.47	0.20	0.14	0.88	0.46
6	76	0.40	0.17	0.11	0.85	0.37
7	81	0.41	0.18	0.11	0.91	0.36
8	78	0.33	0.16	0.07	0.72	0.30
A1	72	0.32	0.18	0.05	0.71	0.27
GO	81	0.34	0.20	0.06	0.88	0.34
A2	79	0.32	0.17	0.05	0.73	0.29

Note. SD = standard deviation; A1 = Algebra I; GO = Geometry; A2 = Algebra II.

Table 5.3 Summary of Pre-Administration Item-Total Polyserial Correlations for ELA/L Operational Items by Grade and Mode

Grade	N of Unique Items	Mean Polyserial	SD Polyserial	Min. Polyserial	Max. Polyserial	Median Polyserial
3	44	0.52	0.13	0.23	0.79	0.52
4	66	0.49	0.15	0.24	0.82	0.46
5	62	0.50	0.16	0.19	0.86	0.47
6	50	0.52	0.15	0.31	0.87	0.50
7	55	0.51	0.18	0.19	0.86	0.45
8	64	0.50	0.17	0.13	0.87	0.49
9	66	0.50	0.18	0.24	0.86	0.47
10	52	0.50	0.19	0.20	0.87	0.44

Note. SD = standard deviation.

Table 5.4 Summary of Item-Total Correlations for Mathematics Operational Items by Grade and Mode

Grade/ Course	N of Unique Items	Mean Polyserial	SD Polyserial	Min. Polyserial	Max. Polyserial	Median Polyserial
3	90	0.51	0.14	0.19	0.79	0.52
4	88	0.53	0.12	0.26	0.80	0.54
5	91	0.50	0.15	0.16	0.85	0.51
6	76	0.53	0.14	0.24	0.82	0.55
7	81	0.50	0.16	0.18	0.84	0.50
8	78	0.46	0.16	0.19	0.86	0.46
A1	72	0.50	0.16	0.18	0.96	0.49
GO	81	0.49	0.16	0.19	0.94	0.48
A2	79	0.46	0.14	0.19	0.73	0.46

Note. SD = standard deviation; A1 = Algebra I; GO = Geometry; A2 = Algebra II.

Section 6: Differential Item Functioning

6.1 Overview

Differential item functioning (DIF) analyses were conducted using the data obtained from the operational items. If an item performs differentially across identifiable subgroups (e.g., gender or ethnicity) when students are matched on ability, the item may be measuring something other than the intended construct (i.e., possible evidence of DIF). It is important, however, to recognize that item performance differences flagged for DIF might be related to actual differences in relevant knowledge or skills (item impact) or statistical Type I error. As a result, DIF statistics are used to identify *potential* item bias. Subsequent reviews by content experts and bias/sensitivity committees are required to determine the source and meaning of performance differences.

In this section, the DIF statistics used at test construction to make decisions about items are provided for all mathematics online and paper and English language arts/literacy (ELA/L) tests. In addition, DIF statistics are presented for the ELA/L online post-equated tests.

6.2 DIF Procedures

Dichotomous Items

The Mantel-Haenszel (MH) DIF statistic was calculated for selected-response items and for dichotomously scored constructed-response items. In this method, students are classified into relevant subgroups of interest (e.g., gender or ethnicity). Using the raw score total as the criteria, students in a certain total score category in the focal group (e.g., females) are compared with students in the same total score category in the reference group (e.g., males). For each item, students in the focal group are also compared to students in the reference group who performed equally well on the test as a whole. The common odds ratio is estimated across all categories of matched student ability using the following formula (Dorans & Holland, 1993), and the resulting estimate is interpreted as the relative likelihood of success on a particular item for members of two groups when matched on ability:

$$\hat{\alpha}_{MH} = \frac{\sum_{s=1}^S \frac{R_{rs} W_{fs}}{N_{ts}}}{\sum_{s=1}^S \frac{R_{fs} W_{rs}}{N_{ts}}}, \quad (6-1)$$

in which:

S = the number of score categories,

R_{rs} = the number of students in the reference group who answer the item correctly,

W_{fs} = the number of students in the focal group who answer the item incorrectly,

R_{fs} = the number of students in the focal group who answer the item correctly,

W_{rs} = the number of students in the reference group who answer the item incorrectly, and

N_{ts} = the total number of students.

To facilitate the interpretation of MH results, the common odds ratio is frequently transformed to the delta scale using the following formula (Holland & Thayer, 1988):

$$MH\ D-DIF = -2.35 \ln(\hat{\alpha}_{MH}) \quad (6-2)$$

Positive values indicate DIF in favor of the focal group (i.e., positive DIF items are differentially easier for the focal group), whereas negative values indicate DIF in favor of the reference group (i.e., negative DIF items are differentially easier for the reference group).

Polytomous Items

For polytomously scored constructed-response items, the MH D-DIF statistic is not calculated; instead, the standardization DIF (Dorans, 2013; Dorans & Schmitt, 1991; Zwick et al., 1997), in conjunction with the Mantel chi-square statistic (Mantel, 1963; Mantel & Haenszel, 1959), is used to identify items with DIF.

The standardization DIF compares the item means of the two groups after adjusting for differences in the distribution of students across the values of the matching variable (i.e., total test score) and is calculated using the following formula:

$$STD-EISDIF = \frac{\sum_{s=1}^S N_{fs} \times E_f(Y | X = s)}{\sum_{s=1}^S N_{fs}} - \frac{\sum_{s=1}^S N_{fs} \times E_r(Y | X = s)}{\sum_{s=1}^S N_{fs}}, \quad (6-3)$$

in which:

X = the total score,

Y = the item score,

S = the number of score categories,

N_{fs} = the number of students in the focal group in score category s ,

E_r = the expected item score for the reference group, and

E_f = the expected item score for the focal group.

A positive *STD-EISDIF* value means that, conditional on the total test score, the focal group has a higher mean item score than the reference group. In contrast, a negative *STD-EISDIF* value means that, conditional on the total test score, the focal group has a lower mean item score than the reference group.

Classification

Based on the DIF statistics and significance tests, items are classified into three categories and assigned values of A, B, or C (Zieky, 1993). Category A items contain negligible DIF, Category B items exhibit slight-to-moderate DIF, and Category C items possess moderate-to-large DIF values. Positive values indicate that, conditional on the total score, the focal group has a higher mean item score than the reference group. In contrast, negative DIF values indicate that, conditional on the total test score, the focal group has a lower mean item score than the reference group. The flagging criteria for dichotomously scored items are presented in Table 6.1; the flagging criteria for polytomously scored constructed-response items are provided in Table 6.2.

Table 6.1 DIF Categories for Dichotomous Selected-Response and Constructed-Response Items

DIF Category	Criteria
A (negligible)	Absolute value of the MH D-DIF is not significantly different from zero, or is less than one.
B (slight to moderate)	1. Absolute value of the MH D-DIF is significantly different from zero but not from one, and is at least one; or 2. Absolute value of the MH D-DIF is significantly different from one, but is less than 1.5. Positive values are classified as “B+” and negative values as “B-.”
C (moderate to large)	Absolute value of the MH D-DIF is significantly different from one, and is at least 1.5. Positive values are classified as “C+” and negative values as “C-.”

Table 6.2 DIF Categories for Polytomous Constructed-Response Items

DIF Category	Criteria
A (negligible)	Mantel Chi-square p-value > 0.05 or $ STD-EISDIF/SD \leq 0.17$
B (slight to moderate)	Mantel Chi-square p-value < 0.05 and $ STD-EISDIF/SD > 0.17$
C (moderate to large)	Mantel Chi-square p-value < 0.05 and $ STD-EISDIF/SD > 0.25$

Note. $STD-EISDIF$ = standardized DIF; SD = total group standard deviation of item score.

6.3 Operational Analysis DIF Comparison Groups

DIF analyses were conducted on each test form for designated comparison groups defined on the basis of demographic variables including gender, race/ethnicity, economic disadvantage, and special instructional needs such as students with disabilities or English learners. Student demographic information was provided by the states, including the District of Columbia, and captured in PearsonAccess^{next} by means of a student data upload. The demographic data was verified by the states prior to score reporting. These comparison groups are specified in Table 6.3.

Table 6.3 Traditional DIF Comparison Groups

Grouping Variable	Focal Group	Reference Group
Gender	Female	Male
Ethnicity	American Indian/Alaska Native (AmerIndian)	White
	Asian	White
	Black or African American	White
	Hispanic/Latino	White
	Native Hawaiian or Pacific Islander	White
	Multiple Race Selected	White
Economic Status*	Economically Disadvantaged (EcnDis)	Not Economically Disadvantaged (NoEcnDis)
Special Instructional Needs	English Learner (ELY)	Non-English Learner (ELN)
	Students with Disabilities (SWDY)	Students without Disabilities (SWDN)

*Economic status was based on participation in National School Lunch Program (receipt of free or reduced-price lunch).

DIF analyses were conducted when the following sample size requirements were met:

- the smaller group, reference or focal, had at least 100 students, and

- the combined group, reference and focal, had at least 400 students.

6.4 Operational Differential Item Functioning Results

Appendix 6 presents tables summarizing the DIF results for the spring pre-administration item DIF results that were used to inform decisions at test construction for both ELA/L and mathematics, as well as the post-administration item DIF results for ELA/L. There is one table prepared for each content and grade level (e.g., ELA/L grade 3).

Spoiled or “do not score” items were excluded from the total test score for each form in DIF analysis. These items were removed from scoring because of item performance, technical scoring issues, content concerns, multiple correct answers, or no correct answers. However, the tables in this section may include items for certain grade levels that were excluded from scoring based on later analyses.

In the DIF results tables, the column “DIF Comparisons” identifies the focal and reference groups for the analysis performed; “Total N of Unique Items” reports the number of unique items included in the analysis. Because DIF analysis is conducted at the parent level for prose constructed responses in ELA/L tests, the total number of unique items reported in the DIF analysis is smaller than the total number of items reported in the classical item analysis (see Tables 5.1 and 5.2) and the item response theory summary statistics (see Tables 7.1 and 7.2 for each ELA/L test. In addition, “0” indicates that the DIF analysis did not classify any items in the particular DIF category, while “n/a” indicates that the DIF analysis was not performed due to insufficient sample sizes.

Table 6.4 Pre-Administration Differential Item Functioning for ELA/L Grade 3

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	39			1	3	38	97				
White vs. Black	39					39	100				
White vs. Hispanic	39			2	5	37	95				
White vs. Asian	39					38	97	1	3		
White vs. AmerIndian	39					39	100				
White vs. Pacific Islander	39			2	5	36	92	1	3		
White vs. Multiracial	39					39	100				
NoEcnDis vs. EcnDis	39					39	100				
ELN vs. ELY	39			4	10	35	90				
SWDN vs. SWDY	39			1	3	38	97				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table 6.5 Pre-Administration Differential Item Functioning for Mathematics Grade 3

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	89			3	3	85	96	1	1		
White vs. Black	89			4	4	84	94	1	1		
White vs. Hispanic	89					89	100				
White vs. Asian	89					85	96	3	3	1	1
White vs. AmerIndian	89			2	2	87	98				
White vs. Pacific Islander	89			1	1	88	99				
White vs. Multiracial	89			1	1	87	98	1	1		
NoEcnDis vs. EcnDis	89					89	100				
ELN vs. ELY	89			1	1	88	99				
SWDN vs. SWDY	89			1	1	88	99				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Section 7: IRT Model and Parameters

7.1 Overview

Multiple operational core forms were administered for each grade in English language arts/literacy (ELA/L) and mathematics assessments. All tests in spring 2022 were pre-equated. This section describes the item response theory (IRT) model used in this assessment program and provides descriptive statistics of the item parameters.

7.2 Two-Parameter Logistic/Generalized Partial Credit Model

The operational items used pre-equated parameters in the context of the 2PL/GPC model, which is denoted as

$$p_{im}(\theta_j) = \frac{\exp\left[\sum_{k=0}^m Da_i(\theta_j - b_i + d_{ik})\right]}{\sum_{v=0}^{M_i-1} \exp\left[\sum_{k=0}^v Da_i(\theta_j - b_i + d_{ik})\right]} \quad (7-1)$$

where $a_i(\theta_j - b_i + d_{i0}) \equiv 0$; $p_{im}(\theta_j)$ is the probability of a student with θ_j getting score m on item i ; D is the IRT scale constant (1.7); a_i is the discrimination parameter of item i ; b_i is the item difficulty parameter of item i ; d_{ik} is the k^{th} step deviation value for item i ; M_i is the number of score categories of item i with possible item scores as consecutive integers from zero to $M_i - 1$; and v indexes the response categories and is iterated from 0 to $M_i - 1$.

7.3 Summary Statistics and Distributions from IRT Analyses

Tables 7.1 through 7.4 present summary statistics for the IRT (b - and a -) parameter estimates, the standard errors of the parameter estimates, and the IRT model fit values (chi-square and adjusted fit) for ELA/L and mathematics assessments. The summary statistics for IRT parameter estimates include all the items administered in the spring administration except the items on the reused forms, if applicable, for which the summary results were reported in the technical reports of the source administrations.

The information is provided by content area (ELA/L and mathematics) for all items at each grade level or course. The summary statistics shown include the total number of items and score points, along with the mean, standard deviation (SD), minimum, and maximum.

7.3.1 IRT Summary Statistics for English Language Arts/Literacy

Table 7.1 shows the pre-equated b - and a -parameter estimates for all ELA/L assessments. Table 7.2 shows the source year for the item statistics for each of the ELA/L assessments that were pre-equated. IRT summary statistics are provided in Appendix 7 for ELA/L for all items, Reading claim items, and Writing claim items.

Table 7.1 Pre-Equated IRT Parameter Estimates Summary for All Items for ELA/L by Grade

Grade	No. of Score Points	No. of Items	<i>b</i> Estimates Summary				<i>a</i> Estimates Summary			
			Mean	SD	Min.	Max.	Mean	SD	Min.	Max.
3	98	44	0.37	0.98	-1.64	2.40	0.57	0.22	0.19	1.04
4	148	66	0.40	0.92	-1.99	2.66	0.47	0.23	0.17	0.99
5	141	62	0.55	0.96	-1.34	3.59	0.48	0.23	0.13	0.99
6	112	50	0.34	0.82	-1.00	2.95	0.52	0.21	0.18	1.16
7	125	55	0.37	0.77	-1.21	1.77	0.50	0.27	0.13	1.23
8	146	64	0.31	0.78	-1.42	2.83	0.47	0.24	0.08	1.06
9	150	66	0.55	0.88	-1.21	2.77	0.50	0.29	0.17	1.22
10	119	52	0.68	0.78	-1.08	2.67	0.49	0.29	0.13	1.18
11	208	92	1.00	0.98	-1.09	5.29	0.47	0.25	0.09	1.13

Note. SD = standard deviation.

Table 7.2 Pre-Equated IRT Parameter Distribution by Year for All Items for ELA/L by Grade

Grade	No. of Items	2014	2015	2016	2017	2018	2019
3	44	0	0	0	3	6	35
4	66	0	4	0	14	3	45
5	62	0	0	0	3	22	37
6	50	0	0	0	11	13	26
7	55	0	6	10	14	15	10
8	64	0	4	1	15	6	38
9	66	0	4	8	17	21	16
10	52	0	4	0	14	2	32

7.3.2 IRT Summary Statistics for Mathematics

Table 7.3 shows the *b*- and *a*-parameter estimates for the mathematics assessments. Table 7.4 shows the source year for the item statistics for each of the assessments. IRT summary statistics are provided in Appendix 7 for mathematics for all items, single-select multiple-choice items, constructed-response items, and subclaims.

Table 7.3 Pre-Equated IRT Parameter Estimates Summary for All Items for Mathematics by Grade/Course

Grade	No. of Score Points	No. of Items	<i>b</i> Estimates Summary				<i>a</i> Estimates Summary			
			Mean	SD	Min.	Max.	Mean	SD	Min.	Max.
3	142	90	-0.38	1.14	-3.29	2.48	0.76	0.26	0.26	1.35
4	148	88	-0.11	0.97	-2.65	1.50	0.74	0.18	0.32	1.35
5	152	91	0.08	1.10	-2.33	2.13	0.65	0.23	0.16	1.50
6	129	76	0.47	0.82	-1.88	2.05	0.70	0.26	0.28	1.34
7	128	81	0.60	1.02	-1.78	4.22	0.66	0.26	0.19	1.49
8	126	78	1.00	0.94	-1.42	2.89	0.58	0.27	0.16	1.40
A1	137	72	1.22	1.11	-1.16	3.57	0.65	0.29	0.17	1.62
G1	144	81	0.95	1.05	-1.60	3.83	0.75	0.38	0.19	1.76
A2	158	79	1.29	1.10	-1.39	3.90	0.61	0.26	0.19	1.20

Note. SD = standard deviation; A1 = Algebra I; GO = Geometry; A2 = Algebra II.

Table 7.4 Pre-Equated IRT Parameter Distribution by Year for All Items for Mathematics by Grade/Course

Grade	ALL	2014	2015	2016	2017	2018	2019
3	90	0	17	17	5	12	39
4	88	0	14	19	11	5	39
5	91	1	20	9	20	7	34
6	76	0	16	14	12	6	28
7	81	0	12	11	8	12	38
8	78	0	16	7	17	4	34
A1	72	0	9	11	7	12	33
G1	81	0	16	15	15	10	25
A2	79	0	14	18	14	10	23

Note. A1 = Algebra I; GO = Geometry; A2 = Algebra II.

Section 8: Performance Level Setting

8.1 Performance Standards

Performance standards relate levels of performance on an assessment directly to what students are expected to learn. This is done by establishing threshold scores that distinguish between performance levels. Performance level setting (PLS) is the process of establishing the threshold scores that define the performance levels for an assessment.

8.2 Performance Levels and Policy Definitions

For the summative assessments, the performance levels are

- Level 5: Exceeded expectations
- Level 4: Met expectations
- Level 3: Approached expectations
- Level 2: Partially met expectations
- Level 1: Did not yet meet expectations

More detailed descriptions of each performance level, known as policy definitions, are provided in the following subsections.

Level 5: Exceeded Expectations

Students performing at this level **exceed academic expectations** for the knowledge, skills, and practices contained in the standards assessed at their grade level or course.

Grades 3–10: Students performing at this level **exceed academic expectations** for the knowledge, skills, and practices contained in the standards for English language arts/literacy (ELA/L) or mathematics assessed at their grade level. They are **academically well prepared** to engage successfully in further studies in this content area.

Algebra II, Integrated Mathematics III, and ELA/L Grade 11: Students performing at this level **exceed academic expectations** for the knowledge, skills, and practices contained in the mathematics and ELA/L standards assessed at grade 11. They are very likely to engage successfully in entry-level, credit-bearing courses in mathematics and ELA/L, as well as technical courses requiring an equivalent command of the content area. Students performing at this level are exempt from having to take and pass placement tests in two- and four-year public institutions of higher education designed to determine whether they are academically prepared for such courses without need for remediation.

Level 4: Met Expectations

Students performing at this level **meet academic expectations** for the knowledge, skills, and practices contained in the standards assessed at their grade level or course.

Grades 3–10: Students performing at this level **meet academic expectations** for the knowledge, skills, and practices contained in the standards for ELA/L or mathematics assessed at their grade level. They are **academically prepared** to engage successfully in further studies in this content area.

Algebra II, Integrated Mathematics III, and ELA/L Grade 11: Students performing at this level **meet academic expectations** for the knowledge, skills, and practices contained in mathematics and ELA/L at grade 11. They are very likely to engage successfully in entry-level, credit-bearing courses in mathematics and ELA/L, as well as technical courses requiring an equivalent command of the content area. Students performing at this level are exempt from having to take and pass placement tests in two- and four-year public institutions of higher education designed to determine whether they are academically prepared for such courses without need for remediation.

Level 3: Approached Expectations

Students performing at this level **approach academic expectations** for the knowledge, skills, and practices contained in the standards assessed at their grade level or course.

Grades 3–10: Students performing at this level **approach academic expectations** for the knowledge, skills, and practices contained in the standards for ELA/L or mathematics assessed at their grade level. They are likely prepared to engage successfully in further studies in this content area.

Algebra II, Integrated Mathematics III, and ELA/L Grade 11: Students performing at this level **approach academic expectations** for the knowledge, skills, and practices contained in the ELA/L and mathematics standards assessed at grade 11. They are likely to engage successfully in entry-level, credit-bearing courses in mathematics and ELA/L, as well as technical courses requiring an equivalent command of the content area. **Students performing at Level 3 are strongly encouraged to continue to take challenging high school coursework in English and mathematics through graduation.** Postsecondary institutions are encouraged to use additional information about students performing at Level 3, such as course completion, course grades, and scores on other assessments to determine whether to place them directly into entry-level courses.

Level 2: Partially Met Expectations

Students performing at this level **partially meet academic expectations** for the knowledge, skills, and practices contained in the standards assessed at their grade level or course.

Grades 3–10: Students performing at this level **partially meet academic expectations** for the knowledge, skills, and practices contained in the standards for ELA/L or mathematics assessed at their grade level. They will likely need academic support to engage successfully in further studies in this content area.

Algebra II, Integrated Mathematics III, and ELA/L Grade 11: Students performing at this level **partially meet academic expectations** for the knowledge, skills, and practices contained in the ELA/L and mathematics standards assessed at grade 11. They will likely need academic support to engage successfully in entry-level, credit-bearing courses, and technical courses requiring an equivalent command of the content area. Students performing at this level are not exempt from having to take and pass placement tests designed to determine whether they are academically prepared for such courses without the need for remediation in two- and four-year public institutions of higher education.

Level 1: Did Not Yet Meet Expectations

Students performing at this level **do not yet meet academic expectations** for the knowledge, skills, and practices contained in the standards assessed at their grade level or course.

Grades 3–10: Students performing at this level **do not yet meet academic expectations** for the knowledge, skills, and practices contained in the standards for ELA/L or mathematics assessed at their grade level. They will need academic support to engage successfully in further studies in this content area.

Algebra II, Integrated Mathematics III, and ELA/L Grade 11: Students performing at this level **do not yet meet academic expectations** for the knowledge, skills, and practices contained in the ELA/L and mathematics standards assessed at grade 11. They will need academic support to engage successfully in entry-level, credit-bearing courses in college algebra, introductory college statistics, and technical courses requiring an equivalent level of mathematics. Students performing at this level are not exempt from having to take and pass placement tests in two- and four-year public institutions of higher education designed to determine whether they are academically prepared for such courses without need for remediation.

8.3 Performance Level Setting Process for the Assessment System

One of the main objectives of the assessment system is to provide information to students, parents, educators, and administrators as to whether students are on track in their learning for success after high school, defined as college- and career-readiness. To set performance levels associated with this objective, participating states and agencies used the evidence-based standard setting (EBSS) method (Beimers et al., 2012) for the PLS process. The EBSS method is a systematic method for combining various considerations into the process for setting performance levels, including policy considerations, content standards, educator judgment about what students should know and be able to demonstrate, and research to support policy goals related to college- and career-readiness. A defined multistep process was used to allow a diverse set of stakeholders to consider the interaction of these elements in recommending performance level threshold scores for each assessment.

The seven steps of the EBSS process that were followed to establish performance standards for the summative assessments are as follows:

- Step 1: Define outcomes of interest and policy goals.
- Step 2: Develop research, data collection, and analysis plans.
- Step 3: Synthesize the research results.
- Step 4: Conduct pre-policy meeting.
- Step 5: PLS meetings with panels.
- Step 6: Conduct reasonableness review with post-policy panel.
- Step 7: Continue to gather evidence in support of standards.

A summary of key components within these steps is provided below. Additional detail about each step in the PLS process is provided in the *Performance Level Setting Technical Report*.

8.3.1 Research Studies

Participating states and agencies conducted two research studies in support of their policy goals—the benchmarking study and the postsecondary educators’ judgment (PEJ) study. The benchmarking study included a review of the literature relative to college- and career-readiness as well as consideration of the percentage of students obtaining a level equivalent to college- and career-readiness on a set of external assessments (e.g., ACT, SAT, NAEP). The PEJ study involved a group of nearly 200 college faculty reviewing items on the Algebra II and ELA/L grade 11 assessments and making judgments about the level of performance needed on each item to be academically ready for an entry-level college-credit bearing course in mathematics or ELA/L. Additional detail² about the benchmarking study can be found in the *Performance Level Setting Technical Report* as well as in the *PARCC Benchmarking Study Report*. Additional detail about the PEJ study can be found in the *Performance Level Setting Technical Report* as well as in the *Postsecondary Educators’ Judgment Study Final Report*.

8.3.2 Pre-Policy Meeting

Prior to the PLS meetings, a pre-policy meeting was convened to determine reasonable ranges that would be shown to panelists during the high school PLS meetings. Pre-policy meeting participants included representatives from both K–12 and higher education who served in roles such as commissioner/superintendent, deputy/assistant commissioner, state board member, director of assessment, director of academic affairs, senior policy associate, and so on. The reasonable ranges recommended by the pre-policy meeting defined the minimum and maximum percentage of students that would be expected to be classified as college- and career-ready. The pre-policy meeting participants reviewed the test purpose, how the performance standards would be used, and the results of the research studies to provide the recommendations for the reasonable ranges without viewing any student performance data.

8.3.3 Performance Level Setting Meetings

The task of the PLS committee was to recommend four threshold scores that would define the five performance levels for each assessment. Participating states and agencies solicited nominations from all states that had administered the assessments in 2014–2015 for panelists to serve on the PLS committees. Nominations were solicited both from state departments of public education (K–12) and higher education (primarily for participation on the high school panels). When selecting panelists, an emphasis was placed on those educators who had content knowledge as well as experience with a variety of student groups and attempted to balance the panels in terms of state representation.

Participating states and agencies used an extended modified Angoff (Yes/No) method to collect educator judgments on the items. This method asked panelists to review each item on a reference form of the assessment and to make the following judgment:

How many points would a borderline student at each performance level likely earn if they answered the question?

² More information is available online from <https://resources.newmeridiancorp.org/research/>.

This extension to the Yes/No standard setting method (Plake et al., 2005) allowed for incorporation of the multipoint items by asking educators to evaluate (Yes or No) whether a borderline student would earn the maximum number of points on an item, a lesser number of points on an item, or no points on the item. In the case of a single-point or multiple-choice item, this task simplifies to the standard Yes/No method.

After receiving training on the PLS procedure, panelists participated in three rounds of judgments for each assessment. Within each round, panelists were asked to consider the items in the test form, starting with the performance-based assessment component and then the end-of-year component. Each panelist made a judgment for the Level 2 performance level, followed by judgments for the Level 3 performance level, the Level 4 performance level, and the Level 5 performance level, in that order. The panelists entered their item judgments for each round by completing an online item judgment survey. Educator judgments were summed across items to create an estimated total score on the reference form for each performance level threshold. Feedback data relative to panelist agreement, student performance on the items, and student performance on the test as a whole were provided in between each of the three rounds of judgment. Panelists were shown the pre-policy reasonable ranges prior to making their Round 1 judgments and again as feedback data following each round of judgment.

A dry run of the PLS meeting process was held for grade 11 ELA/L and Algebra II in order to evaluate the implementation of the PLS method with the innovative characteristics of the summative assessments. These content areas were selected because they combined all the various aspects of the assessments, including the various types of items, scoring rules, and performance level decisions. The dry-run PLS meetings provided the opportunity to implement and evaluate multiple aspects of the operational plan for the actual PLS meeting, including prework, meeting materials, data analysis and feedback, and staff and panelist functions. The results of the dry-run PLS meeting were used to implement improvements in the process for the operational PLS meetings. Additional information about the methods and results of the dry-run PLS meeting is available in the full report in the *Performance Level Setting Dry-Run Meeting Report*.

The PLS meetings for the summative assessments were conducted during three one-week sessions. The dates of the twelve PLS committee meetings that were conducted are shown in Table 8.1.

Additional information about the methods and results of the PLS meetings is available in the *Performance Level Setting Technical Report*.

8.3.4 Post-Policy Reasonableness Review

Performance standards for all summative assessments were recommended by PLS committees and reviewed by the Governing Board and (for the Algebra II, Integrated Mathematics III, and ELA/L grade 11 assessments) the Advisory Committee on College Readiness as part of a post-policy reasonableness review. This group reviewed both the median threshold score recommendations from each committee and the variability in the threshold scores as represented by the standard error of judgment of the committee. Adjustments to the median threshold scores that were within two standard errors of judgment were considered to be consistent with the PLS panels' recommendation.

Table 8.1 Performance Level Setting Committee Meetings and Dates

Dates	Committees by Subjects and Grades
July 27–31, 2015	Algebra I/Integrated Mathematics I
	Geometry/Integrated Mathematics II
	Algebra II/Integrated Mathematics III
	Grade 9 English Language Arts/Literacy
	Grade 10 English Language Arts/Literacy
	Grade 11 English Language Arts/Literacy
August 17–21, 2015	Grades 7 & 8 Mathematics
	Grades 7 & 8 English Language Arts/Literacy
August 24–28, 2015	Grades 3 & 4 Mathematics
	Grades 5 & 6 Mathematics
	Grades 3 & 4 English Language Arts/Literacy
	Grades 5 & 6 English Language Arts/Literacy

In addition to voting to adopt the performance standards based on the committees' recommendations, this group also voted to conduct a shift in the performance levels to better meet the intended inferences about student performance. Holding the college- and career-ready (or on track) expectations (i.e., the current Level 4) constant, performance levels above this expectation were combined and performance levels below this expectation were expanded to create the final system of performance levels with three below and two above the college- and career-ready (or on track) expectation. The shift in performance levels was accomplished using a scale anchoring process that involved two primary steps. In the first step, the top two performance levels, above college- and career-ready (or on track), were combined into a single performance level and an additional performance level below college- and career-ready (or on track) was created by empirically determining the midpoint between the existing two levels. In the second step, the performance level descriptors (PLDs) were updated using items that discriminated student performance well at this level to create a PLD aligned with the new empirically determined performance level. At this same time, PLDs for all performance levels were reviewed for consistency and continuity. Members of the original PLS committees were recruited to participate in this process. Additional information about this process can be found in the *Performance Level Setting Technical Report*.

Section 9: Quality Control Procedures

Quality control in a testing program is a comprehensive and ongoing process. This section describes procedures put into place to monitor the quality of the item bank, test form, and ancillary material development. The quality checks for scanning, image editing, scoring, and data screening during psychometric analyses are also outlined. Additional quality information can be found in the Program Quality Plan document.

9.1 Quality Control of the Item Bank

The summative item bank consists of test passages and items, their associated metadata, and status (e.g., operational-ready, field-test ready, released, etc.). The items on the assessments were developed by Pearson and West Ed and put in the item bank once created.

The Pearson Assessment Banking for Building and Interoperability (ABBI) bank houses the passages and items, art, associated metadata, rubrics, alternate text for use on accommodated forms, and text complexity documentation. It provides an item previewer that allows items to be viewed and interacted with in the same way students see and interact with items and tools, and manages versioning of items with a date/time stamp. It allows reviewers to vote on item acceptance, and to record and retain their review notes for later reconciliation and reference. Item and passage review committee participants conducted their review in the item banking system. The committee members viewed the items as the student would, and could vote to alter the item, accept or reject the item, and record their comments in the system. After each meeting, reports were forwarded to New Meridian. The reports were generated by the item banking system and summarized feedback from the committee reviewers.

All new development for the summative assessments is being created within the ABBI system, which employs templates to control the consistency of the underlying scoring logic and Question and Test Interoperability creation for each item type. The ABBI system incorporates a previewer that allows the reviewers to validate the content of the item and validate the expected scoring of tasks. It supports the full range of review activities, including content review, bias and sensitivity review, expert editorial review, data review, and test construction review. It provides insight into the item edit process through versioning. A series of metadata validations at key points in the development cycle provides support for metadata consistency. The bank can be queried on the full range of metadata values to support bank analysis.

9.2 Quality Control of Test Form Development

Test forms were built based upon targets and the established blueprints set. The construction process started with specification and requirement capture to create the test specification document. From there items were pulled into forms based on the criteria approved in the test specifications document. After forms composition, the forms went through a review process that involved groups from New Meridian, Pearson, and participating states. Quality control steps were conducted on the items and forms that evaluated several item characteristics (e.g., content accuracy, completeness, style guide conformity, tools function). Revisions were incorporated into the forms before final review and approval. Section 2.2 provides more details on the form development process.

The forms quality assurance was performed by Pearson's Assessment and Information Quality (AIQ) organization. AIQ completed a comprehensive review of all *online* forms for the administration cycle. This

group is part of Pearson's larger Organizational Quality group and operates exclusively to validate form operability. The group verifies that the functionality of every online form is working to specifications. The overall functionality and maneuverability of each form is checked, and the behavior of each item within the form is verified. (Quality processes for paper forms are described in Section 9.3.)

The items within each form were tested to verify that they operated as expected for students. As a further aspect of the testing process, AIQ confirmed that forms were loaded correctly and that the audio was correct when compared to text. Sections and overviews were reviewed. Technology-enhanced items also were tested as an additional measure. As enumerated in the *Technology Guidelines for Assessments*, user interfaces were compatible with a range of common computer devices, operating systems, and browsers.

Pearson also performed quality control tests to verify that a standard set of responses was outputted to XML as expected after the final version of the form was approved. These responses were based on the keys provided in the test map or a standard open-ended responses string that contained a valid range of characters. As part of these tests, the test maps also were validated against the form layout and item types for correctness.

Pearson conducted a multifaceted validation of all item layout, rendering, and functionality. Reviewers conducted comparisons between the approved item and the item as it appeared in the field-test form or how it previously appeared; verified that tools and functions in the test delivery system, TestNav, were accurately applied; and verified that the style and layout met all requirements. In addition, answer keys were validated through a formal key review process. More details on the test development procedures are provided in Section 2.

9.3 Quality Control of Test Materials

Pearson provided high-quality materials in a timely and efficient manner to meet the test administration needs. Since the majority of printing work was done in-house, it was possible to fully control the production environment, press schedule, and quality process for print materials. Additionally, strict security requirements were employed to protect secure materials production; Section 3 provides details on the secure handling of test materials. Materials were produced according to the style guide and to the detailed specifications supplied in the materials list.

Pearson Print Service operates within the sanctions of an ISO 9001:2008 Quality Management System, and practices process improvement through Lean principles and employee involvement.

Raw materials (paper and ink) used for scannable forms production were manufactured exclusively for Pearson Print Service using specifications created by Pearson Print Service. Samples of ink and paper were tested by Pearson prior to use in production. Project specialists were the point of contact for incoming production.

Purchase orders and other order information were assessed against manufacturing capabilities and assigned to the optimal production methodology. Expectations, quality requirements, and cost considerations were foremost in these decisions. Prior to release for manufacture, order information was checked against specifications, technical requirements, and other communication that includes expected outcomes. Records of these checks were maintained.

Files for image creation flow through one of two file preparation functions: digital pre-press for digital print methodology, or plateroom for offset print methodology. Both the digital prepress and plateroom functions verify content, file naming, imposition, pagination, numbering stream, registration of technical components, color mapping, workflow, and file integrity. Records of these checks are created and saved.

Offset production requires printing that uses a lithographic process. Offline finishing activities are required to create books and package offset output. Digital output may flow through an inkjet digital production line or a sheet-fed toner application process in the Xpress Center. A battery of quality checks was performed in these areas. The checks included color match, correct file selection, content match to proof, litho-code to serial number synchronization, registration of technical components, ink density controlled by densitometry, inspection for print flaws, perforations, punching, pagination, scanning requirements, and any unique features specified for the order. Records of these checks and samples pulled from planned production points were maintained. Offline finishing included cutting, shrink-wrapping, folding, and collating. The collation process has three robust inline detection systems that inspected each book for:

- Caliper validation that detects too few or too many pages. This detector will stop the collator if an incorrect caliper reading is registered.
- An optical reader that will only accept one sheet. Two or zero sheets will result in a collator stoppage.
- The correct bar code for the signature being assembled. An incorrect or upside down signature will be rejected by the bar code scanner and will result in a collator stoppage.

Pearson's Quality Assurance department personnel inspected print output prior to collation and shipment. Quality Assurance also supported process improvement, work area documentation, audited process adherence, and established training programs for employees.

9.4 Quality Control of Scanning

Establishing and maintaining the accuracy of scanning, editing, and imaging processes is a cornerstone of the Pearson scoring process. While the scanners are designed to perform with great precision, Pearson implements other quality assurance processes to confirm that the data captured from scan processing produces a complete and accurate map to the expected results.

Pearson pioneered optical mark reading and image scanning, and continues to improve in-house scanners for this purpose. Software programs drive the capture of student demographic data and student responses from the test materials during scan processing. Routinely scheduled maintenance and adjustments to the scanner components (e.g., camera) maintain scanner calibration. Test sheets inserted into every batch test scanner accuracy and calibration.

Controlled processes for developing and testing software specifications included a series of validation and verification procedures to confirm the captured data can be mapped accurately and completely to the expected results and that editing application rules are properly applied.

9.5 Quality Control of Image Editing

The final step in producing accurate data for scoring is the editing process. Once information from the documents was captured in the scanning process, the scan program file was executed, comparing the data captured from the student documents to the project specifications. The result of the comparison was a report (or edit listing) of documents needing corrections or validation. Image Editing Services performed the tasks necessary to correct and verify the student data prior to scoring.

Using the report, editors verified that all unscanned documents were scanned, or the data were imported into the system through some other method such as flatbed scan or key entry.

Documents with missing or suspect data were pulled and verified, and corrections or additional data were entered. Standard edits included

- Incorrect or double gridding
- Incorrect dates (including birth year)
- Mismatches between pre-ID label and gridded information
- Incomplete names

When all edits were resolved, corrections were incorporated into the document file containing student records.

Additional quality checks were also performed. These included student n-count checks to make certain that

- students were placed under the correct header,
- all sheets belonged to the appropriate document,
- documents were not scanned twice, and
- no blank documents existed.

Finally, accuracy checks were performed by checking random documents against scanned data to verify the accuracy of the scanning process.

Once all corrections were made, the scan program was tested a second time to verify all data were valid. When the resulting output showed that no fields were flagged as suspect, the file was considered clean and scoring began. Once all scanning was completed, the right/wrong response data were securely handed off.

9.6 Quality Control of Answer Document Processing and Scoring

Quality control of answer document processing and scoring involves all aspects of the scoring procedures, including key-based and rule-based machine scoring and handscoring for constructed-response items and performance tasks.

For the 2015 operational administration, Pearson's validation team prepared test plans used throughout the scoring process. Test plan preparation was organized around detailed specifications.

Based on lessons learned from previous administrations, the following quality steps were implemented:

- Raw score validation (e.g., score key validation; evidence statement, field-test nonscore; double-grid combinations; possible correct combination, if applicable; out-of-range/negative test cases)

- Matching (e.g., validation of high-confidence criteria, low-confidence criteria, cross document, external or forced matching by customer; prior to and after data updates; extract file of matched and unmatched documents)
- Demographic update tests (e.g., verification of data extract against corresponding layout; valid values for updatable fields; invalid values for updatable/nonupdatable fields; negative test for nonexisting record or empty file)

The following components were added to the quality control process specifically for the program. These additional steps were introduced to address issues with item-level scoring that were identified in the 2014 field-test administration:

- XML Validation: A combination of automated validation against 100 percent of item XMLs and human inspection of XML from selected difficult item types or composite items
- Administration/End-to-End Data Validation: An automated generation of response data from approved test maps that have known conditions against the operational scoring systems and data generation systems to verify scoring accuracy
- Psychometric Validation: Verification of data integrity using criteria typically used in psychometric processes (e.g., statistical keychecks) and categorization of identified issues to help inform investigation by other groups
- Content Validation: An examination, by subject matter experts, of all items using a combination of automated tools to generate response and scoring data

In addition to the steps described above, the following quality control process for answer keys and scoring that was implemented for the first operational administration was used:

- Pearson's psychometrics team conducted empirical analyses based on preliminary data files and flagged items based on statistical criteria.
- Pearson content team reviewed the flagged items and provided feedback on the accuracy of content, answer keys, and scoring.
- Items potentially requiring changes were added to the product validation log for further investigation by other Pearson teams.
- Staff was notified of items for which keys or scoring changes were recommended.
- Participating states and agencies approved/rejected scoring changes.
- All approved scoring changes were implemented and validated prior to the generation of the data files used for psychometric processing.

9.7 Quality Control of Psychometric Processes

High-quality psychometric work for the operational administrations was necessary to provide accurate and reliable results of student performance. Pearson was responsible for the psychometric analyses of the operational administration and implemented measures to ensure the quality of work. The psychometric analyses were all conducted according to well-defined specifications. Data cleaning rules were clearly articulated and applied consistently throughout the process. Results from all analyses underwent comprehensive quality checks by a team of psychometricians and data analysts. Detailed checklists were used by members of the team for each statistical procedure.

Described below is an overview of the quality control steps performed at different stages of the psychometric analyses. Greater detail is provided in Sections 5 (“Classical Item Analysis”), 6 (“Differential Item Functioning”), 7 (“IRT Model and Parameters”), and 12 (“Scale Scores”).

Data Screening

Data screening is an important first step to ensure quality data input for meaningful analysis. The Pearson Customer Data Quality team validated all student data files used in the operational psychometric analyses. The data validation for the student data files and item response files included the following steps:

1. Validated variables in the data file for values in acceptable ranges
2. Validated that the test form ID, unique item numbers, and item sequence on the data file were consistent with the test form values on the corresponding test map
3. Computed the composite raw score, claim raw scores, and subclaim raw scores, given the item scores in the student data file
4. Compared computed raw scores to the raw scores in the student data file
5. Compared the student item response block to the item scores
6. Flagged student records with inconsistencies for further investigation

Classical Item Analysis

Classical item analysis produces item-level statistics (e.g., item difficulty and item-total correlations). The item analysis results were reviewed by Pearson psychometricians. Items flagged for unusual statistical properties were reviewed by the content team. If items were identified as having key issues, scoring issues, or content issues, they were presented to the Priority Alert Task Force, whose task was to make decisions on whether to exclude them from the calculation of reported student scores. Refer to Section 5.4 for classical item analysis item flagging criteria.

Conversion Tables

Conversion tables must be accurate because they are used to generate reported scores for students. Comprehensive records were meticulously maintained on item-level decisions, and thorough checks were made to ensure that the correct items were included in the final score. Pre-equated conversion tables were developed independently by two psychometricians and completely matched. A reasonableness check was also conducted by psychometricians for each content and grade level to make sure the results were in alignment with observations during the analyses prior to conversion table creation. Refer to Section 12.3 for the procedure to create conversion tables.

Section 10: Operational Test Forms

Each operational test form is constructed to reflect the alternate New Meridian blueprint. Multiple operational forms are constructed for each grade/subject. The test construction process determined the Common Core State Standards that are assessed in more than one evidence statement when selecting the items for the spring 2022 blueprint. As part of this process, the number of items was reduced in an attempt to keep the proportion of subclaims close to the original, while still maintaining enough points to report at the subclaim level. The process adhered to the Council of Chief State School Officers criteria for procuring and evaluating high-quality assessments.

Core forms are the operational test forms consisting of only those items that will count toward a student's score. Core forms are constructed to meet the blueprint and psychometric properties outlined in the test construction specifications. New Meridian creates multiple core forms for a given assessment to enhance test security and to support opportunity for item release. The number of core operational forms per grade/subject and mode is provided in Table 10.1.

Table 10.1 Number of Core Operational Forms per Grade/Subject and Mode for ELA/L and Mathematics

Grade/Subject	ELA/L		Mathematics	
	CBT	PBT	CBT	PBT
Grade 3	2	1	2	1
Grade 4	2	1	2	1
Grade 5	2	1	2	1
Grade 6	2	1	2	1
Grade 7	2	1	2	1
Grade 8	2	1	2	1
Grade 10	2	1		
Grade 11	2	1		
Algebra I			2	1
Geometry			2	1
Algebra II			2	1

Note. ELA/L = English language arts/literacy; CBT = computer-based test; PBT = paper-based test.

In addition to the operational core forms, appropriate forms were identified as accessibility and accommodated forms. Grades 3 through 8 and 10 through 11 English language arts/literacy (ELA/L) and Integrated Mathematics I and II have two operational accommodated forms, and mathematics grades 3 through 8 and the high school traditional assessments have three accommodated forms. The forms are accommodated to support braille, large print, human reader/human signers, assistive technology, text-to-speech, closed-captioning, and Spanish. Human reader/human signers and Spanish are provided for mathematics assessments only. Closed captioning is provided for ELA/L assessments only.

The summative assessments were administered in either a computer-based test or a paper-based test format. English learner assessments focused on writing effectively when analyzing text. Mathematics assessments focused on applying skills and concepts, and featured multistep problems that require abstract reasoning and modeling of real-world problems. In both content areas, students also demonstrated their acquired skills and knowledge by answering selected-response items and fill-in-the-blank questions. Each assessment was comprised of multiple units; one of the mathematics units was split into calculator and noncalculator sections.

Section 11: Student Characteristics

11.1 Overview of Test Taking Population

Over a million forms were administered in the Department of Defense Education Activity, the District of Columbia, Illinois, and New Jersey during the 2021–2022 school year. Not all participating states and agencies had students testing in all grades. Assessments were administered for English language arts/literacy (ELA/L) in grades 3 through 8 and 10; mathematics assessments were administered in grades 3 through 8, as well as for traditional high school mathematics (Algebra I, Geometry, and Algebra II). New Jersey launched the New Jersey Graduation Proficiency Assessment (NJGPA) as a field test for the class of 2023 in spring 2022. Student characteristics for this group are presented in the NJGPA Addendum of this technical report. The majority of students tested during the spring administration when all grades and content areas were administered mostly online with small numbers of paper testers.

11.2 Rules for Inclusion of Students in Analyses

Criteria for inclusion of students were implemented prior to all operational analyses. These rules were established by Pearson psychometricians in consultation with participating states and agencies to determine which, if any, student records should be removed from analyses. This data screening process resulted in higher-quality, albeit slightly smaller, data sets.

Student response data were included in analyses if they met the following criteria:

1. Valid form numbers were observed for each unit for online assessments or for the full form for paper assessments.
2. Student records were not flagged as “void” (i.e., do not score).
3. The student attempted at least 25 percent of the items in each unit or form.

Additionally, in cases where students had more than one valid record, the record with the higher raw score was chosen. Records for students with administration issues or anomalies were excluded from analyses.

11.3 Students by Grade/Course, Mode, and Gender

Table 11.1 presents, for each grade of ELA/L, the number and percentage of students who took the test in each mode (computer-based test or paper-based test). This information is provided for all participating states combined. Table 11.2 presents the same information for all students who took the mathematics assessments, and Table 11.3 provides this information for students who took the mathematics assessments in Spanish.

Markedly more students tested online than on paper across all grades for both content areas. For ELA/L, the percentages of online students by grade level were greater than 99 percent. For all mathematics students, the percentages of students testing online was greater than 99 percent. The percentage of students taking Spanish-language mathematics online forms was greater than or equal to 99 percent. Overall, fewer students tested at the higher grades for both content areas.

Table 11.1 ELA/L Students by Grade and Mode: All States Combined

Grade	No. of Valid Cases	CBT		PBT	
		N	%	N	%
3	229,185	228,358	99.6	827	0.4
4	232,317	231,743	99.8	574	0.2
5	237,847	237,297	99.8	550	0.2
6	238,260	237,783	99.8	477	0.2
7	245,106	244,677	99.8	429	0.2
8	249,388	248,973	99.8	415	0.2
9	103,744	103,693	100	51	0
10	7,069	7,060	99.9	9	0.1
Grand Total	1,542,916	1,539,584	99.8	3,332	0.2

Note. Includes students taking accommodated forms of English language arts/literacy. CBT = computer-based test; PBT = paper-based test.

Table 11.2 Mathematics Students by Grade/Course and Mode: All States Combined

Grade/Course	No. of Valid Cases	CBT		PBT	
		N	%	N	%
3	230,054	229,434	99.7	620	0.3
4	232,955	232,377	99.8	578	0.2
5	238,545	237,994	99.8	551	0.2
6	238,630	238,150	99.8	480	0.2
7	235,692	235,274	99.8	418	0.2
8	211,563	211,159	99.8	404	0.2
A1	115,174	115,106	99.9	68	0.1
GO	40,404	40,389	100	15	0
A2	13,965	13,963	100	2	0
Grand Total	1,556,982	1,553,846	99.8	3,136	0.2

Note. Includes students taking mathematics in English, students taking Spanish-language forms for mathematics, and students taking accommodated forms. CBT = computer-based test; PBT = paper-based test; A1 = Algebra I; GO = Geometry; A2 = Algebra II; n/a = not applicable.

Table 11.3 Spanish-Language Mathematics Students by Grade/Course and Mode: All States Combined

Grade/Course	No. of Valid Cases	CBT		PBT	
		N	%	N	%
3	5,063	5,048	99.7	15	0.3
4	4,605	4,599	99.9	6	0.1
5	4,220	4,215	99.9	5	0.1
6	3,721	3,718	99.9	3	0.1
7	3,226	3,224	99.9	2	0.1
8	2,844	2,840	99.9	4	0.1
A1	2,563	2,562	100	1	0
GO	425	425	100	0	0
A2	138	138	100	0	0
Grand Total	26,805	26,769	99.9	36	0.1

Note. CBT = computer-based test; PBT = paper-based test.

Tables A.11.1, A.11.2, and A.11.3 in Appendix 11 show the number and percentage of students with valid test scores in each content area (including Spanish-language mathematics), grade/course, and mode of assessment for all states and agencies combined and for each state or agency separately. Tables A.11.4, A.11.5, and A.11.6 present the distribution by content area, grade/course, mode, and gender, for all states combined.

11.4 Demographics

Also presented in Appendix 11 is student demographic information for the following characteristics: economically disadvantaged, students with disabilities, English learners, gender, and race/ethnicity (American Indian/Alaska Native, Asian, Black/African American, Hispanic/Latino, White/Caucasian, Native Hawaiian or Other Pacific Islander, two or more races reported, race not reported). Student demographic information was provided by the states and districts and captured in PearsonAccess^{next} or PearsonAccess 5, depending on which platform was used by the respective state, by means of a student data upload. The demographic data was verified by the states and districts prior to score reporting. Not all demographics were provided for all students. Students missing information on one or more demographic variables were omitted from the corresponding subgroup analyses.

Tables A.11.7 through A.11.14 provide demographic information for students with valid ELA/L scores, and Tables A.11.15 through A.11.23 present demographics for students with valid mathematics scores. All tables of demographic information are organized by grade/course; the results are first aggregated across all participating states and agencies and then presented for each state or agency. Percentages are not reported in instances where fewer than 20 students tested in a grade/course area.

Section 12: Scale Scores

Participating states and agencies report results according to five performance levels that delineate the knowledge, skills, and practices students are able to demonstrate:

- Level 5: Exceeded expectations
- Level 4: Met expectations
- Level 3: Approached expectations
- Level 2: Partially met expectations
- Level 1: Did not yet meet expectations

The assessments are designed to measure and report results in categories called master claims and subclaims. Master claims (or simply “claims”) are at a higher level than subclaims with content representing multiple subclaims contributing to each claim outcome. In addition, four scale scores are reported for the assessments. A summative scale score is reported for each mathematics assessment. A summative scale score and separate claim scores for Reading and Writing are reported for each English language arts/literacy (ELA/L) assessment.

Subclaim outcomes describe student performance for content-specific subsets of the item scores contributing to a particular claim. For example, Written Expression and Knowledge of Conventions subclaim outcomes are reported along with Writing claim scores. Subclaim outcomes are reported as *Below Expectations*, *Nearly Meets Expectations*, or *Meets or Exceeds Expectations*.

12.1 Operational Test Content (Claims and Subclaims)

A claim is a statement about student performance based on how students respond to test questions. The tests are designed to elicit evidence from students that supports valid and reliable claims about the extent to which they are college- and career-ready or on track toward that goal and are making expected academic gains based on the Common Core State Standards.

The number of items associated with each claim and subclaim outcome varies depending on subject and grade. The item types vary in terms of the number of points associated with them, so that both the number of items and the number of points are important in evaluating the quality of a claim or subclaim score.

12.1.1 English Language Arts/Literacy

Table 12.1³ includes the number of items and the number of points by subclaim and claim for ELA/L grade 3. Corresponding information is provided in Appendix 12.1 for all ELA/L grades.

³ Table A.12.1 in Appendix 12.1 is identical to Table 12.1.

Table 12.1 Form Composition for ELA/L Grade 3

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	4–7	8–17
	Reading Informational Text	4–7	11–20
	Vocabulary	4–5	8–10
	Claim Total	12–14	30–31
Writing	Written Expression	1	18
	Knowledge of Conventions	1	6
	Claim Total	2	24
SUMMATIVE TOTAL		14–16	54–55

Note. Each prose constructed-response trait is identified as a separate item in this table for the two writing subclaims and, in some cases, either the Reading Literary Text or the Reading Informational Text subclaim.

Each ELA/L form contains items of varying types. The prose constructed-response (PCR) traits contribute to different claims, and the aggregate of the traits contributes to the summative scale score. ELA/L assessments consist of two PCR tasks. The following details the number of possible points and the associated subclaims for the three PCR tasks:

- Literary Analysis Task
- Research Simulation Task
- Narrative Writing Task

All ELA/L assessments include the Research Simulation Task and either the Literary Analysis Task or the Narrative Writing Task. The Literary Analysis Task and the Research Simulation Task are scored for two traits: (1) Reading Comprehension and Written Expression and (2) Knowledge of Conventions. The Narrative Writing Task is scored for two traits: (1) Written Expression and (2) Knowledge of Conventions. All traits are initially scored as either 0–3 or 0–4; the Written Expression traits are multiplied by 3 (or weighted) to increase their contribution to the total score, making possible subclaim scores 0, 3, 6, and 9, or 0, 3, 6, 9, and 12. The maximum possible points for ELA/L PCR items are provided in Table 12.2.

Table 12.2 Contribution of Prose Constructed-Response Items to ELA/L

Grade	Score	Possible Points		
		Literary Analysis Task	Research Simulation Task	Narrative Writing Task
3	Reading	3	3	0
	Written Expression	9	9	9
	Knowledge of Conventions	3	3	3
	Total	15	15	12
4–5	Reading	4	4	0
	Written Expression	12	12	9
	Knowledge of Conventions	3	3	3
	Total	19	19	12
6–11	Reading	4	4	0
	Written Expression	12	12	12
	Knowledge of Conventions	3	3	3
	Total	19	19	15

Note. English language arts/literacy assessments consist of the Research Simulation Task and either the Literary Analysis Task or the Narrative Writing Task.

12.1.2 Mathematics

Table 12.3⁴ includes the numbers of items and points associated with subclaim scores for mathematics grade 3, as an example of the composition of the mathematics tests.

Table 12.3 Mathematics Form Composition for Grade 3

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	18	20
	Additional & Supporting Content	9	10
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	12
TOTAL		33	52

Because there is substantial variation in the composition of the tests, corresponding information is provided in the tables in Appendix 12.1 for all mathematics grades/courses.

12.2 Establishing the Reporting Scales

Reporting scales designate student performance into one of five performance levels⁵ with Level 1 indicating the lowest level of performance and Level 5 indicating the highest level of performance. Threshold or cut scores associated with performance levels were initially expressed as raw scores on the performance level setting (PLS) forms approved by the Governing Board. A scale score task force was assembled, which made recommendations about how threshold levels would be represented on the reporting scale.

12.2.1 Summative Score Scale and Performance Levels

There are 201 defined summative scale score points for both ELA/L and mathematics, ranging from 650 to 850. The lowest obtainable scale score is 650 and the highest obtainable scale score is 850. The thresholds for summative performance levels on the scale score metric recommended by the scale score task force are Level 2 and Level 4. The cuts are the anchors for establishing the linear transformation between the theta scale and the reported scale score. A scale score of 700 is associated with minimum Level 2 performance, and a scale score of 750 is associated with minimum Level 4 performance. Not all possible scale scores may be realized in a scoring table.

For spring 2015, scale scores were defined for each test as a linear transformation of the theta (θ_{2015}) scale. The theta values associated with the Level 2 and Level 4 performance levels were identified using the test characteristic curve (TCC) associated with the PLS form. With Levels 2 and 4 scale scores fixed at 700 and 750, respectively, the relationship between theta (θ_{2015}) and scale scores ($ScaleScore_{2015}$) was established as

$$ScaleScore_{2015} = A_{2015} \times \theta_{2015} + B_{2015} \quad (12-1)$$

⁴ Table A.12.9 in Appendix 12.1 is identical to Table 12.3.

⁵ Section 8 provides an overview of the PLS process, and detailed information can be found in the *Performance Level Setting Technical Report*.

where A_{2015} is the slope and B_{2015} is the intercept. The slope and intercept were established as

$$A_{2015} = \frac{750 - 700}{\theta_{2015_{Level4}} - \theta_{2015_{Level2}}} \quad (12-2)$$

and

$$B_{2015} = 750 - A_{2015} \times \theta_{2015_{Level4}} \quad (12-3)$$

As indicated by these formulas, the slope and intercept for the summative scale scores were based on the theta scale, and by default the item response theory (IRT) parameter scale, established in 2015. Since the spring 2016 IRT parameter scale is the base scale for the IRT parameters, the scaling constants A_{2015} and B_{2015} were updated in order to continue reporting performance levels, summative scale scores, claim scores, and subclaim performance levels on the same scale as 2015. Maintaining the 2015 scale allows for prior year scores to be compared to current and future scores, and it maintains the performance levels cut scores.

New scaling constants for the summative scale score were needed for the linear transformation of the theta scale θ_{2016} to the 2015 reporting scale ($ScaleScore_{2015}$):

$$ScaleScore_{2015} = SA_{2016} \times \theta_{2016} + SB_{2016} \quad (12-4)$$

The slope ($slope_{2015_to_2016}$) and intercept ($intercept_{2015_to_2016}$) generated during the year-to-year linking defined the linear relationship between the 2015 theta scale (θ_{2015}) and the 2016 theta scale (θ_{2016}). These values were included in the scale score formula, and the formulas were used to solve for the slope (SA_{2016}) and (SB_{2016}) intercept for 2016.

The slope (A_{2016}) was updated using the following formula:

$$SA_{2016} = \frac{A_{2015}}{slope_{2015_to_2016}} \quad (12-5)$$

where A_{2015} is the current scale score multiplicative constant, $slope_{2015_to_2016}$ is the multiplicative coefficient from the year-to-year linking, and SA_{2016} is the scale score slope constant for 2016 and beyond.

The intercept (B_{2016}) was updated using the following formula:

$$SB_{2016} = B_{2015} - A_{2016} \times intercept_{2015_to_2016} \quad (12-6)$$

where B_{2015} is the current scale score additive constant, A_{2016} is the updated scale score slope, and (SB_{2016}) is the scale score intercept constant for 2016 and beyond.

In addition, new scaling constants for the Reading and Writing claim scales were needed. The same formulas were applied by replacing the slope (A_{2015}) and intercept (B_{2015}) with the Reading claim slope and intercept and the Writing claim slope and intercept.

A and B values resulting from these calculations as well as the theta values associated with the threshold performance levels are included in Appendix 12.2. Also, the 2015–2016 technical report includes raw to scale score conversion tables for the PLS forms.

12.2.2 ELA/L Reading and Writing Claim Scale

There are 81 defined scale score points possible for Reading, ranging from 10 to 90. The threshold Reading and Writing performance levels on the scale score metric recommended by the scale score task force are Level 2 and Level 4. A scale score of 30 is associated with minimum Level 2 performance, and a scale score of 50 is associated with minimum Level 4 performance. There are 51 defined scale score points possible for Writing, ranging from 10 to 60. A scale score of 25 is associated with minimum Level 2 performance, and a scale score of 35 is associated with minimum Level 4 performance. Not all possible scale scores may be realized in a scoring table.

As with the summative scale scores, scale scores for Reading and Writing were defined for each test as a linear transformation of the IRT theta (θ) scale. The same IRT theta scale was used for Reading and Writing as was used for the ELA/L summative scores. The theta values associated with the Level 2 and Level 4 performance levels were identified using the TCC associated with the PLS form. As with the summative scores, the relationship between theta and scale scores was established with Level 2 and Level 4 theta scores and the corresponding predefined scale scores. The formulas used for this are provided in Table 12.4.

Table 12.4 Calculating Scaling Constants for Reading and Writing Claim Scores

Reading	Writing
$Scale = A_R \times \theta + B_R$	$Scale = A_W \times \theta + B_W$
$A_R = \frac{50 - 30}{\theta_{Level4} - \theta_{Level2}}$	$A_W = \frac{35 - 25}{\theta_{Level4} - \theta_{Level2}}$
$B_R = 50 - A \times \theta_{Level4}$	$B_W = 35 - A \times \theta_{Level4}$

Note. A and B values resulting from these calculations are included in Appendix 12.2.

12.2.3 Subclaims Scale

The Level 4 cut is defined as *Meets or Exceeds Expectations* because high school students at Level 4 or above are likely to have the skills and knowledge to meet the definition of career- and college-readiness. The Level 3 cut is defined as *Nearly Meets Expectations*. Subclaim outcomes center on the Level 3 and Level 4 performance levels and are reported at three levels:

- Below Expectations
- Nearly Meets Expectations
- Meets or Exceeds Expectations

The subclaim performance levels are designated through the IRT theta (θ) scale for the items associated with a particular subclaim. The theta values and corresponding raw scores associated with the Level 3 and Level 4 performance levels were identified using the TCC. Students earning a raw subclaim score equal to or greater than the Level 4 threshold were designated as *Meets or Exceeds Expectations*. Students not earning a raw

subclaim score equal to or greater than the Level 3 threshold were designated as *Below Expectations*. Students whose raw subclaim score fell between the Level 3 and 4 thresholds were designated as *Nearly Meets Expectations*.

12.3 Creating Conversion Tables

A conversion table relates the number of points earned by a student on the ELA/L summative score, the mathematics summative score, the Reading claim score, or the Writing claim score to the corresponding scale score for the test form administered to that student. An IRT inverse TCC approach is used to develop the relationship between point scores and theta, θ_s (IRT ability estimates). In carrying out the calculations, estimates of item parameters and thetas are substituted for parameters in the formulas in each step.

Step 1: Calculate the expected item score (i.e., estimated item true score) for every theta in the selected range (between -15 and +15, in 0.0001 increments) based on the generalized partial credit model for both dichotomous and polytomous items:

$$s_i(\theta_j) = \sum_{m=0}^{M_i-1} mp_{im}(\theta_j) \quad (12-7)$$

$$p_{im}(\theta_j) = \frac{\exp\left[\sum_{k=0}^m Da_i(\theta_j - b_i + d_{ik})\right]}{\sum_{v=0}^{M_i-1} \exp\left[\sum_{k=0}^v Da_i(\theta_j - b_i + d_{iv})\right]} \quad (12-8)$$

where $a_i(\theta_j - b_i + d_{i0}) \equiv 0$; $s_i(\theta_j)$ is the expected item score for item i on theta, θ_j ; $p_{im}(\theta_j)$ is the probability of a student, j , with θ_j getting score m on item i ; m_i is the number of score categories of item i ; with possible item scores as consecutive integers from 0 to $m_i - 1$; D is the IRT scale constant (1.7); a_i is a slope parameter; b_i is a location parameter reflecting overall item difficulty; d_{ik} is a location parameter incrementing the overall item difficulty to reflect the difficulty of earning score category k ; and v is the number of score categories.

Step 2: Calculate the expected (weighted) test score for every theta in the selected range:

$$T_j = \sum_{i=1}^I w_i s_i(\theta_j) \quad (12-9)$$

where T_j is the expected (weighted) test score on theta, θ_j ; w_i is the item weight for item i (e.g., with $w_i = 2$, a dichotomous item is scored as 0 or 2, and a three-category item is scored as 0, 2, or 4); and I is the total number of items in a test form.

Step 3: Calculate the estimated conditional standard error of measurement (CSEM) for each theta in the selected range:

$$CSEM_j = \sqrt{\frac{1}{\sum_{i=1}^I L_i(\theta_j)}} \quad (12-10)$$

$$L_i(\theta_j) = (Da_i)^2 [s_{i2}(\theta_j) - s_i^2(\theta_j)] \quad (12-11)$$

$$s_{i2}(\theta_j) = \sum_{m=0}^{M_i-1} m^2 p_{im}(\theta_j) \quad (12-12)$$

where $L_i(\theta_j)$ is the estimated item information function for item i on theta, θ_j .

Step 4: Match every raw score with a theta. θ_j is the theta for a raw score r_h , if $T_j - r_h$ is minimum across all T_j .

Step 5: Calculate the reported scale score. Using the A and B scaling constants in Appendix 12.2, convert each theta value to a scale score and each theta CSEM to a scale score CSEM:

$$ScaleScore = A \times \theta + B \quad (12-13)$$

$$CSEM = CSEM_\theta \times A \quad (12-14)$$

The scale scores are rounded to the nearest whole number, and CSEMs are rounded to the tenths place. Furthermore, the scale scores are truncated with the lowest obtainable scale score (LOSS) of 650 and highest obtainable scale score (HOSS) of 850.

Figure 12.1 contains TCCs, estimated CSEM curves, and estimated information (INF) curves for ELA/L grade 3.⁶ The curves in each figure are for the two core online forms (O1 and O2), one core paper form (P1), and one or more accommodated forms A(O). The curves are reported on the theta scale. Vertical dotted lines indicate the performance level cuts on the theta scale. For ELA/L grade 3, all forms had similar TCCs. CSEM and INF curves were also similar.

Appendix 12.3 contains TCC, CSEM, and INF curves for all ELA/L grades and all mathematics grades/courses. The curves are based on IRT parameters from a prior operational or field-test administration.

⁶ Grade 3 TCC, CSEM, and INF curves are also included in Appendix Figure A.12.1.

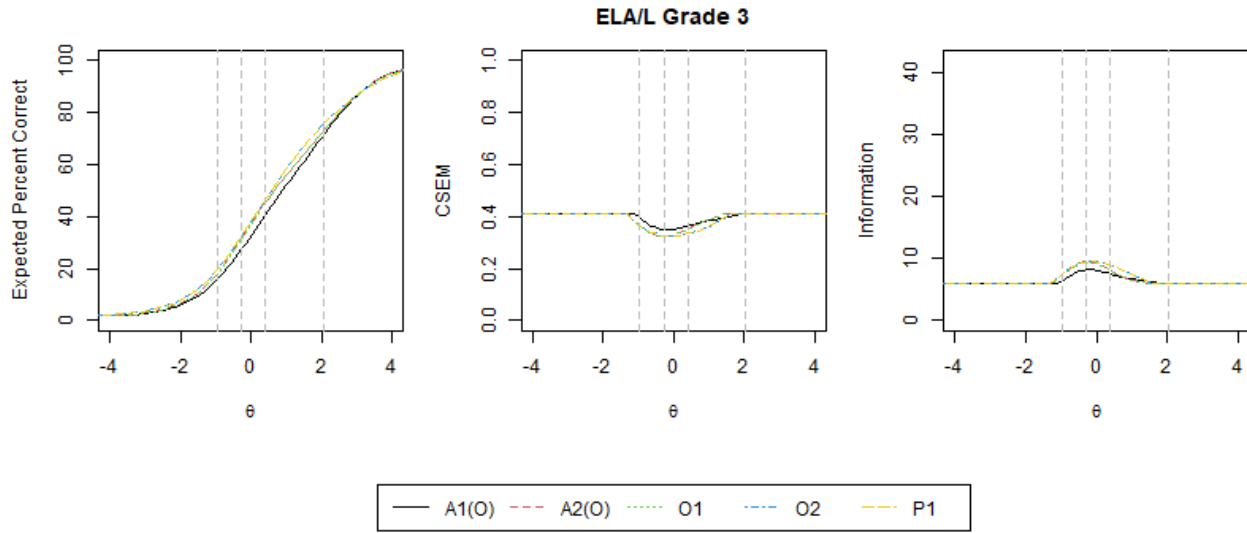


Figure 12.1 Test Characteristic Curves, Conditional Standard Error of Measurement Curves, and Information Curves for ELA/L Grade 3

12.4 Score Distributions

12.4.1 Score Distributions for English Language Arts/Literacy

Figures 12.2 through 12.4 graphically represent the distributions of scale scores for grades 3 through 10 ELA/L summative, Reading, and Writing, respectively. The vertical axis of each graph, labeled “Density,” represents the proportion of students earning the scale score point indicated along the horizontal axis. For the summative distributions, the y-axis ranges from 0 to .02 and the x-axis from 650 to 850. For the Reading distributions, the y-axis ranges from 0 to .05 and the x-axis from 10 to 90. For the Writing distributions, the y-axis ranges from 0 to .10 and the x-axis from 10 to 60.

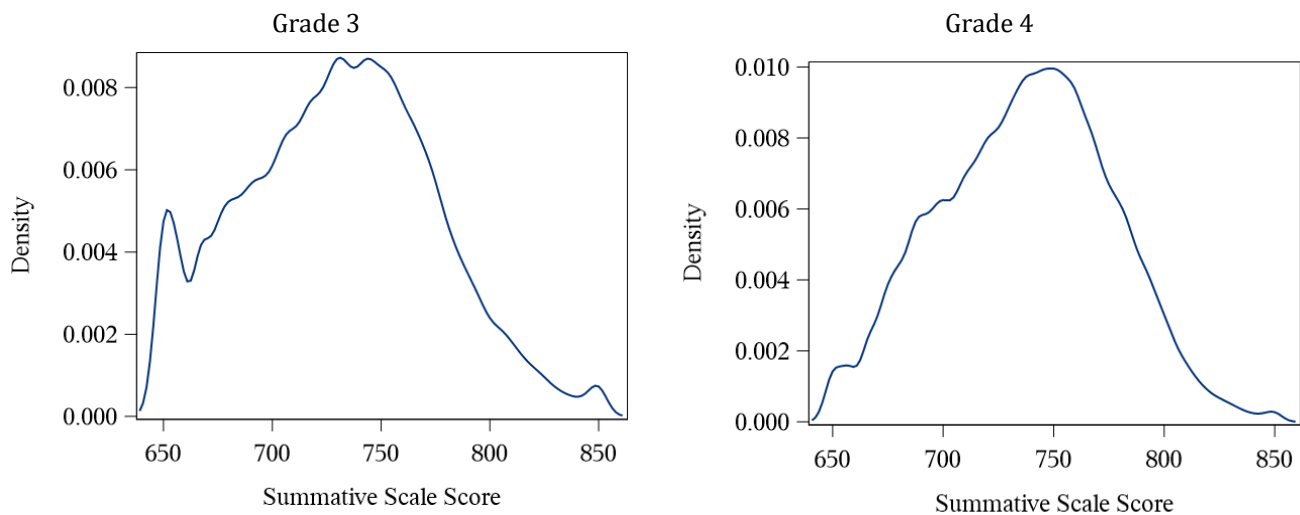
The distributions of the ELA/L summative scale scores were fairly symmetrical and centered around the Level 4 cut score (750) or slightly below, except for grade 11, which was centered closer to 700.

Reading scale scores tended to be centered around or slightly below the Level 4 cut score of 50 and were slightly more irregular than the summative scale scores. Distributions tended to be fairly symmetric, except for grade 11, which was skewed right.

Writing scale score distributions were noticeably less smooth than Reading or ELA/L summative distributions due to peaks related to the weighting of the Written Expression portion of the PCR tasks and a noticeable proportion of students at the LOSS. Due to the weighting of the Written Expression trait, multiple Writing scale score values are not likely to be obtained resulting in multiple peaks across the range of the Writing scale score. A noticeable proportion of students earned the LOSS of 10 in Writing across all ELA/L grades. Students with zero raw score points on the written portion of the assessment are automatically assigned the LOSS value of a

scale. Writing items are embedded exclusively in PCR tasks, which tended to be difficult. The Written Expression trait also tended to be the most difficult of the PCR traits.

Across the ELA/L grades, there are relatively few students between 11 and about 20, depending on the grade.⁷ As noted in Section 12.2.2, the scale score task force selected 10 as the LOSS. This value was selected to be consistent with the Reading LOSS and reduce truncation at the lower ends of the scale. However, the scale is defined by the theta values associated with the Level 2 and Level 4 performance levels. All other scale score values are identified through a theta-to-scale score linear transformation applying the scaling constants (Table 12.4). For Writing, the lowest theta estimate associated with raw scores ranging from one to two are linearly transformed to scale score values generally between 15 and 20, meaning that there may be multiple scale scores between 11 and 20 that are not assigned to a raw score. In contrast, the Reading lowest theta estimates associated with raw scores ranging from one to two are linearly transformed to scale score values closer to the LOSS. The gap in the proportion of students at the scale scores between the LOSS value of 10 and the scale score values around 17 to 19 is an artifact of the scale score task force selecting the LOSS value of 10.



⁷ Due to smoothing of the kernel density function, in some figures, particularly those with small sample sizes, the line representing the distribution may appear to remain above zero near the region.

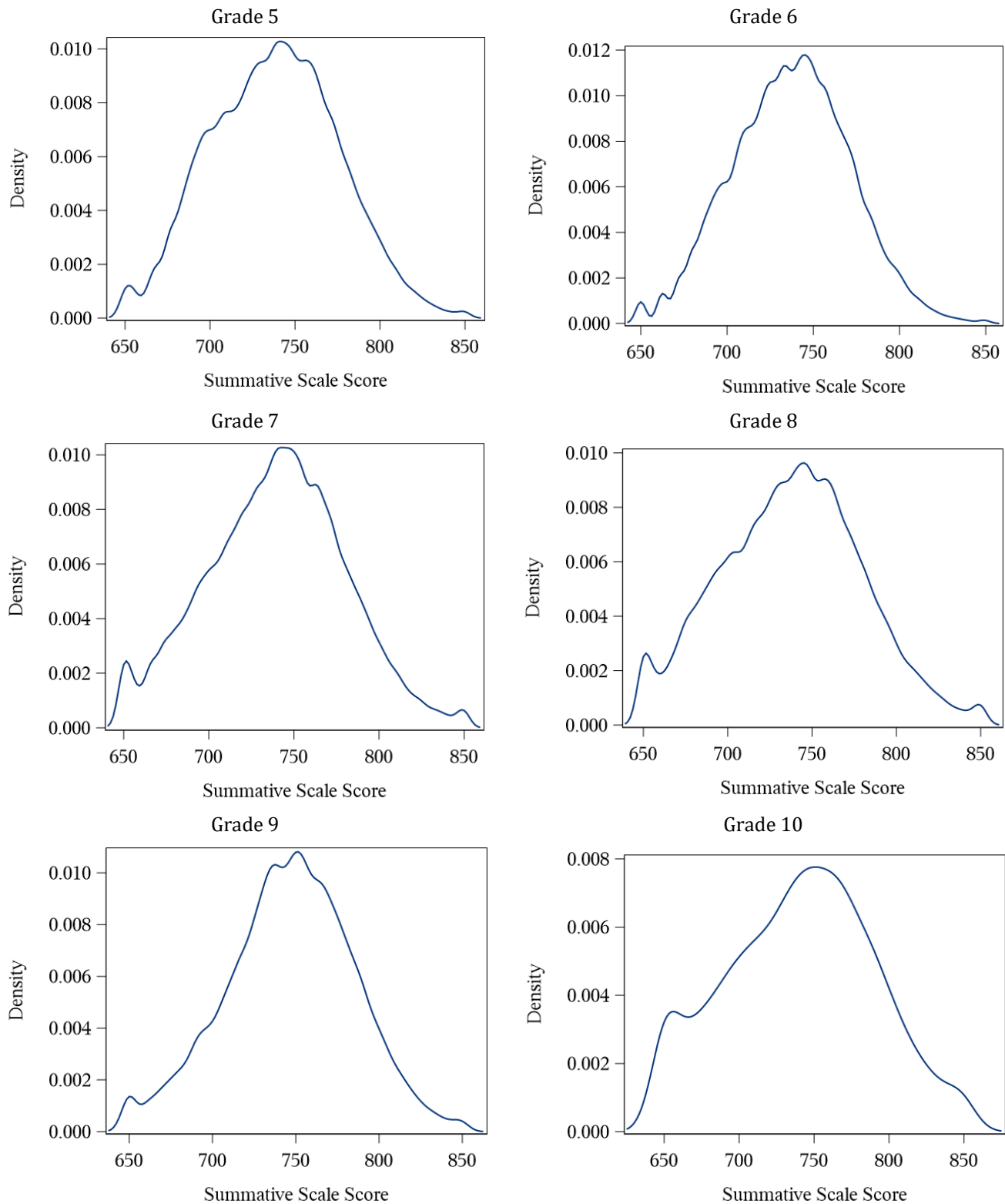
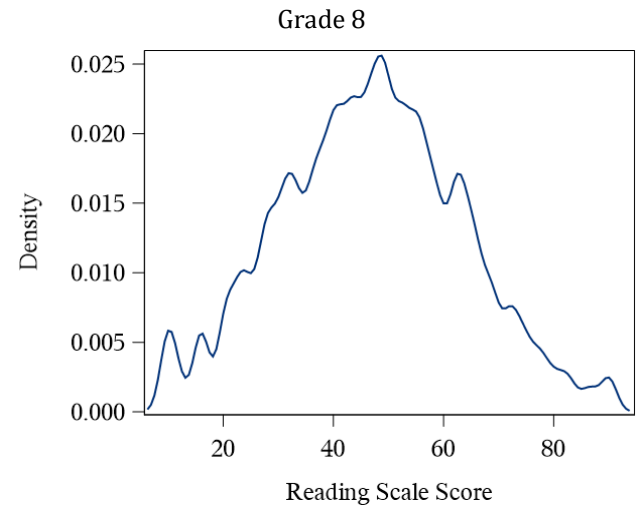
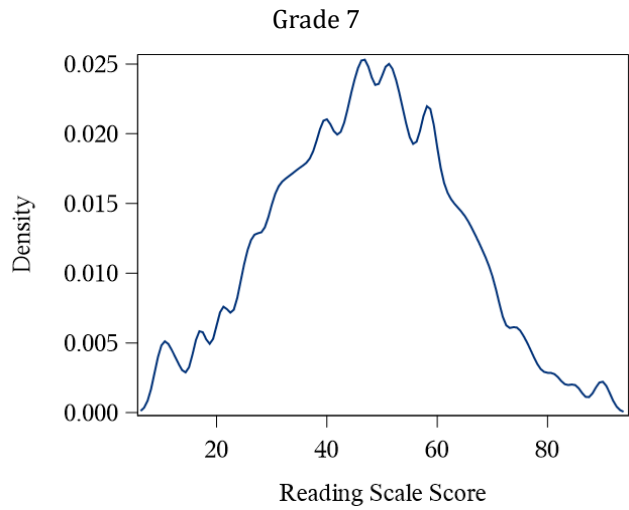
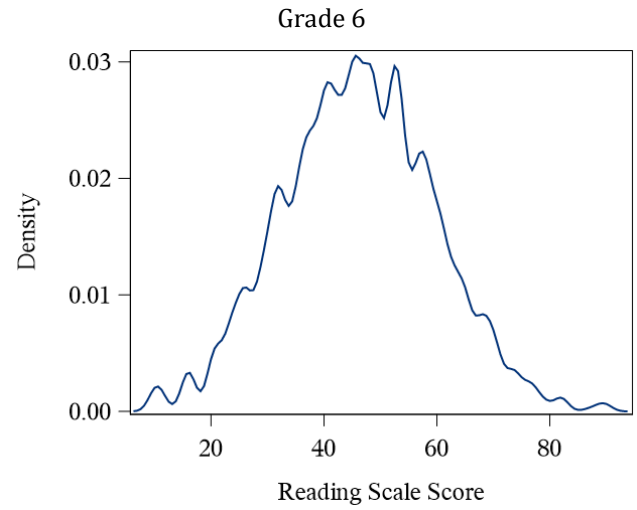
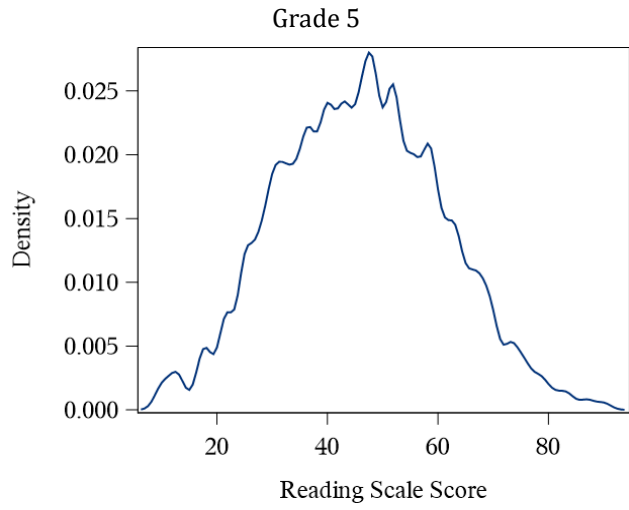
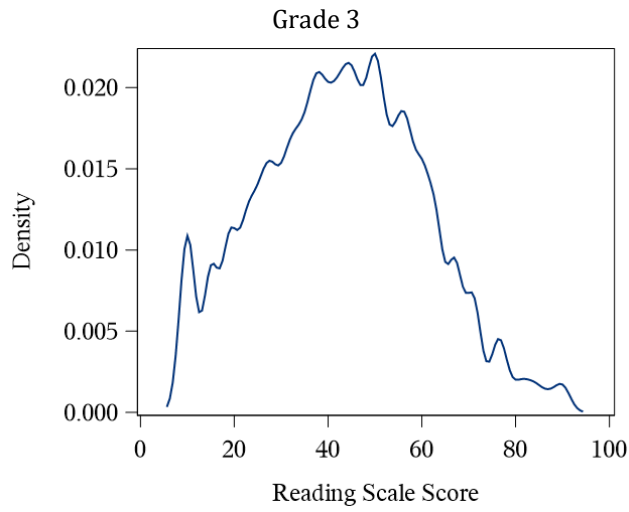


Figure 12.2 Distributions of ELA/L Scale Scores: Grades 3–10



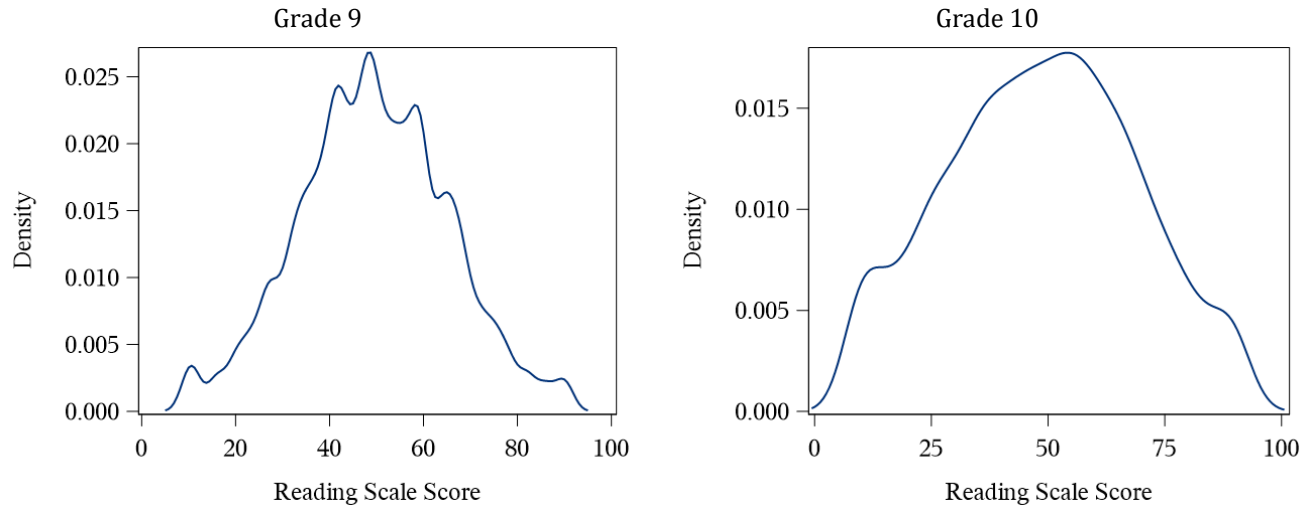


Figure 12.3 Distributions of Reading Scale Scores: Grades 3–10

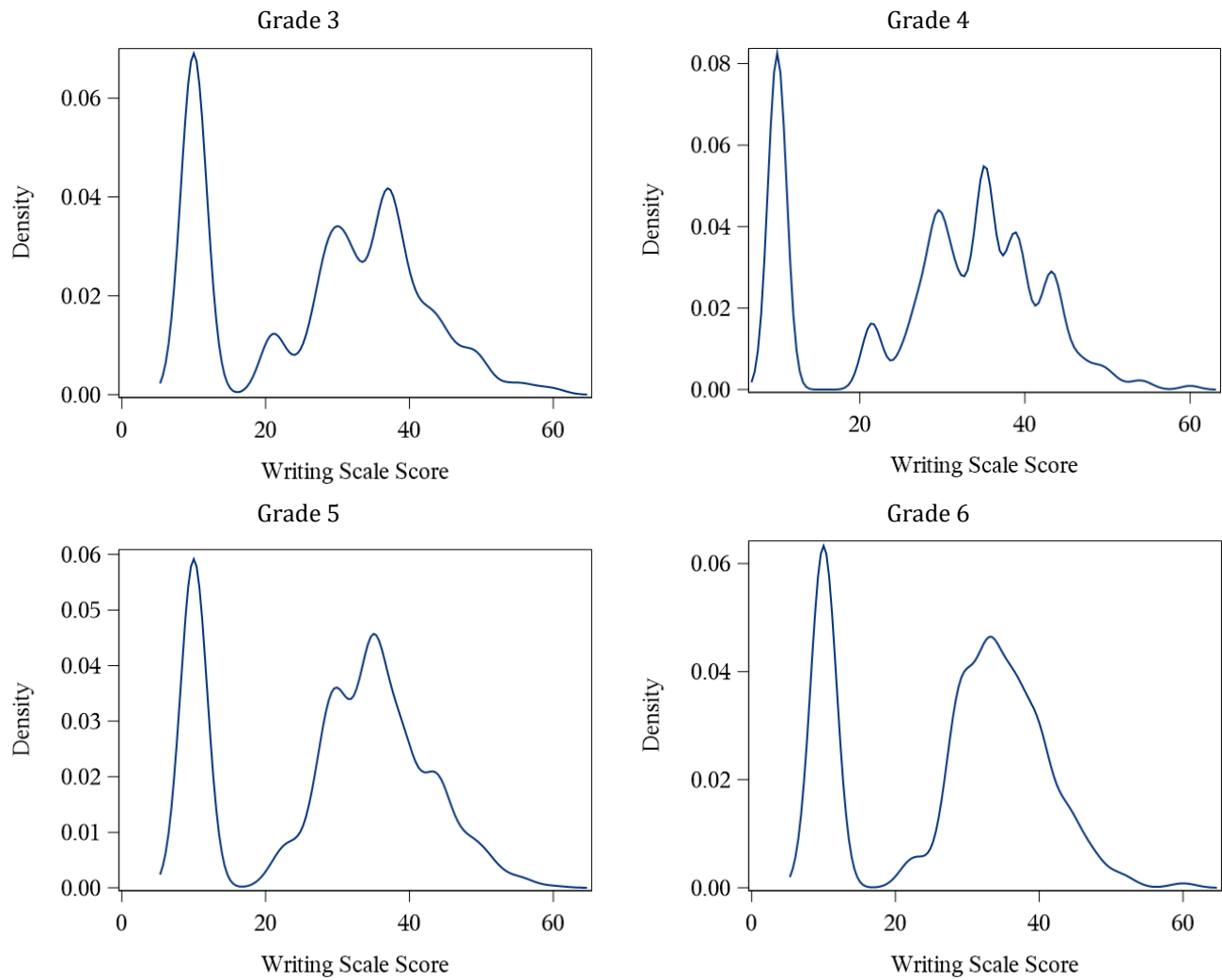


Figure 12.4 Distributions of Writing Scale Scores: Grades 3–10

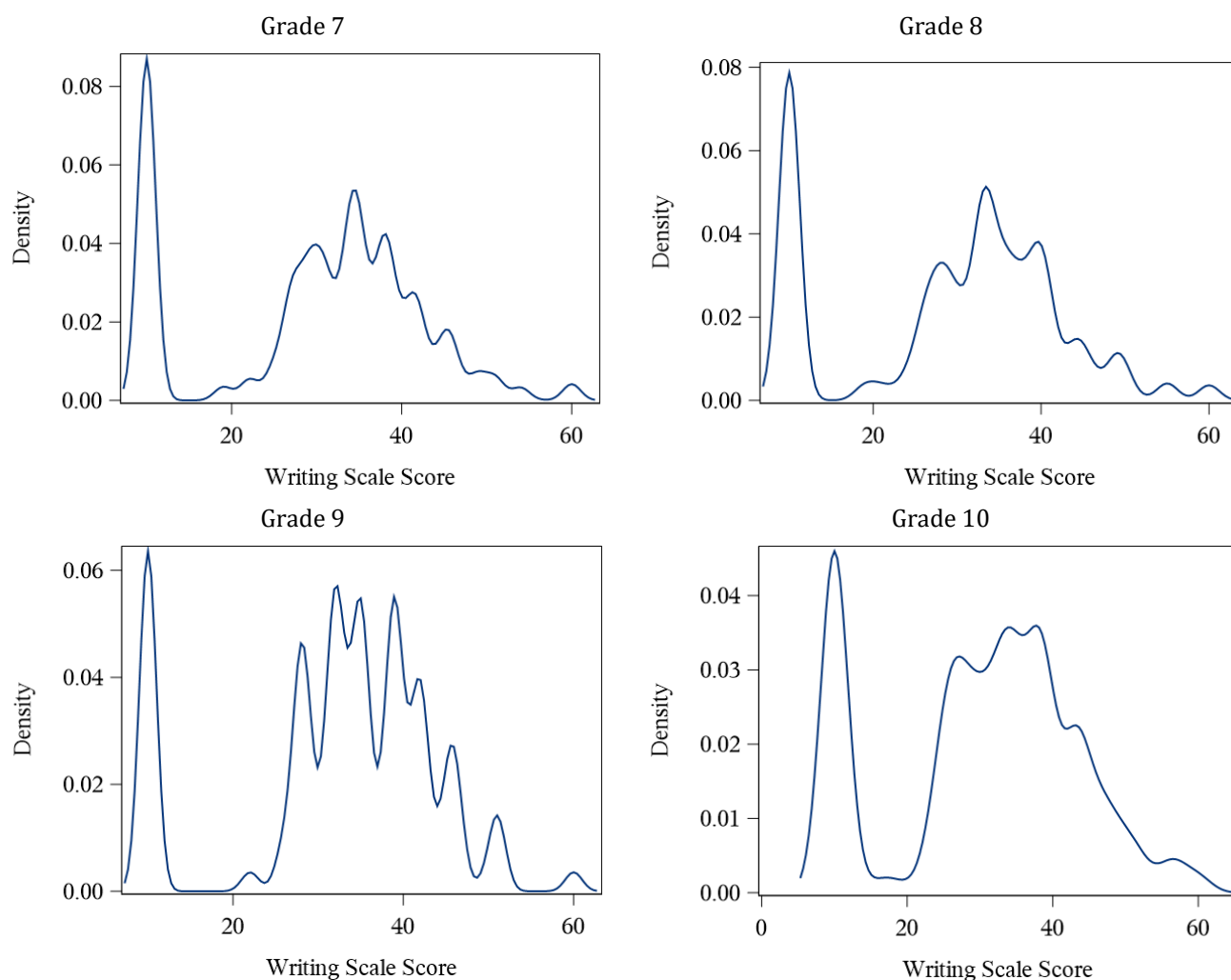


Figure 12.4 (continued) Distributions of Writing Scale Scores: Grades 3–10

12.4.2 Scale Score Cumulative Frequencies for ELA/L

The cumulative frequency distribution for the summative scale score is presented in Appendix 12.4 for ELA/L assessments.

12.4.3 Summary Scale Score Statistics for ELA/L Groups

Subgroup statistics for ELA/L full summative, Reading, and Writing scale scores are presented in Tables 12.5 and 12.6⁸ for ELA/L grades 3 and 10, respectively. The results for all ELA/L grades are provided in Appendix 12.5. Grade 3 ELA/L subgroup statistics are presented in Table 12.5.⁹

⁸ Due to omitted demographic values, subgroup sample sizes may not sum to the total sample size.

⁹ Table A.12.40 in Appendix 12.5 is identical to Table 12.5.

Table 12.5 Subgroup Performance for ELA/L Scale Scores: Grade 3

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		229,184	730.45	44.30	650	850
Gender	Female	112,531	735.44	44.88	650	850
	Male	116,647	725.63	43.19	650	850
Ethnicity	American Indian/Alaska Native	507	728.13	46.08	650	850
	Asian	17,564	762.49	42.61	650	850
	Black or African American	38,247	709.72	40.39	650	850
	Hispanic/Latino	65,085	716.97	41.59	650	850
	Native Hawaiian or Pacific Islander	358	741.97	45.24	650	850
	Two or More Races	9,558	736.35	44.73	650	850
	White	97,143	741.20	40.94	650	850
Economic Status*	Not Economically Disadvantaged	123,355	744.71	42.33	650	850
	Economically Disadvantaged	96,131	711.28	39.65	650	850
English Learner Status	Non-English Learner	183,158	735.44	43.84	650	850
	English Learner	35,250	704.19	37.04	650	850
Disabilities	Students without Disabilities	188,159	735.84	43.24	650	850
	Student with Disability (SWD)	38,991	705.46	40.61	650	850
Reading Summative Score		229,184	42.98	17.62	10	90
Gender	Female	112,531	44.55	17.70	10	90
	Male	116,647	41.46	17.41	10	90
Ethnicity	American Indian/Alaska Native	507	41.67	17.75	10	84
	Asian	17,564	54.70	16.71	10	90
	Black or African American	38,247	35.07	16.16	10	90
	Hispanic/Latino	65,085	37.41	16.36	10	90
	Native Hawaiian or Pacific Islander	358	46.67	17.17	10	90
	Two or More Races	9,558	45.53	17.79	10	90
	White	97,143	47.43	16.53	10	90
Economic Status*	Not Economically Disadvantaged	123,355	48.57	16.89	10	90
	Economically Disadvantaged	96,131	35.40	15.69	10	90
English Learner Status	Non-English Learner	183,158	44.97	17.41	10	90
	English Learner	35,250	32.24	14.47	10	90
Disabilities	Students without Disabilities	188,159	45.06	17.20	10	90
	Student with Disability (SWD)	38,991	33.33	16.37	10	90
Writing Summative Score		229,184	27.46	13.47	10	60
Gender	Female	112,531	29.19	13.45	10	60
	Male	116,647	25.79	13.27	10	60
Ethnicity	American Indian/Alaska Native	507	27.17	13.81	10	60

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Asian	17,564	36.43	12.14	10	60
	Black or African American	38,247	21.63	12.54	10	60
	Hispanic/Latino	65,085	24.21	13.03	10	60
	Native Hawaiian or Pacific Islander	358	31.01	13.45	10	60
	Two or More Races	9,558	28.69	13.57	10	60
	White	97,143	30.18	12.73	10	60
Economic Status*	Not Economically Disadvantaged	123,355	31.24	12.90	10	60
	Economically Disadvantaged	96,131	22.42	12.59	10	60
English Learner Status	Non-English Learner	183,158	28.70	13.40	10	60
	English Learner	35,250	21.15	12.06	10	60
Disabilities	Students without Disabilities	188,159	28.95	13.22	10	60
	Student with Disability (SWD)	38,991	20.60	12.44	10	60

Note.

SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Grade 10 subgroup statistics for ELA/L, Reading, and Writing scale scores are presented in Table 12.6.¹⁰ Mean scores were very similar to what was observed for grades 3 through 8. Corresponding tables for grades 9 and 10 are presented in Appendix 12.5.

Table 12.6 Subgroup Performance for ELA/L Scale Scores: Grade 10

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		7,069	740.86	49.34	650	850
Gender	Female	3,446	748.54	47.92	650	850
	Male	3,617	733.48	49.52	650	850
Ethnicity	American Indian/Alaska Native	n/r	n/r	n/r	n/r	n/r
	Asian	283	764.70	46.52	650	850
	Black or African American	3,016	723.24	46.41	650	850
	Hispanic/Latino	1,472	739.63	46.64	650	850
	Native Hawaiian or Pacific Islander	54	756.54	43.47	650	850
	Two or More Races	477	752.41	44.60	650	850
	White	1,590	765.13	45.09	650	850
Economic Status*	Not Economically Disadvantaged	n/r	n/r	n/r	n/r	n/r
	Economically Disadvantaged	2,162	714.49	45.39	650	850
English Learner Status	Non-English Learner	n/r	n/r	n/r	n/r	n/r
	English Learner	468	702.08	39.55	650	819
Disabilities	Students without Disabilities	5,456	749.49	47.48	650	850

¹⁰ Table A.12.47 in Appendix 12.5 is identical to Table 12.6.

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Student with Disability (SWD)	1,324	709.22	44.94	650	850
Reading Summative Score		7,069	48.52	20.41	10	90
Gender	Female	3,446	50.67	19.66	10	90
	Male	3,617	46.44	20.88	10	90
Ethnicity	American Indian/Alaska Native	n/r	n/r	n/r	n/r	n/r
	Asian	283	58.36	19.56	10	90
	Black or African American	3,016	40.89	18.85	10	90
	Hispanic/Latino	1,472	47.50	19.31	10	90
	Native Hawaiian or Pacific Islander	54	53.43	17.58	16	90
	Two or More Races	477	54.47	18.75	10	90
	White	1,590	59.29	18.40	10	90
Economic Status*	Not Economically Disadvantaged	n/r	n/r	n/r	n/r	n/r
	Economically Disadvantaged	2,162	37.16	18.15	10	90
English Learner Status	Non-English Learner	n/r	n/r	n/r	n/r	n/r
	English Learner	468	31.72	15.52	10	78
Disabilities	Students without Disabilities	5,456	51.91	19.62	10	90
	Student with Disability (SWD)	1,324	36.09	19.21	10	90
Writing Summative Score		7,069	30.42	13.09	10	60
Gender	Female	3,446	32.86	12.68	10	60
	Male	3,617	28.07	13.04	10	60
Ethnicity	American Indian/Alaska Native	n/r	n/r	n/r	n/r	n/r
	Asian	283	36.08	11.42	10	60
	Black or African American	3,016	26.63	12.85	10	60
	Hispanic/Latino	1,472	30.53	12.58	10	60
	Native Hawaiian or Pacific Islander	54	35.06	11.33	10	60
	Two or More Races	477	32.35	12.09	10	60
	White	1,590	35.32	12.34	10	60
Economic Status*	Not Economically Disadvantaged	n/r	n/r	n/r	n/r	n/r
	Economically Disadvantaged	2,162	24.74	12.70	10	60
English Learner Status	Non-English Learner	n/r	n/r	n/r	n/r	n/r
	English Learner	468	22.57	11.57	10	51
Disabilities	Students without Disabilities	5,456	32.53	12.59	10	60
	Student with Disability (SWD)	1,324	22.56	12.40	10	60

Note. SD = standard deviation; n/r = not reported due to n<20.

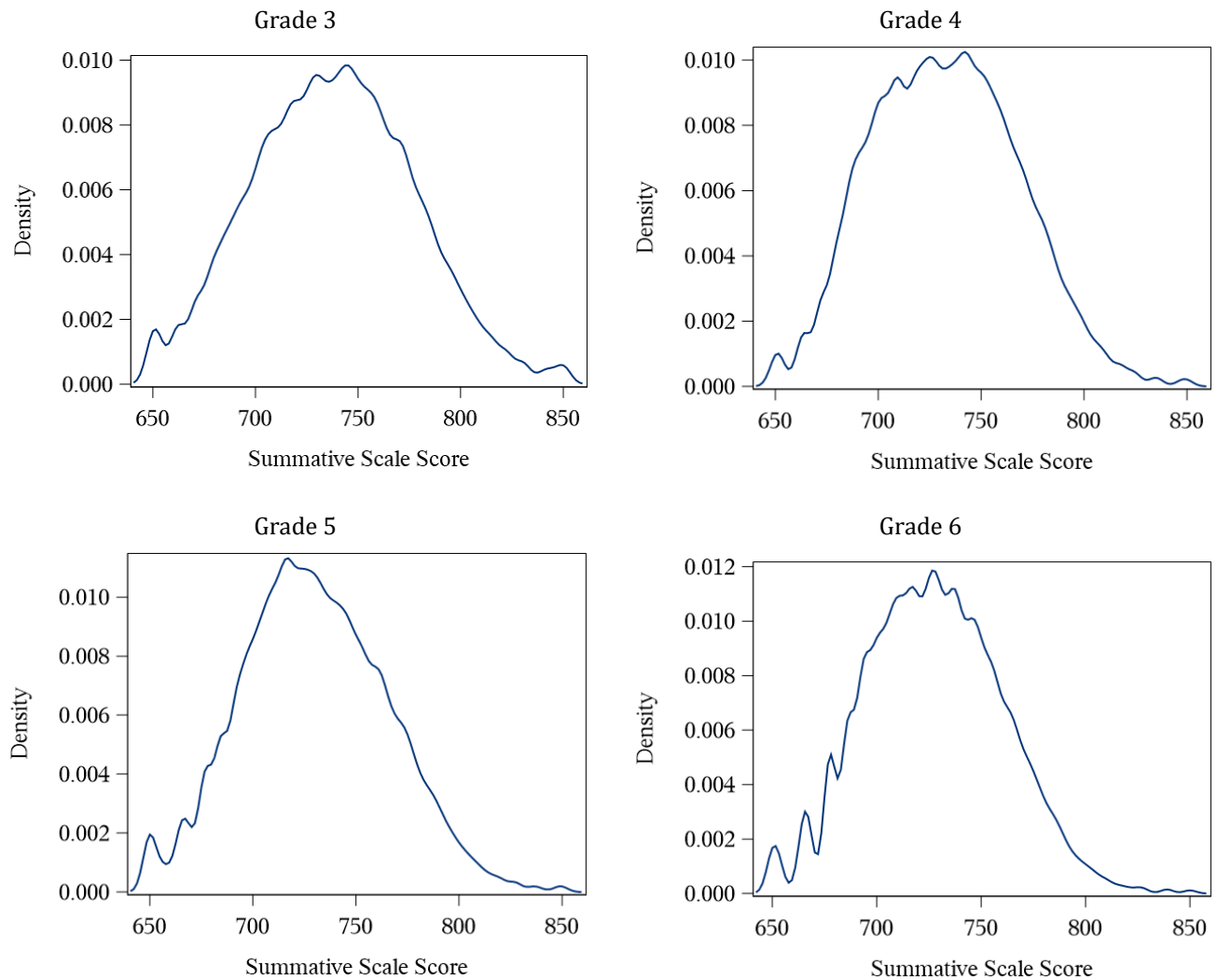
*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

12.4.4 Score Distributions for Mathematics

Figure 12.5 graphically represents the distributions of scale scores for grades 3 through 8 mathematics. The y -axis for these distributions ranges from 0 to .02 and the x -axis from 650 to 850. Scale score distributions generally peaked between approximately 700 and the Level 4 performance level cut of 750. Figure 12.6 graphically represents the distributions of scale scores for Algebra I, Geometry, Algebra II, and Integrated Mathematics I and II. Scale score distributions generally peaked between approximately 700 and the 750 Level 4 performance level cut score for Algebra I and Geometry. Integrated Mathematics results are omitted from this section due to low sample size.

12.4.5 Scale Score Cumulative Frequencies for Mathematics

The cumulative frequency distribution for the summative scale score is presented in Appendix 12.4 for mathematics assessments.



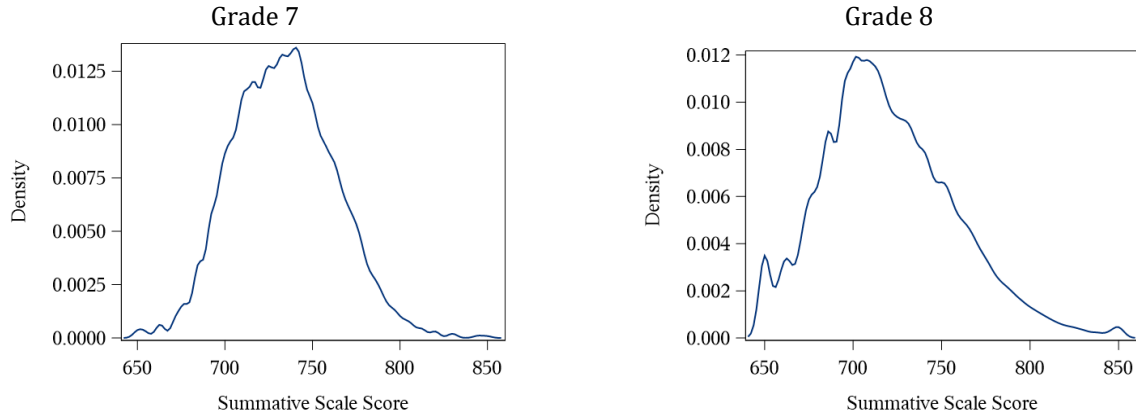


Figure 12.5 Distributions of Mathematics Scale Scores: Grades 3–8

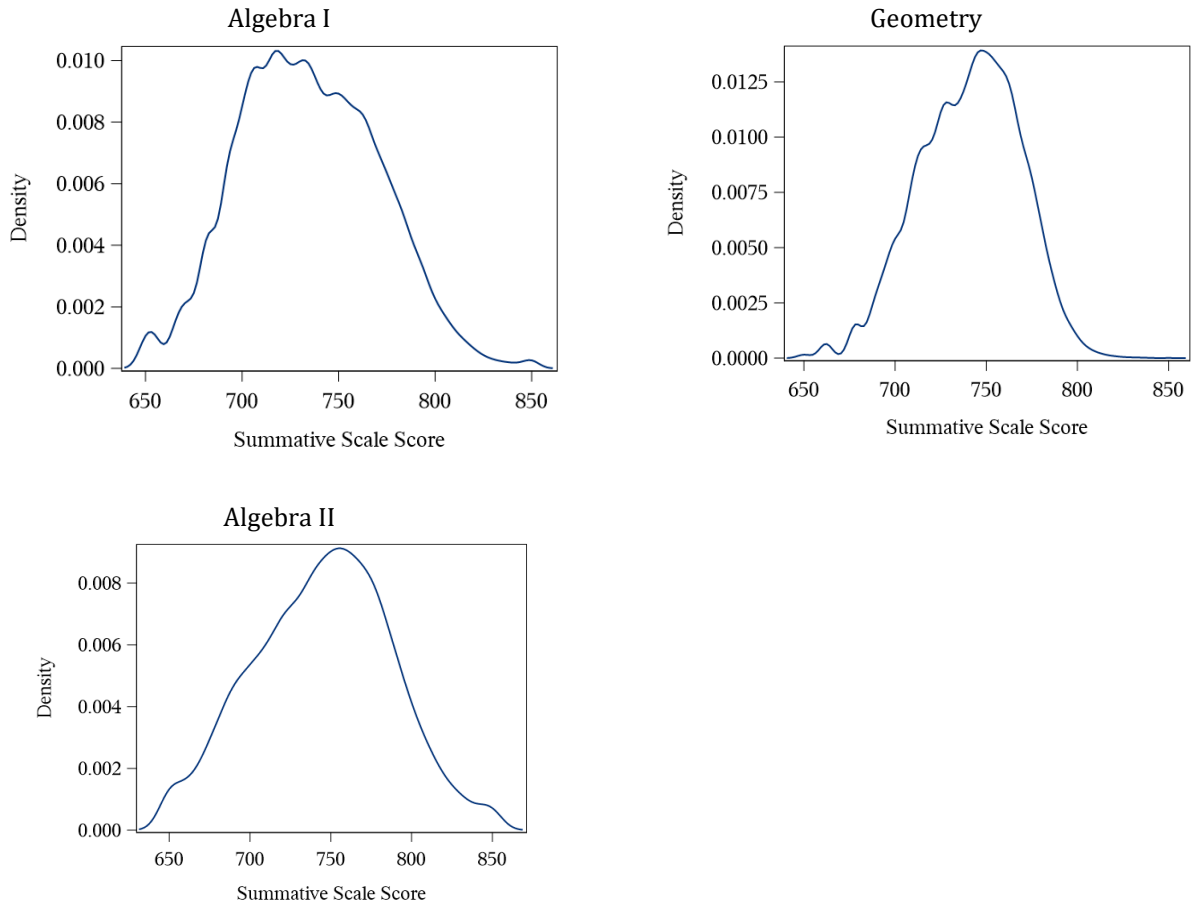


Figure 12.6 Distributions of Mathematics Scale Scores: High School

12.4.6 Summary Scale Score Statistics for Mathematics Groups

Subgroup statistics for mathematics scale scores are presented in Tables 12.7 and 12.8¹¹ for grade 3 and Algebra I, respectively. Grade 3 subgroup statistics are presented in Table 12.7.¹² Students using the Spanish language form tended to have lower mean scores. Corresponding tables for all grades/courses are presented in Appendix 12.5.

Table 12.7 Subgroup Performance for Mathematics Scale Scores: Grade 3

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		230,053	737.99	39.38	650	850
Gender	Female	112,938	736.26	37.99	650	850
	Male	117,109	739.67	40.61	650	850
Ethnicity	American Indian/Alaska Native	508	738.03	40.49	650	850
	Asian	17,645	771.70	37.65	650	850
	Black or African American	38,055	714.02	34.27	650	850
	Hispanic/Latino	66,053	724.77	34.31	650	850
	Native Hawaiian or Pacific Islander	359	746.36	38.97	650	850
	Two or More Races	9,537	743.19	40.24	650	850
	White	97,179	749.69	35.75	650	850
Economic Status*	Not Economically Disadvantaged	124,021	751.91	37.40	650	850
	Economically Disadvantaged	96,292	719.20	34.00	650	850
English Learner Status	Non-English Learner	182,726	741.86	39.36	650	850
	English Learner	36,588	718.55	33.08	650	850
Disabilities	Students without Disabilities	189,138	742.03	38.44	650	850
	Student with Disability (SWD)	38,894	719.38	38.40	650	850
Language Form	Spanish	5,063	708.73	29.38	650	850

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table 12.8 Subgroup Performance for Mathematics Scale Scores: Algebra I

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		115,174	734.83	36.12	650	850
Gender	Female	56,085	735.00	34.88	650	850
	Male	58,986	734.66	37.27	650	850
Ethnicity	American Indian/Alaska Native	167	734.59	37.89	650	850
	Asian	11,431	767.76	34.96	650	850

¹¹ Due to omitted demographic values, subgroup sample sizes in these tables may not sum to total sample size.

¹² Table A.12.48 in Appendix 12.5 is identical to Table 12.7. Table A.12.54 in Appendix 12.5 is identical to Table 12.8.

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Black or African American	19,311	716.48	30.43	650	850
	Hispanic/Latino	37,136	721.06	31.33	650	850
	Native Hawaiian or Pacific Islander	269	745.09	34.38	650	850
	Two or More Races	3,019	743.77	35.26	650	850
	White	43,602	745.36	32.61	650	850
Economic Status*	Not Economically Disadvantaged	74,152	742.29	36.10	650	850
	Economically Disadvantaged	34,995	718.79	31.10	650	850
English Learner Status	Non-English Learner	99,804	737.72	35.72	650	850
	English Learner	7,893	702.97	25.80	650	850
Disabilities	Students without Disabilities	92,916	739.65	35.44	650	850
	Student with Disability (SWD)	21,857	714.74	31.81	650	850
Language Form	Spanish	2,563	698.87	22.15	650	786

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

12.5 Interpreting Claim Scores and Subclaim Scores

12.5.1 Interpreting Claim Scores

ELA/L assessments provide separate claim scale scores for both Reading and Writing. The claim scale scores and the summative scale score are on different scales; therefore, the sum of the scale scores for each claim will not equal the summative scale score. Reading scale scores range from 10 to 90 and Writing scale scores range from 10 to 60.

The claim scores can be interpreted by comparing a student's claim scale score to the average performance for the school, district, and state. The Individual Student Report provides the student scale score results and the average scale score results for the school, district, and state.

12.5.2 Interpreting Subclaim Scores

Within each reporting category are specific skill sets (subclaims) students demonstrate on the summative assessments. Subclaim categories are not reported using scale scores or performance levels. Subclaim performance for the assessments is reported using graphical representations that indicate how the student performed relative to the Level 3 and Level 4 performance levels for the content area.

Subclaim indicators represent how well students performed in a subclaim category relative to Level 3 and Level 4 thresholds for the items associated with the subclaim category. To determine a student's subclaim performance, the Level 3 and Level 4 thresholds corresponding to the IRT-based performance for the items for a given subclaim determined the reference points for *Approached Expectations* and *Did Not Yet Meet Expectations* or *Partially Met Expectations*, respectively.

Student performance for each subclaim is marked with a subclaim performance indicator:

- An up arrow for the specified subclaim indicates that the student *Met or Exceeded Expectations*, meaning that the student's subclaim performance reflects a level of proficiency consistent with Performance Level 4 or 5. Students in this subclaim category are likely academically well prepared to engage successfully in further studies in the subclaim content area and may need instructional enrichment.
- A bidirectional arrow for the specified subclaim indicates that the student *Approached Expectations*, meaning that the student's subclaim performance reflects a level of proficiency consistent with Performance Level 3. Students in this subclaim category likely need academic support to engage successfully in further studies in the subclaim content area.
- A down arrow for the specified subclaim indicates that the student *Did Not Yet Meet or Partially Met Expectations*, meaning that the student's subclaim performance reflects a level of proficiency consistent with Performance Level 1 or 2. Students in this subclaim category are likely not academically well prepared to engage successfully in further studies in the subclaim content area. Such students likely need instructional interventions to increase achievement in the subclaim content area.

Section 13: Reliability

13.1 Overview

Reliability focuses on the extent to which differences in test scores reflect true differences in the knowledge, ability, or skill being tested rather than fluctuations due to chance. Thus, reliability measures the consistency of the scores across conditions that can be assumed to differ at random, especially which form of the test the student is administered and which persons are assigned to score responses to constructed-response questions. In statistical terms, the variance in the distributions of test scores, essentially the differences among individuals, is partly due to real differences in the knowledge, skill, or ability being tested (true variance) and partly due to random errors in the measurement process (error variance). Reliability is an estimate of the proportion of the total variance that is true variance.

There are several different ways of estimating reliability. The type of raw score reliability estimate reported here is an internal-consistency measure, which is derived from analysis of the consistency of the performance of individuals across items within a test. It is used because it serves as a good estimate of alternate forms reliability, but it does not take into account form-to-form variation due to lack of test form parallelism, nor is it responsive to day-to-day variation due to, for example, the student's state of health or the testing environment. The scale score reliability results use a modified measure of internal consistency that accounts for the conversions between raw scores and scale scores.

Reliability coefficients range from 0 to 1. The higher the reliability coefficient for a set of scores, the more likely students would be to obtain very similar scores upon repeated testing occasions, if the students do not change in their level of the knowledge or skills measured by the test. The reliability estimates in the tables to follow attempt to answer the question, "How consistent would the scores of these students be over replications of the entire testing process?"

Reliability of classification estimates the proportion of students who are accurately classified into proficiency levels. There are two kinds of classification reliability statistics: decision accuracy and decision consistency. Decision accuracy is the agreement between the classifications actually made and the classifications that would be made if the test scores were perfectly reliable. Decision consistency is the agreement between the classifications that would be made on two independent forms of the test.

Another index is inter-rater reliability for the human-scored constructed-response items, which measures the agreement between individual raters (scorers). The inter-rater reliability coefficient answers the question, "How consistent is the scoring such that a set of similarly trained raters would produce similar scores to those obtained?"

Standard error of measurement (SEM) quantifies the amount of error in the test scores. SEM is the extent by which students' scores tend to differ from the scores they would receive if the test were perfectly reliable. As the SEM increases, the variability of students' observed scores is likely to increase across repeated testing. Observed scores with large SEMs pose a challenge to the valid interpretation of a single test score.

Reliability and SEM estimates were calculated at the full assessment level, and at the claim and subclaim levels. In addition, conditional SEMs were calculated and reported in Appendix 13, Tables A.13.1 through A.13.17.

13.2 Reliability and SEM Estimation

13.2.1 Raw Score Reliability Estimation

Coefficient alpha (Cronbach, 1951), which measures internal consistency reliability, is the most commonly used measure of reliability. Coefficient alpha is estimated by substituting sample estimates for the parameters in the formula below:

$$\alpha = \frac{n}{n-1} \left[1 - \frac{\sum_{i=1}^n \sigma_i^2}{\sigma_X^2} \right] \quad (13-1)$$

where n is the number of items, σ_i^2 is the variance of scores on the i th item, and σ_X^2 is the variance of the total score (sum of scores on the individual items). Other things being equal, the more items a test includes, the higher the internal consistency reliability.

Since the test forms have mixed item types (dichotomous and polytomous items), it is more appropriate to report stratified alpha (Feldt & Brennan, 1989). Stratified alpha is a weighted average of coefficient alphas for item sets with different maximum score points or “strata.” Stratified alpha is a reliability estimate computed by dividing the test into parts (strata), computing alpha separately for each part, and using the results to estimate a reliability coefficient for the total score. Stratified alpha is used here because different parts of the test consist of different item types and may measure different skills. The formula for the stratified alpha is

$$\rho_{strata} = 1 - \frac{\sum_{h=1}^H \sigma_{x_h}^2 (1 - \alpha_h)}{\sigma_X^2} \quad (13-2)$$

where $\sigma_{X_h}^2$ is the variance for part h of the test, σ_X^2 is the variance of the total scores, and α_h is coefficient alpha for part h of the test. Estimates of stratified alpha are computed by substituting sample estimates for the parameters in the formula. The average stratified alpha is a weighted average of the stratified alphas across the test forms.

The formula for the standard error of measurement is

$$\sigma_E = \sigma_X \sqrt{1 - \rho_{XX'}} \quad (13-3)$$

where σ_X is the standard deviation of the test raw score and $\rho_{xx'}$ is the reliability estimated by substitution of appropriate statistics for the parameters in equation 13-1 or 13-2.

In this section, reliability estimates are reported for overall summative scores, claim scores, and subclaim scores. Estimates are also reported for subgroups for summative scores. Cronbach’s alpha and stratified alpha coefficients are influenced by test length, test characteristics, and sample characteristics (Cortina, 1993; Lord & Novick, 1968; Tavakol & Dennick, 2011). As test length decreases and samples become smaller and more homogeneous, lower estimates of alpha are obtained (Pike & Hudson, 1998; Tavakol & Dennick, 2011). A decrease in the number of items may result in a decrease in stratified alpha estimates. The decrease in sample

size and the homogeneity of the samples is likely to result in lower stratified alpha estimates. A smaller, more homogenous sample will likely result in lower stratified alpha estimates. Moderate-to-acceptable ranges of reliability tend to exceed .5 (Cortina, 1993; Schmitt, 1996). Estimates lower than .5 may indicate a lack of internal consistency. Additional analyses investigate whether lower estimates of alpha are due to restriction in range of the sample. In these cases, the alpha estimates are not appropriate measures of internal consistency. As a result, sample-free reliability estimates are also provided, such as scale score reliability (Kolen et al., 1996).

13.2.2 Scale Score Reliability Estimation

Like the stratified alpha coefficients, scale score reliability coefficients range from 0 to 1. The higher the reliability coefficient for a set of scores, the more likely individuals would be to obtain similar scores upon repeated testing occasions, if the students do not change in their level of the knowledge or skills measured by the test. Because the scale scores are computed from a total score and do not have an item-level component, a stratified alpha coefficient cannot be computed for scale scores. Instead, Kolen et al.'s (1996) method for scale score reliability was used.

The general formula for a reliability coefficient,

$$\rho = 1 - \frac{\sigma^2(E)}{\sigma^2(X)}, \quad (13-4)$$

involves the error variance, $\sigma^2(E)$, and the total score variance, $\sigma^2(X)$. Using Kolen et al.'s (1996) method, conditional raw score distributions are estimated using Lord and Wingersky's (1984) recursion formula. The conditional raw score distributions are transformed into conditional scale score distributions. Denote X as the raw sum score ranging from 0 to X , and s as a resulting scale score after transformation. The conditional distribution of scale scores is written as $P(X = x | \theta)$. The mean and variance, $\sigma^2[s(X)]$, of this distribution can be computed using these scores and their associated probabilities.

The average error variance of the scale scores is computed as

$$\sigma^2(Error_{scale}) = \int_{\theta} \sigma^2(s(X) | \theta) g(\theta) d\theta \quad (13-5)$$

where $g(\theta)$ is the ability distribution. The square root of the error variance is the conditional standard error of measurement of the scale scores.

Just as the reliability of raw scores is one minus the ratio of error variance to total variance, the reliability of scale scores is one minus the ratio of the average variance of measurement error for scale scores to the total variance of scale scores,

$$\rho_{scale} = 1 - \frac{\sigma^2(Error_{scale})}{\sigma^2[s(X)]} \quad (13-6)$$

The Windows program POLYSEM (Kolen, 2004) was used to estimate scale score error variance and reliability.

13.3 Reliability Results for Total Group

13.3.1 Raw Score Reliability Results

Tables 13.1 and 13.2 summarize test reliability estimates for the total testing group for English language arts/literacy (ELA/L) and mathematics, respectively. The tables provide the average reliability, which is estimated by averaging the internal consistency estimates computed for all the individual forms of the test and the raw score SEMs. In addition, the number of forms, the sample size of the minimum reliability, sample size of the maximum reliability, and the average maximum possible score for each set of tests are provided. Estimates were calculated only for groups of 100 or more students administered a specific test form.

English Language Arts/Literacy

The average reliability estimates for grades 3 through 10 ELA/L range from a low of .87 to a high of .91; note that grade 9 had a low sample size. The average raw score SEM is consistently between about 6 percent and 8 percent of the maximum possible score.

Table 13.1 Summary of ELA/L Test Reliability Estimates for Total Group

Grade Level	Number of Forms	Avg. Max Possible Score	Avg. Raw Score SEM	Average Reliability	Minimum Reliability N	Alpha	Maximum Reliability N	Alpha
3	5	54	3.69	0.89	603	0.81	117,193	0.90
4	5	70	4.32	0.89	879	0.84	114,919	0.91
5	5	73	4.37	0.90	9,534	0.84	129,428	0.90
6	5	73	4.58	0.90	684	0.85	118,609	0.92
7	5	71	4.51	0.91	704	0.82	114,605	0.91
8	5	71	4.63	0.90	705	0.84	116,341	0.92
9	4	71	5.08	0.87	4,554	0.81	344	0.91
10	5	73	4.98	0.88	167	0.79	3,188	0.90

Mathematics

The average reliability estimates for mathematics assessments range from .89 to .93. The raw score SEM consistently ranges from about 5 to 7 percent of the maximum score.

Table 13.2 Summary of Mathematics Test Reliability Estimates for Total Group

Grade Level	Number of Forms	Avg. Max Possible Score	Avg. Raw Score SEM	Average Reliability	Minimum Reliability N	Alpha	Maximum Reliability N	Alpha
3	6	52	3.13	0.93	1,043	0.91	78,760	0.93
4	6	52	3.02	0.93	411	0.90	76,069	0.94
5	6	52	3.01	0.92	363	0.87	80,635	0.93
6	6	52	3.03	0.92	350	0.87	862	0.94
7	5	52	3.14	0.91	29,369	0.91	779	0.95
8	6	52	3.00	0.89	564	0.77	73,879	0.89
A1	6	55	2.82	0.91	326	0.90	46,782	0.92
GO	6	55	3.12	0.90	160	0.88	1,472	0.90
A2	6	55	3.30	0.89	6,087	0.88	697	0.93

Note. A1=Algebra I; GO=Geometry; A2=Algebra II.

13.3.2 Scale Score Reliability Results

Tables 13.3 and 13.4 summarize scale score reliability estimates for the total testing group for ELA/L and mathematics for spring 2022. The tables provide average reliabilities by grade/course, which are estimated by averaging the reliability estimates computed for all forms of the test within the grade/course level. In addition, the number of forms, the total sample size across all forms, and the average maximum possible score for each set of tests are provided. Scale score reliability requires an ability distribution, which is not reasonable to assume for Integrated Mathematics due to the low sample sizes.

English Language Arts/Literacy

Reliability estimates for ELA/L are presented in Table 13.3. Average reliabilities range from .87 to .89. The average SEM ranges from 10.18 to 14.8.

Table 13.3 Summary of ELA/L Test Pre-Equated Scale Score Reliability Estimates for Total Group

Grade Level	Number of Forms	Avg. Scale Score SEM	Avg. Scale Score Reliability	Min. Scale Score Reliability	Max. Scale Score Reliability
3	5	14.19	0.87	0.87	0.88
4	5	12.41	0.87	0.86	0.88
5	5	11.48	0.87	0.85	0.88
6	5	10.18	0.89	0.89	0.9
7	5	12.43	0.89	0.88	0.89
8	5	12.62	0.88	0.87	0.9
9	5	12.55	0.89	0.87	0.9
10	5	14.8	0.89	0.87	0.9

Mathematics

The scale score reliability estimates for the mathematics assessments are presented in Table 13.4. Average scale score reliability estimates for the grades 3 through 8 mathematics assessments range from .85 to .92. For the high school assessments, these quantities range from .85 to .87. For grades 3 through 8, the average scale score SEM ranges from 9.24 to 14.19. For high school tests, the average scale score SEM ranges from 10.82 to 15.37.

Table 13.4 Summary of Mathematics Test Scale Score Reliability Estimates for Total Group

Grade Level	Number of Forms	Avg. Scale Score SEM	Avg. Raw Score Reliability	Min. Scale Score Reliability	Max. Scale Score Reliability
3	6	10.09	0.91	0.91	0.91
4	6	9.28	0.92	0.91	0.92
5	6	9.75	0.9	0.9	0.91
6	6	9.66	0.9	0.89	0.9
7	6	9.24	0.89	0.88	0.89
8	6	14.19	0.85	0.84	0.85
A1	6	12.39	0.87	0.86	0.88
GO	6	10.82	0.86	0.85	0.87
A2	6	15.37	0.85	0.84	0.86

Note. A1=Algebra I, GO=Geometry, A2=Algebra II.

13.4 Reliability Results for Subgroups of Interest

When the sample size was sufficiently large, raw score reliability and SEM were estimated for the groups identified for differential item functioning analysis. Estimates were calculated only for groups of 100 or more students administered a specific test form.

Tables 13.5 and 13.6 summarize test reliability for groups of interest for ELA/L grade 3 and mathematics grade 3, respectively. Corresponding information is provided in Appendix 13.1 for all ELA/L and mathematics grades. For each group, the average, minimum, and maximum reliability estimates are listed, as well as the sample sizes of the reported minimum and maximum reliabilities. Note that reliability estimates are dependent on score variance, and subgroups with smaller variance are likely to have lower reliability estimates than the total group.

13.4.1 Reliability Results for Gender

English Language Arts/Literacy

The average reliability estimates and the average SEMs for males and females reflect the corresponding reliabilities for the total group. For most tests, the reliabilities between males and females are equal or within .02. The SEMs for females were slightly higher than for males for ELA/L assessments.

Mathematics

As with the ELA/L test components, the average reliability estimates and SEMs for males and females reflect the corresponding reliabilities for the total group. For most tests, the reliabilities between males and females are equal or within .03. The SEMs for females are slightly higher than for males for the majority of tests.

13.4.2 Reliability Results for Ethnicity

English Language Arts/Literacy

The majority of the average reliabilities for the ethnicity groups are .01 to .03 lower than for the total group. There is not a consistent difference among the average reliabilities for white, Black/African American, Asian/Pacific Islander, Hispanic/Latino, and multiple-ethnicity students, with the majority of the reliabilities between .84 and .89. Average SEMs were generally slightly higher for white and Asian/Pacific Islander students than for Black/African American and Hispanic/Latino students.

Mathematics

As with the ELA/L reliabilities, the reliabilities for ethnicity groups are marginally lower than for the total group of students. While there is variation across tests, the average reliabilities are often highest for multiple-ethnicity students. The average SEMs reflect the total group SEMs. Average SEMs were generally higher for white, Asian/Pacific Islander, and multiple-ethnicity students than for Hispanic, Black/African American, and American Indian/Alaska Native students.

13.4.3 Reliability Results for Special Education Needs

English Language Arts/Literacy

The average reliabilities for five groups of students (economically disadvantaged, not economically disadvantaged, non-English learner, students with disabilities, and students without disabilities) are generally equal to or .01 to .02 less than the average reliability for the total group of students. Average reliabilities for English learner students are lower, ranging from .81 to .85. The SEMs are generally higher for the larger student groups (not economically disadvantaged students, non-English learner students, and students without disabilities).

Mathematics

The average reliabilities for the larger student groups (not economically disadvantaged, non-English learner, and students without disabilities) are generally equal to or .01 to .04 less than the average reliability for the total group of students. For economically disadvantaged, English learner, and students with disabilities, the average reliabilities are lower than those for the total group. The SEMs are generally higher for the larger student groups (not economically disadvantaged students, non-English learner students, and students without disabilities).

13.4.4 Reliability Results for Students Taking Accommodated Forms

English Language Arts/Literacy

Reliability information for accommodated forms is sparse due to small sample sizes or because the form was not administered. Reliabilities for text-to-speech forms tended to be lower than the overall reliabilities, while those for closed-caption forms tended to be higher.

Mathematics

The text-to-speech forms had sufficient sample sizes for reliability and SEM estimation across grades/subjects, except for high school courses where the sample was not sufficient. For almost all tests, text-to-speech reliabilities are similar to the total group reliabilities, with SEMs slightly lower than the total group SEMs.

13.4.5 Reliability Results of Students Taking Translated Forms

Mathematics

There were sufficient numbers of students taking the Spanish-language form for reliability and SEM estimation for grades 3 through 8. The average reliability ranged from .67 to .87. The SEMs are generally lower for the students administered the Spanish-language forms.

Table 13.5 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 3

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	54	3.69	0.89	603	0.81	117,193	0.90
Gender							
Male	54	3.52	0.88	348	0.79	36,316	0.88
Female	54	3.66	0.89	739	0.79	36,154	0.89
Ethnicity							
Black/African American	54	3.47	0.87	176	0.76	18,749	0.89
Asian/Pacific Islander	54	3.98	0.87	8,053	0.85	9,349	0.87
Hispanic/Latino	54	3.57	0.88	238	0.76	32,478	0.89
American Indian/Alaska Native	54	3.79	0.89	265	0.89	206	0.90
Multiple	54	3.71	0.89	271	0.84	4,879	0.90
White	54	3.78	0.87	153	0.82	51,277	0.88
Special Instruction Needs							
Economically Disadvantaged	54	3.50	0.87	401	0.75	47,431	0.88
Not Economically Disadvantaged	54	3.82	0.87	166	0.81	64,135	0.88
English Learner	54	3.40	0.85	157	0.75	18,143	0.86
Non-English Learner	54	3.74	0.88	412	0.81	94,294	0.89
Students with Disabilities	54	3.31	0.88	540	0.80	15,577	0.90
Students without Disabilities	54	3.77	0.89	82,449	0.88	100,768	0.89
Students Taking Accommodated Forms							
ASL	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	54	2.90	0.83	7,524	0.83	7,524	0.83

Note. n/a = not applicable.

Table 13.6 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 3

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	52	3.13	0.93	1,043	0.91	78,760	0.93
Gender							
Male	52	3.11	0.93	388	0.90	22,171	0.93
Female	52	3.12	0.92	271	0.91	155	0.93
Ethnicity							
Black/African American	52	2.88	0.91	210	0.89	10,627	0.92
Asian/Pacific Islander	52	3.15	0.91	7,001	0.91	1,764	0.93
Hispanic/Latino	52	3.02	0.91	289	0.89	18,729	0.92
American Indian/Alaska Native	52	3.18	0.93	154	0.93	166	0.93
Multiple	52	3.16	0.93	3,447	0.93	955	0.93
White	52	3.23	0.91	35,969	0.91	9,280	0.93
Special Instruction Needs							
Economically Disadvantaged	52	2.96	0.91	243	0.89	27,047	0.91
Not Economically Disadvantaged	52	3.21	0.92	45,291	0.91	141	0.94
English Learner	52	2.93	0.90	219	0.89	8,067	0.91
Non-English Learner	52	3.16	0.93	719	0.92	21,108	0.93
Students with Disabilities	52	2.94	0.92	9,423	0.90	9,358	0.93
Students without Disabilities	52	3.16	0.92	244	0.91	68,814	0.93
Students Taking Accommodated Forms							
ASL	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	52	3.00	0.93	27,682	0.93	30,959	0.93
Students Taking Translated Forms							
Spanish Language	52	2.72	0.87	4,930	0.87	4,930	0.87

Note. n/a = not applicable.

13.5 Reliability Results for ELA/L Claims and Subclaims

Participating states and agencies developed subclaims in addition to major claims based on the Common Core State Standards. ELA/L has two major claims relating to Reading and Writing. The major claim for Reading is that students read and comprehend a range of sufficiently complex texts independently. The major claim for Writing is that students write effectively when using and/or analyzing sources. Refer to Table 13.7 for a summary of the ELA/L claims and subclaims.

Table 13.7 Descriptions of ELA/L Claims and Subclaims

English Language Arts/Literacy		
Major Claim Subclaim		Description
Reading	Reading Literature	Students demonstrate comprehension and draw evidence from readings of grade-level, complex literary text.
Reading	Reading Information	Students demonstrate comprehension and draw evidence from readings of grade-level, complex informational text.
Reading	Reading Vocabulary	Students use context to determine the meaning of words and phrases.
Writing	Writing Written Expression	Students produce clear and coherent writing in which the development, organization, and style are appropriate to the task, purpose, and audience.
Writing	Writing Knowledge Language and Conventions	Students demonstrate knowledge of conventions and other important elements of language.

Reliability indices were calculated for each major claim and subclaim. Table 13.8 presents the average reliability estimates for all forms of the test at the specified grade and testing mode for the ELA/L tests. In order to assist in understanding the reliability estimates, the range of maximum number of points for each major claim and subclaim is also provided. Reliabilities from grade 11 tended to be lower than the other grades, so they are omitted from the descriptions in the following paragraphs. However, they can be found in Table 13.8.

The average reliabilities for the Reading claim for grades 3 through 8 and 10 range from .8 to .87. They are based on maximum scores of 40–44 points per form, except for grade 3 (30–31 points), grade 7 (40–42 points), and grade 8 (36–44 points). The Writing claim average reliabilities are based on a lower number of points than those for the Reading claim, and are slightly lower, ranging from .79 to .84. The reliabilities for the Writing claim for grade 3 are based on a maximum raw score of 24 points, and the average reliabilities for grades 4 and 5 are based on between 27 and 30 points per form. The average reliabilities for the grades 5 through 10 Writing claims are based on a maximum score of 30 points.

The average reliabilities of the Reading Literature subclaim scores vary from .64 to .75. The maximum number of points per form ranges from 12 to 18. The average reliabilities of the Reading Information subclaim scores vary from .60 to .71, with 11–22 points per form. The average reliabilities of the Reading Vocabulary subclaim scores vary from .41 to .66. The maximum number of points per form for this subclaim ranges from 8 to 14.

The Writing Written Expression subclaim is based on 18 points for grade 3 and 21–24 points for grades 4 and 5. Grades 6 through 10 are based on 24 points for all forms. The average reliabilities range from .74 to .85. The Writing Knowledge of Language and Conventions subclaims are all based on six points. The reliabilities range from .81 to .88.

Table 13.8 Average ELA/L Reliability Estimates for Subscores

	Reading: Total		Reading: Literature		Reading: Information		Reading: Vocabulary		Writing: Total		Writing Expression		Writing: Knowledge Language and Conventions	
Grade Level	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability
3	30–31	0.85	11–12	0.72	11–11	0.67	8–8	0.59	24–24	0.79	18–18	0.74	6–6	0.83
4	40–44	0.85	16–18	0.72	12–16	0.62	8–12	0.58	27–30	0.80	21–24	0.79	6–6	0.84
5	40–44	0.86	16–18	0.68	14–16	0.65	8–12	0.66	27–30	0.82	21–24	0.80	6–6	0.83
6	40–44	0.87	16–18	0.75	14–16	0.71	8–14	0.60	30–30	0.83	24–24	0.82	6–6	0.85
7	40–42	0.84	16–16	0.64	14–16	0.69	8–10	0.59	30–30	0.84	24–24	0.85	6–6	0.88
8	36–44	0.83	12–18	0.68	14–16	0.69	8–10	0.46	30–30	0.84	24–24	0.85	6–6	0.87
9	40–44	0.80	12–18	0.67	14–20	0.60	8–14	0.41	30–30	0.82	24–24	0.80	6–6	0.82
10	40–44	0.83	12–14	0.68	16–22	0.68	10–12	0.44	30–30	0.82	24–24	0.81	6–6	0.81

13.6 Reliability Results for Mathematics Subclaims

For mathematics, there are four subclaims related to whether students are on track or ready for college and careers:

- Subclaim A: Students solve problems involving the major content for their grade/course level with connections to the Standards for Mathematical Practice.
- Subclaim B: Students solve problems involving the additional and supporting content for their grade/course level with connections to the Standards for Mathematical Practice.
- Subclaim C: Students express grade/course-level appropriate mathematical reasoning by constructing viable mathematical arguments and critiquing the reasoning of others, and/or attending to precision when making mathematical statements.
- Subclaim D: Students solve real-world problems with a degree of difficulty appropriate to the grade/course by applying knowledge and skills articulated in the standards and by engaging particularly in the modeling practice.

Reliability estimates were calculated for each subclaim for mathematics. Table 13.9 presents the average reliability estimates for mathematics subclaims.

Subclaims with greater numbers of points tend to have greater reliability estimates. The Major Content subclaim has the largest number of points for each assessment and, accordingly, has higher average reliabilities than the other three subclaims. For grades 3 through 8, Algebra I, Geometry, and Algebra II, the median of the average reliabilities for the Major Content range from .73 to .86. The maximum number of points per form range from 16 to 21.

The median of the average reliabilities for the Additional and Supporting Content subclaim for grades 3 through 8, Algebra I, Geometry, and Algebra II ranges from .52 to .73. The maximum number of points per form for this subclaim ranges from 9 to 12.

The average reliabilities for Mathematics Reasoning range from .57 to .73 for grades 3 through 8, Algebra I, Geometry, and Algebra II. The maximum number of points for this subclaim is 10 for all grades and forms.

For the Modeling Practice subclaim, the average reliabilities for grades 3 through 8, Algebra I, Geometry, and Algebra II range from .54 to .74. The number of points is 12 for grades 3 through 8 and 15 for all high school courses except Algebra I (9–15 points).

Table 13.9 Average Mathematics Reliability Estimates for Subscores

	Major Content		Additional & Supporting Content		Mathematics Reasoning		Modeling Practice	
Grade Level	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability	Range of Raw Score	Average Reliability
3	20–20	0.86	10–10	0.73	10–10	0.62	12–12	0.74
4	21–21	0.86	9–9	0.68	10–10	0.73	12–12	0.69
5	20–20	0.83	10–10	0.65	10–10	0.68	12–12	0.67
6	20–20	0.79	10–10	0.69	10–10	0.69	12–12	0.67
7	20–20	0.80	10–10	0.61	10–10	0.67	12–12	0.71
8	20–20	0.73	10–10	0.52	10–10	0.57	12–12	0.69
A1	17–23	0.77	9–9	0.68	10–10	0.73	9–15	0.59
GO	18–18	0.81	12–12	0.57	10–10	0.62	15–15	0.63
A2	16–18	0.79	12–12	0.68	10–10	0.58	15–15	0.54

Note. Integrated Mathematics I, II, and III had insufficient sample sizes. A1 = Algebra I; GO = Geometry; A2 = Algebra II.

13.7 Reliability of Classification

The reliability of the classifications for the students was calculated using the computer program BB-CLASS (Brennan, 2004), which operationalizes a statistical method developed by Livingston and Lewis (1993, 1995). As Livingston and Lewis (1993, 1995) explain, this method uses information from the administration of one test form (i.e., distribution of scores, the minimum and maximum possible scores, the cut points used for classification, and the reliability coefficient) to estimate two kinds of statistics, decision accuracy and decision consistency. Decision accuracy refers to the extent to which the classifications of students based on their scores on the test form agree with the classifications made on the basis of the classifications that would be made if the test scores were perfectly reliable. Decision consistency refers to the agreement between these classifications based on two nonoverlapping, equally difficult forms of the test.

Decision consistency values are always lower than the corresponding decision accuracy values, because in decision consistency, both of the classifications are subject to measurement error. In decision accuracy, only one of the classifications is based on a score that contains an error(s). It is not possible to know which students were accurately classified, but it is possible to estimate the proportion of the students who were accurately classified. Similarly, it is not possible to know which students would be consistently classified if they were retested with another form, but it is possible to estimate the proportion of the students who would be consistently classified.

13.7.1 English Language Arts/Literacy

Table 13.11 provides information about the accuracy and the consistency of two types of classifications made on the basis of the summative scale scores on the grades 3 through 10 ELA/L assessments. The columns labeled “Exact Level” provide the estimates of the indices based on classifications of students into one of five performance levels. The columns labeled “Level 4 or Higher vs. 3 or Lower” provide the estimates of the indices based on classifications of students as being either in one of the upper two levels (Levels 4 and 5) or in one of the lower three levels (Levels 1, 2, and 3). Performance Level 4 is considered the college- and career-readiness standard on the summative assessments.

The table shows that for classifying each student into one of the five performance levels, the proportion accurately classified ranges from .64 to .72; the proportion who would be consistently classified on two different test forms ranges from .54 to .62. For classifying each student as being at Level 4 or higher vs. being at Level 3 or lower, the proportion accurately classified ranges from .89 to .91; the proportion who would be consistently classified this way on two different test forms ranges from .85 to .88.

Table 13.10 Reliability of Classification: Summary for ELA/L

Level	Decision Accuracy: Proportion Accurately Classified		Decision Consistency: Proportion Consistently Classified	
	Exact Level	Level 4 or Higher vs. 3 or Lower	Exact Level	Level 4 or Higher vs. 3 or Lower
3	0.69	0.91	0.60	0.87
4	0.69	0.90	0.59	0.86
5	0.72	0.91	0.62	0.87
6	0.72	0.91	0.62	0.87
7	0.70	0.91	0.60	0.88
8	0.70	0.91	0.59	0.87
9	0.67	0.89	0.56	0.85
10	0.64	0.90	0.54	0.86

Table 13.11 provides more detailed information about the accuracy and the consistency of the classification of students into performance levels for ELA/L grade 3. Each cell in the 5-by-5 table shows the estimated proportion of students who would be classified into a particular combination of performance levels. The sum of the five bold values on the diagonal is approximately equal to the level of decision accuracy or consistency presented in Table 13.10. For “Level 4 and Higher vs. 3 and Lower” found in Table 13.10, the sum of the shaded values in Table 13.11 is approximately equal to the level of decision accuracy or consistency presented in Table 13.10. Note that the sums based on values in Table 13.11 may not match exactly to the values in Table 13.10 due to truncation and rounding.

Detailed information for all ELA/L spring results are provided in Appendix 13, Tables A.13.18 through A.13.25. The structure of these tables is the same as that of Table 13.11 and the values in the tables should be interpreted in the same manner. Table 13.11 includes the same information as Table A.13.18.

Table 13.11 Reliability of Classification: Grade 3 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.22	0.03	0.00	0.00	0.00	0.25
	700–724	0.04	0.10	0.05	0.00	0.00	0.19
	725–749	0.00	0.04	0.12	0.04	0.00	0.21
	750–809	0.00	0.00	0.05	0.25	0.03	0.33
	810–850	0.00	0.00	0.00	0.01	0.01	0.02
Decision Consistency	650–699	0.21	0.05	0.01	0.00	0.00	0.27
	700–724	0.04	0.07	0.06	0.01	0.00	0.19
	725–749	0.01	0.04	0.09	0.05	0.00	0.19
	750–809	0.00	0.01	0.07	0.22	0.02	0.32
	810–850	0.00	0.00	0.00	0.02	0.01	0.03

13.7.2 Mathematics

Table 13.12 provides information about the accuracy and the consistency of two types of classifications made on the basis of the summative scale scores on the mathematics assessments. For the grades 3 through 8 mathematics tests, the table shows that for classifying each student into one of the five performance levels,

the proportion accurately classified ranges from .68 to .76; the proportion who would be consistently classified on two different test forms ranges from .58 to .67. For the six high school mathematics courses, the table shows that for classifying each student into one of the five performance levels, the proportion accurately classified ranges from .69 to .73; the proportion who would be consistently classified on two different test forms ranges from .59 to .63.

For classifying each student as being at Level 4 or higher vs. being at Level 3 or lower, for the grades 3 through 8 mathematics tests, the proportion accurately classified ranges from .92 to .93; the proportion who would be consistently classified on two different test forms is .89 to .90 for grades 3 and 8. For high school mathematics courses, the proportion accurately classified as being at Level 4 or higher vs. being at Level 3 or lower ranges from .89 to .90; the proportion who would be consistently classified on two different test forms ranges from .85 to .87.

Appendix 13, Tables A.13.26 through A.13.34, provide more detailed information about the accuracy and the consistency of the classification of students into performance levels for mathematics. Each cell in the 5-by-5 table shows the estimated proportion of students who would be classified into a particular combination of performance levels.

Table 13.12 Reliability of Classification: Summary for Mathematics

Level	Decision Accuracy: Proportion Accurately Classified		Decision Consistency: Proportion Consistently Classified	
	Exact Level	Level 4 or Higher vs. 3 or Lower	Exact Level	Level 4 or Higher vs. 3 or Lower
3	0.74	0.92	0.65	0.89
4	0.76	0.93	0.67	0.90
5	0.74	0.92	0.64	0.89
6	0.75	0.93	0.66	0.90
7	0.75	0.92	0.66	0.89
8	0.68	0.93	0.58	0.89
A1	0.72	0.90	0.62	0.87
GO	0.73	0.89	0.63	0.85
A2	0.69	0.89	0.59	0.85

Note. A1 = Algebra I; GO = Geometry; A2 = Algebra II.

13.8 Inter-Rater Agreement

Inter-rater agreement is the agreement between the first and second scores assigned to student responses. Inter-rater agreement measurements include exact, adjacent, and nonadjacent agreement. Pearson scoring staff used these statistics as one factor in determining the needs for continuing training and intervention on both individual and group levels. Table 13.13 displays both the expectations and the actual agreement percentages for perfect agreement and perfect plus adjacent agreement.

Table 13.13 Inter-Rater Agreement Expectations and Results

Subject	Score Point Range	Perfect Agreement Expectation	Perfect Agreement Result	Within One Point Expectation	Within One Point Result
Mathematics	0–1	90%	98%	96%	100%
Mathematics	0–2	80%	97%	96%	100%
Mathematics	0–3	70%	96%	96%	99%
Mathematics	0–4	65%	94%	95%	99%
Mathematics	0–5	65%	91%	95%	98%
Mathematics	0–6	65%	95%	95%	98%
ELA/L	Multi-trait	65%	83%	96%	100%

Note. A 0 or 1 score compared to a blank score will have a disagreement greater than 1 point.

Pearson's ePEN2 scoring system included comprehensive inter-rater agreement reports that allowed supervisory personnel to monitor both individual and group performance. Based on reviews of these reports, scoring experts targeted individuals for increased backreading and feedback and, if necessary, retraining. Table 13.13 shows that the actual percentages for perfect reader agreement were higher than the inter-rater agreement expectations, and the percentages for within one point were very close. Refer to Section 4 for more information on handscoring.

Section 14: Validity

14.1 Overview

The *Standards for Educational and Psychological Testing*, issued jointly by the American Educational Research Association [AERA], American Psychological Association [APA], and National Council on Measurement in Education [NCME] (2014), reports:

Validity refers to the degree to which evidence and theory support the interpretations of test scores for proposed uses of tests. Validity is, therefore, the most fundamental consideration in developing tests and evaluating tests. The process of validation involves accumulating relevant evidence to provide a sound scientific basis for the proposed score interpretations (p. 11).

The purpose of test validation is not to validate the test itself but to validate interpretations of the test scores for particular uses. Test validation is not a quantifiable property but an ongoing process, beginning at initial conceptualization and continuing throughout the lifetime of an assessment. Every aspect of an assessment provides evidence in support of its validity (or evidence of lack of validity), including design, content specifications, item development, and psychometric characteristics. The 2021–2022 operational assessments provided an opportunity to gather evidence of validity based on both test content and on the internal structure of the tests.

Pearson applies the principles of universal design, as articulated in materials developed by the National Center for Educational Outcomes at the University of Minnesota (Thompson et al., 2002).

14.2 Evidence Based on Test Content

Evidence based on content of achievement tests is supported by the degree of correspondence between test items and content standards. The degree to which the test measures what it claims to measure is known as construct validity. The summative assessments adhere to the principles of evidence-centered design, in which the standards to be measured (the Common Core State Standards) are identified, and the performance a student needs to achieve to meet those standards is delineated in the evidence statements. Test items are reviewed for adherence to universal design principles, which maximize the participation of the widest possible range of students.

Pearson and New Meridian built spreadsheets at the evidence statement level that incorporate the probability statements from the test blueprints and attrition rates at committee review and data review. The basis of our entire item development is driven by the use of these item development target spreadsheets. Before beginning item development, Pearson uses these target spreadsheets to develop an internal item development plan to correlate with the expectations of the test design. These are reviewed and approved by state or agency leads and New Meridian. All parties acknowledge that each assessment has multiple parts and each part specifies the types of tasks and standards eligible for assessment.

In addition to the evidence statements, content is aligned through the articulation of performance in the performance level descriptors (PLDs). At the policy level, the PLDs include policy claims about the educational achievement of students who attain a particular performance level, and a broad description of the

grade-level knowledge, skills, and practices students performing at a particular achievement level are able to demonstrate. Those policy-level descriptors are the foundation for the subject- and grade-specific PLDs, which, along with the evidence frameworks, guide the development of the items and tasks.

The college- and career-ready determinations in English language arts/literacy (ELA/L) and mathematics describe the academic knowledge, skills, and practices students must demonstrate to show readiness for success in entry-level, credit-bearing college courses and relevant technical courses. The states and agencies determined that this level means graduating from high school and having at least a 75 percent likelihood of earning a grade of “C” or better in credit-bearing courses without the need for remedial coursework. After reviewing the standards and assessment design, the Governing Board (made up of the K–12 education chiefs in participating states or agencies) in conjunction with the Advisory Committee on College Readiness (composed of higher education chiefs in the participating states or agencies), determined that students who achieve at Levels 4 and 5 on the final high school assessments are likely to have acquired the skills and knowledge to meet the definition of college- and career-readiness. To validate the determinations, a postsecondary educator judgment study and a benchmark study of the SAT, ACT, National Assessment of Educational Progress, Trends in International Mathematics and Science Study, Programme of International Student Assessment, and Progress in International Reading Literacy Study tests were conducted (McClarty et al., 2015).

Gathering construct validity evidence for the assessments is embedded in the process by which the assessment content is developed and validated. At each step in the assessment development process, participating states or agencies involved hundreds of educators, assessment experts, and bias and sensitivity experts in review of text, items, and tasks for accuracy, appropriateness, and freedom from bias. See Section 2 for an overview of the content development process. In the early stages of development, Pearson conducted research studies to validate the item and task development approach. One such study was a student task interaction study designed to collect data on the student’s experience with the assessment tasks and technological functionalities, as well as the amount of time needed for answering each task. Pearson also conducted a rubric choice study that compared the functioning of two rubrics developed to score the prose constructed-response (PCR) tasks in ELA/L. Quantitative and qualitative evidence was collected to support the use of a condensed or expanded trait scoring rubric in scoring student responses.

The items and tasks were field-tested prior to their use on an assessment. During the initial field-test administration in 2014, participating states and agencies collected feedback from students, test administrators, test coordinators, and classroom teachers on their experience with the assessments, including the quality of test items and student experience. Information pertaining to this process can be found at <https://resources.newmeridiancorp.org/research/>. The feedback from that survey was used to inform test directions, test timing, and the function of online task interactions. Performance data from the field test also informed the future development of additional items and tasks.

All item developers and item writers are provided with an electronic version of the accessibility guidelines and the linguistic complexity rubric. Items and passages are reviewed internally by accessibility and fairness experts trained in the principles of universal design and who become well versed in the accessibility guidelines. Items received internal review for alignment to evidence tables, task generation model, item selection guidelines, and accessibility and fairness reviews.

An important consideration when constructing test forms is recognition of items that may introduce construct-irrelevant variance. Such items should not be included on test forms to help ensure fairness to all subgroups of students. New Meridian convened bias and sensitivity committees to review all items. Additionally, content experts facilitated reviews of all items. All reviewers were trained using the bias and sensitivity guidelines, and the guidelines were used to review items and ELA/L passages. Accommodations were made available based on individual need documented in the student's approved Individualized Education Program (IEP), 504 Plan, or if required by the participating state or agency, an English Learner (EL) Plan. An accessibility specialist worked in consultation with the accessibility specialist to review forms and determine which forms should be used for students with accommodations.

The ELA/L and mathematics operational test forms, as described in Section 2, were carefully constructed to align with the test blueprints and specifications that are based on the Common Core State Standards (CCSS). During fall 2016, content experts representing various participating states and agencies, along with other content experts, held a series of meetings to review the operational forms for ELA/L and mathematics. These meetings provided an opportunity to evaluate test forms in their entirety and recommend changes. Requested item replacements were accommodated to the extent possible while a concerted effort was made maintain the integrity of the various linking designs required for the operational test analyses. Psychometricians were available throughout this process to provide guidance with regard to implications of item replacements for the linking and statistical requirements.

Further information regarding the college- and career-ready content standards, PLDs s, and accessibility features and accommodations is provided at <http://resources.newmeridiancorp.org/>.

14.3 Evidence Based on Internal Structure

Analyses of the internal structure of a test typically involve studies of the relationships among test items and/or test components (i.e., subclaims) in the interest of establishing the degree to which the items or components appear to reflect the construct on which a test score interpretation is based (AERA, APA, & NCME, 2014, p. 16). The term *construct* is used here to refer to the characteristics that a test is intended to measure; in the case of the operational tests, the characteristics of interest are the knowledge and skills defined by the test blueprint for ELA/L and for mathematics.

The summative assessments provide a full summative test score, Reading claim score, and Writing claim score as well as ELA/L subclaim and mathematics subclaim scores. The goal of reporting at this level is to provide criterion-referenced data to assess the strengths and weaknesses of a student's achievement in specific components of each content area. This information can then be used by teachers to plan for further instruction, to plan for curriculum development, and to report progress to parents. The results can also be used as one factor in making administrative decisions about program effectiveness, teacher effectiveness, class grouping, and needs assessment.

14.3.1 Intercorrelations

The ELA/L full summative tests comprise two claim scores, Reading (RD) and Writing (WR), and five subclaim scores—Reading Literature (RL), Reading Information (RI), Reading Vocabulary (RV), Writing Written Expression (WE), and Writing Knowledge Language and Conventions (WKL). The RD claim score is a composite of RL, RI, and RV. The Writing claim score, a composite of WE and WKL, comprises only PCR items, and the same PCR items are in each subclaim. The ELA/L operational test analyses were performed by

evaluating the separate trait scores of WE and WKL, and for some PCR items, also RL or RI; therefore, the trait scores were used for the intercorrelations.

The mathematics full summative tests have four subclaim scores: Major Content (MC), Mathematical Reasoning (MR), Modeling Practice (MP), and Additional and Supporting Content (ASC).

High total group internal consistencies as well as similar reliabilities across subgroups provide additional evidence of validity. High reliability of test scores implies that the test items within a domain are measuring a single construct, which is a necessary condition for validity when the intention is to measure a single construct. Refer to Section 13 for reliability estimates for the overall population, subgroups of interest, as well as for claims and subclaims for ELA/L and subclaims for mathematics.

Another way to assess the internal structure of a test is through the evaluation of correlations among scores. These analyses were conducted between the ELA/L Reading and Writing claim scores and the ELA/L subclaims (RL, RI, RV, WE, and WKL) and between the mathematics subclaims. If these components within a content area are strongly related to each other, this is evidence of unidimensionality.

A series of tables are provided to summarize the results for the spring 2022 administration. Tables 14.1 through 14.8 present the Pearson correlations observed between the ELA/L Reading and Writing claim scores and subclaim scores for each grade. The tables provide the weighted average intercorrelations by averaging the intercorrelations computed for all the core operational forms of the test within each grade level. The total sample size across all forms is provided in the upper triangle portion of the tables. The subclaim reliabilities (from Section 13) are reported along the diagonal. The WR, WE, and WKL scores tended to be highly correlated; this is expected given that these three intercorrelations are based on the trait scores from the same Writing items. RL, RI, and RV, all subclaims of Reading, are moderately to highly correlated. Additionally, the WR claim and the WE and WKL subclaims are moderately correlated with RD subclaims (of RL, RI, and RV). These moderate-to-high ELA/L intercorrelations among the subclaims are sufficiently high to provide evidence that the ELA/L tests are unidimensional. The moderate intercorrelations among the subclaims and claims suggest the claims may be sufficient for individual student reporting.

The intercorrelations and reliability estimates for mathematics are provided in Tables 14.9 through 14.17. The shaded values along the diagonal are the reliabilities as reported in Section 13. The average intercorrelations are provided in the lower portion of the table, and the total sample sizes are provided in the upper portion of the table. Please refer to Appendix 12.1 ("Form Composition") for information about the number of items and number of score points in each claim and subclaim.

The mathematics intercorrelations are moderate. The main observable pattern in the mathematics intercorrelations is that the MC subclaim generally has slightly higher correlations with the ASC, MR, and MP subclaims; the intercorrelations among the ASC, MR, and MP subclaims are usually slightly lower. The mathematics intercorrelations are sufficiently high to suggest that the mathematics tests are likely to be unidimensional with some minor secondary dimensions.

Table 14.1 Average Intercorrelations and Reliability between Grade 3 ELA/L Subclaims

	RD	RL	RI	RV	WR	WE	WKL
RD	0.86	229,122	229,122	229,122	229,122	229,122	229,122
RL	0.90	0.71	229,122	229,122	229,122	229,122	229,122
RI	0.88	0.68	0.67	229,122	229,122	229,122	229,122
RV	0.85	0.65	0.63	0.60	229,122	229,122	229,122
WR	0.71	0.64	0.70	0.53	0.79	229,122	229,122
WE	0.70	0.63	0.69	0.51	0.99	0.75	229,122
WKL	0.66	0.59	0.63	0.50	0.90	0.82	0.83

Note. RD = Reading; RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WR = Writing; WE = Written Expression; WKL = Writing Knowledge and Conventions.

Table 14.2 Average Intercorrelations and Reliability between Grade 4 ELA/L Subclaims

	RD	RL	RI	RV	WR	WE	WKL
RD	0.84	232,286	232,286	232,286	232,286	232,286	232,286
RL	0.91	0.71	232,286	232,286	232,286	232,286	232,286
RI	0.86	0.65	0.63	232,286	232,286	232,286	232,286
RV	0.83	0.65	0.60	0.58	232,286	232,286	232,286
WR	0.74	0.66	0.69	0.54	0.80	232,286	232,286
WE	0.73	0.66	0.69	0.54	0.99	0.79	232,286
WKL	0.69	0.62	0.64	0.52	0.92	0.87	0.84

Note. RD = Reading; RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WR = Writing; WE = Written Expression; WKL = Writing Knowledge and Conventions.

Table 14.3 Average Intercorrelations and Reliability between Grade 5 ELA/L Subclaims

	RD	RL	RI	RV	WR	WE	WKL
RD	0.86	237,781	237,781	237,781	237,781	237,781	237,781
RL	0.90	0.68	237,781	237,781	237,781	237,781	237,781
RI	0.86	0.65	0.66	237,781	237,781	237,781	237,781
RV	0.86	0.67	0.61	0.66	237,781	237,781	237,781
WR	0.71	0.63	0.69	0.53	0.82	237,781	237,781
WE	0.70	0.63	0.68	0.52	0.99	0.80	237,781
WKL	0.68	0.61	0.66	0.51	0.95	0.91	0.83

Note. RD = Reading; RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WR = Writing; WE = Written Expression; WKL = Writing Knowledge and Conventions.

Table 14.4 Average Intercorrelations and Reliability between Grade 6 ELA/L Subclaims

	RD	RL	RI	RV	WR	WE	WKL
RD	0.87	238,200	238,200	238,200	238,200	238,200	238,200
RL	0.91	0.75	238,200	238,200	238,200	238,200	238,200
RI	0.89	0.69	0.71	238,200	238,200	238,200	238,200
RV	0.83	0.67	0.63	0.60	238,200	238,200	238,200
WR	0.75	0.65	0.75	0.54	0.83	238,200	238,200
WE	0.74	0.64	0.75	0.53	1	0.82	238,200
WKL	0.73	0.63	0.73	0.53	0.97	0.94	0.85

Note. RD = Reading; RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WR = Writing; WE = Written Expression; WKL = Writing Knowledge and Conventions.

Table 14.5 Average Intercorrelations and Reliability between Grade 7 ELA/L Subclaims

	RD	RL	RI	RV	WR	WE	WKL
RD	0.86	245,076	245,076	245,076	245,076	245,076	245,076
RL	0.88	0.68	245,076	245,076	245,076	245,076	245,076
RI	0.89	0.65	0.72	245,076	245,076	245,076	245,076
RV	0.82	0.60	0.62	0.63	245,076	245,076	245,076
WR	0.74	0.63	0.74	0.52	0.84	245,076	245,076
WE	0.73	0.62	0.73	0.52	1	0.85	245,076
WKL	0.72	0.62	0.72	0.52	0.96	0.93	0.87

Note. RD = Reading; RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WR = Writing; WE = Written Expression; WKL = Writing Knowledge and Conventions.

Table 14.6 Average Intercorrelations and Reliability between Grade 8 ELA/L Subclaims

	RD	RL	RI	RV	WR	WE	WKL
RD	0.84	249,310	249,310	249,310	249,310	249,310	249,310
RL	0.89	0.70	249,310	249,310	249,310	249,310	249,310
RI	0.89	0.67	0.70	249,310	249,310	249,310	249,310
RV	0.75	0.55	0.53	0.51	249,310	249,310	249,310
WR	0.76	0.67	0.74	0.47	0.83	249,310	249,310
WE	0.75	0.66	0.73	0.47	1	0.84	249,310
WKL	0.75	0.67	0.73	0.48	0.97	0.95	0.86

Note. RD = Reading; RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WR = Writing; WE = Written Expression; WKL = Writing Knowledge and Conventions.

Table 14.7 Average Intercorrelations and Reliability between Grade 9 ELA/L Subclaims

	RD	RL	RI	RV	WR	WE	WKL
RD	0.83	103,678	103,678	103,678	103,678	103,678	103,678
RL	0.90	0.73	103,678	103,678	103,678	103,678	103,678
RI	0.86	0.63	0.64	103,678	103,678	103,678	103,678
RV	0.75	0.53	0.51	0.50	103,678	103,678	103,678
WR	0.72	0.60	0.73	0.46	0.80	103,678	103,678
WE	0.71	0.59	0.73	0.45	1	0.78	103,678
WKL	0.72	0.60	0.72	0.46	0.97	0.96	0.81

Note. RD = Reading; RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WR = Writing; WE = Written Expression; WKL = Writing Knowledge and Conventions.

Table 14.8 Average Intercorrelations and Reliability between Grade 10 ELA/L Subclaims

	RD	RL	RI	RV	WR	WE	WKL
RD	0.82	7,069	7,069	7,069	7,069	7,069	7,069
RL	0.87	0.66	7,069	7,069	7,069	7,069	7,069
RI	0.90	0.66	0.66	7,069	7,069	7,069	7,069
RV	0.77	0.57	0.56	0.42	7,069	7,069	7,069
WR	0.73	0.65	0.71	0.47	0.81	7,069	7,069
WE	0.72	0.65	0.71	0.47	1	0.79	7,069
WKL	0.73	0.66	0.71	0.48	0.98	0.97	0.80

Note. RD = Reading; RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WR = Writing; WE = Written Expression; WKL = Writing Knowledge and Conventions.

Table 14.9 Average Intercorrelations and Reliability between Grade 3 Mathematics Subclaims

	MC	ASC	MR	MP
MC	0.85	229,942	229,942	229,942
ASC	0.80	0.73	229,942	229,942
MR	0.71	0.65	0.62	229,942
MP	0.77	0.72	0.69	0.74

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

Table 14.10 Average Intercorrelations and Reliability between Grade 4 Mathematics Subclaims

	MC	ASC	MR	MP
MC	0.86	232,885	232,885	232,885
ASC	0.74	0.67	232,885	232,885
MR	0.77	0.69	0.73	232,885
MP	0.73	0.68	0.72	0.69

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

Table 14.11 Average Intercorrelations and Reliability between Grade 5 Mathematics Subclaims

	MC	ASC	MR	MP
MC	0.82	238,435	238,435	238,435
ASC	0.71	0.64	238,435	238,435
MR	0.76	0.66	0.68	238,435
MP	0.76	0.66	0.72	0.66

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

Table 14.12 Average Intercorrelations and Reliability between Grade 6 Mathematics Subclaims

	MC	ASC	MR	MP
MC	0.79	238,564	238,564	238,564
ASC	0.72	0.69	238,564	238,564
MR	0.74	0.67	0.68	238,564
MP	0.75	0.67	0.70	0.67

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

Table 14.13 Average Intercorrelations and Reliability between Grade 7 Mathematics Subclaims

	MC	ASC	MR	MP
MC	0.80	235,663	235,663	235,663
ASC	0.69	0.60	235,663	235,663
MR	0.75	0.64	0.66	235,663
MP	0.76	0.67	0.71	0.71

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

Table 14.14 Average Intercorrelations and Reliability between Grade 8 Mathematics Subclaims

	MC	ASC	MR	MP
MC	0.71	211,478	211,478	211,478
ASC	0.66	0.49	211,478	211,478
MR	0.66	0.59	0.53	211,478
MP	0.69	0.63	0.64	0.67

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

Table 14.15 Average Intercorrelations and Reliability between Algebra I Subclaims

	MC	ASC	MR	MP
MC	0.74	115,095	115,095	115,095
ASC	0.75	0.63	115,095	115,095
MR	0.73	0.69	0.69	115,095
MP	0.68	0.62	0.63	0.57

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

Table 14.16 Average Intercorrelations and Reliability between Geometry Subclaims

	MC	ASC	MR	MP
MC	0.79	40,403	40,403	40,403
ASC	0.69	0.55	40,403	40,403
MR	0.68	0.57	0.56	40,403
MP	0.74	0.62	0.70	0.58

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

Table 14.17 Average Intercorrelations and Reliability between Algebra II Subclaims

	MC	ASC	MR	MP
MC	0.74	13,941	13,941	13,941
ASC	0.71	0.60	13,941	13,941
MR	0.70	0.66	0.52	13,941
MP	0.66	0.62	0.63	0.52

Note. MC = Major Content; ASC = Additional and Supporting Content; MR = Mathematical Reasoning; MP = Modeling Practice.

14.3.2 Reliability

Additionally, the reliability analyses presented in Section 13 of this technical report provide information about the internal consistency of the summative assessments. Internal consistency is typically measured via correlations among the items on an assessment and provides an indication of how much the items measure the same general construct. The reliability estimates, computed using coefficient alpha (Cronbach, 1951), are presented in Tables 13.1 and 13.2 and are along the diagonals of Tables 14.1 through 14.17¹³ The average reliabilities for ELA/L and mathematics summative assessments range from .87 up to .93 Appendix Tables A.13.1 through A.13.8 summarize test reliability for groups of interest for ELA/L grades 3 through 10, and Appendix Tables A.13.9 through A.13.17 summarize test reliability for groups of interest for mathematics grades/courses. Along with the subclaim intercorrelations, the reliability estimates indicate that the items within each assessment are measuring the same construct and provide further evidence of unidimensionality.

14.3.3 Local Item Dependence

In addition to the intercorrelations for ELA/L and mathematics, local item independence was evaluated. Local independence is one of the primary assumptions of item response theory (IRT) that states the probability of success on one item is not influenced by performance on other items, when controlling for ability level. This implies that ability or theta accounts for the associations among the observed items. Local item dependence (LID) when present essentially overstates the amount of information predicted by the IRT model. It can exert other undesirable psychometric effects and represents a threat to validity since other factors besides the construct of interest are present. Classical statistics are also affected when LID is present since estimates of test reliability like IRT information can be inflated (Zenisky et al., 2003).

The LID issue affects the choice of item scoring in IRT calibrations. Specifically, if evidence suggests these items indeed have local dependence, then it might be preferable to sum the item scores into clusters or testlets as a method of minimizing LID. However, if these items do not appear to have strong local item dependence, then retaining the scores as individual item scores in an IRT calibration is preferred since more information concerning item properties is retained. During the initial operational administration of the summative assessments in spring 2015, a study that included two methods of investigating the presence of LID was conducted. A description of the methods and study findings is summarized below.

First, analyses of the internal consistency in items and testlets were conducted under classical test theory (Wainer & Thissen, 2001) as a way to evaluate the degree of LID. Two estimates of Cronbach's alpha (Cronbach, 1951) were compared based on individual items in a test and those clustered into testlets. Cronbach's alpha is formulated as:

$$\alpha = \frac{l}{l-1} \frac{\sum_{i \neq i'} \sigma_{ii'}}{\sigma_X^2} \quad (14-1)$$

where l is the total number of items, $\sigma_{ii'}$ is the covariance of items i and i' ($i \neq i'$), and σ_X^2 is the variance of total scores. To compute an alpha coefficient, sample standard deviations and variances are substituted for the $\sigma_{ii'}$ and σ_X^2 . The alpha for the total test based on individual items is compared with those that form

¹³ Section 13 provides information on the computations of the reliability estimates.

testlets based on larger subparts. If the item-level configuration has appreciably higher levels of internal consistency compared with the testlets, LID may be present.

For IRT-based methods, local dependence can be evaluated using statistics such as Q_3 (Yen, 1984). The item residual is the difference between observed and expected performance. The Q_3 index is the correlation between residuals of each item pair defined as

$$d_i = (O - \hat{E}), \quad (14-2)$$

$$Q_3 = r(d_i, d_j) \quad (14-3)$$

where O is the observed score and \hat{E} is the expected value of O under a proposed IRT model and the index is defined as the correlation between the two item residuals.

LID manifests itself as a residual correlation that is nonzero and large. For Q_3 , LID can be either positive or negative. Positive (negative) LID indicates that performance is higher (lower) than expectation. The residual Q_3 correlation matrix can be inspected to determine if there are any blocks of locally dependent items (e.g., perhaps blocks of items belonging to the same reading passage). For Q_3 , the null hypothesis is that local independence holds. The expected value of Q_3 is $-1/n-1$ where n is the number of items such that the statistic shows a small negative bias. As a rule of thumb, item pairs with moderate levels of LID for Q_3 are $|.2|$ or greater. Significant levels of LID are present when the statistic is greater than $|.4|$. An alternative is to use the Fisher r to z transformation and evaluate the resulting p -values.

For the LID comparisons, the following eight test levels administered in spring 2015 were selected:

- Grade 4 for span 3–5 in ELA/L
- Grade 4 for span 3–5 in mathematics
- Grade 7 for span 6–8 in ELA/L
- Grade 7 for span 6–8 in mathematics
- Grade 10 for span 9–11 in ELA/L
- Integrated Mathematics II for Integrated Mathematics I–III
- Algebra I
- Algebra II

One spring 2015 computer-based test (CBT) form for each of the eight tests was selected that was roughly at the median in terms of test difficulty. For ELA/L, reading items were summed according to passage assignment. For mathematics, items were summed according to subclaims. Cronbach's alpha was computed for the entire sets of forms using the two different approaches as described above, one involving calculations at the item level and the second utilizing scores on summed items (i.e., testlets). A further description of the data is given in Table 14.18.

To cross-validate the internal consistency analysis, the Q_3 statistic was computed from spring CBT data based on grade 4 ELA/L and Integrated Mathematics II items. All items in the pool at that test level were included. The CBT item pool for grade 4 ELA/L contained 125 items while Integrated Mathematics II had 77 items.

The results for the internal consistency analysis are shown in Figure 14.1. In every instance, the item-level Cronbach's alpha is higher than in the testlet configuration. The greatest difference was for Algebra II, which showed a difference of .07. Although this was not unexpected, the magnitude of the differences in the respective alpha coefficients in general do not suggest a concerning level of LID. Table 14.19 shows the summary for the Q3 values. Figures 14.2 and 14.3 show graphs of the distribution of Q3 values. Most of the Q3 values were small and negative, again suggesting that LID is not at a level of concern. For these two test levels, the difference in the alpha coefficients was .03 and was consistent with the low values of Q3.

In summary, this investigation did not find evidence for the existence of pervasive LID. The results of both the internal consistency analyses and Q3 methods support a claim of minimal LID. For a multiple-choice-only test containing four reading passages with 5 to 12 items associated with a reading passage, Sireci et al. (1991) reported that testlet alpha was approximately 10 percent lower than the item-level coefficient. In comparison, the tests have complex test structures and exhibited smaller differences in alpha coefficients. In addition, the median Q3 values presented in Table 14.19 centered around the expectation of $-1/n-1$.

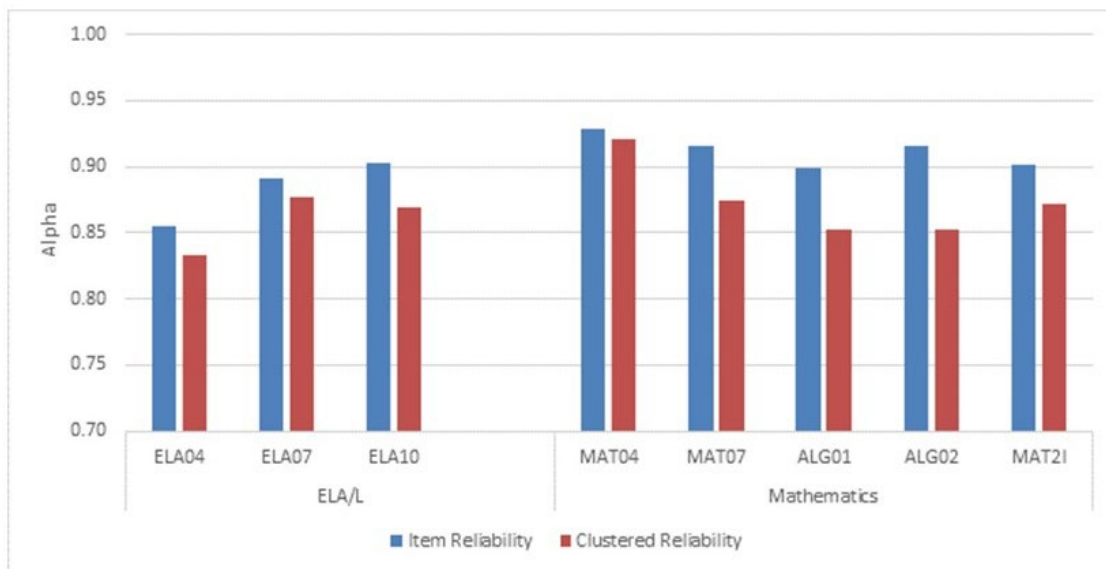


Figure 14.1 Comparison of Internal Consistency by Item and Cluster (Testlet)

Table 14.18 Conditions Used in LID Investigation and Results

Content	Grade/ Course	N Valid	N Complete	Percent Incomplete	No. Items	No. Tasks	Item Rel.	Task Rel.
ELA/L								
ELA/L	4	13,660	13,518	1.04	31	5	0.86	0.83
ELA/L	7	12,757	12,685	0.56	41	7	0.89	0.88
ELA/L	10	3,097	3,033	2.07	41	7	0.90	0.87
Mathematics								
Math	4	10,332	10,255	0.75	53	4	0.93	0.92
Math	7	10,295	10,188	1.04	50	6	0.92	0.87
Math	A1	5,072	4,885	3.69	52	6	0.90	0.85
Math	A2	4,982	4,769	4.28	54	6	0.92	0.85
Math	M2	2,708	2,645	2.33	51	6	0.90	0.87

Note. A1 = Algebra I; A2 = Algebra II; M2 = Integrated Mathematics II.

Table 14.19 Summary of Q3 Values for ELA/L Grade 4 and Integrated Mathematics II (Spring 2015)

Min.	Q1	Median	Mean	Q3	Max.	SD
ELA/L Grade 4						
-0.138	-0.047	-0.031	-0.031	-0.017	0.279	0.030
Integrated Mathematics II						
-0.160	-0.038	-0.017	-0.019	0.001	0.280	0.032

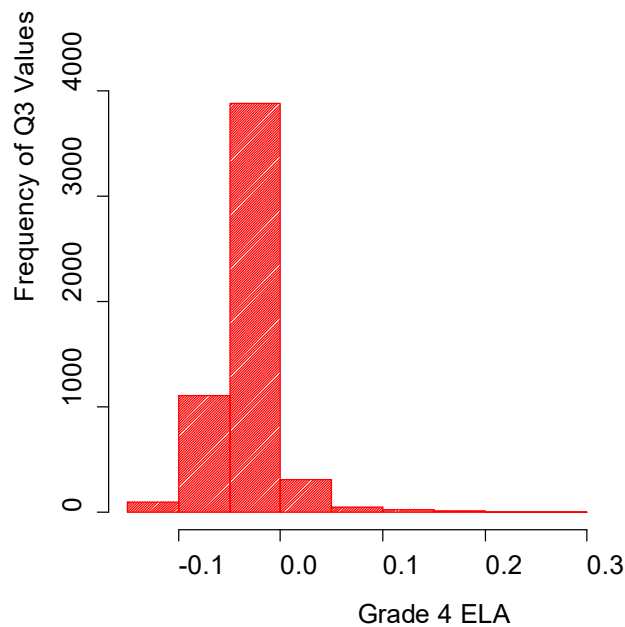


Figure 14.2 Distribution of Q3 Values for Grade 4 ELA/L (Spring 2015)

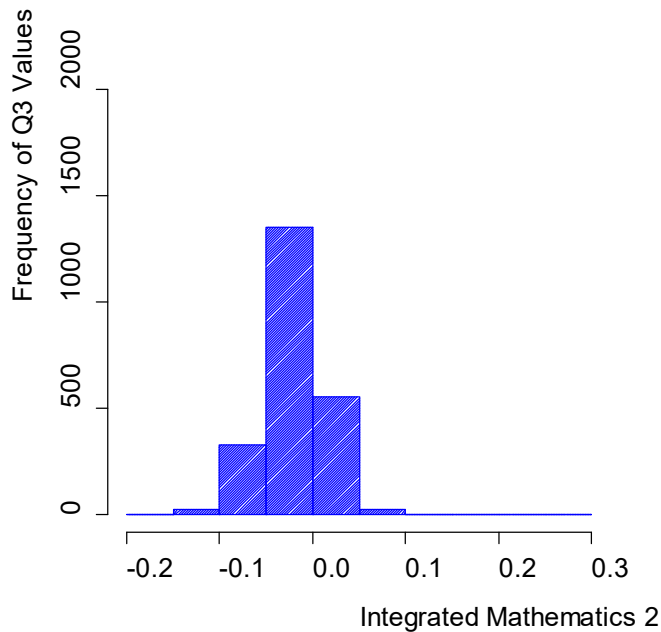


Figure 14.3 Distribution of Q3 Values for Integrated Mathematics II (Spring 2015)

14.4 Evidence Based on Relationships to Other Variables

Empirical results concerning the relationships between scores on a test and measures of other variables external to the test can also provide evidence of validity when these relationships are found to be consistent with the definition of the construct that the test is intended to measure. As indicated in the AERA, APA, and NCME standards (2014), the variables investigated can include other tests that measure the same construct and different constructs, criterion measures that scores on the test are expected to predict, as well as demographic characteristics of students that are expected to be related and unrelated to test performance.

The relationship of the scores across the ELA/L and mathematics assessments was evaluated using correlational analyses. Tables 14.20 through 14.25 present the Pearson correlations observed between the ELA/L scale scores and the mathematics scale scores for each grade. For grades 3 through 8, students must have a valid test score for both ELA/L and mathematics at the same grade level to be included in the tables. These tables provide the correlation in the lower triangle, and the sample size is provided in the upper triangle. In computing the correlations between a particular pair of ELA/L and mathematics tests, students must have taken both tests in spring 2022. ELA/L, Reading (RD), and Writing (WR) are moderately to highly correlated with mathematics; the correlations range from .68 up to .79 for grades 3 through 8. These correlations suggest that the ELA/L and mathematics tests are assessing different content. The higher intercorrelations between the ELA/L, Reading (RD), and Writing (WR) scores suggest stronger internal relationships when compared to the correlations with the mathematics content area.

The ELA/L and mathematics correlations for the high school tests are presented in Tables 14.26 through 14.28. Because students in high school can take the mathematics courses in different years (e.g., one student may take Algebra I in grade 9 while another student may take Algebra I in grade 10), the high school mathematics scores were correlated with several of the ELA/L grades (e.g., Algebra I correlated with both grades 9 and 10). Only correlations for pairings with total sample sizes of at least 100 are shown in the tables.

Blank cells indicate pairings with sample sizes less than 100. Across grades 8 through 11, ELA/L, Reading (RD), and Writing (WR) scores have correlations with high school mathematics tests that range from .48 to .71. Correlations between high school mathematics scores and corresponding ELA/L scores demonstrate low to moderate correlations.

Table 14.20 Correlations between ELA/L and Mathematics for Grade 3

	ELA/L	RD	WR	MA
ELA/L		227,746	227,746	227,746
RD	0.96		227,746	227,746
WR	0.89	0.74		227,746
MA	0.78	0.76	0.68	

Note. ELA/L = English language arts/literacy; RD = Reading; WR = Writing; MA = Mathematics.

Table 14.21 Correlations between ELA/L and Mathematics for Grade 4

	ELA/L	RD	WR	MA
ELA/L		230,777	230,777	230,777
RD	0.96		230,777	230,777
WR	0.88	0.74		230,777
MA	0.79	0.77	0.69	

Note. ELA/L = English language arts/literacy; RD = Reading; WR = Writing; MA = Mathematics.

Table 14.22 Correlations between ELA/L and Mathematics for Grade 5

	ELA/L	RD	WR	MA
ELA/L		236,312	236,312	236,312
RD	0.95		236,312	236,312
WR	0.87	0.71		236,312
MA	0.76	0.75	0.65	

Note. ELA/L = English language arts/literacy; RD = Reading; WR = Writing; MA = Mathematics.

Table 14.23 Correlations between ELA/L and Mathematics for Grade 6

	ELA/L	RD	WR	MA
ELA/L		236,112	236,112	236,112
RD	0.96		236,112	236,112
WR	0.88	0.74		236,112
MA	0.77	0.77	0.66	

Note. ELA/L = English language arts/literacy; RD = Reading; WR = Writing; MA = Mathematics.

Table 14.24 Correlations between ELA/L and Mathematics for Grade 7

	ELA/L	RD	WR	MA
ELA/L		232,991	232,991	232,991
RD	0.95		232,991	232,991
WR	0.89	0.73		232,991
MA	0.75	0.75	0.64	

Note. ELA/L = English language arts/literacy; RD = Reading; WR = Writing; MA = Mathematics.

Table 14.25 Correlations between ELA/L and Mathematics for Grade 8

	ELA/L	RD	WR	MA
ELA/L		210,023	210,023	210,023
RD	0.94		210,023	210,023
WR	0.89	0.72		210,023
MA	0.68	0.69	0.57	

Note. ELA/L = English language arts/literacy; RD = Reading; WR = Writing; MA = Mathematics.

Table 14.26 Correlations between ELA/L and Mathematics for High School

ELA/L	Mathematics Courses		
	A1	GO	A2
8	0.72 2,321	0.52 255	
9	0.68 8,120	0.67 1,915	0.60 382
10		0.56 261	

Note. ELA/L = English language arts/literacy; A1 = Algebra I; GO = Geometry; A2 = Algebra II.

Table 14.27 Correlations between ELA/L Reading and Mathematics for High School

RD	Mathematics		
	A1	GO	A2
8	0.71 2,321	0.57 255	
9	0.67 8,120	0.68 1,915	0.59 382
10		0.59 261	

Note. RD = Reading; A1 = Algebra I; GO = Geometry; A2 = Algebra II.

Table 14.28 Correlations between ELA/L Writing and Mathematics for High School

WR	Mathematics		
	A1	GO	A2
8	0.64 2,321	0.38 255	
9	0.59 8,120	0.60 1,915	0.48 382
10		0.42 261	

Note. WR = Writing; A1 = Algebra I; GO = Geometry; A2 = Algebra II.

14.5 Evidence from Special Studies

Several research studies were conducted to provide additional validity evidence for the participating state and agencies' goals of assessing more rigorous academic expectations, helping to prepare students for college and careers, and providing information back to teachers and parents about their students' progress toward college- and career-readiness. Some of the special studies conducted include the following:

- content alignment studies
- a benchmarking study
- a longitudinal study of external validity
- a mode comparability study
- a device comparability study
- Quality Testing Standards study

The following paragraphs briefly describe each of these studies.

14.5.1 Content Alignment Studies

In 2016, content of the ELA/L assessments at grades 5, 8, and 11 and the Algebra II and Integrated Mathematics II assessments were evaluated to determine how well the assessments were aligned to the CCSS (Doorey & Polikoff, 2016; Schultz et al., 2016). These content alignment studies were conducted by the Fordham Institute for grades 5 and 8 and by Human Resources Research Organization (HumRRO) for the high school assessments. Both of these studies used the same methodology by having content experts review the assessment items and answers (for the constructed-response items the rubrics were reviewed). The content experts then judged how well the items aligned to the CCSS, the depth of knowledge of the items, and the accessibility of the items to all students, including English learners and students with disabilities. The authors of both studies noted that the content experts reviewing the assessments were required to be familiar with the CCSS but could not be employed by participating organizations or be the writers of the CCSS. Therefore, an effort was made to eliminate any potential conflicts of interest.

The content studies had the individual content experts review and rate each item; then as a group the content experts came to a consensus on the final ratings for the content alignment, depth of knowledge, and accessibility to all students. In addition to the ratings, the content experts were asked to make comments that provided an explanation of their ratings; these comments were then used by the full group of content experts to provide narrative comments regarding the overall ratings and to provide feedback and recommendation about the assessment programs.

The assessment program was rated as Excellent Match for ELA/L content and depth and Good Match for mathematics content and depth for grades 5 and 8. However, for grade 11 ELA/L content was rated as Excellent Match but depth was rated as Limited/Uneven Match. The high school mathematics assessments were rated at Excellent Match for content and Good Match for depth.

The content studies noted some weaknesses and strengths of the assessments. For ELA/L, it was noted that the assessments include complex texts, a range of cognitive demands, and have a variety of item types. Furthermore, the ELA/L "assessments require close reading, assess writing to sources, research, and inquiry, and emphasize vocabulary and language skills" (Doorey & Polikoff, 2016). The grade 11 ELA/L assessment had a smaller range of depth and included items assessing the higher-demand cognitive level. A weakness of

the ELA/L assessments is the lack of a listening and speaking component. It was also suggested that the ELA/L assessments could be enhanced by the inclusion of a research task that requires the use of two or more sources of information.

The strengths of the mathematics assessments include assessments that are aligned to the major work for each grade level. While the grade 5 assessment includes a range of cognitive demand items, the grade 8 assessment includes a number of higher-demand items and may not fully assess the standards at the lowest level of cognitive demand. It was suggested that the grade 5 assessment could include more focus on the major work and the grade 8 assessment could include items at the lowest cognitive demand level. Additionally, the reviewers noted that some of the mathematics items should be carefully reviewed for editorial and mathematical accuracy.

The high school report noted that the assessment program incorporates a number of accessibility features and test accommodations for students with disabilities and for English learners. Furthermore, the assessments included items designed to accommodate the needs of students with disabilities.

In 2017, HumRRO conducted a study to evaluate the quality and alignment of ELA/L and mathematics assessments for grades 3, 4, 6, and 7 (Schultz et al., 2017). This alignment study followed a similar methodology as the 2016 study. For the study, cognitive complexity was consistent with the current assessments' definition. An item's cognitive complexity is a measure of the rigor of an individual item based on the amount of text a student must process from the corresponding passage to answer the item correctly, the way in which students are expected to interact with the item's functionality, and the linguistic demands and reading load that exists within the components of the item itself. Reviewers were asked to determine the extent to which items were aligned to the CCSS, using "fully," "partially," or "not aligned" as the rating categories. Ratings were averaged to determine overall alignment. For ELA/L, 99.6 percent of grade 3 and 4 items, 95.5 percent of grade 6 items, and 94.6 percent of grade 7 items were fully aligned. For mathematics, 92.0 percent of grade 3, 91.1 percent of grade 4 items, 83.1 percent of grade 6 items, and 94.0 percent of grade 7 items were fully aligned. The majority of the items that did not fall into fully aligned were considered partially aligned to the standards. CCSS are designed to be measured by multiple items, so items that aligned to multiple CCSS received a partially aligned rating. The overall item-to-CCSS alignment was captured by a holistic alignment rating that indicated if an item captured the identified standards as a set. Holistic ratings (either yes or no) were found by averaging review ratings across clusters for items that included more than one standard. For ELA, for all four grades, at least 93 percent of items had a holistic alignment rating of yes to indicate that the identified standards captured the skills or knowledge required. For mathematics, grade 6 had the lowest percentage for the holistic alignment rating of yes (84.8 percent), and grade 7 had the highest (96.3 percent). Overall the alignment study suggests that the identified CCSS capture the knowledge and skills required in the items.

In addition to the alignment study, HumRRO also evaluated the CCSSO criteria for content and depth for ELA/L and mathematics grades 3, 4, 6, and 7, as well as the cognitive complexity levels of these same grades (Schultz et al., 2017). There are five criteria for ELA/L content: close reading, writing, vocabulary and language skills, research and inquiry, and speaking and listening. Reviewers were asked to rate the content as Excellent, Good, Limited/Uneven, or Weak Match. For grades 3, 4, 6, and 7, the ELA/L assessments received a composite rating of Excellent Match for assessing the content needed for college- and career-readiness. There are four criteria for ELA/L depth: text quality and types, complexity of texts, cognitive demand, and high-quality items and item variety. All grades in this study received a composite rating of Good Match for depth. For mathematics content, the composite rating is based on two criteria: focus and concepts, procedures and applications. Grades 3, 4, and 6 received a composite content rating of Good Match, and grade 7 received a

composite content rating of Excellent Match. The mathematics composite depth rating is based on three criteria: connecting practice to content, cognitive demand, and high-quality items and item variety. All grades in the study were rated as Excellent Match at assessing the depth needed to successfully meet college and career readiness.

Finally, the 2017 HumRRO study looked at cognitive complexity of the items on ELA/L and mathematics at grades 3, 4, 6, and 7 (Schultz et al., 2017). Reviewers indicated their agreement with the intended cognitive complexity ratings provided by participating states and agencies of low, medium, or high. The results indicated that the reviewers generally agreed with the distribution of complexity levels. There were differences in agreements in ELA/L language cluster and a few exceptions to agreement in math, particularly at grade 6, where there was disagreement in the ratings at the medium complexity level for two domains and the high complexity level for one domain. For grade 7, there was agreement across low, medium, and high in all domains.

14.5.2 Benchmarking Study

The purpose of the benchmarking study (McClarty et al., 2015) was to provide information that would inform the performance level setting process. An evidence-based standard setting approach (EBSS; McClarty et al., 2013) was used to establish the performance levels for its assessments. In EBSS, the threshold scores for performance levels are set based on a combination of empirical research evidence and expert judgment. This benchmarking study provided one source of empirical evidence to inform the college- and career-readiness performance level (i.e., Level 4). The study findings were provided to a pre-policy standard setting committee. The charge of this committee was to suggest a reasonable range for the percentage of students meeting or exceeding the Level 4 threshold score and therefore considered college- and career-ready. Section 8.3.2 of this report provides more information about the pre-policy meeting.

For the benchmarking study, external information was analyzed to provide information about the Level 4 threshold scores for the grade 11 ELA/L, Algebra II, and Integrated Mathematics III assessments, the grade 8 ELA/L and mathematics assessments, and the grade 4 ELA/L and mathematics assessments. The assessments and Level 4 expectations were compared with comparable assessments and expectations for the Programme of International Student Assessment, Trends in International Mathematics and Science Study, Progress in International Reading Literacy Study, National Assessment of Educational Progress, ACT, SAT, the Michigan Merit Exam, and the Virginia end-of-course exams. For each external assessment, the best-matched performance level was determined and the percentage of students reaching that level across the nation and in the participating states and agencies was determined. Across all grades and subjects, the data indicated approximately 25 to 50 percent of students were college- and career-ready or on track to readiness based on the Level 4 expectations.

For details on how the benchmarking study was used during the standard setting process, refer to Section 8 of this technical report.

14.5.3 Longitudinal Study of External Validity of Performance Levels (Phase 1)

In 2016–2017, the first phase of a two-part external validity study of claims about the alignment of Level 4 to college readiness was completed (Steedle et al., 2017) using the summative assessment scores from the 2014–2015 and 2015–2016 academic years. Associations between the performance levels and college-readiness benchmarks established by the College Board and ACT were used to study the claim that students who achieve Level 4 have a .75 probability of attaining at least a C in entry-level, credit-bearing, postsecondary coursework. Regression estimates measured the relationship between the summative assessment scores and external test scores. The Level 4 benchmark was used to estimate the expected score on an external test, and vice versa. Assessment scores were dichotomized for additional analyses. Cross-tabulation tables provided classification agreement among tests. Logistic regression modeled the relationship between students' summative scores and their probabilities of meeting the external assessment benchmark, and vice versa.

These methods were used to make the following comparisons in mathematics: Algebra I and PSAT10 Math, Geometry and PSAT10 Math, Algebra II and PSAT10 Math, Algebra II and PSAT/NMSQT Math, Algebra II and SAT Math, and Algebra II and ACT Math. The classification agreement (meeting the benchmark on both tests or not meeting the benchmark on both tests) ranged from 62.5 percent to 86.5 percent. The overall trend indicated that students who met the benchmark on a mathematics assessment were likely to meet or exceed the benchmark on an external test (probabilities ranged from .509 to .886). However, students who met the benchmark on the external test had relatively low probabilities of meeting the mathematics benchmark (.097 to .310).

The following comparisons were made in ELA/L: grade 9 and PSAT10 evidence-based reading and writing (EBRW), grade 10 and PSAT10 EBRW, grade 10 and PSAT/NMSQT EBRW, grade 10 and SAT EBRW, grade 11 and PSAT/NMSQT EBRW, grade 11 and SAT EBRW, grade 11 and ACT English, and grade 11 and ACT reading. In the majority of comparisons, the trend in ELA/L results was similar to mathematics. The classification agreements ranged from 67.3 percent to 79.7 percent. Students meeting the ELA/L benchmark had probabilities between .667 and .825 of meeting the benchmark on the external assessment. However, a student taking the external test had lower probabilities of meeting the benchmark on the ELA/L assessments (.326 to .513).

Overall, results indicated that a student meeting the benchmark on the summative assessment had a high probability of making the benchmark on the external test, but the converse did not hold for students meeting the benchmark on the external test, for the majority of comparisons. These results suggest that meeting the summative benchmark is an indicator of academic readiness for college. However, it may be that students who meet the summative benchmark have a greater than .75 probability of earning a C or higher in first-year college courses.

Phase 1 is a preliminary study using indirect comparisons; therefore, there are limitations to interpretations. Phase 2 of this study was to occur in 2018 and use longitudinal data including academic performance in entry-level college courses for students who took the summative assessments during high school. Currently, this study is on hold due to challenges obtaining student academic data from entry-level college courses and/or matching the data to the student summative scores.

14.5.4 Mode and Device Comparability Studies

The summative assessments have been operational since the 2014–2015 school year. In addition to the traditional paper format, the assessments were available for online administration via a variety of electronic devices, including desktop computers, laptop computers, and tablets. The research agenda includes several studies evaluating the interchangeability of scale scores across modes and devices.

This report describes a two-pronged study consisting of a mode comparability analysis and a device comparability analysis. In the mode comparability analysis, scores arising from the paper administration were compared to those arising from any type of online administration. In the device comparability analysis, online scores arising from tests administered using a tablet are compared with online scores arising from any other type of electronic administration where a tablet was not present (i.e., laptops, desktops, Chromebooks).

The goal of this study was threefold: (1) to investigate whether assessment items were of similar difficulty across the levels of conditions for each analysis (i.e., paper and online for the mode comparability analysis and tablet and non-tablet for the device comparability analysis); (2) to determine whether the psychometric properties of test scores were similar across the levels of conditions for each analysis; and (3) to determine whether overall test performance was similar across the levels of conditions for each analysis.

This study examined performance on 12 assessments, split evenly between mathematics and ELA/L. Students were matched on demographic variables as well as on the score from the summative assessment in the same content area in the prior year, creating comparable samples that allowed for an unbiased comparison of performance across different conditions.

The results of the mode comparability analysis were mixed and found to be consistent with prior research. The item means suggested that items were of similar difficulty on paper and online modes. Only two items were flagged for mode effects, both of which were on the mathematics assessments. C-level differential item functioning (DIF) was present in both analyses. All the items flagged for C-level DIF in the mathematics assessments favored the online students, whereas the majority of items flagged for C-level DIF in the ELA/L assessments favored the paper students. An examination of test reliability displayed comparable reliability values between the two modes; none of the test forms were flagged for mode effects with respect to test reliability. The test-level adjustment analysis as well as the change of the paper students' performance levels after the adjustment constants were applied to the paper students' scores indicated that more scale scores were adjusted downward than were adjusted upward on the paper test form for each assessment except grades 5 and 7 mathematics. However, all adjustments were less than the minimum standard error of theta except for grade 11 ELA/L, which was the same as the minimum standard error of theta. Therefore, the adjustments are within measurement precision for each assessment.

The results of the device comparability study revealed consistent evidence supporting the comparability between the tablet condition (TC) and the non-tablet condition (NTC). Specifically, the item means suggested that items were similarly difficult for the TC and NTC, and none of the items were flagged for device effects. The DIF analysis revealed that none of the items had C-level DIF. Consistent with the findings at the item level, an examination of test reliability indicated that the TC and NTC test forms were similarly reliable and that none of the test forms were flagged for device effects. Furthermore, the test-level adjustment analysis as well as the change of the students' performance levels after the adjustment constants were applied did not indicate strong evidence of device effects.

The generalizability of the findings from this study may be limited due to the small sample size of both the paper students (for mode comparability) and the tablet students (for device comparability) at the high school grades; however, it appears that high-quality matching supports the internal validity of this study's findings. For mode and device comparability, there were few to no items flagged for mode or device effects, the psychometric properties of test scores were similar across assessment conditions, and any adjustments to student performance for the paper or tablet condition were within measurement precision.

14.5.5 Quality Testing Standards

New Meridian, in coordination with multiple states and vendors, developed an alternate form of the summative assessment to meet the provision for shorter testing times requested by several states. Research conducted using 2017 (Boyd et al., 2018) and 2018 (Minchen et al., 2018) student data evaluated the effects of removing items from the original assessments to determine if scores arising from the two versions would be comparable. Research was conducted in several steps. First, subject matter experts identified item subsets from the original forms that maintained the integrity of the assessment and were approximately 65 to 80 percent of the original test length. Then, students were rescored on the item subsets, producing a set of hypothetical scores, as if the students had only taken the subset of items. Finally, a series of analyses were conducted. While the research generally supported the comparability of the two versions, a limitation of the methodology was that the alternate blueprints were not actually administered as such. In this report, the shorter version of the blueprint is referred to as the alternate assessment and the original blueprint is referred to as the original assessment.

Through extensive research and guidance from the Technical Advisory Committee, the alternate blueprint was available in spring 2019 in addition to the original blueprint. In 2019, the option to administer either blueprint was made at the state or agency level. Since some states administered the alternate blueprint and some states administered the original blueprint, the following research evaluated the comparability between the two blueprints with respect to scale score comparability and performance level comparability.

The goal was to determine additional evidence to support scale score comparability and performance level comparability, according to the guidelines outlined in the Quality Testing Standards (Center for Assessment, 2018). For the purpose of this work, scale score and performance level comparability have formal definitions. Scale score comparability is defined by the Center for Assessment (2018) as follows: If a student taking the alternate assessments with New Meridian content took the original assessment, would the student obtain a similar scale score? Performance level comparability is defined by the Center for Assessment (2018) as follows: If a student taking the alternate assessment with New Meridian content took the original assessment, would the student receive a similar designation in terms of college- and career-readiness or performance level 4 on the original blueprint?

For the spring 2019 assessments, the mathematics items on the alternate forms also appeared on the corresponding original forms; however, for ELA/L assessments, a small number of items were unique to the alternate forms. The scale scores were reported on the same scale regardless of the form and used the same performance level cut scores.

Three sets of analyses were conducted. Most of the analyses were conducted on a set of matched samples from the 2019 alternate and original forms, allowing for direct comparisons of assessment characteristics and outcomes to be made. Such samples were obtained through coarsened exact matching (Iacus et al., 2012), which used demographic information and prior achievement scores, where possible. Prior achievement scores were grouped into bands within each performance level, and students taking the alternate forms were

matched with students who took the original forms who had identical information on all demographic and prior achievement variables. The prior assessments used in the matching process can be found in Tables 14.29 and 14.30. For grade 3 assessments, only demographic information is used in the matching process due to the lack of prior assessment data. Due to differences in high school assessment requirements across states and agencies, multiple prior assessments may have been used. For ELA/L grade 10, the prior assessment was ELA/L grade 8 for the matching process.

Table 14.29 Prior Grades Used in ELA/L Matching

Current Grade	Prior Grade	Prior Test Year
Grade 3	N/A	N/A
Grade 4	Grade 3	2018
Grade 5	Grade 4	2018
Grade 6	Grade 5	2018
Grade 7	Grade 6	2018
Grade 8	Grade 7	2018
Grade 10	Grade 8	2017

Table 14.30 Prior Grades/Courses Used in Mathematics Matching

Current Grade/ Course	Prior Grade /Course	Prior Test Year
Grade 3	N/A	N/A
Grade 4	Grade 3	2018
Grade 5	Grade 4	2018
Grade 6	Grade 5	2018
Grade 7	Grade 6	2018
Grade 8	Grade 7	2018
Algebra I	Grade 7 (44%), Grade 8 (56%)	2018
Geometry	Algebra I	2018
Algebra II	Algebra I (10%), Geometry (90%)	2018

Sample sizes before and after the matching process are listed in Table 14.31 for ELA/L and Table 14.32 for mathematics. ELA/L grade 9, Geometry, and Algebra II, matched samples were fairly small, ranging from 75 to 1,540. Due to the small sample for ELA/L grade 9, the comparability analyses were not conducted. Geometry and Algebra II were included in the comparability analyses; however, the results should be interpreted with caution given the small samples.

Table 14.31 ELA/L Matching Sample Size Results

ELA/L	Form	Unmatched		Matched	
		Current Forms N	Original Forms N	Current Forms N	Original Forms N
Grade 3	1	105,482	32,034	31,481	31,481
	2	105,309	31,861	31,272	31,272
Grade 4	1	105,826	28,153	27,695	27,695
	2	126,875	34,071	33,444	33,444
Grade 5	1	136,148	36,313	35,742	35,742
	2	101,869	27,272	26,721	26,721
Grade 6	1	119,838	31,031	30,667	30,667
	2	120,218	30,802	30,506	30,506
Grade 7	1	116,933	29,877	29,544	29,544
	2	117,757	29,835	29,593	29,593
Grade 8	1	118,198	29,638	29,312	29,312
	2	119,059	29,248	28,898	28,898
Grade 9	1	30,648	86	75	75
	2	71,029	116	102	102
Grade 10	1	55,046	27,951	22,970	22,970
	2	41,439	20,758	17,193	17,193

Table 14.32 Mathematics Matching Sample Size Results

	Form	Unmatched		Matched	
		Current Forms N	Original Forms N	Current Forms N	Original Forms N
Grade 3	1	88,858	26,531	25,970	25,970
	2	88,919	26,595	25,987	25,987
Grade 4	1	87,291	25,941	25,070	25,070
	2	87,488	26,192	25,207	25,207
Grade 5	1	91,136	27,333	26,377	26,377
	2	91,739	27,611	26,754	26,754
Grade 6	1	95,174	28,514	27,677	27,677
	2	94,800	28,342	27,665	27,665
Grade 7	1	93,777	24,547	23,855	23,855
	2	93,265	24,141	23,485	23,485
Grade 8	1	83,289	15,293	14,962	14,962
	2	76,135	13,973	13,695	13,695
Algebra I	1	43,232	21,530	16,926	16,926
	2	46,482	23,036	18,157	18,157
Geometry	1	40,673	3,252	1,540	1,540
	2	40,918	3,360	1,514	1,514
Algebra II	1	27,568	1,037	823	823
	2	27,527	1,066	753	753

Detailed matching results for select assessments can be found in the Appendix, Tables A.14.1 through A.14.3. ELA/L and mathematics for grade 6 and ELA/L grade 10 matching results are presented. Other grade levels had very similar results to grade 6, except for ELA/L grade 10.

The remaining analyses were conducted on assessment data from 2018 and 2019, rather than the matched samples. The second set of analyses was conducted at the grade level, using all available data from both 2018 and 2019, examining grade-level statistics over the course of two years, ensuring state participation was similar within each grade for both years. Finally, the last set of analyses used two-year student cohorts, examining students' scores over two years. Only students who completed assessments in both 2018 and 2019 were included; therefore, grade 3 student data from 2019 were not included.

Effect sizes were used throughout the research to determine the degree to which differences were practically significant. For differences between continuous distributions, such as scale score and claim score means, Cohen's (1988) D was used, and is calculated as

$$D = \frac{\bar{x}_1 - \bar{x}_2}{S_p} \quad (14-4)$$

where \bar{x}_1 and \bar{x}_2 are the means of interest, and S_p is the pooled standard deviation of the scores in both distributions. For differences in proportions, Cohen's (1988) h was used, and is given by

$$h = 2 \left(\sin^{-1} \sqrt{p_1} - \sin^{-1} \sqrt{p_2} \right) \quad (14-5)$$

where p_1 and p_2 are the proportions of interest. And for differences in ordinal distributions, Cramer's (1946) V was used, which is given as

$$V = \sqrt{\frac{\chi^2}{n \times \min(r-1, c-1)}} \quad (14-6)$$

where χ^2 is the chi-squared value from the contingency table calculation, n is the total sample size, r is the number of rows in the contingency table, and c is the number of columns in the contingency table. Cohen (1988) defined effect sizes as .25, .5, and .8 as constituting small, medium, and large effects, respectively. A number of regression analyses are also performed, and the change in R^2 between the full and reduced models is examined; R^2 values of .01, .06, and .15 constitute the small, medium, and large effect sizes (Cohen, 1988).

Scale Score Comparability: Item-Level Analysis

Item-level evaluations (i.e., p-values, polyserial correlations, and DIF) were conducted separately for alternate and original forms on the matched sample for items that were common to both forms for each grade/course. First, p-values were compared. Scatterplots for the alternate form p-values and original form p-values for ELA/L grades 3 to 6 and mathematics grades 3 to 6 are presented in Figures 14.4 and 14.5, respectively.

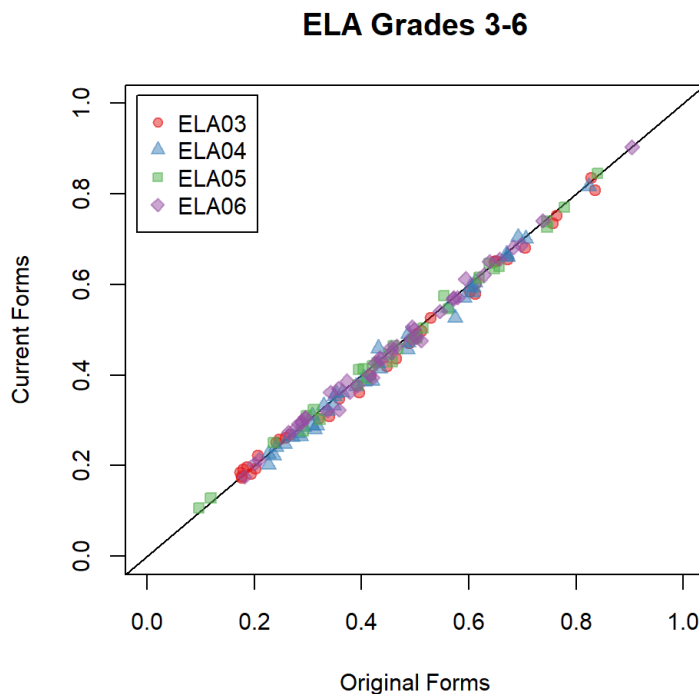


Figure 14.4 ELA/L Grades 3–6 P-Values

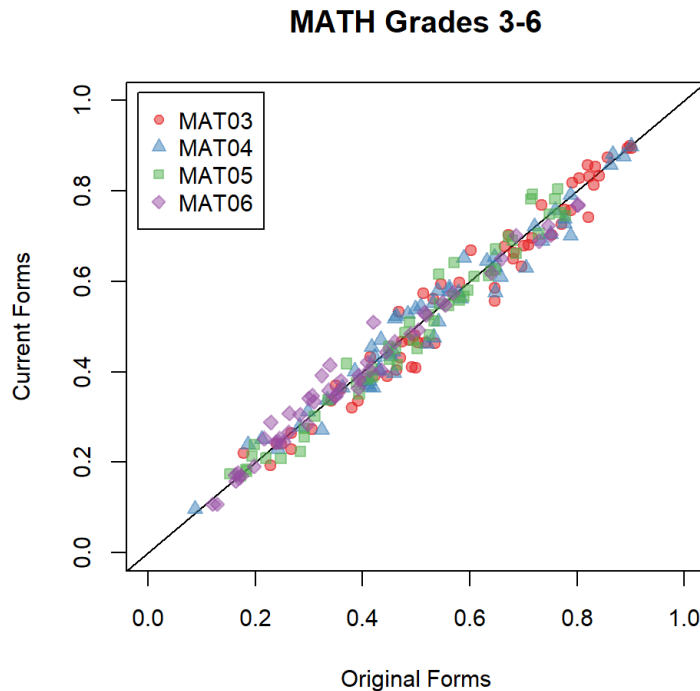


Figure 14.5 Mathematics Grades 3–6 P-Values

The scatterplots for all grades and courses are presented in Figures A.14.1 through A.14.6. Scatterplots show that most points cluster closely and evenly around the $y = x$ line, showing that items perform similarly on both forms with the matched samples, with the exception of ELA/L grade 10, Algebra II, and Geometry.

The distributions of p-value differences for all grades are presented in Tables A.14.4 and A.14.5. Differences tend to be small and center around zero, except for ELA/L grade 10, Algebra II, and Geometry. For ELA/L grades 3 through 8, differences in item difficulties range from $-.049$ to $.070$. For mathematics grades 3 through 8 and Algebra I, differences in item difficulties range from $-.105$ to $.090$. The high school assessments show larger differences. P-values for ELA/L grade 10 on the alternate forms were lower than on the original forms.

The polyserial correlations of common items on the alternate and original forms using the matched sample were also analyzed. Scatterplots, which are presented in Figures A.14.7 through A.14.12, show that most points cluster closely and evenly around the $y = x$ line, showing that items perform similarly on both forms with the matched sample, with the exception of Algebra I, Algebra II, and Geometry. The distributions of these differences, which are presented in Tables A.14.6 and A.14.7, tend to be small and center around 0, except for ELA/L grade 10, Algebra II, and Geometry. For ELA grades 3 through 8, differences in polyserial values range from $-.058$ to $.043$. For Mathematics grades 3 through 8, differences in polyserial values range from $-.090$ to 0.125 . The high school assessments show larger differences.

Common items were checked for DIF on several categories separately for the alternate and original forms, using the matched samples. The resulting cross tabulation of DIF categories was examined. Percentages were computed for each possible combination of DIF categories and represented the total number of cross

tabulations divided by the total number of DIF calculations (items multiplied by categories for which the sample size was sufficient for DIF calculations) within a grade. For most tests, at least 90 percent of calculations displayed no DIF on both the alternate and original forms. DIF results summaries can be found in Tables A.14.8 through A.14.10.

Scale Score Comparability: Test-Level Analysis

Test-level evaluations included analyzing reliability, scale score distributions, ELA/L claim score distributions, and subclaim distributions. Analyses showed that reliability, calculated as the stratified alpha, was slightly lower for alternate forms compared to their original form counterparts, as expected. For each assessment, the Spearman Brown Prophecy Formula was used to predict the alternate form reliabilities based on the reduction in items. The alternate form reliability estimates tended to be generally similar to the Spearman-Brown Prophecy values based on the corresponding reduction in points. This indicated that the loss of precision was approximately commensurate with the reduction in length. Similar results were found at the claim and subclaim levels.

Both raw score (RS) and scale score (SS) standard error of measurement (SEMs) are presented, as well as an adjusted raw score SEM that is simply the proportion of total points represented by the raw score SEM. The scale score and adjusted raw score SEMs were always slightly larger for the alternate forms, as expected. Reliability and SEM results at the summative level are available in Tables A.14.11 through A.14.16, while results for the claim and subclaim levels are available in Tables A.14.19 through A.14.24 and A.14.42 through A.14.52.

Scale score and subclaim distributions between the alternate and original forms tended to be similar, as evidenced by small effect sizes with respect to the difference in the means of the scale scores and distributions of the performance levels, except for ELA/L grade 10. The effect sizes, computed as Cohen's D, of the differences between the summative scale score alternate and original means were less than .20 in magnitude for all ELA/L and mathematics grades except ELA/L grade 10. Results are available in Tables A.14.17 and A.14.18. The effect sizes of the differences between the alternate and original Writing claim scale score means were less than .20 in magnitude for all ELA/L grades except ELA/L grade 10. Results are available in Table A.14.19. The effect sizes of the differences between the alternate and original Reading claim scale score means were also less than .20 in magnitude for all ELA/L grades except ELA/L grade 10. Results are presented in Table A.14.20. Subclaim distributions for alternate and original forms using the matched sample were compared using Cramer's V effect size. All effect sizes were .20 or lower. Detailed results for ELA/L and mathematics grade 6 assessments are presented in Tables A.14.21 and A.14.22, respectively, while results summaries for all grades and courses can be found in Tables A.14.23 and A.14.24.

Scale Score Comparability: Longitudinal Analysis

Longitudinal analyses generally revealed stability in scale score means when controlling for state participation. Effect sizes ranged in magnitude from 0 to .16, with all but two being smaller than .10. No clear directional pattern emerged. Detailed results can be found in Tables A.14.25 through A.14.28. Additionally, a regression analysis approach was used to examine the relationship between students' 2018 and 2019 scale scores. The full and reduced models are given below.

Full Model:

$$SS_{2019} = \beta_0 + \beta_1 \times SS_{2018} + \beta_2 \times C + \beta_3 \times SS_{2018} \times C \quad (14-7)$$

Reduced Model:

$$SS_{2019} = \beta_0 + \beta_1 \times SS_{2018} \quad (14-8)$$

where SS_{2019} is the scale score on the 2019 assessment, SS_{2018} is the scale score on the 2018 assessment, C is a categorical variable in which students taking the alternate assessment are indicated with a one, and students taking the original assessment are indicated with a zero.

The changes in R^2 ranged from less than .0001 to .0260, demonstrating that the form choice for 2019 did not explain much additional variance in the 2019 scale scores. Regression results can be found in Tables A.14.29 and A.14.30.

As an additional component of the research, student growth percentiles (SGPs) were compared for students in the matched samples for grades 4 and higher who have prior achievement scores. Section 15 describes the SGP analyses conducted for the spring 2019 administration. SGPs can be computed using either each individual state or the entire consortium as the peer group. For these analyses, SGPs are computed based on the consortium peer group.

The mean SGPs for students in the matched sample who were administered the alternate forms were compared with those in the sample who were administered the original forms. Means were computed across all students in the sample as well as for various subgroups. Similar means indicated that student growth can be measured similarly regardless of the type of form, providing additional evidence of comparability. SGP mean differences greater than 5 percentile points in magnitude, which corresponds to an effect size of approximately 0.18 (D. Betebenner, personal communication, September 10, 2019), may warrant further investigation.

For ELA/L and mathematics grades 4 through 8, differences between the mean SGPs were generally less than 5 percentile points in magnitude. At the overall level, mean differences (measured in percentile points and computed as the alternate form mean SGP minus original form mean SGP) ranged from -3.0 to 1.3 for ELA/L and from -2.7 to 3.5 for mathematics. Subgroups evaluated were African American or Black, Asian, Hispanic, multiple races, Native American, white, economically disadvantaged, English learners, and students with disabilities. Except the Asian and Native American subgroups, the differences in the means were less than 5 in magnitude. For Asian students in mathematics grade 8, the difference in the means was 5.2. For Native American students, the differences for ELA/L grade 4 and mathematics grades 4, 6, and 8 were -5.3, -8.4, -9.1, and -6.5, respectively. Of note is that each of these exceptions occurs when the sample size is relatively small. For mathematics grade 8, there were only 730 Asian students administered each type of form; all Native American grades contained fewer than 200 students for each type of form. SGP mean differences for all students as well as for each of the subgroups for Algebra I tended to be slightly higher than 5 in absolute value, but always less than 10. Results for Geometry and Algebra 2 are not included due to small sample sizes.

These results provide additional evidence in support of comparability between the alternate and original scale scores at grades 4 through 8. For high school analyses, small samples, potential differences in course progressions, and possible differences in administration characteristics (e.g., graduation requirements) within each state complicate the interpretation of the results.

Performance Level Comparability: Test-Level Analyses

The performance level distributions for the alternate and original forms were compared using Cramer's V as the effect size measure. For summative performance level and college- and career-readiness (CCR), which is defined as students who attained performance levels 4 or 5, distributions tended to be similar across the alternate and original forms, with effect sizes of less than .10 in magnitude relative to the differences in their distributions, except for ELA/L grade 10. Detailed results for ELA/L and mathematics grade 3 can be found in Tables A.14.31 and A.14.32, respectively. A summary of the effect sizes for all assessments can be found in Table A.14.33. Additionally, the percentage of students attaining or exceeding the CCR indicator for Alternate and Original forms was calculated and compared using Cohen's h as the measure of effect size. All effect sizes were less than .10 in magnitude, except for ELA/L grade 10. These results can be found in Table A.14.34.

Performance Level Comparability: Classification Analyses

Classification accuracy and consistency were also computed using BB-Class (Brennan, 2004) in two ways: using all five performance levels and using only the CCR indicator. Both classification accuracy and consistency were always lower for alternate forms compared to the original forms, as expected, as there are differences in measurement precision discussed above. Effect sizes, as computed by Cohen's h , measuring the differences were small to moderate in magnitude, and ranged from $-.04$ to $-.23$ for performance level classification accuracy (Tables A.14.35 and A.14.37), from $-.05$ to $-.25$ for performance level classification consistency (Tables A.14.36 and A.14.38), from $-.02$ to $-.10$ for CCR classification accuracy (Tables A.14.35 and A.14.37), and from $-.02$ to $-.12$ for CCR classification consistency Tables (A.14.36 and A.14.38).

Performance Level Comparability: Longitudinal Analyses

Finally, a longitudinal evaluation of performance levels was conducted using all available data, rather than the matched samples. Performance level and CCR distributions were examined for each grade in 2018 and 2019, ensuring that data from both years represented the same states. Cramer's V and Cohen's h were used as the measures of effect size for the performance level and CCR comparisons, respectively. All effect sizes were .10 or less in magnitude. Detailed results for ELA/L and mathematics grade 6 can be found in Tables A.14.39 and A.14.40, while a summary of results across all assessments can be found in Table A.14.41.

Quality Testing Standards Summary

The purpose of the Quality Testing Standards study was to compare the results from the alternate and original assessments. Because states only administered one type, comparable samples were extracted from the data using coarsened exact matching. Using this data, a variety of analyses demonstrated that there appears to be broad comparability between the alternate and original scale scores and performance levels, that the alternate forms have less measurement precision than the original forms, and that the results from many of the high school tests were slightly less clear. Several factors limited the analysis of high school results. First, for ELA/L grade 10, the prior assessment used was ELA/L grade 8 from 2017. A test and results that are two years removed may be less than ideal. Second, high school tests tended to have smaller samples and were obtained from fewer states. Third, high school curriculum and course progressions may vary from state to state.

Additionally, several longitudinal analyses were conducted using assessment data from 2018 and 2019 rather than the matched sample. Although the analyses were limited in scope, the results support the findings from the matched analyses.

14.6 Evidence Based on Response Processes

As noted in the AERA, APA, and NCME Standards (2014), additional support for a particular score interpretation or use can be provided by theoretical and empirical evidence indicating that students are using the intended response processes when responding to the items in a test. This type of evidence may be gathered from interacting with students in order to understand what processes underlie their item responses. Evidence may also be derived from feedback provided by test proctors/teachers involved in the administration of the test and raters involved in the scoring of constructed-response items. Evidence may also be gathered by evaluating the correct and incorrect responses to short constructed-response items (e.g., items requiring a few words to respond) or by evaluating the response patterns to multi-part items.

New Meridian has undertaken research investigating the quality of the items, tasks, and stimuli, focusing on whether students interact with items/tasks as intended, whether they were given enough time to complete the assessments, and the degree to which scoring rubrics allow accurate and reliable scoring. In addition, the accessibility of the test for students with disabilities and English learners has been examined. This research has included examining students' understanding of the format of the assessments and the use of technology.

One such study involved a series of four component studies that were conducted to evaluate the usability and effect of a drawing tool for online mathematics items. The purpose of these studies was to determine if results could support the use of the drawing tool, which is a way to expand students' ability to demonstrate their understanding and reasoning, thereby enhancing accessibility and construct validity of the assessment. This goal is in keeping with guidance from the CCSS and the National Council of Teachers of Mathematics that students should have multiple paths and tools available to express their responses. Additionally, the drawing tool was intended to boost comparability across modes.

The first two studies (Brandt, Bercovitz, McNally, & Zimmerman, 2015; Brandt, Bercovitz, & Zimmerman, 2015) focused on evaluating the usability of the tool itself both in the general population and among students with low vision and fine motor impairment disabilities. During these studies, detailed information regarding the functionality of the tool was collected, and it was determined that the items should be tested operationally.

The third and fourth studies (Minchen et al., 2018b; Steedle & LaSalle, 2016) involved evaluating the effect of the tool in the context of the operational assessments. The third study was conducted in grade 3, and the fourth study was conducted in grades 4 and 5. To evaluate the drawing tool in context, a set of items was studied by field-testing the items with and without the drawing tool. The drawing tool version of each item was randomly assigned to students so that comparisons could be made. The goal was to explore the impact of the drawing tool on item performance. In general, the results showed that the drawing tool usually did not have a significant impact on performance or item statistics. Items that included access to the drawing tool, however, did show longer response times for grades 4 and 5, prompting a limitation to be placed on the number of drawing tool items in each unit.

Several other research efforts have investigated questions relevant to response processes evidence. Descriptions of the research conducted can be found online.¹⁴

¹⁴ Descriptions of the research are available at <http://resources.newmeridiancorp.org/>.

14.7 Interpretations of Test Scores

The summative assessment scores are expressed as scale scores (both total scores and claim scores), along with performance levels, to describe how well students met the academic standards for their grade level. Additionally, information on specific skills (the subclaims) is also provided and is reported as *Below Expectations*, *Nearly Meets Expectations*, and *Meets or Exceeds Expectations*. On the basis of a student's total score, an inference is drawn about how much knowledge and skill in the content area the student has acquired. The total score is also used to classify students in terms of their level of knowledge and skill in the content area as students progress in their K–12 education. These levels are called performance levels and are reported as follows:

- Level 5: Exceeded expectations
- Level 4: Met expectations
- Level 3: Approached expectations
- Level 2: Partially met expectations
- Level 1: Did not yet meet expectations

Students classified as either Level 4 or Level 5 are meeting or exceeding the grade level expectations. PLDs assist with the understanding and interpretation of the ELA/L scores (<https://resources.newmeridiancorp.org/ela-test-design/>) and mathematics scores (<https://resources.newmeridiancorp.org/math-test-design/>). Additionally, resource information is available online to educators, parents, and students (<http://resources.newmeridiancorp.org/>). Section 12 of this technical report provides more information on the scale scores and the subclaim scores.

14.8 Evidence Based on the Consequences to Testing

The consequences of testing should also be investigated to support the validity evidence for the use of the summative assessments as the standards note that tests are usually administered “with the expectation that some benefit will be realized from the intended use of the scores” (AERA, APA, & NCME, 2014). When this is the case, evidence that the expected benefits accrue will provide support for the intended use of the scores. Evidence of the consequence of testing will also accrue with the continued implementation of the CCSS and the continued administration of the assessments.

The consequences of the tests may vary by state or by school district. For example, some states may require “passing” the assessments as one of several criteria for high school graduation, while other states/districts may not require students to “pass” the assessments for high school graduation. Additionally, some school districts may use the scores along with other information, such as school grades and teacher recommendations, for placing students into special programs (e.g., remedial support, gifted and talented program) or for course placement (e.g., Algebra I in grade 8). Because the consequences for the assessments can vary by each state, it is suggested that each member state provide school districts, teachers, parents, and students with information on how to interpret and use the scores. Additionally, the states should monitor how scores are used to ensure that the scores are being used as intended.

14.9 Summary

In this section of the technical report, several aspects of validity were included, such as validity evidence based on content, the internal structure of the assessments, relationships across the content assessments, and evidence from special studies.

The item development process involved educators, assessment experts, and bias and sensitivity experts in review of text, items, and tasks for accuracy, appropriateness, and freedom from bias. Several studies were conducted during the item development process to evaluate the item development process (e.g., technological functionalities, answer time required, and student experiences). Additionally, items were field-tested prior to the initial operational administration, and data and feedback from students, test administrators, and classroom teachers were used to improve the operational administration of the items and to inform future item development. The multiple item and form reviews conducted by educators and studies to evaluate item administration help to ensure the integrity of the assessments.

The intercorrelations of the subclaims, the reliability analyses, and the local item dependence analyses indicated that the ELA/L and the mathematics assessments are both essentially unidimensional. Furthermore, the correlations between ELA/L and mathematics indicated that the two assessments are measuring different content.

Several studies were conducted as part of the assessment program (e.g., benchmarking study, content evaluation/alignment studies, longitudinal study, and mode and device comparability studies). The benchmarking study was conducted in support of the standard setting meeting. This study indicated students performing at or above Level 4 could be considered to be college- and career-ready or on track to readiness.

The content evaluation/alignment studies performed by the Fordham Institute and HumRRO indicate that the assessments are good-to-excellent matches to the CCSS in terms of content and depth of knowledge. Thus, the assessments are assessing the college- and career-readiness standards. However, the reports noted that the program could improve by adding a wider range of depth of knowledge to some of the assessments. The reports also suggested enhancing the ELA/L assessments by including a research task that requires the use of two or more sources of information.

In the longitudinal study of external validity, associations between the performance levels and college-readiness benchmarks established by the College Board and ACT were used to study the claim that students who achieve Level 4 have a .75 probability of attaining at least a C in entry-level, credit-bearing, postsecondary coursework. In the first phase of the study, the relationship between the summative assessment and external tests was studied. Overall, results indicated that a student meeting the benchmark on the summative assessment had a high probability of making the benchmark on the external test, but the converse did not hold for students meeting the benchmark on the external test, for the majority of comparisons. These results suggest that meeting the benchmark is an indicator of academic readiness for college. In the next phase of the study, the relationship between scores and performance in first-year college courses will be explored.

The mode comparability study indicated that the comparability across modes was inconsistent across content domains and grade levels. The results of the mode comparability analysis were mixed and found to be consistent with prior research. The results of the device comparability study revealed consistent evidence supporting the comparability between the tablet condition and the non-tablet condition. In both the mode and device comparability studies, there were few to no items flagged for mode or device effects, the

psychometric properties of test scores were similar across assessment conditions, and any adjustments to student performance for the paper or tablet condition were within measurement precision.

In addition to the validity information presented in this section of the technical report, other information in support of the uses and interpretations of the scores appear in the following sections:

- Section 5 provides information concerning the test characteristics based on classical test theory.
- Section 6 provides information regarding the differential item functioning (DIF) analyses.
- Section 11 presents information regarding student characteristics for the spring administration of the ELA/L and mathematics administration.
- Section 12 provides detailed information concerning the scores that were reported and the cut scores for ELA/L and mathematics.
- Section 13 provides information on the test reliability (total test score and for subclaims) and includes information on the inter-rater reliability/agreement.

Section 15: Student Growth Measures

Student growth percentiles (SGPs) are normative measures of annual progress. Normative measures are useful in answering questions like “How does my academic progress compare with the academic progress of my peers?” In contrast to criterion-referenced measures of growth, which describe academic growth toward a particular goal, norm-referenced measures of growth describe students’ growth relative to that of students who performed similarly in the past (Betebenner, 2009).

SGPs measure individual student progress by tracking student scores from one year to the next. SGPs compare a student’s performance to that of his or her academic peers both within the state and across the consortium. Academic peers are defined as students in the norm group who took the same assessment as the student in prior years and achieved a similar score.

Some participating states or agencies chose to implement norm groups based on their respective student data. State-specific SGP results are not reported in this technical report. As a result, SGPs were only summarized for states using norm groups based on the consortium. The following sections describe the norm groups, the estimation procedure, and the results for SGPs based on consortium norm groups.

The SGP describes a student’s location in the distribution of current test scores for all students who performed similarly in the past. SGPs indicate the percentage of academic peers above whom the student scored. With a range of 1 to 99, higher numbers represent higher growth and lower numbers represent lower growth. For example, an SGP of 60 on grade 7 English language arts/literacy (ELA/L) means that the student scored better than 60 percent of the students in the state or consortium who took grade 7 ELA/L in spring 2019 *and* who had achieved a similar score as this student on the grade 6 ELA/L assessment in spring 2018 and the grade 5 ELA/L assessment in spring 2017.¹⁵ A SGP of 50 represents typical (median) student growth for the state or consortium. Because students are only compared with other students who performed similarly in the past, all students, regardless of starting point, can demonstrate high or low growth.

The 2020–2021 academic year is the seventh year of test administration, including an abbreviated administration to a small number of students in one state in 2020. Data from 2020 was not used in SGP calculations. Students in states that participated in spring 2018 and spring 2019 generally received SGPs based on two prior scores. Students in states that participated in spring 2019 received SGPs based on one prior score. Students who do not have a previous test score, which include any new students and all grade 3 and 4 students, do not receive an SGP.

15.1 Norm Groups

The norm groups consisted of students with the same prior scores based on grade or content area progressions (academic peers). SGPs were based on up to two years of prior test scores from spring 2018 and spring 2019 administrations. States administering traditional mathematics assessments in fall 2018 or fall 2019 may also have SGPs based on these prior scores. Tables 15.1 through 15.8 list the grade or content area progressions required for SGPs based on one prior or two prior test scores for ELA/L grades 3 through 11, mathematics grades 3 through 8, Algebra I, Geometry, Algebra II, and Integrated Mathematics I, II, and III,

¹⁵ Because regression modeling is used to establish the relationship between prior and current scores, the SGP is for students with the exact same prior scores. This often leads to confusion among nontechnical stakeholders who often ask, “How many students are there with exactly the same prior scores?” To avoid explaining regression to nontechnical stakeholders, the “similar scores” is often used to finesse the idea of regression without mentioning it.

respectively. In general, the progressions of grade levels and content areas are consecutive. The traditional and integrated mathematics courses have progressions that are not consecutive but reflect student progression for high school mathematics courses. SGPs were calculated for all norm groups with at least 1,000 students. Some progressions did not meet the minimum sample size for SGP calculations.

Table 15.1 ELA/L Grade-Level Progressions for One- and Two-Year Prior Test Scores

Two Prior Year Test Scores	One Prior Year Test Score	Current Year Test Score
N/A	N/A	Grade 3*
N/A	Grade 3	Grade 4
Grades 3 and 4	Grade 4	Grade 5
Grades 4 and 5	Grade 5	Grade 6
Grades 5 and 6	Grade 6	Grade 7
Grades 6 and 7	Grade 7	Grade 8
Grades 7 and 8	Grade 8	Grade 9
Grades 8 and 9	Grade 9	Grade 10
Grades 9 and 10	Grade 10	Grade 11

*SGP not calculated for grade 3 since there are no prior scores.

Table 15.2 Mathematics Grade-Level Progressions for One- and Two-Year Prior Test Scores

Two Prior Year Test Scores	One Prior Year Test Score	Current Year Test Score
N/A	N/A	Grade 3*
N/A	Grade 3	Grade 4
Grades 3 and 4	Grade 4	Grade 5
Grades 4 and 5	Grade 5	Grade 6
Grades 5 and 6	Grade 6	Grade 7
Grades 6 and 7	Grade 7	Grade 8

*SGP not calculated for grade 3 since there are no prior scores.

Table 15.3 Algebra I Grade/Content Area Progressions for One- and Two-Year Prior Test Scores

Two Prior Year Test Scores	One Prior Year Test Score	Current Year Test Score
Grades 5 and 6	Grade 6	Algebra I
Grades 6 and 7	Grade 7	Algebra I
Grades 6 or 7 and 8	Grade 8	Algebra I
Grades 6, 7, or 8 and Geometry	Geometry	Algebra I
Grade 8 and Integrated Mathematics I	Integrated Mathematics I	Algebra I
Integrated Mathematics I and Integrated Mathematics II	Integrated Mathematics II	Algebra I

Table 15.4 Geometry Grade/Content Area Progressions for One- and Two-Year Prior Test Scores

Two Prior Year Test Scores	One Prior Year Test Score	Current Year Test Score
Grades 5 and 6	Grade 6	Geometry
Grades 6 and 7	Grade 7	Geometry
Grades 6 or 7 and 8	Grade 8	Geometry
Grades 6, 7, or 8 and Algebra I	Algebra I	Geometry
Grade 8 and Integrated Mathematics I	Integrated Mathematics I	Geometry
Integrated Mathematics I and Integrated Mathematics II	Integrated Mathematics II	Geometry

Table 15.5 Algebra II Grade/Content Area Progressions for One- and Two-Year Prior Test Scores

Two Prior Year Test Scores	One Prior Year Test Score	Current Year Test Score
Grades 6 and 7	Grade 7	Algebra II
Grades 7 and 8	Grade 8	Algebra II
Grades 7 or 8 and Algebra I	Algebra I	Algebra II
Grade 8 or Algebra I and Geometry	Geometry	Algebra II
Grade 8 and Integrated Mathematics I	Integrated Mathematics I	Algebra II
Integrated Mathematics I and Integrated Mathematics II	Integrated Mathematics II	Algebra II

Table 15.6 Integrated Mathematics I Grade/Content Area Progressions for One- and Two-Year Prior Test Scores

Two Prior Year Test Scores	One Prior Year Test Score	Current Year Test Score
Grades 5 and 6	Grade 6	Integrated Mathematics I
Grades 6 and 7	Grade 7	Integrated Mathematics I
Grades 6 or 7 and 8	Grade 8	Integrated Mathematics I
Grades 7 or 8 and Algebra I	Algebra I	Integrated Mathematics I
Grade 8 or Algebra I and Geometry	Geometry	Integrated Mathematics I

Table 15.7 Integrated Mathematics II Grade/Content Area Progressions for One- and Two-Year Prior Test Scores

Two Prior Year Test Scores	One Prior Year Test Score	Current Year Test Score
Grades 6 and 7	Grade 7	Integrated Mathematics II
Grades 7 and 8	Grade 8	Integrated Mathematics II
Grades 7 or 8 and Integrated Mathematics I	Algebra I	Integrated Mathematics II

Table 15.8 Integrated Mathematics III Grade/Content Area Progressions for One- and Two-Year Prior Test Scores

Two Prior Year Test Scores	One Prior Year Test Score	Current Year Test Score
Grades 6 and 7	Grade 7	Integrated Mathematics III
Grades 7 and 8	Grade 8	Integrated Mathematics III
Grades 7 or 8 and Integrated Mathematics I	Algebra I	Integrated Mathematics III
Integrated Mathematics I and Integrated Mathematics II	Integrated Mathematics II	Integrated Mathematics III

In addition to the above progressions, in 2018 the State Leads approved a state-specific SGP progression for one state. In this state, grade 9 students are not required to take the test. Therefore, grade 10 students were not receiving an SGP. For this state, both mathematics and ELA/L progressions were adjusted (see Table 15.9) such that the grade 10 students would receive growth estimates. Other states were not affected by this change.

Table 15.9 State-Specific SGP Progressions

Two Prior Test Scores	One Prior Test Score	Current Test Score
ELA/L Grades 7 and 8	ELA/L Grade 8	ELA/L Grade 10
Mathematics Grade 7 and 8	Mathematics Grade 8	Geometry
Mathematics Grade 7 and Algebra I	Algebra I	Geometry

15.2 Student Growth Percentile Estimation

SGPs are calculated using quantile regression, which describes the conditional distribution of the response variable with greater precision than traditional linear regression, which describes only the conditional mean (Betebenner, 2009). This application of quantile regression uses B-spline smoothing to fit a curvilinear relationship between a norm group's prior and current scores. Cubic B-spline basis functions are used when calculating SGPs to better model the heteroscedasticity, nonlinearity, and skewness in assessment data.

For each group, the quantile regression fits 100 relationships (one for each percentile) between students' prior and current scores. The result is a single coefficient matrix that relates students' prior achievement to their current achievement at each percentile. The National Center for the Improvement of Educational Assessment performed the analyses using Betebenner's (2009) nonlinear quantile-regression based SGP. The analysis was done in the SGP package in R (Betebenner et al., 2017). For details on student growth percentiles, see Betebenner's *A Technical Overview of the Student Growth Percentile Methodology: Student Growth Percentiles and Percentile Growth Projections/Trajectories* (2011).

Betebenner's (2009) SGP model uses Koenker's (2005) quantile regression approach to estimate the conditional density associated with a student's score at administration t conditioned on the student's prior score(s). Quantile regression functions represent the solution to a loss function much in the way that least squares regression represents the solution to a minimization of squared deviations. The conditional quantile functions are parametrized as a linear combination of B-spline basis functions (Wei & He, 2006) to smooth irregularities found in the data. For scores from administration t (where $t \geq 2$), the τ th quantile function for y_t conditional on prior scores (Y_{t-1}, \dots, Y_1) is

$$Q_{Y_t}(\tau | Y_{t-1}, \dots, Y_1) = \sum_{u=1}^{t-1} \sum_{j=1}^n \phi_{ju}(Y_u) \beta_{ju}(\tau) \quad (15-1)$$

where ϕ_{ju} ($j=1,2,\dots, n$ students; $u=1, \dots, t-1$ administrations) represent the B-spline basis functions. The SGP of each student i is the midpoint between the two consecutive τ whose quantile scores capture the student's current score, multiplied by 100. For example, a student with a current score that lies between the fitted value for $\tau = .595$ and $\tau = .605$ would receive an SGP of 60.

SGPs are assumed to be uniformly distributed and uncorrelated with prior achievement. Scale score conditional standard errors of measurement were incorporated for calculation of SGP standard errors of measurement. Goodness of fit results were checked (i.e., uniform distribution of SGPs by prior achievement) for indications of ceiling/floor effects for each SGP norm-group analysis.

15.3 Student Growth Percentile Results/Model Fit for Total Group

The estimation of SGPs was conducted for each student who had at least one prior score. Each analysis is defined by the norm cohort group (grade/sequence). A goodness of fit plot is produced for each analysis run. A ceiling/floor effects test identifies potential problems at the highest obtainable scale scores and lowest obtainable scale scores. Other fit plots compare the observed conditional density of SGP estimates with the theoretical uniform density. If there is perfect model fit, 10 percent of the estimated growth percentiles are expected within each decile band. A Q-Q plot compares the observed distribution with the theoretical distribution; ideally, the step function lines do not deviate much from the ideal line of perfect fit.

Tables 15.10 and 15.11 summarize SGP estimates for the total testing group for ELA/L and mathematics, respectively. SGPs were calculated at the consortium level and, if sample size was sufficient, the state level. Median SGPs were all 50. If the model is a perfect fit, the median is expected to be 50 with norm-referenced data. The minimum SGP is 1 and the maximum SGP is 99. The average standard error for the SGPs is within expectations for these models.

In general, SGPs can be divided into three categories: (1) below 30, indicating that a student is not meeting a year's worth of growth; (2) an SGP of 30–70, indicating that a student did achieve a year's worth of growth; and (3) an SGP over 70, indicating that the student surpassed a year's worth of growth. It is important to note that definitions such as these are not inherent to the SGP method, but rather require expert judgment (Betebenner, 2009). The observed standard errors, ranging from 12.99 to 16.10, support these interpretations (Betebenner et al., 2017).

Table 15.10 Summary of ELA/L SGP Estimates for Total Group

Grade	Sample Size	Average SGP	Average Standard Error	Median SGP
4	96,655	49.99	14.09	50
5	98,944	49.99	14.13	50
6	99,392	50.00	13.66	50
7	99,326	50.00	13.89	50
8	98,201	49.98	13.89	50
9	—	—	—	—
10	—	—	—	—

Note. “—” indicates insufficient sample for SGP calculation for these tests.

Table 15.11 Summary of Mathematics SGP Estimates for Total Group

Grade	Sample Size	Average SGP	Average Standard Error	Median SGP
4	94,459	50.08	13.48	50
5	96,780	50.01	14.15	50
6	97,424	50.01	14.63	50
7	94,704	49.96	15.34	50
8	93,869	49.90	16.75	50
A1	1,455	49.54	15.74	49
GO	1,861	49.63	15.24	50
A2	1,558	49.58	15.89	49

Note. “—” indicates insufficient sample for SGP calculation for these tests.

A1 = Algebra I; GO = Geometry; A2 = Algebra II.

15.4 Student Growth Percentile Results for Subgroups of Interest

Median SGPs are provided for subgroups of interest. With norm-referenced data, the median of all SGPs is expected to be close to 50. Median subgroup growth percentiles below 50 represent growth lower than the median, and median growth percentiles above 50 represent growth higher than the median. Table 15.12 summarizes SGPs for groups of interest for ELA/L grade 4. The ELA/L tables for grades 4 through 8 are provided in the Appendix (Tables A.15.1 through A.15.5). Table 15.13 summarizes SGPs for groups of interest for mathematics grade 4; complete mathematics subgroup results are provided in the Appendix (Tables A.15.6 through A.15.13). Median SGPs for subgroups of interest fell within the band of 30–70, which is considered to be adequate growth. ELA/L grades 10 and 11 had insufficient sample size for SGP subgroup results to be reported.

15.4.1 SGP Results for Gender

English Language Arts/Literacy

The median SGPs for females tend to be higher than the median SGPs for males. The median SGP for females ranges from 48 to 54, whereas the median SGP for males ranges from 46 to 50.5. The standard error for males and females is comparable to the total group.

Mathematics

There was no consistent pattern between median SGPs for females and males. The median SGP for females ranges from 48 to 51, and the median SGP for males ranges from 49 to 51. The standard errors for both are similar to the total group.

15.4.2 SGP Results for Ethnicity

English Language Arts/Literacy

The African American group median SGP ranges from 34 to 47, with students in higher grades at the higher range. Asian/Pacific Islanders tend to have the highest median SGPs, over 60 for all tests but grade 10. American Indian/Alaska Native students had median SGPs ranging from 43 to 52 in grades 5 through 8. The median SGP for Hispanics ranges from 43 to 51. For all ethnicity groups, standard errors are similar to that of the total group.

Mathematics

The median SGP for African Americans ranges from 33 to 41, with the highest growth in mathematics grade 8 and Algebra II. Asian/Pacific Islanders tend to have the highest SGPs across all tests, with a minimum of 51 and a maximum of 66. American Indian/Alaska Native had median SGPs ranging from 31 to 46. The median SGP for Hispanics ranges from 42 to 48. For all ethnicities, the standard errors for all groups are under 20 points.

15.4.3 SGP Results for Special Instructional Needs

English Language Arts/Literacy

Economically disadvantaged and English language learner students tended to have moderate median SGPs. The median SGP ranges from 41 to 48 for economically disadvantaged students and from 40 to 49 for English language learners. Students with disabilities had an observed median SGP of 40 to 44. The standard errors for special instructional needs subgroups are similar to those observed for the total group.

Mathematics

Economically disadvantaged and English language learner students tend to have lower median SGPs than the general population. The median SGP ranges from 39 to 45 for economically disadvantaged students and from 42 to 47 for English language learners. The median SGP for students with disabilities ranges from 34.5 to 47, whereas for students without disabilities the median SGP ranges from 51 to 52. The standard errors for special education students are similar to the total group.

Table 15.12 Summary of SGP Estimates for Subgroups: Grade 4 ELA/L

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	49,196	49.67	14.11	50
Female	47,458	50.31	14.06	50
Ethnicity				
White	51,770	51.44	13.99	52
African American	12,578	42.71	14.41	40
Asian/Pacific Islander	5,241	59.02	13.78	63
American Indian/Alaska Native	207	45.45	14.43	44
Hispanic	21,886	48.48	14.24	48
Multiple	4,703	50.78	13.95	51
Special Instruction Needs				
Economically Disadvantaged	40,641	45.49	14.24	44
Not-economically Disadvantaged	56,014	53.25	13.98	55
English Learner (EL)	14,992	47.00	14.37	46
Non-English Learner	81,663	50.54	14.03	51
Students with Disabilities (SWD)	16,915	41.49	14.39	38
Students without Disabilities	79,740	51.79	14.02	53

15.4.4 SGP Results for Students Taking Spanish Forms

Mathematics

There is a wide range of median growth percentiles for students taking Spanish forms. The sample size is less than 50 for all grade levels. These forms had a slightly higher standard error on average, likely due to lower sample sizes.

Table 15.13 Summary of SGP Estimates for Subgroups: Grade 4 Mathematics

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	48,093	50.44	13.48	51
Female	46,365	49.71	13.49	50
Ethnicity				
White	51,416	50.45	13.05	51
African American	12,328	44.39	14.79	42
Asian/Pacific Islander	5,187	59.13	13.14	63
American Indian/Alaska Native	197	49.16	13.55	48
Hispanic	20,401	50.15	13.88	50
Multiple	4,661	51.37	13.49	51
Special Instruction Needs				
Economically Disadvantaged	39,044	46.91	14.05	46
Not-economically Disadvantaged	55,415	52.32	13.08	53
English Learner (EL)	13,590	49.62	14.13	50
Non-English Learner	80,869	50.16	13.37	50
Students with Disabilities (SWD)	16,526	43.36	14.15	41
Students without Disabilities	77,933	51.51	13.34	52
Spanish Language Form	1,268	42.63	14.57	41

References

- American Educational Research Association [AERA], American Psychological Association [APA], & National Council on Measurement in Education [NCME] (2014). *Standards for educational and psychological testing*. American Educational Research Association.
- Beimers, J. N., Way, W. D., McClarty, K. L., & Miles, J. A. (2012). Evidence based standard setting: Establishing cut scores by integrating research evidence with expert content judgments. *Bulletin*, Issue 21. http://images.pearsonassessments.com/images/tmrs/bulletin21_evidence_based_standard_setting.pdf
- Betebenner, D. (2009). Norm-and criterion-referenced student growth. *Educational Measurement: Issues and Practice*, 28(4), 42–51.
- Betebenner, D. W. (2011). *A technical overview of the student growth percentile methodology: Student growth percentiles and percentile growth projections/trajectories*. National Center for the Improvement of Educational Assessment.
- Betebenner, D. W., Van Iwaarden, A., Domingue, B., & Shang, Y. (2017). *SGP: Student growth percentiles & percentile growth trajectories* (R package version, 1-7 [Computer software]).
- Boyd, A., Minchen, N., & McBride, M. (2018). *Alternative blueprinting options research report*. Pearson.
- Brandt, R., Bercovitz, E., McNally, S., & Zimmerman, L. (2015). *Drawing response interaction usability study for PARCC* (July 28–July 30, 2015). Partnership for Assessment of Readiness for College and Careers. <https://files.eric.ed.gov/fulltext/ED599260.pdf>
- Brandt, R., Bercovitz, E., & Zimmerman, L. (2015). *Drawing response interaction usability study for PARCC November 16–19, 2015*. Pearson. <https://eric.ed.gov/?id=ED599261>
- Brennan, R. L. (2004). *Manual for BB-CLASS: A computer program that uses the beta-binomial model for classification consistency and accuracy (Version 1.0)*. (CASMA Research Report No. 9). Center for Advanced Studies in Measurement, University of Iowa.
- Center for Assessment. (2018). *PARCC comparability review guidelines*.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Lawrence Erlbaum.
- Cortina, J. M. (1993). What is coefficient alpha? An examination of theory and applications. *Journal of Applied Psychology*, 78(1), 98–104.
- Cramer, H. (1946). *Mathematical methods of statistics*. Princeton University Press.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16(3), 297–334.
- Doorey, N., & Polikoff, M. (2016). *Evaluating the content and quality of next generation assessments*. Thomas B. Fordham Institute.
- Dorans, N. J. (2013). *ETS contributions to the quantitative assessment of item, test and score fairness* (ETS R&D Science and Policy Contributions Series, ETS SPC-13-04). Educational Testing Service.
- Dorans, N. J., & Holland, P. W. (1993). DIF detection and description: Mantel-Haenszel and standardization. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 35–66). Lawrence Erlbaum.
- Dorans, N. J., & Schmitt, A. P. (1991). *Constructed response and differential item functioning: A pragmatic approach* (ETS Research Report No. RR-91-47). Educational Testing Service.
- Feldt, L. S., & Brennan, R. (1989). Reliability. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 105–146). Macmillan.
- Holland, P. W., & Thayer, D. T. (1988). Differential item performances and the Mantel-Haenszel procedure. In H. Wainer & H. I. Braun (Eds.), *Test validity* (pp. 129–145). Lawrence Erlbaum.
- Iacus, S. M., King, G., & Porro, G. (2012). Causal inference without balance checking: Coarsened exact matching. *Political Analysis*, 20(1), 1–24. <https://doi.org/10.1093/pan/mpr013>
- Koenker, R. (2005). *Quantile regression*. Cambridge University Press.

- Kolen, M. J. (2004). *POLYCEM windows console version* [Computer software]. The Center for Advanced Studies in Measurement and Assessment (CASMA), University of Iowa.
- Kolen, M. J., Zeng, L., & Hanson, B. A. (1996). Conditional standard errors of measurement for scale scores using IRT. *Journal of Educational Measurement*, 33(2), 129–140.
- Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical data. *Biometrics*, 33(1), 159–174.
- Livingston, S. A., & Lewis, C. (1993). *Estimating the consistency and accuracy of classifications based on test scores* (ETS Research Report No. RR-93-48). Educational Testing Service.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32(2), 179–197.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theories of mental test scores*. Addison-Wesley.
- Lord, F. M., & Wingersky, M. S. (1984). Comparison of IRT true score and equipercentile observed-score “equatings.” *Applied Psychological Measurement*, 8(4), 453–461.
- Mantel, N. (1963). Chi-square tests with one degree of freedom: Extensions of the Mantel-Haenszel procedure. *Journal of the American Statistical Association*, 58(303), 690–700.
- Mantel, N., & Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute*, 22(4), 719–748.
- McClarty, K. L., Korbin, J. L., Moyer, E., Griffin, S., Huth, K., Carey, S., & Medberry, S. (2015). *PARCC benchmarking study*. Pearson Educational Measurement, Pearson.
- McClarty, K. L., Way, W. D., Porter, A. C., Beimers, J. N., & Miles, J. A. (2013). Evidence-based standard setting: establishing a validity framework for cut scores. *Educational Researcher*, 42(2), 78–88.
- Minchen, N., Boyd, A., & McBride, M. (2018a). *Alternative blueprinting options 2018 research report*. Pearson.
- Minchen, N., LaSalle, A., & Boyd, A. (2018b). *Operational study 4: Accessibility of new items/functionality component 4 report*. Pearson.
- Olsson, U., Drasgow, F., & Dorans, N. J. (1982). The polyserial correlation coefficient, *Biometrika*, 47, 337–347.
- Pike, C. K., & Hudson, W. W. (1998). Reliability and measurement error in the presence of homogeneity. *Journal of Social Service Research*, 24(1–2), 149–163.
- Plake, B. S., Ferdous, A. A., Impara, J. C., & Buckendahl, C. W. (2005). Setting multiple performance standards using the Yes/No method: An alternative item mapping method [Paper presentation]. Annual meeting of the National Council on Measurement in Education, Montreal, Quebec, Canada.
- Schmitt, N. (1996). Uses and abuses of coefficient alpha. *Psychological Assessment*, 8(4), 350–353.
- Schultz, S. R., Michaels, H. R., Norman Dvorak, R., & Wiley, C. R. H. (2016). *Evaluating the content and quality of next generation high school assessments*. (HumRRO Report 2016 No. 001). Human Resources Research Organization.
- Schultz, S. R., Norman Dvorak, R., & Chen, J. (2017). *Evaluating the quality and alignment of PARCC ELA/literacy and mathematics assessments: Grades 3, 4, 6, and 7*. (HumRRO Report 2017 No. 040). Human Resources Research Organization.
- Sireci, S. G., Thissen, D., & Wainer, H. (1991). On the reliability of testlet-based tests. *Journal of Educational Measurement*, 28(3), 237–247.
- Steedle, J., & LaSalle, A. (2016). *Operational study 4: Accessibility of new items/functionality component 3 report*. Pearson.
- Steedle, J., Quesen, S., & Boyd, A. (2017). *Longitudinal study of external validity of the PARCC performance levels: Phase I report*. Pearson.
- Tavakol, M., & Dennick, R. (2011). Making sense of Cronbach’s alpha. *International Journal of Medical Education*, 2, 53–55. <https://doi.org/10.5116/ijme.4dfb.8dfd>
- Thompson, S. J., Johnstone, C. J., & Thurlow, M. L. (2002). *Universal design applied to large scale assessments* (Synthesis Report 44). University of Minnesota, National Center on Educational Outcomes.
- Wainer, H., & Thissen, D. (2001). *Test scoring*. Lawrence Erlbaum.

- Wei, Y., & He, X. (2006). Conditional growth charts. *Annals of Statistics*, 34(5), 2069–2097.
- Williamson, D. M., Xi, X., & Breyer, F. J. (2012). A framework for evaluation and use of automated scoring. *Educational Measurement: Issues and Practices*, 31(1), 2–13.
- 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.
- Zenisky, A. L., Hambleton, R. K., & Sireci, S.C. (2003). *Effects of local dependence on the validity of IRT item test, and ability statistics* (Technical Report). American College Admissions Test.
- Zieky, M. (1993). Practical questions in the use of DIF statistics in test development. In P. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 337–348). Lawrence Erlbaum
- Zwick, R., Thayer, D. T., & Mazzeo, J. (1997). *Describing and categorizing DIF in polytomous items* (ETS Research Report RR-97-05). Educational Testing Service.

Appendices

Appendix 6: Summary of Differential Item Function (DIF) Results

Table A.6.1 Pre-Administration Differential Item Functioning for ELA/L Grade 3

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	39			1	3	38	97				
White vs. Black	39					39	100				
White vs. Hispanic	39			2	5	37	95				
White vs. Asian	39					38	97	1	3		
White vs. AmerIndian	39					39	100				
White vs. Pacific Islander	39			2	5	36	92	1	3		
White vs. Multiracial	39					39	100				
NoEcnDis vs. EcnDis	39					39	100				
ELN vs. ELY	39			4	10	35	90				
SWDN vs. SWDY	39			1	3	38	97				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability, SWDY = student with disability.

Table A.6.2 Pre-Administration Differential Item Functioning for ELA/L Grade 4

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	60			2	3	53	88	4	7	1	2
White vs. Black	60	2	3	4	7	54	90				
White vs. Hispanic	60	1	2	3	5	56	93				
White vs. Asian	60			1	2	59	98				
White vs. AmerIndian	60			1	2	59	98				
White vs. Pacific Islander	60			2	3	57	95	1	2		
White vs. Multiracial	60					60	100				
NoEcnDis vs. EcnDis	60	1	2	1	2	58	97				
ELN vs. ELY	60	2	3	3	5	55	92				
SWDN vs. SWDY	60			2	3	58	97				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.3 Pre-Administration Differential Item Functioning for ELA/L Grade 5

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	56			3	5	50	89	3	5		
White vs. Black	56			4	7	52	93				
White vs. Hispanic	56			3	5	53	95				
White vs. Asian	56					54	96	2	4		
White vs. AmerIndian	56			1	2	55	98				
White vs. Pacific Islander	56			1	2	55	98				
White vs. Multiracial	56					56	100				
NoEcnDis vs. EcnDis	56					56	100				
ELN vs. ELY	56			6	11	50	89				
SWDN vs. SWDY	56	1	2	2	4	53	95				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.4 Pre-Administration Differential Item Functioning for ELA/L Grade 6

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	46	1	2	3	7	40	87	2	4		
White vs. Black	46	2	4	3	7	41	89				
White vs. Hispanic	46	1	2	3	7	42	91				
White vs. Asian	46			1	2	45	98				
White vs. AmerIndian	46	2	4	3	7	40	87	1	2		
White vs. Pacific Islander	46	1	2	3	7	42	91				
White vs. Multiracial	46					46	100				
NoEcnDis vs. EcnDis	46			3	7	43	93				
ELN vs. ELY	46			9	20	37	80				
SWDN vs. SWDY	46	1	2	2	4	43	93				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.5 Pre-Administration Differential Item Functioning for ELA/L Grade 7

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	50			7	14	42	84	1	2		
White vs Black	50			2	4	48	96				
White vs. Hispanic	50	1	2	2	4	47	94				
White vs. Asian	50					49	98			1	2
White vs. AmerIndian	50			1	2	49	98				
White vs. Pacific Islander	50	1	2	2	4	47	94				
White vs. Multiracial	50					50	100				
NoEcnDis vs. EcnDis	50			1	2	49	98				
ELN vs. ELY	50	2	4	4	8	44	88				
SWDN vs. SWDY	50					50	100				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.6 Pre-Administration Differential Item Functioning for ELA/L Grade 8

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	58	1	2	2	3	53	91	2	3		
White vs. Black	58	1	2	1	2	56	97				
White vs. Hispanic	58	1	2	2	3	55	95				
White vs. Asian	58					56	97			2	3
White vs. AmerIndian	58					58	100				
White vs. Pacific Islander	58					58	100				
White vs. Multiracial	58					58	100				
NoEcnDis vs. EcnDis	58					58	100				
ELN vs. ELY	58	1	2	5	9	52	90				
SWDN vs. SWDY	58					58	100				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.7 Pre-Administration Differential Item Functioning for ELA/L Grade 9

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	60	1	2	5	8	52	87	2	3		
White vs. Black	60			2	3	58	97				
White vs. Hispanic	60			1	2	59	98				
White vs. Asian	60					60	100				
White vs. AmerIndian	60			2	3	58	97				
White vs. Pacific Islander	60			1	2	59	98				
White vs. Multiracial	60					60	100				
NoEcnDis vs. EcnDis	60					60	100				
ELN vs. ELY	60	4	7	7	12	48	80	1	2		
SWDN vs. SWDY	60			2	3	58	97				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.8 Pre-Administration Differential Item Functioning for ELA/L Grade 10

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	47			4	9	41	87	2	4		
White vs. Black	47					47	100				
White vs. Hispanic	47			1	2	46	98				
White vs. Asian	47					47	100				
White vs. AmerIndian	47	1	2	1	2	45	96				
White vs. Pacific Islander	47					47	100				
White vs. Multiracial	47					47	100				
NoEcnDis vs. EcnDis	47					47	100				
ELN vs. ELY	47	3	6	3	6	41	87				
SWDN vs. SWDY	47	1	2			46	98				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.9 Pre-Administration Differential Item Functioning for Mathematics Grade 3

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	89			3	3	85	96	1	1		
White vs. Black	89			4	4	84	94	1	1		
White vs. Hispanic	89					89	100				
White vs. Asian	89					85	96	3	3	1	1
White vs. AmerIndian	89			2	2	87	98				
White vs. Pacific Islander	89			1	1	88	99				
White vs. Multiracial	89			1	1	87	98	1	1		
NoEcnDis vs. EcnDis	89					89	100				
ELN vs. ELY	89			1	1	88	99				
SWDN vs. SWDY	89			1	1	88	99				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.10 Pre-Administration Differential Item Functioning for Mathematics Grade 4

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	84	1	1	4	5	78	93	1	1		
White vs. Black	84	1	1	5	6	78	93				
White vs. Hispanic	84					84	100				
White vs. Asian	84					82	98	2	2		
White vs. AmerIndian	84			3	4	79	94	2	2		
White vs. Pacific Islander	84					82	98	2	2		
White vs. Multiracial	84			1	1	83	99				
NoEcnDis vs. EcnDis	84					84	100				
ELN vs. ELY	84	1	1	1	1	82	98				
SWDN vs. SWDY	84			1	1	83	99				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.11 Pre-Administration Differential Item Functioning for Mathematics Grade 5

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	85	1	1			84	99				
White vs. Black	85			5	6	79	93	1	1		
White vs. Hispanic	85					85	100				
White vs. Asian	85			1	1	79	93	5	6		
White vs. AmerIndian	85	1	1	3	4	81	95				
White vs. Pacific Islander	85					85	100				
White vs. Multiracial	85			1	1	84	99				
NoEcnDis vs. EcnDis	85					85	100				
ELN vs. ELY	85	1	1	3	4	81	95				
SWDN vs. SWDY	85			2	2	81	95	2	2		

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.12 Pre-Administration Differential Item Functioning for Mathematics Grade 6

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	76	1	1	5	7	69	91	1	1		
White vs. Black	76			4	5	71	93	1	1		
White vs. Hispanic	76					76	100				
White vs. Asian	76					73	96	3	4		
White vs. AmerIndian	76			1	1	74	97	1	1		
White vs. Pacific Islander	76			1	1	75	99				
White vs. Multiracial	76					76	100				
NoEcnDis vs. EcnDis	76					76	100				
ELN vs. ELY	76			4	5	71	93	1	1		
SWDN vs. SWDY	76	1	1	1	1	70	92	2	3	2	3

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.13 Pre-Administration Differential Item Functioning for Mathematics Grade 7

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	80					78	98	2	3		
White vs. Black	80					80	100				
White vs. Hispanic	80					80	100				
White vs. Asian	80					79	99	1	1		
White vs. AmerIndian	80					80	100				
White vs. Pacific Islander	80					80	100				
White vs. Multiracial	80					80	100				
NoEcnDis vs. EcnDis	80					80	100				
ELN vs. ELY	80	2	3	2	3	76	95				
SWDN vs. SWDY	80			1	1	79	99				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.14 Pre-Administration Differential Item Functioning for Mathematics Grade 8

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	77			1	1	76	99				
White vs. Black	77			4	5	71	92	2	3		
White vs. Hispanic	77					77	100				
White vs. Asian	77			1	1	74	96	2	3		
White vs. AmerIndian	77			4	5	73	95				
White vs. Pacific Islander	77			1	1	75	97	1	1		
White vs. Multiracial	77			1	1	76	99				
NoEcnDis vs. EcnDis	77					77	100				
ELN vs. ELY	77	1	1	4	5	70	91	2	3		
SWDN vs. SWDY	77			1	1	74	96	2	3		

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.15 Pre-Administration Differential Item Functioning for Algebra I

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	71			2	3	69	97				
White vs. Black	71			1	1	69	97	1	1		
White vs. Hispanic	71					71	100				
White vs. Asian	71					61	86	9	13	1	1
White vs. AmerIndian	71			4	6	67	94				
White vs. Pacific Islander	71			2	3	69	97				
White vs. Multiracial	71			1	1	69	97	1	1		
NoEcnDis vs. EcnDis	71					71	100				
ELN vs. ELY	71	1	1	3	4	65	92	2	3		
SWDN vs. SWDY	71					71	100				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.16 Pre-Administration Differential Item Functioning for Geometry

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	77			2	3	75	97				
White vs. Black	77			4	5	73	95				
White vs. Hispanic	77			2	3	75	97				
White vs. Asian	77					71	92	6	8		
White vs. AmerIndian	77	1	1	9	12	67	87				
White vs. Pacific Islander	77					77	100				
White vs. Multiracial	77			1	1	75	97	1	1		
NoEcnDis vs. EcnDis	77			2	3	75	97				
ELN vs. ELY	77	3	4	7	9	65	84	2	3		
SWDN vs. SWDY	77	1	1	2	3	73	95	1	1		

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Table A.6.17 Pre-Administration Differential Item Functioning for Algebra II

DIF Comparison	Total N of Unique Items	C- DIF		B- DIF		A DIF		B+ DIF		C+ DIF	
		N	% of Total	N	% of Total	N	% of Total	N	% of Total	N	% of Total
Female vs. Male	76			1	1	75	99				
White vs. Black	76			2	3	74	97				
White vs. Hispanic	76			2	3	74	97				
White vs. Asian	76			1	1	63	83	10	13	2	3
White vs. AmerIndian	76	1	1	4	5	71	93				
White vs. Pacific Islander	76					75	99	1	1		
White vs. Multiracial	76					75	99	1	1		
NoEcnDis vs. EcnDis	76					76	100				
ELN vs. ELY	76	1	1	3	4	71	93	1	1		
SWDN vs. SWDY	76			2	3	74	97				

Note. AmerIndian = American Indian/Alaska Native; Black = Black/African American; Hispanic = Hispanic/Latino; Pacific Islander = Native Hawaiian or Pacific Islander; Multiracial = Multiple Race Selected; NoEcnDis = not economically disadvantaged; EcnDis = economically disadvantaged; ELN = not an English learner; ELY = English learner; SWDN = not student with disability; SWDY = student with disability.

Appendix 7.1: Pre-Equated IRT Results for Spring 2022 English Language Arts/Literacy (ELA/L)

Table A.7.1 Pre-Equated IRT Summary Parameter Estimates for All Items for ELA/L by Grade

Grade	Item Grouping	No. of Score Points	No. of Items	<i>b</i> Estimates Summary				<i>a</i> Estimates Summary			
				Mean	SD	Min.	Max.	Mean	SD	Min.	Max.
3	All Items	98	44	0.37	0.98	-1.64	2.40	0.57	0.22	0.19	1.04
	Reading	68	34	0.05	0.88	-1.64	2.40	0.48	0.17	0.19	0.84
	Writing	30	10	1.45	0.15	1.28	1.69	0.85	0.10	0.72	1.04
4	All Items	148	66	0.40	0.92	-1.99	2.66	0.47	0.23	0.17	0.99
	Reading	108	54	0.23	0.93	-1.99	2.66	0.37	0.13	0.17	0.77
	Writing	40	12	1.13	0.35	0.71	1.83	0.88	0.05	0.81	0.99
5	All Items	141	62	0.55	0.96	-1.34	3.59	0.48	0.23	0.13	0.99
	Reading	100	50	0.37	0.96	-1.34	3.59	0.39	0.15	0.13	0.77
	Writing	41	12	1.29	0.51	0.60	2.12	0.86	0.09	0.72	0.99
6	All Items	112	50	0.34	0.82	-1.00	2.95	0.52	0.21	0.18	1.16
	Reading	84	42	0.19	0.79	-1.00	2.95	0.45	0.15	0.18	0.79
	Writing	28	8	1.14	0.39	0.68	1.86	0.88	0.15	0.75	1.16
7	All Items	125	55	0.37	0.77	-1.21	1.77	0.50	0.27	0.13	1.23
	Reading	90	45	0.31	0.82	-1.21	1.77	0.39	0.15	0.13	0.75
	Writing	35	10	0.61	0.40	0.18	1.46	0.98	0.17	0.71	1.23
8	All Items	146	64	0.31	0.78	-1.42	2.83	0.47	0.24	0.08	1.06
	Reading	104	52	0.24	0.84	-1.42	2.83	0.38	0.13	0.08	0.70
	Writing	42	12	0.59	0.42	-0.17	1.19	0.90	0.12	0.68	1.06
9	All Items	150	66	0.55	0.88	-1.21	2.77	0.50	0.29	0.17	1.22
	Reading	108	54	0.50	0.96	-1.21	2.77	0.38	0.15	0.17	0.68
	Writing	42	12	0.74	0.30	0.32	1.21	1.04	0.09	0.92	1.22
10	All Items	119	52	0.68	0.78	-1.08	2.67	0.49	0.29	0.13	1.18
	Reading	84	42	0.66	0.84	-1.08	2.67	0.38	0.17	0.13	0.94
	Writing	35	10	0.79	0.38	0.27	1.35	0.97	0.11	0.77	1.18
11	All Items	208	92	1.00	0.98	-1.09	5.29	0.47	0.25	0.09	1.13
	Reading	152	76	0.99	1.08	-1.09	5.29	0.38	0.16	0.09	0.84
	Writing	56	16	1.08	0.29	0.61	1.53	0.89	0.17	0.63	1.13

Note. SD = standard deviation.

Appendix 7.2: Pre-Equated IRT Results for Spring 2022 Mathematics

Table A.7.2 Pre-Equated IRT Summary Parameter Estimates for All Items for Mathematics by Grade/Subject

Grade	Item Grouping	No. of Score Points	No. of Items	<i>b</i> Estimates Summary				<i>a</i> Estimates Summary			
				Mean	SD	Min.	Max.	Mean	SD	Min.	Max.
3	All Items	142	90	-0.38	1.14	-3.29	2.48	0.76	0.26	0.26	1.35
	SSMC	29	29	-0.88	1.05	-2.86	1.72	0.75	0.21	0.26	1.15
	CR	113	61	-0.14	1.12	-3.29	2.48	0.77	0.28	0.29	1.35
	Type I	82	74	-0.67	1.01	-3.29	1.72	0.81	0.25	0.26	1.35
	Type II	27	8	1.14	0.71	0.35	2.48	0.47	0.09	0.37	0.62
	Type III	33	8	0.84	0.48	0.10	1.58	0.63	0.15	0.46	0.85
4	All Items	148	88	-0.11	0.97	-2.65	1.50	0.74	0.18	0.32	1.35
	SSMC	27	27	-1.01	0.84	-2.65	0.64	0.76	0.16	0.48	1.13
	CR	121	61	0.28	0.74	-1.79	1.50	0.74	0.19	0.32	1.35
	Type I	86	71	-0.35	0.93	-2.65	1.50	0.77	0.19	0.48	1.35
	Type II	26	8	0.96	0.24	0.57	1.38	0.68	0.07	0.57	0.78
	Type III	36	9	0.80	0.32	0.20	1.09	0.64	0.18	0.32	0.92
5	All Items	152	91	0.08	1.10	-2.33	2.13	0.65	0.23	0.16	1.50
	SSMC	33	33	-0.44	1.15	-2.33	2.13	0.59	0.29	0.16	1.50
	CR	119	58	0.38	0.95	-2.16	1.90	0.69	0.18	0.36	1.11
	Type I	89	74	-0.12	1.10	-2.33	2.13	0.66	0.25	0.16	1.50
	Type II	27	8	0.82	0.65	-0.15	1.89	0.59	0.11	0.47	0.73
	Type III	36	9	1.07	0.35	0.69	1.54	0.65	0.17	0.44	0.92
6	All Items	129	76	0.47	0.82	-1.88	2.05	0.70	0.26	0.28	1.34
	SSMC	18	18	0.19	0.80	-0.94	1.64	0.60	0.28	0.28	1.16
	CR	111	58	0.56	0.81	-1.88	2.05	0.73	0.25	0.35	1.34
	Type I	79	62	0.34	0.81	-1.88	2.05	0.73	0.28	0.28	1.34
	Type II	23	7	0.78	0.74	-0.49	1.59	0.57	0.14	0.43	0.83
	Type III	27	7	1.27	0.34	0.76	1.69	0.61	0.09	0.49	0.76
7	All Items	128	81	0.60	1.02	-1.78	4.22	0.66	0.26	0.19	1.49
	SSMC	29	29	0.23	1.06	-1.56	2.57	0.56	0.30	0.19	1.49
	CR	99	52	0.81	0.95	-1.78	4.22	0.71	0.23	0.30	1.26
	Type I	83	68	0.53	1.08	-1.78	4.22	0.67	0.28	0.19	1.49
	Type II	27	8	0.82	0.55	-0.44	1.40	0.60	0.11	0.43	0.80

Grade	Item Grouping	No. of Score Points	No. of Items	<i>b</i> Estimates Summary				<i>a</i> Estimates Summary			
				Mean	SD	Min.	Max.	Mean	SD	Min.	Max.
	Type III	18	5	1.16	0.35	0.70	1.64	0.55	0.11	0.38	0.69
8	All Items	126	78	1.00	0.94	-1.42	2.89	0.58	0.27	0.16	1.40
	SSMC	31	31	0.65	1.05	-1.42	2.89	0.42	0.17	0.16	0.74
	CR	95	47	1.23	0.78	-1.26	2.66	0.68	0.27	0.18	1.40
	Type I	82	66	0.92	0.98	-1.42	2.89	0.59	0.28	0.16	1.40
	Type II	20	6	1.55	0.61	1.08	2.66	0.60	0.17	0.49	0.91
	Type III	24	6	1.31	0.36	0.95	1.91	0.45	0.08	0.35	0.56
A1	All Items	137	72	1.22	1.11	-1.16	3.57	0.65	0.29	0.17	1.62
	SSMC	23	23	0.87	1.38	-1.16	3.57	0.48	0.21	0.17	0.85
	CR	114	49	1.38	0.93	-0.96	3.34	0.73	0.29	0.24	1.62
	Type I	77	57	1.07	1.19	-1.16	3.57	0.65	0.32	0.17	1.62
	Type II	27	8	1.71	0.39	0.95	2.28	0.73	0.15	0.55	0.91
	Type III	33	7	1.86	0.45	1.45	2.60	0.61	0.12	0.41	0.76
GO	All Items	144	81	0.95	1.05	-1.60	3.83	0.75	0.38	0.19	1.76
	SSMC	19	19	0.21	1.27	-1.60	3.83	0.55	0.33	0.26	1.41
	CR	125	62	1.18	0.86	-0.67	2.67	0.81	0.38	0.19	1.76
	Type I	85	67	0.80	1.08	-1.60	3.83	0.74	0.41	0.19	1.76
	Type II	23	7	1.68	0.53	0.96	2.39	0.78	0.20	0.48	1.04
	Type III	36	7	1.72	0.33	1.29	2.07	0.80	0.21	0.61	1.13
A2	All Items	158	79	1.29	1.10	-1.39	3.90	0.61	0.26	0.19	1.20
	SSMC	23	23	0.46	1.04	-1.39	2.48	0.59	0.21	0.26	1.12
	CR	135	56	1.63	0.94	-0.13	3.90	0.61	0.28	0.19	1.20
	Type I	83	62	1.09	1.14	-1.39	3.90	0.62	0.27	0.19	1.20
	Type II	27	8	1.88	0.39	1.36	2.57	0.66	0.19	0.46	1.07
	Type III	48	9	2.13	0.46	1.52	2.90	0.44	0.16	0.28	0.80
M1	All Items	62	34	1.10	1.09	-0.95	4.02	0.61	0.31	0.11	1.61
	SSMC	13	13	0.90	1.16	-0.06	4.02	0.45	0.20	0.11	0.77
	CR	49	21	1.23	1.05	-0.95	2.85	0.71	0.33	0.18	1.61
	Type I	37	28	0.86	1.02	-0.95	4.02	0.60	0.34	0.11	1.61
	Type II	10	3	1.80	0.48	1.31	2.26	0.60	0.04	0.57	0.65
	Type III	15	3	2.66	0.20	2.46	2.85	0.75	0.10	0.68	0.87

Grade	Item Grouping	No. of Score Points	No. of Items	<i>b</i> Estimates Summary				<i>a</i> Estimates Summary			
				Mean	SD	Min.	Max.	Mean	SD	Min.	Max.
M2	All Items	55	30	1.46	1.23	-0.97	3.96	0.56	0.27	0.06	1.41
	SSMC	13	13	1.06	1.20	-0.97	3.75	0.42	0.18	0.06	0.71
	CR	42	17	1.77	1.20	-0.50	3.96	0.67	0.28	0.31	1.41
	Type I	30	24	1.27	1.10	-0.97	3.75	0.57	0.29	0.06	1.41
	Type II	10	3	2.78	1.03	2.04	3.96	0.53	0.14	0.41	0.68
	Type III	15	3	1.70	1.91	-0.50	2.98	0.50	0.17	0.31	0.64
M3	All Items	55	29	1.52	1.44	-1.02	4.30	0.52	0.26	0.17	1.24
	SSMC	12	12	1.49	1.83	-1.02	4.30	0.47	0.32	0.17	1.24
	CR	43	17	1.54	1.16	-0.35	2.80	0.56	0.21	0.25	1.08
	Type I	30	23	1.44	1.56	-1.02	4.30	0.52	0.28	0.17	1.24
	Type II	10	3	1.90	1.01	0.73	2.53	0.50	0.11	0.42	0.63
	Type III	15	3	1.81	0.99	0.71	2.62	0.52	0.25	0.25	0.76

Note. SD = standard deviation; SSMC = single select multiple choice; CR =constructed response; A1 = Algebra I; GO = Geometry; A2 = Algebra II; M1 = Integrated Mathematics I; M2 = Integrated Mathematics II; M3 = Integrated Mathematics III.

Appendix 11: Students by Grade/Subject and Mode, for Each State

Table A.11.1 All ELA/L Test Takers, by State and Grade

State	Category	Total	English Language Arts/Literacy							
			Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	Grade 9	Grade 10
All States	N of Students	1,542,916	229,185	232,317	237,847	238,260	245,106	249,388	103,744	7,069
	N of CBT	1,539,584	228,358	231,743	237,297	237,783	244,677	248,973	103,693	7,060
	% of CBT	99.8	99.6	99.8	99.8	99.8	99.8	99.8	100.0	99.9
	N of PBT	3,332	827	574	550	477	429	415	51	n/r
	% of PBT	0.2	0.4	0.2	0.2	0.2	0.2	0.2	0.0	n/r
DC	% of All Data	2.8	0.4	0.4	0.4	0.4	0.3	0.3	0.3	0.3
	N of Students	42,838	6,234	6,163	5,885	5,419	5,259	5,110	4,680	4,088
	N of CBT	42,794	6,227	6,161	5,879	5,412	5,257	5,100	4,676	4,082
	% of CBT	99.9	99.9	100.0	99.9	99.9	100.0	99.8	99.9	99.9
	N of PBT	44	n/r	n/r	n/r	n/r	n/r	n/r	n/r	n/r
	% of PBT	0.1	n/r	n/r	n/r	n/r	n/r	n/r	n/r	n/r
DD	% of All Data	2.3	0.4	0.4	0.4	0.3	0.3	0.3	n/a	0.2
	N of Students	33,016	5,803	5,559	5,428	4,779	4,458	4,008	n/a	2,981
	N of CBT	32,904	5,773	5,525	5,399	4,773	4,453	4,003	n/a	2,978
	% of CBT	99.7	99.5	99.4	99.5	99.9	99.9	99.9	n/a	99.9
	N of PBT	112	30	34	29	n/r	n/r	n/r	n/a	n/r
	% of PBT	0.3	0.5	0.6	0.5	n/r	n/r	n/r	n/a	n/r

Table A.11.1 All ELA/L Test Takers, by State and Grade

State	Category	Total	English Language Arts/Literacy							
			Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	Grade 9	Grade 10
IL	% of All Data	51.3	8.1	8.2	8.5	8.6	8.8	9.1	n/a	n/a
	N of Students	792,057	125,191	127,131	131,381	132,030	136,346	139,978	n/a	n/a
	N of CBT	789,166	124,430	126,628	130,911	131,600	135,972	139,625	n/a	n/a
	% of CBT	99.6	99.4	99.6	99.6	99.7	99.7	99.7	n/a	n/a
	N of PBT	2,891	761	503	470	430	374	353	n/a	n/a
	% of PBT	0.4	0.6	0.4	0.4	0.3	0.3	0.3	n/a	n/a
NJ	% of All Data	43.8	6.0	6.1	6.2	6.2	6.4	6.5	6.4	n/a
	N of Students	675,005	91,957	93,464	95,153	96,032	99,043	100,292	99,064	n/a
	N of CBT	674,720	91,928	93,429	95,108	95,998	98,995	100,245	99,017	n/a
	% of CBT	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	n/a
	N of PBT	285	29	35	45	34	48	47	47	n/a
	% of PBT	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	n/a

Note. DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; CBT = computer-based test; PBT = paper-based test; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.2 All Mathematics Test Takers, by State and Grade

State	Category	Mathematics									
		Total	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	A1	G0	A2
All States	N of Students	1,556,982	230,054	232,955	238,545	238,630	235,692	211,563	115,174	40,404	13,965
	N of CBT	1,553,846	229,434	232,377	237,994	238,150	235,274	211,159	115,106	40,389	13,963
	% of CBT	99.8	99.7	99.8	99.8	99.8	99.8	99.8	99.9	100.0	100.0
	N of PBT	3,136	620	578	551	480	418	404	68	n/r	n/r
	% of PBT	0.2	0.3	0.2	0.2	0.2	0.2	0.2	0.1	n/r	n/r
DC	% of All Data	2.7	0.4	0.4	0.4	0.3	0.3	0.3	0.3	0.3	0.0
	N of Students	41,950	6,261	6,183	5,877	5,411	5,197	3,951	4,825	4,019	226
	N of CBT	41,899	6,250	6,181	5,871	5,404	5,195	3,942	4,817	4,013	226
	% of CBT	99.9	99.8	100.0	99.9	99.9	100.0	99.8	99.8	99.9	100.0
	N of PBT	51	n/r	n/r	n/r	n/r	n/r	n/r	n/r	n/r	n/r
	% of PBT	0.1	n/r	n/r	n/r	n/r	n/r	n/r	n/r	n/r	n/r
DD	% of All Data	2.2	0.4	0.4	0.3	0.3	n/a	0.2	0.2	0.2	0.2
	N of Students	33,602	5,810	5,513	5,405	4,726	n/a	2,773	3,651	3,053	2,671
	N of CBT	33,482	5,777	5,481	5,372	4,716	n/a	2,770	3,644	3,052	2,670
	% of CBT	99.6	99.4	99.4	99.4	99.8	n/a	99.9	99.8	100.0	100.0
	N of PBT	120	33	32	33	n/r	n/a	n/r	n/r	n/r	n/r
	% of PBT	0.4	0.6	0.6	0.6	n/r	n/a	n/r	n/r	n/r	n/r

Table A.11.2 All Mathematics Test Takers, by State and Grade

Mathematics											
State	Category	Total	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	A1	GO	A2
IL	% of All Data	50.5	8.0	8.1	8.4	8.4	8.7	8.9	n/a	n/a	n/a
	N of Students	788,132	124,674	126,558	130,894	131,359	135,578	139,069	n/a	n/a	n/a
	N of CBT	785,464	124,129	126,048	130,426	130,930	135,211	138,720	n/a	n/a	n/a
	% of CBT	99.7	99.6	99.6	99.6	99.7	99.7	99.7	n/a	n/a	n/a
	N of PBT	2,668	545	510	468	429	367	349	n/a	n/a	n/a
	% of PBT	0.3	0.4	0.4	0.4	0.3	0.3	0.3	n/a	n/a	n/a
NJ	% of All Data	44.5	6.0	6.1	6.2	6.2	6.1	4.2	6.9	2.1	0.7
	N of Students	693,298	93,309	94,701	96,369	97,134	94,917	65,770	106,698	33,332	11,068
	N of CBT	693,001	93,278	94,667	96,325	97,100	94,868	65,727	106,645	33,324	11,067
	% of CBT	100.0	100.0	100.0	100.0	100.0	99.9	99.9	100.0	100.0	100.0
	N of PBT	297	31	34	44	34	49	43	53	n/r	n/r
	% of PBT	0.0	0.0	0.0	0.0	0.0	0.1	0.1	0.0	n/r	n/r

Note. DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; A1 = Algebra I; GO = Geometry; A2 = Algebra II; CBT = computer-based test; PBT = paper-based test; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.3 All Spanish-Language Mathematics Test Takers, by State and Grade

State	Category	Mathematics									
		Total	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	A1	GO	A2
All States	N of Students	26,805	5,063	4,605	4,220	3,721	3,226	2,844	2,563	425	138
	N of CBT	26,769	5,048	4,599	4,215	3,718	3,224	2,840	2,562	425	138
	% of CBT	99.9	99.7	99.9	99.9	99.9	99.9	99.9	100.0	100.0	100.0
	N of PBT	36	n/r	n/r	n/r	n/r	n/r	n/r	n/r	n/a	n/a
	% of PBT	0.1	n/r	n/r	n/r	n/r	n/r	n/r	n/r	n/a	n/a
DC	% of All Data	2.2	0.4	0.3	0.3	0.2	0.2	0.3	0.4	0.1	n/a
	N of Students	609	103	79	83	62	65	75	103	39	n/a
	N of CBT	606	100	79	83	62	65	75	103	39	n/a
	% of CBT	99.5	97.1	100.0	100.0	100.0	100.0	100.0	100.0	100.0	n/a
	N of PBT	n/r	n/r	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
	% of PBT	n/r	n/r	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
DD	% of All Data	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
	N of Students	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
	N of CBT	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
	% of CBT	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
	N of PBT	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a
	% of PBT	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a

Table A.11.3 All Spanish-Language Mathematics Test Takers, by State and Grade

Mathematics											
State	Category	Total	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	A1	GO	A2
IL	% of All Data	42.8	10.7	9.3	8.3	6.3	4.5	3.7	n/a	n/a	n/a
	N of Students	11,497	2,859	2,500	2,230	1,689	1,218	1,001	n/a	n/a	n/a
	N of CBT	11,467	2,847	2,494	2,226	1,686	1,217	997	n/a	n/a	n/a
	% of CBT	99.7	99.6	99.8	99.8	99.8	99.9	99.6	n/a	n/a	n/a
	N of PBT	30	n/r	n/r	n/r	n/r	n/r	n/r	n/a	n/a	n/a
	% of PBT	0.3	n/r	n/r	n/r	n/r	n/r	n/r	n/a	n/a	n/a
NJ	% of All Data	54.7	7.8	7.6	7.1	7.3	7.2	6.6	9.2	1.4	0.5
	N of Students	14,699	2,101	2,026	1,907	1,970	1,943	1,768	2,460	386	138
	N of CBT	14,696	2,101	2,026	1,906	1,970	1,942	1,768	2,459	386	138
	% of CBT	100.0	100.0	100.0	99.9	100.0	99.9	100.0	100.0	100.0	100.0
	N of PBT	n/r	n/a	n/a	n/r	n/a	n/r	n/a	n/r	n/a	n/a
	% of PBT	n/r	n/a	n/a	n/r	n/a	n/r	n/a	n/r	n/a	n/a

Note. DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; A1 = Algebra I; GO = Geometry; A2 = Algebra II; CBT = computer-based test; PBT = paper-based test; n/a = not applicable; n/r = not reported due to n<20.

*No students in DD tested in mathematics using Spanish-language forms.

Table A.11.4 All States Combined: ELA/L Test Takers, by Grade, Mode, and Gender

Grade	Mode	Valid Cases	Female		Male	
			N	%	N	%
3	All	229,178	112,531	49.1	116,647	50.9
	CBT	228,352	112,199	49.1	116,153	50.9
	PBT	826	332	40.2	494	59.8
4	All	232,314	114,052	49.1	118,262	50.9
	CBT	231,740	113,826	49.1	117,914	50.9
	PBT	574	226	39.4	348	60.6
5	All	237,830	116,510	49.0	121,320	51.0
	CBT	237,280	116,252	49.0	121,028	51.0
	PBT	550	258	46.9	292	53.1
6	All	238,216	116,268	48.8	121,948	51.2
	CBT	237,739	116,076	48.8	121,663	51.2
	PBT	477	192	40.3	285	59.7
7	All	245,047	119,492	48.8	125,555	51.2
	CBT	244,619	119,326	48.8	125,293	51.2
	PBT	428	166	38.7	262	61.1
8	All	249,309	122,108	49.0	127,201	51.0
	CBT	248,895	121,935	49.0	126,960	51.0
	PBT	414	173	41.7	241	58.1
9	All	103,624	50,420	48.6	53,204	51.3
	CBT	103,575	50,397	48.6	53,178	51.3
	PBT	49	23	45.1	26	51.0
10	All	7,063	3,446	48.7	3,617	51.2
	CBT	7,054	3,443	48.8	3,611	51.1
	PBT	n/r	n/r	n/r	n/r	n/r

Note. DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; CBT = computer-based test; PBT = paper-based test; n/r = not reported due to n<20.

Table A.11.5 All States Combined: Mathematics Test Takers, by Grade, Mode, and Gender

Grade	Mode	Valid Cases	Female		Male	
			N	%	N	%
3	All	230,047	112,938	49.1	117,109	50.9
	CBT	229,428	112,680	49.1	116,748	50.9
	PBT	619	258	41.7	361	58.3
4	All	232,949	114,345	49.1	118,604	50.9
	CBT	232,371	114,112	49.1	118,259	50.9
	PBT	578	233	40.3	345	59.7
5	All	238,527	116,876	49.0	121,651	51.0
	CBT	237,976	116,618	49.0	121,358	51.0
	PBT	551	258	46.8	293	53.2
6	All	238,585	116,445	48.8	122,140	51.2
	CBT	238,105	116,249	48.8	121,856	51.2
	PBT	480	196	40.8	284	59.2
7	All	235,638	115,051	48.8	120,587	51.2
	CBT	235,220	114,894	48.8	120,326	51.1
	PBT	418	157	37.6	261	62.4
8	All	211,502	102,954	48.7	108,548	51.3
	CBT	211,099	102,783	48.7	108,316	51.3
	PBT	403	171	42.3	232	57.4
A1	All	115,071	56,085	48.7	58,986	51.2
	CBT	115,004	56,053	48.7	58,951	51.2
	PBT	67	32	47.1	35	51.5
GO	All	40,365	20,199	50.0	20,166	49.9
	CBT	40,351	20,195	50.0	20,156	49.9
	PBT	n/r	n/r	n/r	n/r	n/r
A2	All	13,949	6,903	49.4	7,046	50.5
	CBT	13,947	6,902	49.4	7,045	50.5
	PBT	n/r	n/r	n/r	n/r	n/r

Note. DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; A1 = Algebra I; GO = Geometry; A2 = Algebra II; n/r = not reported due to n<20.

Table A.11.6 All States Combined: Spanish-Language Mathematics Test Takers, by Grade, Mode, and Gender

Grade	Mode	Valid Cases	Female		Male	
			N	%	N	%
3	All	5,063	2,562	50.6	2,501	49.4
	CBT	5,048	2,557	50.7	2,491	49.3
	PBT	n/r	n/r	n/r	n/r	n/r
4	All	4,602	2,303	50.0	2,299	49.9
	CBT	4,596	2,299	50.0	2,297	49.9
	PBT	n/r	n/r	n/r	n/r	n/r
5	All	4,219	2,113	50.1	2,106	49.9
	CBT	4,214	2,111	50.1	2,103	49.9
	PBT	n/r	n/r	n/r	n/r	n/r
6	All	3,721	1,802	48.4	1,919	51.6
	CBT	3,718	1,801	48.4	1,917	51.6
	PBT	n/r	n/r	n/r	n/r	n/r
7	All	3,226	1,530	47.4	1,696	52.6
	CBT	3,224	1,530	47.5	1,694	52.5
	PBT	n/r	n/a	n/a	n/r	n/r
8	All	2,844	1,406	49.4	1,438	50.6
	CBT	2,840	1,405	49.5	1,435	50.5
	PBT	n/r	n/r	n/r	n/r	n/r
A1	All	2,562	1,119	43.7	1,443	56.3
	CBT	2,561	1,118	43.6	1,443	56.3
	PBT	n/r	n/r	n/r	n/a	n/a
GO	All	425	181	42.6	244	57.4
	CBT	425	181	42.6	244	57.4
	PBT	n/a	n/a	n/a	n/a	n/a
A2	All	138	64	46.4	74	53.6
	CBT	138	64	46.4	74	53.6
	PBT	n/a	n/a	n/a	n/a	n/a

Note. DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; A1 = Algebra I; GO = Geometry; A2 = Algebra II; CBT = computer-based test; PBT = paper-based test; n/a = not applicable; n/r = not reported due to n<20.

*No students in DD tested in mathematics using Spanish-language forms.

Table A.11.7 Demographic Information: Grade 3 ELA/L, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.9	48.5	n/r	48.8	34.7
SWD (%)	17.0	18.6	17.0	15.5	18.9
EL (%)	15.4	16.8	12.6	18.9	10.6
Male (%)	50.9	52.0	50.4	51.0	50.7
Female (%)	49.1	48.0	49.6	49.0	49.3
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.2
Asian (%)	7.7	1.8	5.3	5.7	10.9
Black/AA (%)	16.7	62.7	9.5	16.4	14.4
Hisp/Lat (%)	28.4	17.4	22.1	26.4	32.3
Wh/Caus (%)	42.4	14.5	44.5	46.4	38.6
NtvHawaii/Pacific (%)	0.2	n/r	1.3	0.1	0.2
Two or More (%)	4.2	n/a	n/r	4.5	3.4
Unknown (%)	0.3	3.5	2.4	0.3	0.0

Note. All States = data from all participating states combined; DC = District of Columbia, DD = Department of Defense Education Activity, IL = Illinois; NJ = New Jersey; Econ Dis = economically disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.8 Demographic Information: Grade 4 ELA/L, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.8	49.7	n/r	48.5	34.8
SWD (%)	18.4	21.0	17.3	16.8	20.5
EL (%)	15.1	15.9	10.3	19.0	9.9
Male (%)	50.9	51.5	51.0	50.9	50.9
Female (%)	49.1	48.5	49.0	49.1	49.1
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.2
Asian (%)	7.8	1.7	5.4	5.7	11.1
Black/AA (%)	16.8	63.3	9.7	16.2	15.0
Hisp/Lat (%)	28.5	17.1	22.8	27.0	31.7
Wh/Caus (%)	42.2	13.9	42.8	46.2	38.7
NtvHawaii/Pacific (%)	0.2	n/r	1.5	0.1	0.2
Two or More (%)	4.0	n/r	15.4	4.3	3.2
Unknown (%)	0.3	3.8	2.1	0.2	0.0

Note. All states = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = economically disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/r = not reported due to n<20.

Table A.11.9 Demographic Information: Grade 5 ELA/L, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.4	48.2	n/r	48.4	33.7
SWD (%)	18.8	21.9	17.4	17.1	21.1
EL (%)	12.8	15.7	9.4	16.6	7.5
Male (%)	51.0	51.3	50.8	51.0	51.0
Female (%)	49.0	48.7	49.2	49.0	48.9
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.2
Asian (%)	7.6	1.4	5.0	5.6	10.9
Black/AA (%)	16.7	66.0	9.9	16.3	14.6
Hisp/Lat (%)	28.8	17.3	23.6	27.4	31.7
Wh/Caus (%)	42.6	12.5	42.4	46.1	39.5
NtvHawaii/Pacific (%)	0.2	n/r	1.6	0.1	0.2
Two or More (%)	3.8	n/r	14.8	4.1	2.9
Unknown (%)	0.3	2.6	2.4	0.3	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = economically disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/r = not reported due to n<20.

Table A.11.10 Demographic Information: Grade 6 ELA/L, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.3	48.2	n/r	48.4	33.1
SWD (%)	18.7	21.7	15.8	17.2	20.9
EL (%)	10.2	13.8	8.5	13.4	5.8
Male (%)	51.2	49.7	49.6	51.4	51.1
Female (%)	48.8	50.3	50.4	48.6	48.9
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.2
Asian (%)	7.6	1.0	5.8	5.5	10.9
Black/AA (%)	16.9	67.3	9.3	16.5	14.8
Hisp/Lat (%)	28.6	17.4	23.5	27.4	31.1
Wh/Caus (%)	42.7	11.0	41.7	46.1	39.9
NtvHawaii/Pacific (%)	0.2	n/r	1.3	0.1	0.2
Two or More (%)	3.7	n/a	n/r	4.0	2.9
Unknown (%)	0.3	3.0	2.1	0.2	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = economically disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.11 Demographic Information: Grade 7 ELA/L, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.0	47.7	n/r	48.2	32.5
SWD (%)	18.6	23.0	15.7	17.0	20.7
EL (%)	9.3	12.0	8.6	12.4	4.9
Male (%)	51.2	50.7	51.4	51.1	51.4
Female (%)	48.8	49.3	48.6	48.9	48.5
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.1
Asian (%)	7.4	1.3	5.7	5.4	10.6
Black/AA (%)	16.7	66.3	9.8	16.1	15.1
Hisp/Lat (%)	29.1	18.9	24.7	28.2	31.0
Wh/Caus (%)	42.8	10.6	40.4	45.8	40.4
NtvHawaii/Pacific (%)	0.2	n/r	2.0	0.1	0.2
Two or More (%)	3.5	n/r	14.7	3.9	2.6
Unknown (%)	0.2	2.8	2.4	0.2	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = economically disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/r = not reported due to n<20.

Table A.11.12 Demographic Information: Grade 8 ELA/L, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	40.7	47.7	n/r	48.2	31.4
SWD (%)	18.3	23.7	14.1	16.7	20.5
EL (%)	8.9	10.8	7.2	12.1	4.5
Male (%)	51.0	49.7	50.5	50.9	51.2
Female (%)	49.0	50.3	49.5	49.1	48.7
AmInd/ANat (%)	0.2	n/r	n/r	0.3	0.1
Asian (%)	7.4	1.2	6.4	5.3	10.7
Black/AA (%)	16.8	67.1	9.8	16.5	15.0
Hisp/Lat (%)	29.0	19.8	24.0	28.3	30.6
Wh/Caus (%)	42.8	9.3	41.0	45.6	40.8
NtvHawaii/Pacific (%)	0.2	n/r	1.5	0.1	0.2
Two or More (%)	3.3	n/a	n/r	3.8	2.4
Unknown (%)	0.2	2.3	2.1	0.2	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = economically disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.13 Demographic Information: Grade 9 ELA/L, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	30.5	50.3	n/a	n/a	29.6
SWD (%)	20.1	23.2	n/a	n/a	19.9
EL (%)	5.0	9.0	n/a	n/a	4.8
Male (%)	51.3	50.3	n/a	n/a	51.3
Female (%)	48.6	49.6	n/a	n/a	48.6
AmInd/ANat (%)	0.1	n/r	n/a	n/a	0.1
Asian (%)	10.0	1.4	n/a	n/a	10.4
Black/AA (%)	17.2	68.7	n/a	n/a	14.7
Hisp/Lat (%)	31.2	18.5	n/a	n/a	31.8
Wh/Caus (%)	38.9	9.1	n/a	n/a	40.4
NtvHawaii/Pacific (%)	0.2	n/r	n/a	n/a	0.2
Two or More (%)	2.2	n/r	n/a	n/a	2.3
Unknown (%)	0.1	2.2	n/a	n/a	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = economically disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.14 Demographic Information: Grade 10 ELA/L, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	30.6	52.9	n/r	n/a	n/a
SWD (%)	18.7	22.9	13.0	n/a	n/a
EL (%)	6.6	7.3	5.7	n/a	n/a
Male (%)	51.2	50.4	52.3	n/a	n/a
Female (%)	48.7	49.5	47.7	n/a	n/a
AmInd/ANat (%)	n/r	n/r	n/r	n/a	n/a
Asian (%)	4.0	1.4	7.6	n/a	n/a
Black/AA (%)	42.7	66.4	10.1	n/a	n/a
Hisp/Lat (%)	20.8	20.7	21.0	n/a	n/a
Wh/Caus (%)	22.5	9.0	41.0	n/a	n/a
NtvHawaii/Pacific (%)	0.8	n/r	1.7	n/a	n/a
Two or More (%)	6.7	n/a	n/r	n/a	n/a
Unknown (%)	2.3	2.4	2.2	n/a	n/a

Note. All States = data from all participating states combined; DC = District of Columbia, DD = Department of Defense Education Activity, IL = Illinois, NJ = New Jersey; Econ Dis = economically disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.15 Demographic Information: Grade 3 Mathematics, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.9	48.0	n/r	48.8	34.8
SWD (%)	16.9	18.4	17.1	15.5	18.7
EL (%)	15.9	17.6	12.7	18.9	11.9
Male (%)	50.9	52.0	50.5	51.0	50.7
Female (%)	49.1	48.0	49.5	49.0	49.3
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.2
Asian (%)	7.7	1.9	5.4	5.7	10.8
Black/AA (%)	16.5	62.0	9.4	16.3	14.3
Hisp/Lat (%)	28.7	17.8	22.1	26.4	33.0
Wh/Caus (%)	42.2	14.7	44.6	46.5	38.2
NtvHawaii/Pacific (%)	0.2	n/r	1.3	0.1	0.2
Two or More (%)	4.1	n/a	n/r	4.5	3.4
Unknown (%)	0.3	3.5	2.3	0.3	0.0

Note. All States = data from all participating states combined; DC = District of Columbia, DD = Department of Defense Education Activity, IL = Illinois; NJ = New Jersey; Econ Dis = Economically Disadvantaged; SWD = student with disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.16 Demographic Information: Grade 4 Mathematics, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.7	49.4	n/r	48.4	34.8
SWD (%)	18.3	20.8	17.2	16.8	20.2
EL (%)	15.6	16.6	10.4	19.0	11.3
Male (%)	50.9	51.4	51.0	50.9	50.9
Female (%)	49.1	48.6	49.0	49.1	49.1
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.2
Asian (%)	7.8	1.7	5.4	5.7	11.1
Black/AA (%)	16.7	62.8	9.8	16.1	14.8
Hisp/Lat (%)	28.8	17.5	22.6	27.0	32.3
Wh/Caus (%)	42.1	14.0	42.7	46.3	38.3
NtvHawaii/Pacific (%)	0.2	n/r	1.5	0.1	0.2
Two or More (%)	4.0	n/r	15.6	4.3	3.1
Unknown (%)	0.3	3.8	2.1	0.2	0.0

Note. All States = data from all participating states combined; DC = District of Columbia, DD = Department of Defense Education Activity, IL = Illinois; NJ = New Jersey; Econ Dis = Economically Disadvantaged; SWD = Student with Disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/r = not reported due to n<20.

Table A.11.17 Demographic Information: Grade 5 Mathematics, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.3	47.6	n/r	48.3	33.8
SWD (%)	18.7	21.7	17.4	17.1	20.8
EL (%)	13.3	16.5	9.5	16.6	8.9
Male (%)	51.0	51.2	50.8	51.0	51.0
Female (%)	49.0	48.8	49.2	49.0	49.0
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.2
Asian (%)	7.6	1.5	5.0	5.6	10.8
Black/AA (%)	16.6	65.2	10.0	16.2	14.4
Hisp/Lat (%)	29.0	17.9	23.7	27.3	32.4
Wh/Caus (%)	42.4	12.6	42.3	46.2	39.1
NtvHawaii/Pacific (%)	0.2	n/r	1.6	0.1	0.2
Two or More (%)	3.7	n/r	14.8	4.1	2.9
Unknown (%)	0.3	2.6	2.3	0.3	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey. Econ Dis = Economically Disadvantaged; SWD = Student with Disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/r = not reported due to n<20.

Table A.11.18 Demographic Information: Grade 6 Mathematics, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	41.1	47.6	n/r	48.3	33.1
SWD (%)	18.6	21.4	16.0	17.1	20.6
EL (%)	10.8	14.9	8.7	13.4	7.1
Male (%)	51.2	49.8	49.7	51.4	51.1
Female (%)	48.8	50.2	50.3	48.6	48.9
AmInd/ANat (%)	0.2	n/r	n/r	0.2	0.2
Asian (%)	7.6	1.2	5.9	5.5	10.8
Black/AA (%)	16.7	66.4	9.4	16.4	14.6
Hisp/Lat (%)	28.9	18.1	23.4	27.4	31.8
Wh/Caus (%)	42.6	11.1	41.6	46.2	39.6
NtvHawaii/Pacific (%)	0.2	n/r	1.3	0.1	0.2
Two or More (%)	3.6	n/a	n/r	4.0	2.8
Unknown (%)	0.2	3.0	2.1	0.2	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey. Econ Dis = Economically Disadvantaged; SWD = Student with Disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/r = not reported due to n<20.

Table A.11.19 Demographic Information: Grade 7 Mathematics, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	42.3	47.8	n/a	48.1	33.8
SWD (%)	18.9	23.1	n/a	17.0	21.3
EL (%)	10.1	13.0	n/a	12.4	6.5
Male (%)	51.2	50.5	n/a	51.1	51.3
Female (%)	48.8	49.5	n/a	48.9	48.7
AmInd/ANat (%)	0.2	n/r	n/a	0.2	0.1
Asian (%)	6.6	1.3	n/a	5.4	8.6
Black/AA (%)	16.9	66.1	n/a	16.1	15.4
Hisp/Lat (%)	29.8	19.6	n/a	28.2	32.7
Wh/Caus (%)	42.9	10.2	n/a	45.9	40.4
NtvHawaii/Pacific (%)	0.1	n/r	n/a	0.1	0.2
Two or More (%)	3.3	n/r	n/a	3.9	2.5
Unknown (%)	0.2	2.6	n/a	0.2	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = Economically Disadvantaged; SWD = Student with Disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.20 Demographic Information: Grade 8 Mathematics, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	44.3	52.7	n/r	48.1	37.5
SWD (%)	19.9	27.2	17.9	16.7	26.4
EL (%)	10.7	13.3	9.2	12.1	7.6
Male (%)	51.3	51.4	52.1	50.9	52.1
Female (%)	48.7	48.6	47.9	49.1	47.8
AmInd/ANat (%)	0.3	n/r	n/r	0.3	0.2
Asian (%)	5.1	0.9	5.3	5.3	5.1
Black/AA (%)	17.9	72.6	10.8	16.4	18.0
Hisp/Lat (%)	30.7	19.7	26.8	28.3	36.7
Wh/Caus (%)	42.3	4.9	37.9	45.7	37.7
NtvHawaii/Pacific (%)	0.1	n/r	1.6	0.1	0.2
Two or More (%)	3.4	n/a	n/r	3.8	2.2
Unknown (%)	0.2	1.7	2.0	0.2	n/r

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = Economically Disadvantaged; SWD = Student with Disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.21 Demographic Information: Algebra I, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	30.4	50.7	n/r	n/a	30.5
SWD (%)	19.0	21.7	14.0	n/a	19.0
EL (%)	6.9	12.7	6.9	n/a	6.6
Male (%)	51.2	48.9	50.3	n/a	51.3
Female (%)	48.7	51.0	49.7	n/a	48.6
AmInd/ANat (%)	0.1	n/r	n/r	n/a	0.1
Asian (%)	9.9	1.3	7.3	n/a	10.4
Black/AA (%)	16.8	66.3	9.4	n/a	14.8
Hisp/Lat (%)	32.2	22.5	23.3	n/a	33.0
Wh/Caus (%)	37.9	7.5	40.3	n/a	39.1
NtvHawaii/Pacific (%)	0.2	n/r	1.7	n/a	0.2
Two or More (%)	2.6	n/a	n/r	n/a	2.3
Unknown (%)	0.2	2.3	2.4	n/a	0.0

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = Economically Disadvantaged; SWD = Student with Disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.22 Demographic Information: Geometry, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	20.3	51.5	n/r	n/a	18.4
SWD (%)	11.5	22.0	12.7	n/a	10.1
EL (%)	3.5	7.3	5.4	n/a	2.8
Male (%)	49.9	49.9	51.5	n/a	49.8
Female (%)	50.0	49.9	48.5	n/a	50.1
AmInd/ANat (%)	0.1	n/r	n/r	n/a	0.1
Asian (%)	16.4	1.5	7.4	n/a	19.0
Black/AA (%)	15.0	66.7	9.9	n/a	9.2
Hisp/Lat (%)	21.4	19.3	21.2	n/a	21.7
Wh/Caus (%)	42.8	10.2	40.2	n/a	47.0
NtvHawaii/Pacific (%)	0.3	n/r	1.7	n/a	0.2
Two or More (%)	3.5	n/r	16.8	n/a	2.7
Unknown (%)	0.4	2.0	2.4	n/a	n/r

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey. Econ Dis = Economically Disadvantaged; SWD = Student with Disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Table A.11.23 Demographic Information: Algebra II, Overall and by State

Demographic	All States	DC	DD	IL	NJ
Econ Dis (%)	15.1	13.3	n/r	n/a	18.8
SWD (%)	10.3	15.5	12.4	n/a	9.6
EL (%)	3.0	n/r	5.0	n/a	2.6
Male (%)	50.5	41.6	50.4	n/a	50.6
Female (%)	49.4	58.4	49.6	n/a	49.2
AmInd/ANat (%)	0.2	n/a	n/r	n/a	0.2
Asian (%)	25.3	n/r	7.6	n/a	29.9
Black/AA (%)	7.5	27.9	9.7	n/a	6.6
Hisp/Lat (%)	21.5	9.3	21.6	n/a	21.7
Wh/Caus (%)	39.5	46.0	42.0	n/a	38.8
NtvHawaii/Pacific (%)	0.6	n/a	n/r	n/a	0.3
Two or More (%)	4.8	n/a	n/r	n/a	2.5
Unknown (%)	0.5	8.8	1.9	n/a	n/r

Note. All States = data from all participating states combined; DC = District of Columbia; DD = Department of Defense Education Activity; IL = Illinois; NJ = New Jersey; Econ Dis = Economically Disadvantaged; SWD = Student with Disabilities; EL = English learner; AmInd/ANat = American Indian/Alaska Native; Black/AA = Black/African American; Hisp/Lat = Hispanic/Latino; Wh/Caus = White/Caucasian; NtvHawaii/Pacific = Native Hawaiian or Other Pacific Islander; Two or More = two or more races reported; n/a = not applicable; n/r = not reported due to n<20.

Appendix 12.1: Form Composition

Table A.12.1 Form Composition for ELA/L Grade 3

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	4–7	8–17
	Reading Informational Text	4–7	11–20
	Vocabulary	4–5	8–10
	Claim Total	12–14	30–31
Writing	Written Expression	1	18
	Knowledge of Conventions	1	6
	Claim Total	2	24
SUMMATIVE TOTAL		14–16	54–55

Note. This table is identical to Table 12.1 in Section 12.

Table A.12.2 Form Composition for ELA/L Grade 4

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	5–8	14–20
	Reading Informational Text	5–9	18–22
	Vocabulary	4–7	8–14
	Claim Total	18	40–44
Writing	Written Expression	1	21–24
	Knowledge of Conventions	1	6
	Claim Total	2	27–30
SUMMATIVE TOTAL		20	67–74

Table A.12.3 Form Composition for ELA/L Grade 5

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	5–8	14–20
	Reading Informational Text	5–9	14–22
	Vocabulary	4–7	8–14
	Claim Total	18	40–44
Writing	Written Expression	1	21–24
	Knowledge of Conventions	1	6
	Claim Total	2	27–30
SUMMATIVE TOTAL		20	67–74

Table A.12.4 Form Composition for ELA/L Grade 6

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	5–9	14–22
	Reading Informational Text	5–11	14–26
	Vocabulary	4–7	8–14
	Claim Total	18	40–44
Writing	Written Expression	1	24
	Knowledge of Conventions	1	6
	Claim Total	2	30
SUMMATIVE TOTAL		20	70–74

Table A.12.5 Form Composition for ELA/L Grade 7

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	5–9	14–22
	Reading Informational Text	5–11	14–26
	Vocabulary	4–7	8–14
	Claim Total	18	40–44
Writing	Written Expression	1	24
	Knowledge of Conventions	1	6
	Claim Total	2	30
SUMMATIVE TOTAL		20	70–74

Table A.12.6 Form Composition for ELA/L Grade 8

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	5–9	14–22
	Reading Informational Text	5–11	14–26
	Vocabulary	4–7	8–14
	Claim Total	18	40–44
Writing	Written Expression	1	24
	Knowledge of Conventions	1	6
	Claim Total	2	30
SUMMATIVE TOTAL		20	70–74

Table A.12.7 Form Composition for ELA/L Grade 9

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	5–9	14–22
	Reading Informational Text	5–11	14–26
	Vocabulary	4–7	8–14
	Claim Total	18	40–44
Writing	Written Expression	1	24
	Knowledge of Conventions	1	6
	Claim Total	2	30
SUMMATIVE TOTAL		20	70–74

Table A.12.8 Form Composition for ELA/L Grade 10

Claims	Subclaims	Number of Items	Number of Points
Reading	Reading Literary Text	5–9	14–22
	Reading Informational Text	5–11	14–26
	Vocabulary	4–7	8–14
	Claim Total	18	40–44
Writing	Written Expression	1	24
	Knowledge of Conventions	1	6
	Claim Total	2	30
SUMMATIVE TOTAL		20	70–74

Table A.12.9 Form Composition for Mathematics Grade 3

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	18	20
	Additional & Supporting Content	9	10
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	12
TOTAL		33	52

Note. This table is identical to Table 12.3 in Section 12.

Table A.12.10 Form Composition for Mathematics Grade 4

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	17	21
	Additional & Supporting Content	8	9
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	12
TOTAL		31	52

Table A.12.11 Form Composition for Mathematics Grade 5

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	17	20
	Additional & Supporting Content	8	10
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	12
TOTAL		31	52

Table A.12.12 Form Composition for Mathematics Grade 6

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	15	20
	Additional & Supporting Content	8	10
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	12
TOTAL		29	52

Table A.12.13 Form Composition for Mathematics Grade 7

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	18	20
	Additional & Supporting Content	7	10
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	12
TOTAL		31	52

Table A.12.14 Form Composition for Mathematics Grade 8

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	18	20
	Additional & Supporting Content	6	10
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	12
TOTAL		30	52

Table A.12.15 Form Composition for Algebra I

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	12	17
	Additional & Supporting Content	8–9	9–11
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	15
	Integrated (Ψ^*)	1–2	2–4
TOTAL		28	55

Table A.12.16 Form Composition for Geometry

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	15	18
	Additional & Supporting Content	9	12
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	15
TOTAL		30	55

Table A.12.17 Form Composition for Algebra II

	Subclaims	Number of Items	Number of Points
Mathematics			
	Major Content	13–14	16–18
	Additional & Supporting Content	9	12
	Expressing Mathematical Reasoning	3	10
	Modeling and Applications	3	15
	Integrated (Ψ^*)	0–2	0–2
TOTAL		29	55

Appendix 12.2: Threshold Scores and Scaling Constants

Table A.12.18 Threshold Scores and Scaling Constants for ELA/L Grades 3 to 8

Assessment	Threshold Cut	Theta	Scale Score	A	B
Grade 3 ELA	Level 2 Cut	-0.9648	700	36.7227	735.4297
	Level 3 Cut	-0.2840	725		
	Level 4 Cut	0.3968	750		
	Level 5 Cut	2.0360	810		
Grade 4 ELA	Level 2 Cut	-1.3004	700	31.5462	741.0214
	Level 3 Cut	-0.5079	725		
	Level 4 Cut	0.2846	750		
	Level 5 Cut	1.5578	790		
Grade 5 ELA	Level 2 Cut	-1.3411	700	29.4580	739.5050
	Level 3 Cut	-0.4924	725		
	Level 4 Cut	0.3563	750		
	Level 5 Cut	2.0224	799		
Grade 6 ELA	Level 2 Cut	-1.3656	700	28.3160	738.6673
	Level 3 Cut	-0.4827	725		
	Level 4 Cut	0.4002	750		
	Level 5 Cut	1.8133	790		
Grade 7 ELA	Level 2 Cut	-1.2488	700	33.9161	742.3542
	Level 3 Cut	-0.5117	725		
	Level 4 Cut	0.2254	750		
	Level 5 Cut	1.2614	785		
Grade 8 ELA	Level 2 Cut	-1.2730	700	34.1183	743.4330
	Level 3 Cut	-0.5402	725		
	Level 4 Cut	0.1925	750		
	Level 5 Cut	1.4696	794		

Table A.12.19 Threshold Scores and Scaling Constants for Mathematics Grades 3 to 8

Assessment	Threshold Cut	Theta	Scale Score	A	B
Grade 3 Mathematics	Level 2 Cut	-1.4141	700	32.1135	745.4119
	Level 3 Cut	-0.6356	725		
	Level 4 Cut	0.1429	750		
	Level 5 Cut	1.3931	790		
Grade 4 Mathematics	Level 2 Cut	-1.3840	700	29.9167	741.4049
	Level 3 Cut	-0.5484	725		
	Level 4 Cut	0.2873	750		
	Level 5 Cut	1.8323	796		
Grade 5 Mathematics	Level 2 Cut	-1.4571	700	29.0301	742.2997
	Level 3 Cut	-0.5959	725		
	Level 4 Cut	0.2653	750		
	Level 5 Cut	1.6262	790		
Grade 6 Mathematics	Level 2 Cut	-1.3829	700	28.1465	738.9252
	Level 3 Cut	-0.4948	725		
	Level 4 Cut	0.3935	750		
	Level 5 Cut	1.7567	788		
Grade 7 Mathematics	Level 2 Cut	-1.4464	700	25.1033	736.3102
	Level 3 Cut	-0.4505	725		
	Level 4 Cut	0.5453	750		
	Level 5 Cut	1.9919	786		
Grade 8 Mathematics	Level 2 Cut	-0.8851	700	32.9505	729.1640
	Level 3 Cut	-0.1264	725		
	Level 4 Cut	0.6323	750		
	Level 5 Cut	2.1896	801		

Table A.12.20 Threshold Scores and Scaling Constants for High School ELA/L

Assessment	Threshold Cut	Theta	Scale Score	A	B
Grade 10 ELA/L	Level 2 Cut	-0.8909	700	43.1280	738.4223
	Level 3 Cut	-0.3112	725		
	Level 4 Cut	0.2684	750		
	Level 5 Cut	1.2858	794		
Grade 11 ELA/L	Level 2 Cut	-1.1017	700	34.9278	738.4801
	Level 3 Cut	-0.3859	725		
	Level 4 Cut	0.3298	750		
	Level 5 Cut	1.5206	792		

Table A.12.21 Threshold Scores and Scaling Constants for High School Mathematics

Assessment	Threshold Cut	Theta	Scale Score	A	B
Algebra I	Level 2 Cut	-1.1781	700	31.5325	737.1490
	Level 3 Cut	-0.3853	725		
	Level 4 Cut	0.4075	750		
	Level 5 Cut	2.1651	805		
Algebra II	Level 2 Cut	-0.5759	700	37.7676	721.7509
	Level 3 Cut	0.0860	725		
	Level 4 Cut	0.7480	750		
	Level 5 Cut	2.2728	808		
Geometry	Level 2 Cut	-1.3013	700	25.9775	733.8039
	Level 3 Cut	-0.3389	725		
	Level 4 Cut	0.6235	750		
	Level 5 Cut	1.8940	783		
Integrated Mathematics I	Level 2 Cut	-1.0919	700	32.0043	734.9446
	Level 3 Cut	-0.3107	725		
	Level 4 Cut	0.4704	750		
	Level 5 Cut	1.9934	799		
Integrated Mathematics II	Level 2 Cut	-0.9175	700	29.2865	726.8695
	Level 3 Cut	-0.0638	725		
	Level 4 Cut	0.7898	750		
	Level 5 Cut	1.9817	785		

Table A.12.22 Scaling Constants for Reading and Writing Grades 3 to 11

	Reading		Writing	
	AR	BR	AW	BW
Grade 3 ELA/L	14.6891	44.1719	7.3445	32.0859
Grade 4 ELA/L	12.6184	46.4086	6.3093	33.2043
Grade 5 ELA/L	11.7832	45.8019	5.8916	32.9010
Grade 6 ELA/L	11.3264	45.4669	5.6632	32.7335
Grade 7 ELA/L	13.5664	46.9416	6.7832	33.4708
Grade 8 ELA/L	13.6472	47.3732	6.8237	33.6866
Grade 10 ELA/L	17.2512	45.3690	8.6256	32.6845
Grade 11 ELA/L	13.9712	45.3920	6.9856	32.6961

Appendix 12.3: IRT Test Characteristic Curves, CSEM Curves, and Information Curves

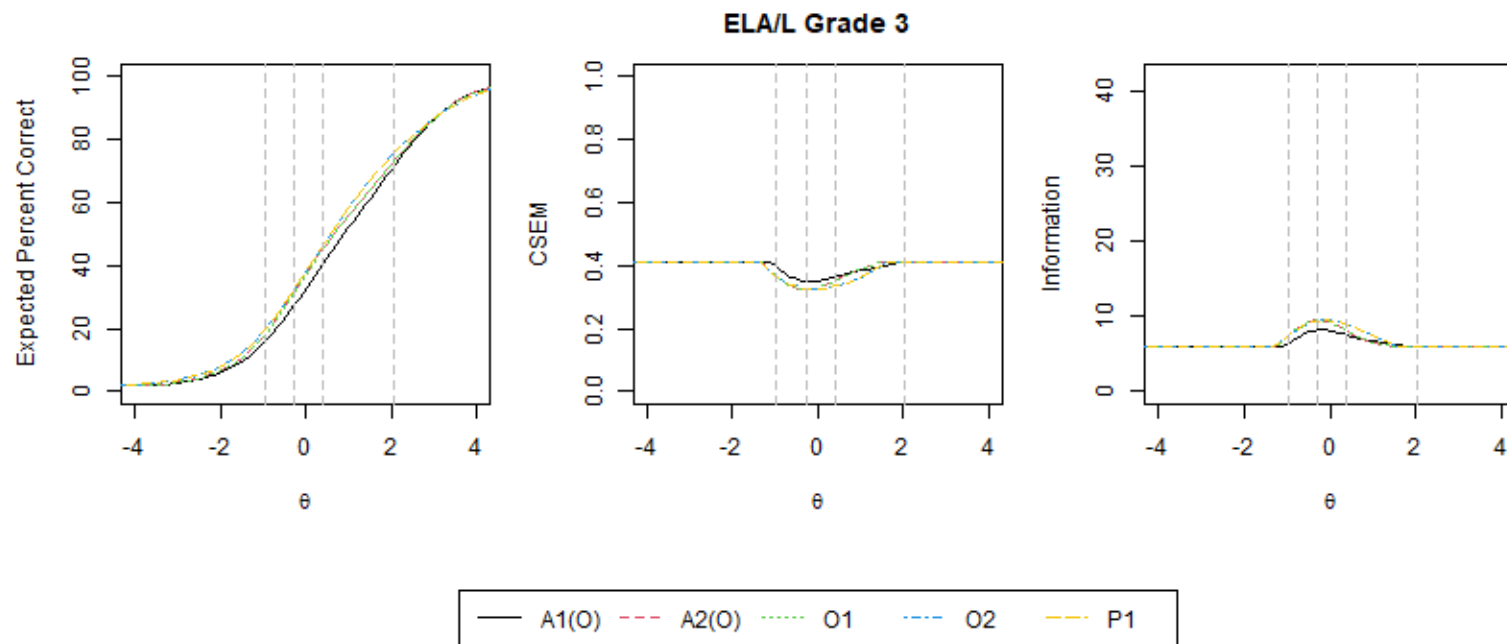


Figure A.12.1 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 3

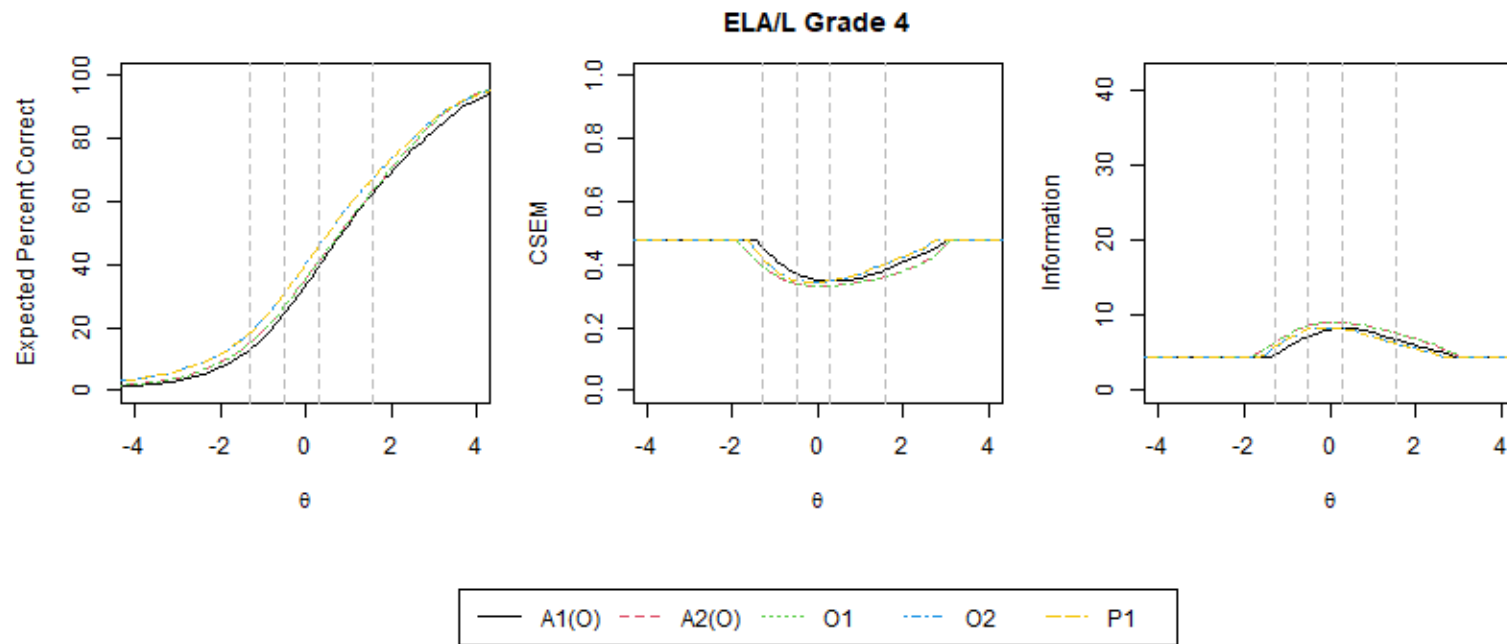


Figure A.12.2 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 4

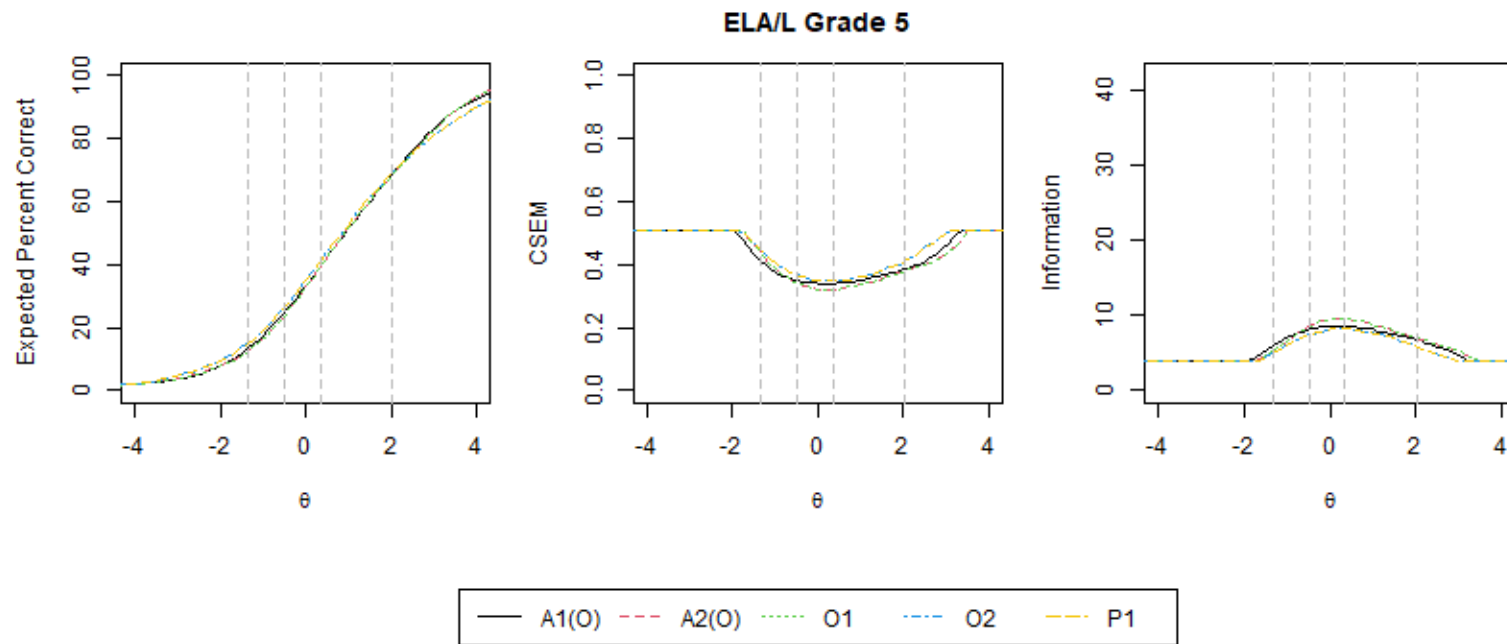


Figure A.12.3 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 5

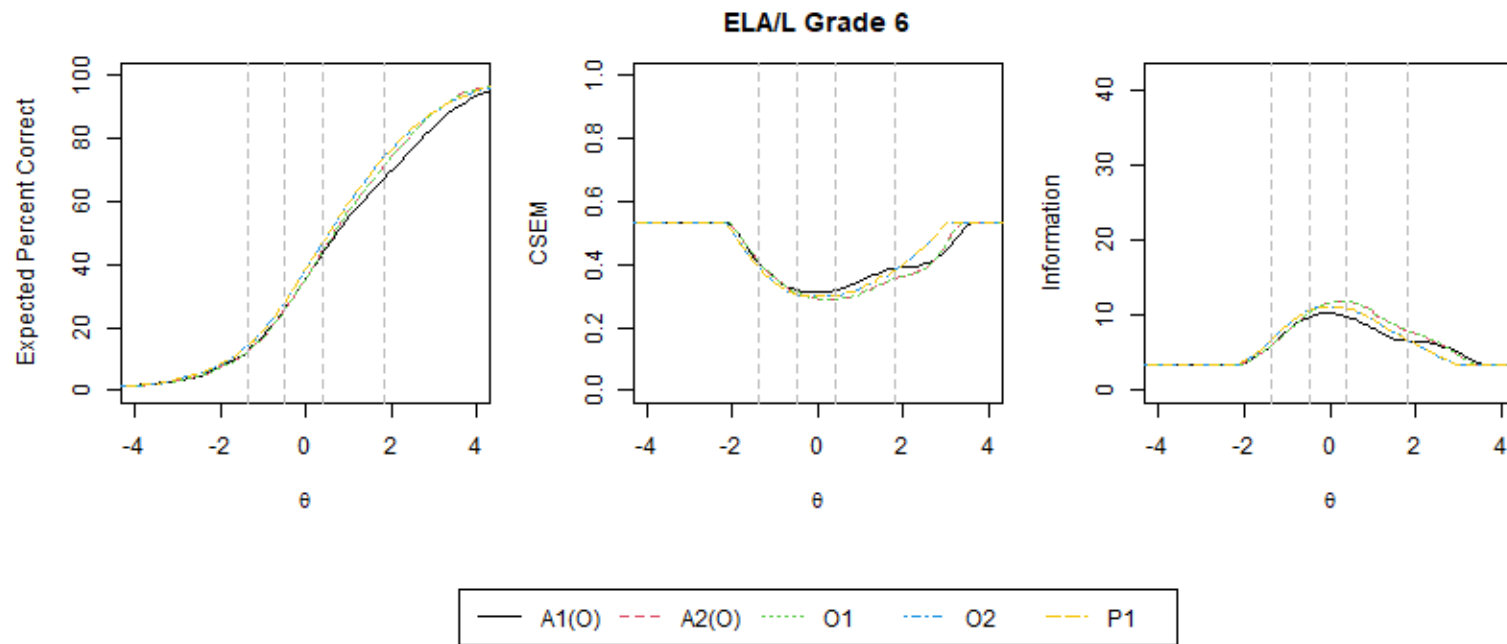


Figure A.12.4 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 6

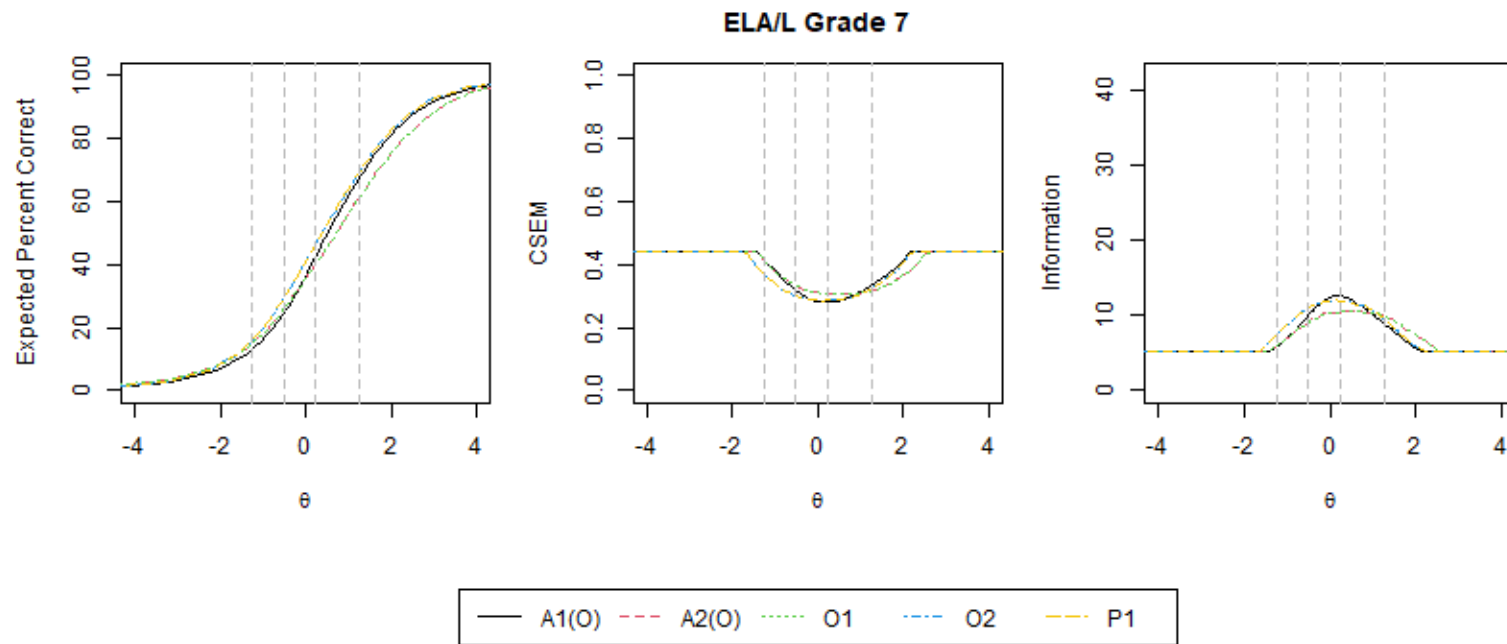


Figure A.12.5 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 7

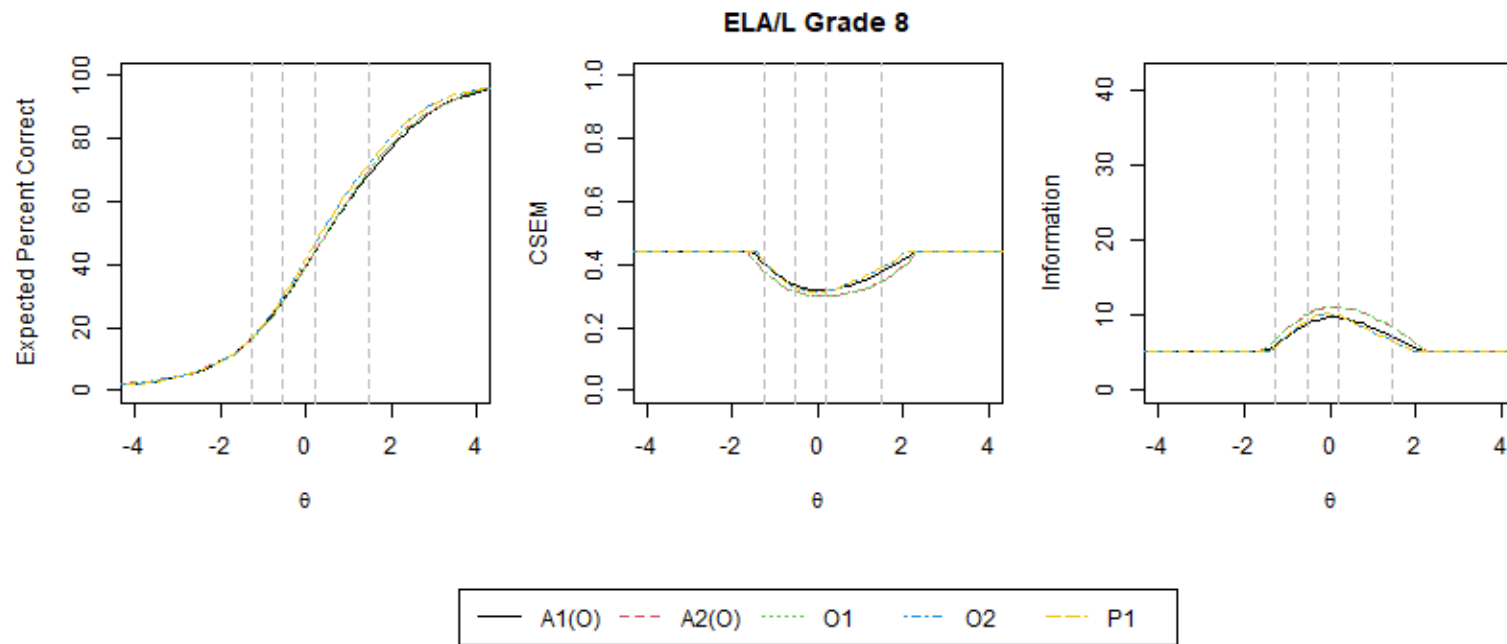


Figure A.12.6 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 8

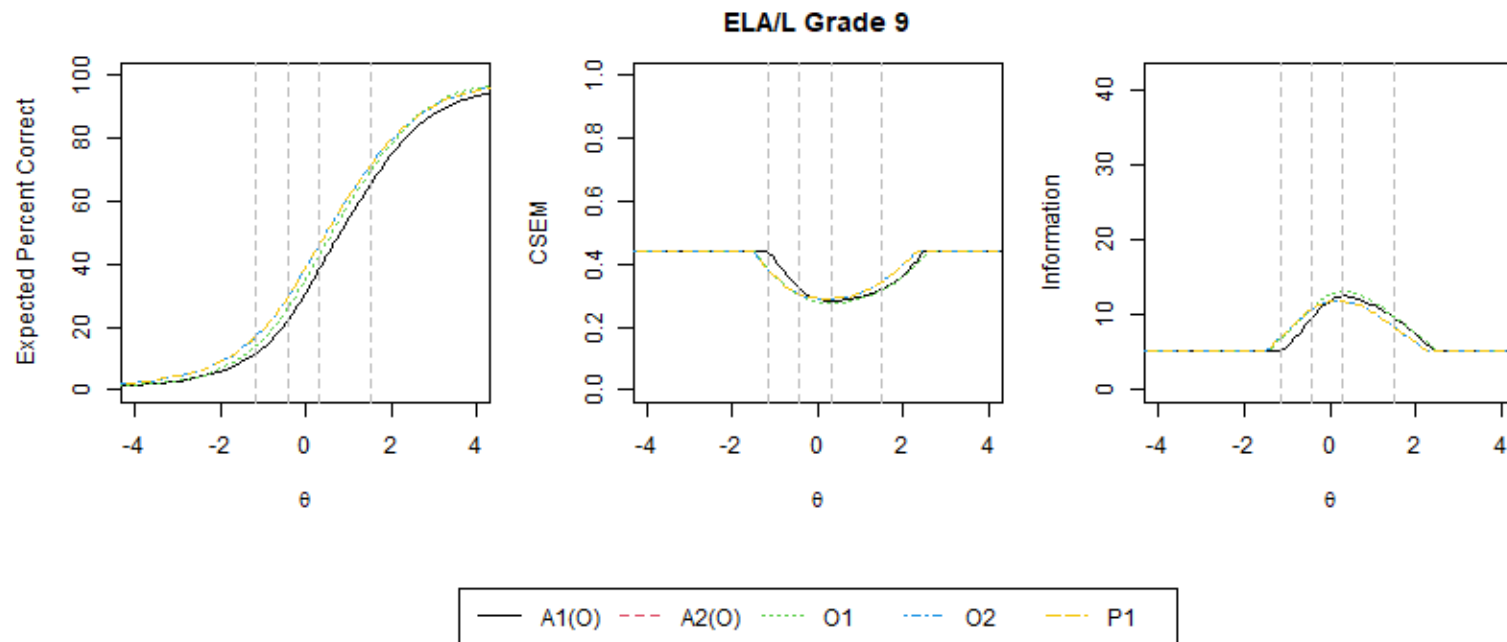


Figure A.12.7 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 9

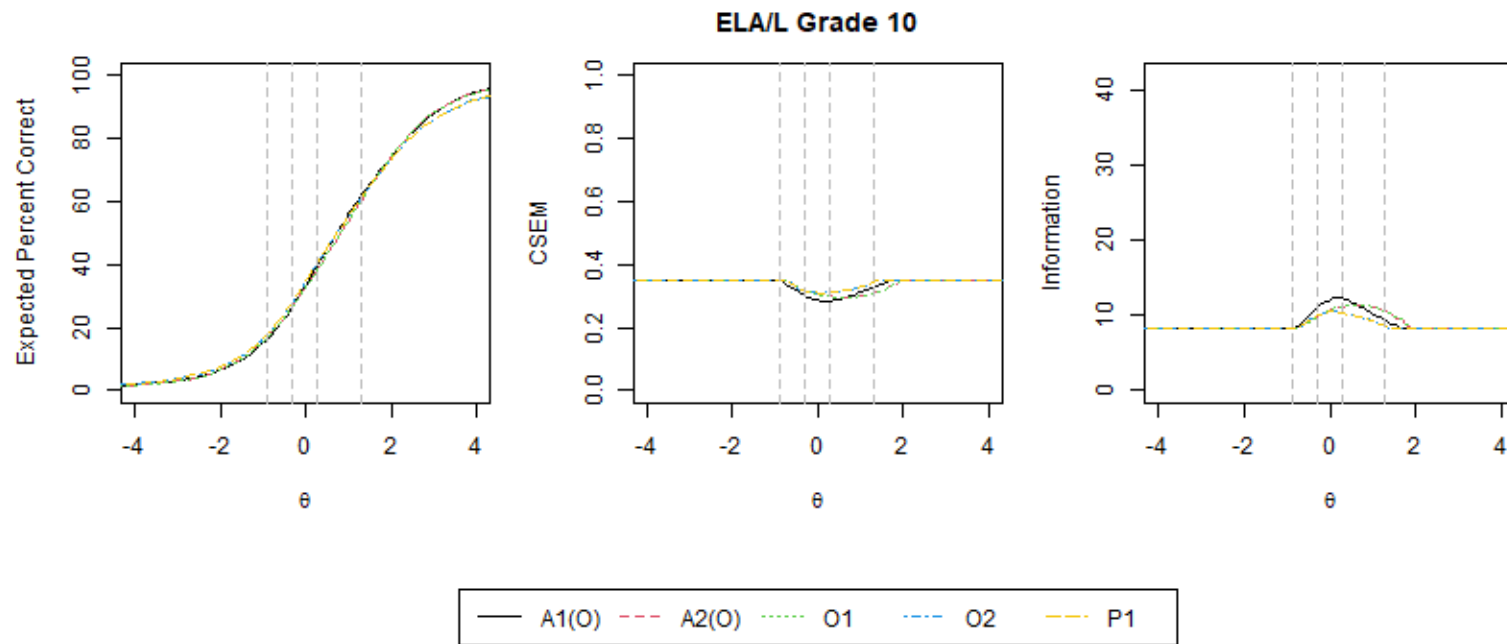


Figure A.12.8 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves ELA/L Grade 10

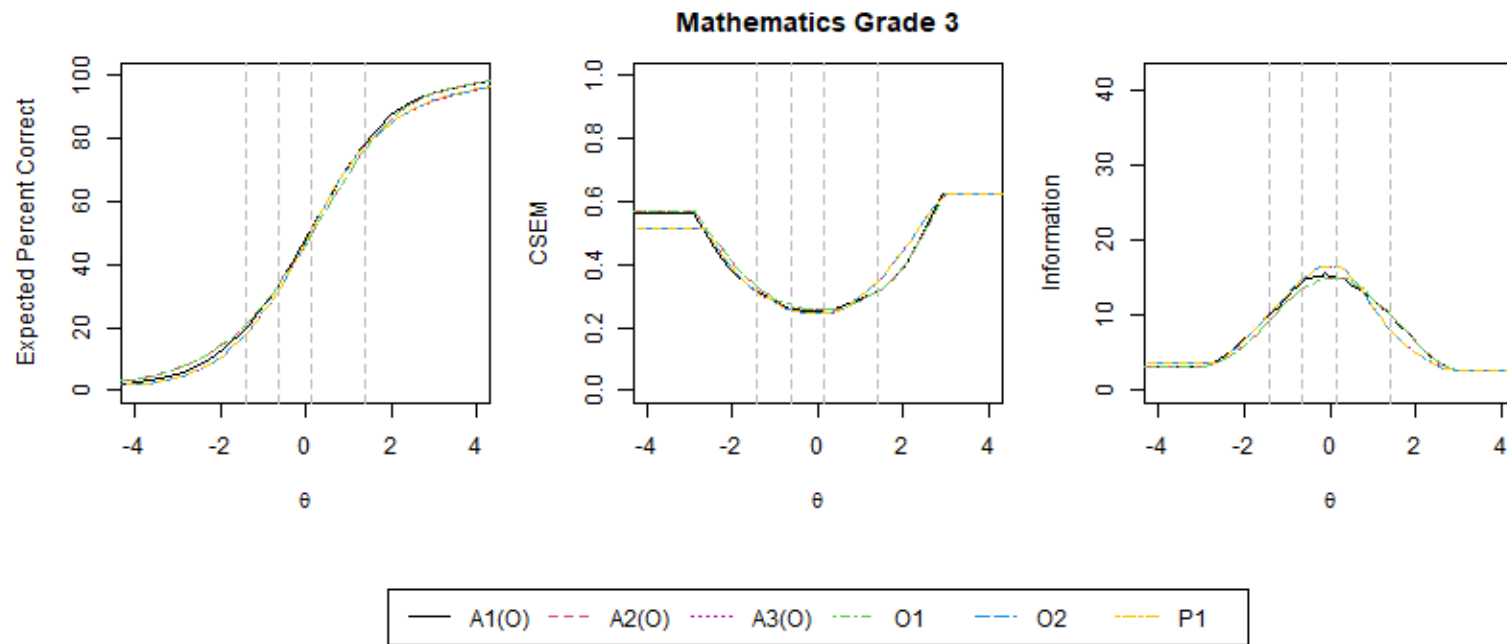


Figure A.12.9 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 3

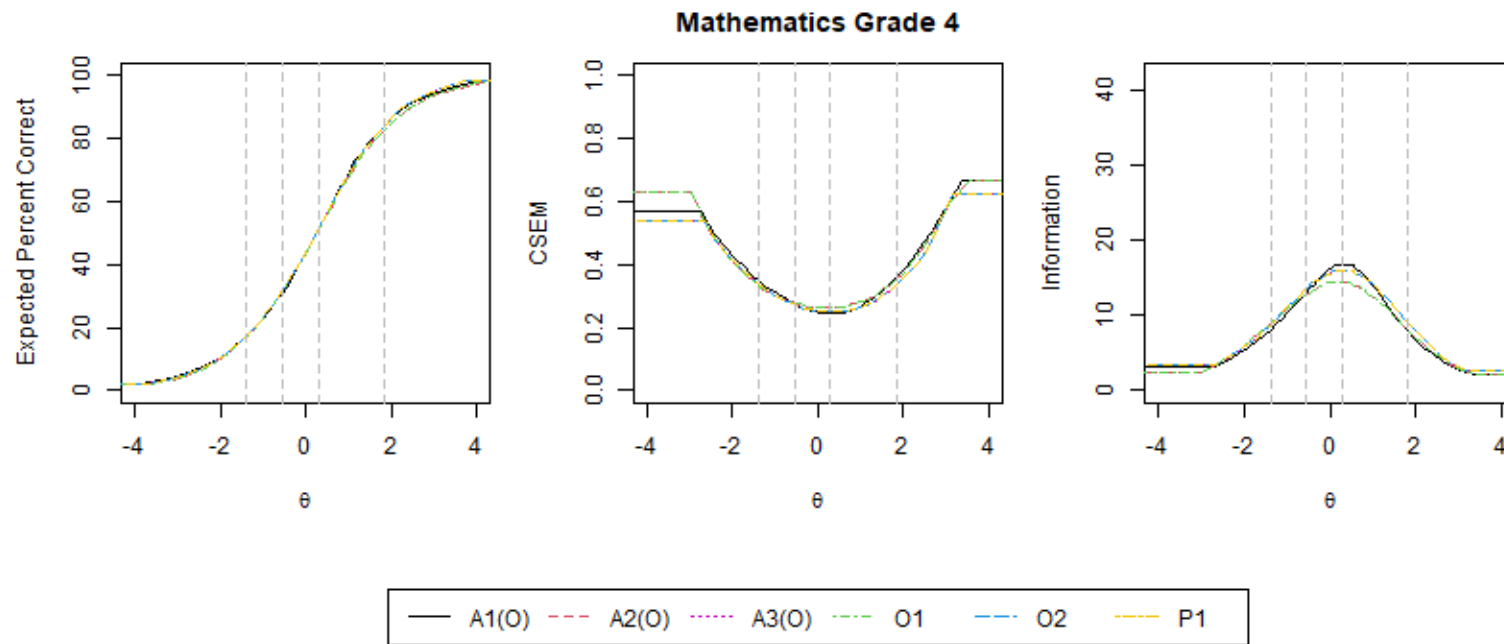


Figure A.12.10 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 4

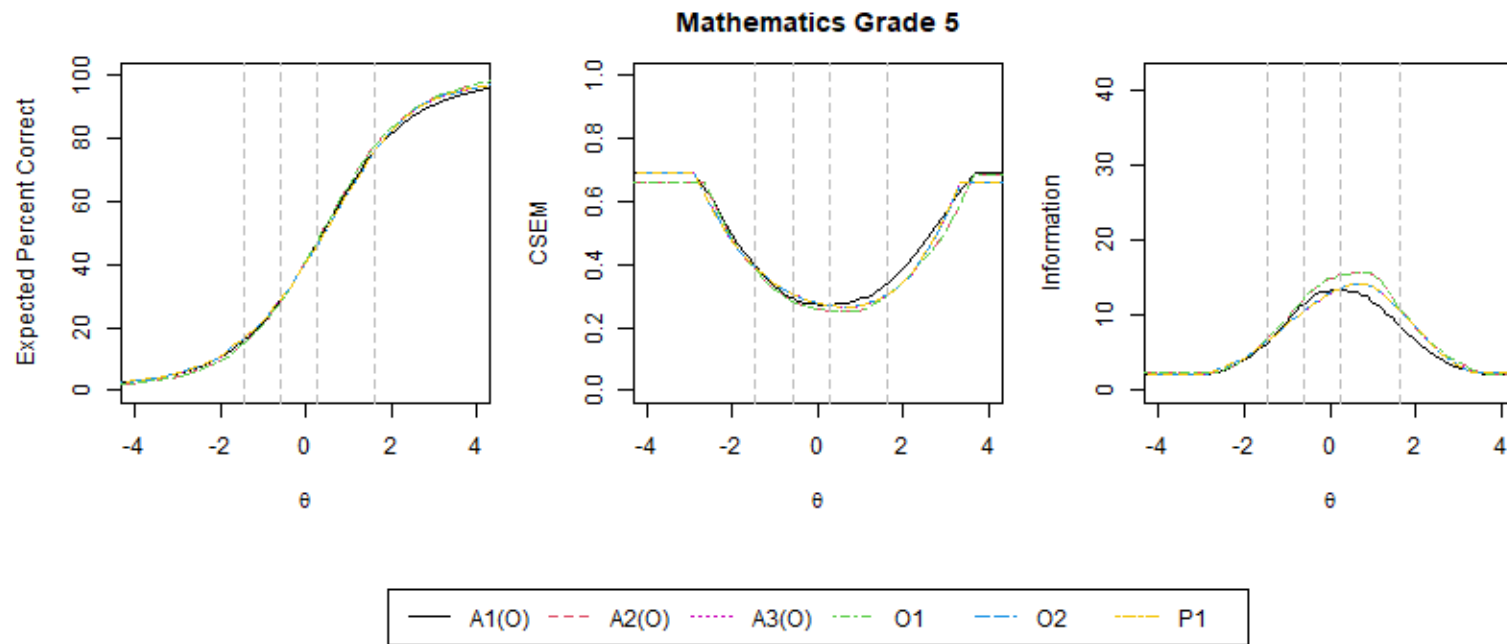


Figure A.12.11 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 5

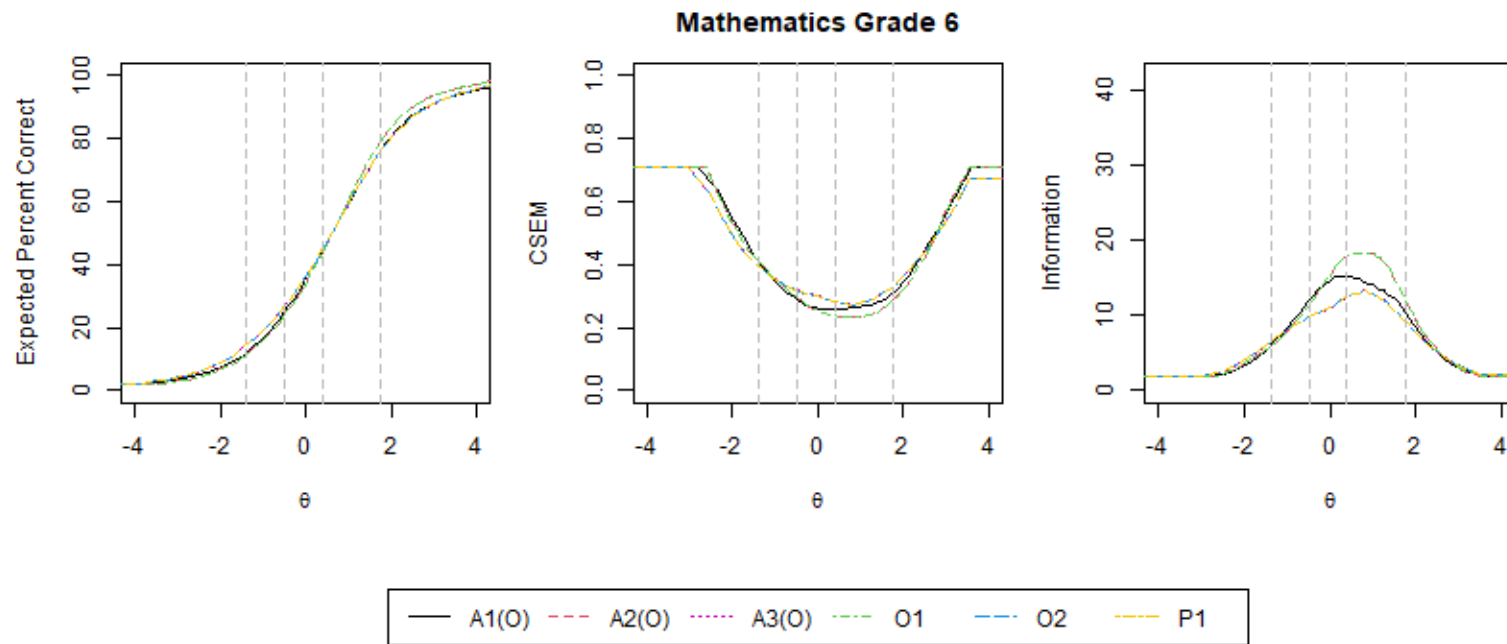


Figure A.12.12 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 6

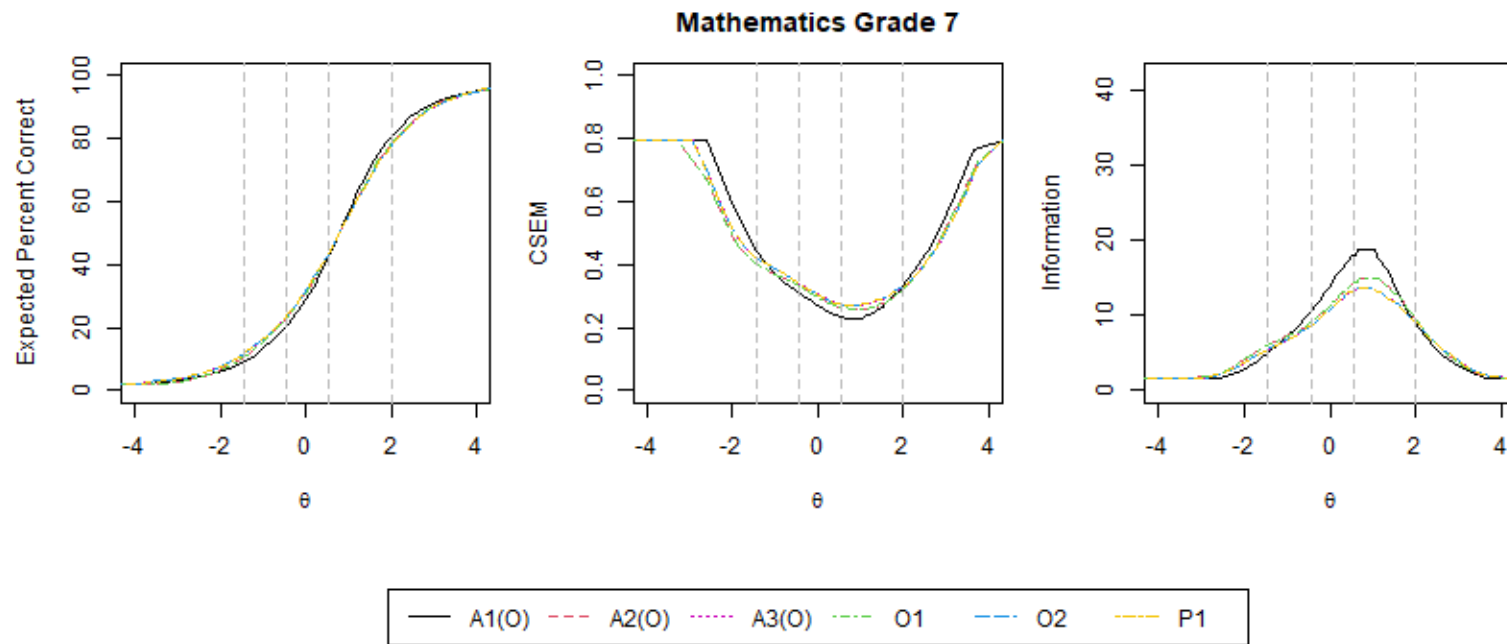


Figure A.12.13 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 7

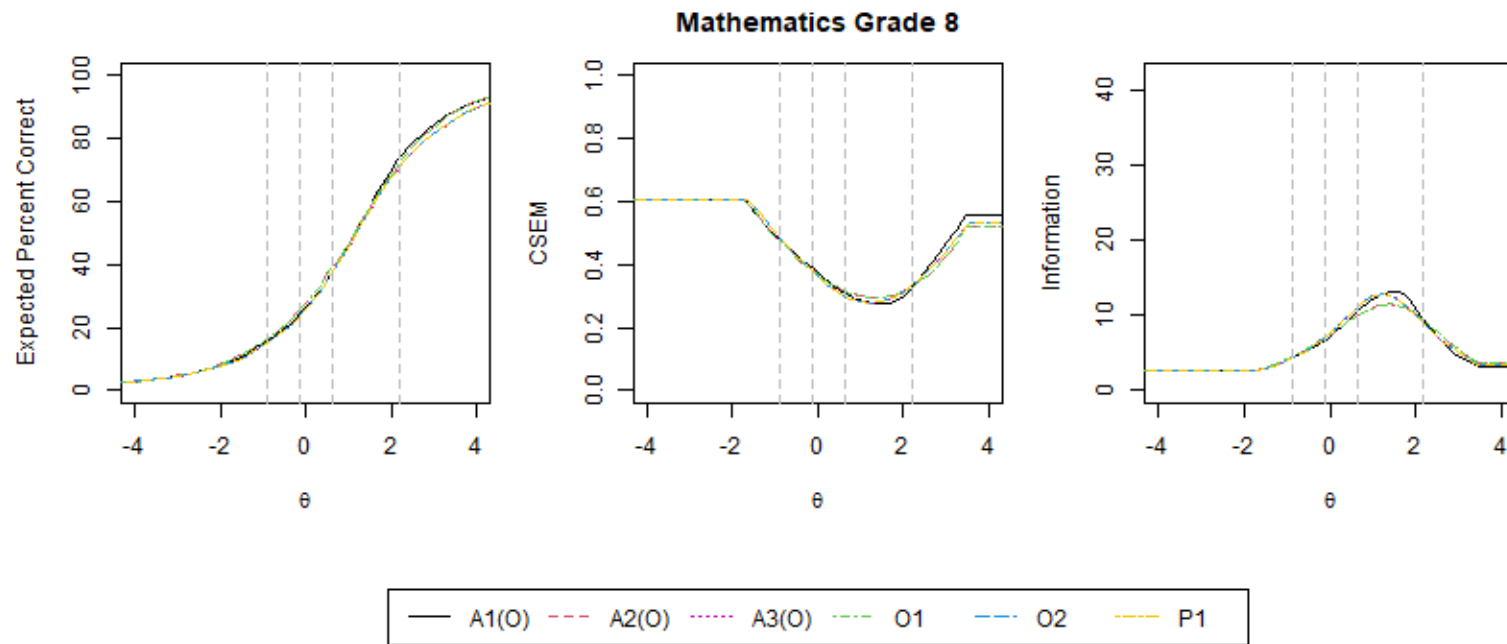


Figure A.12.14 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Mathematics Grade 8

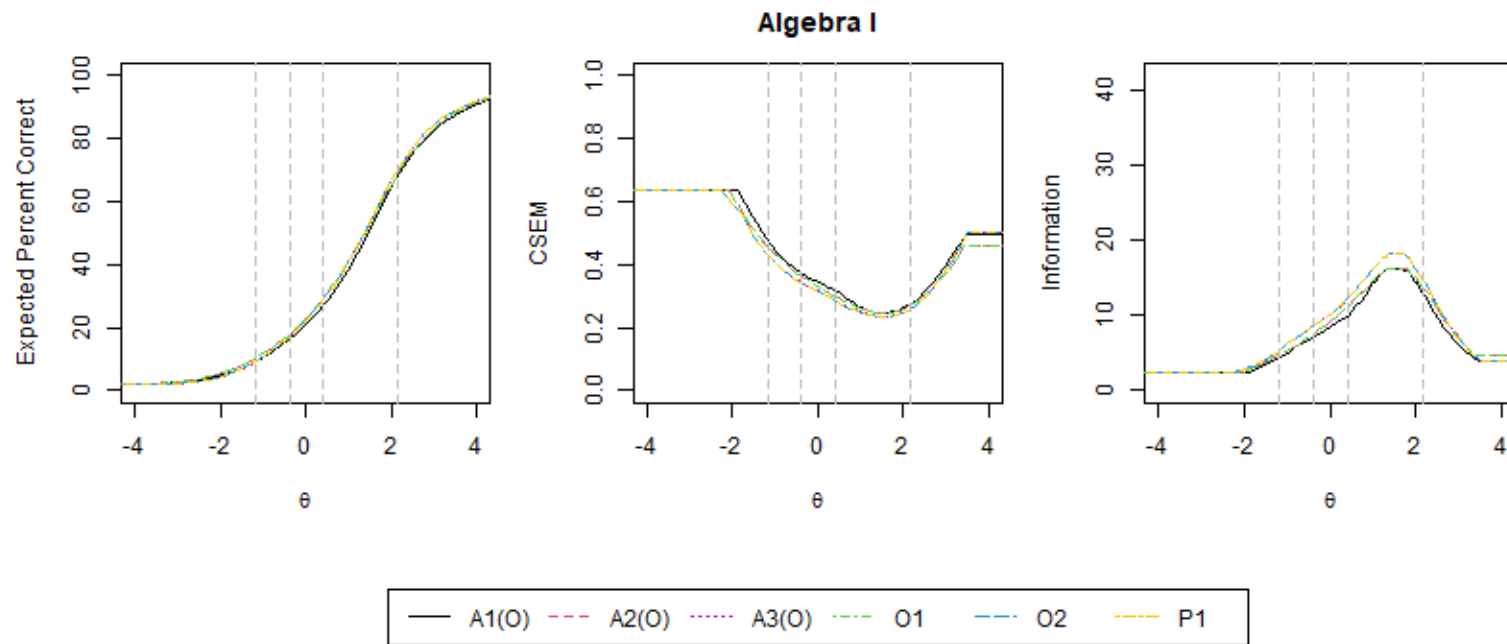


Figure A.12.15 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Algebra I

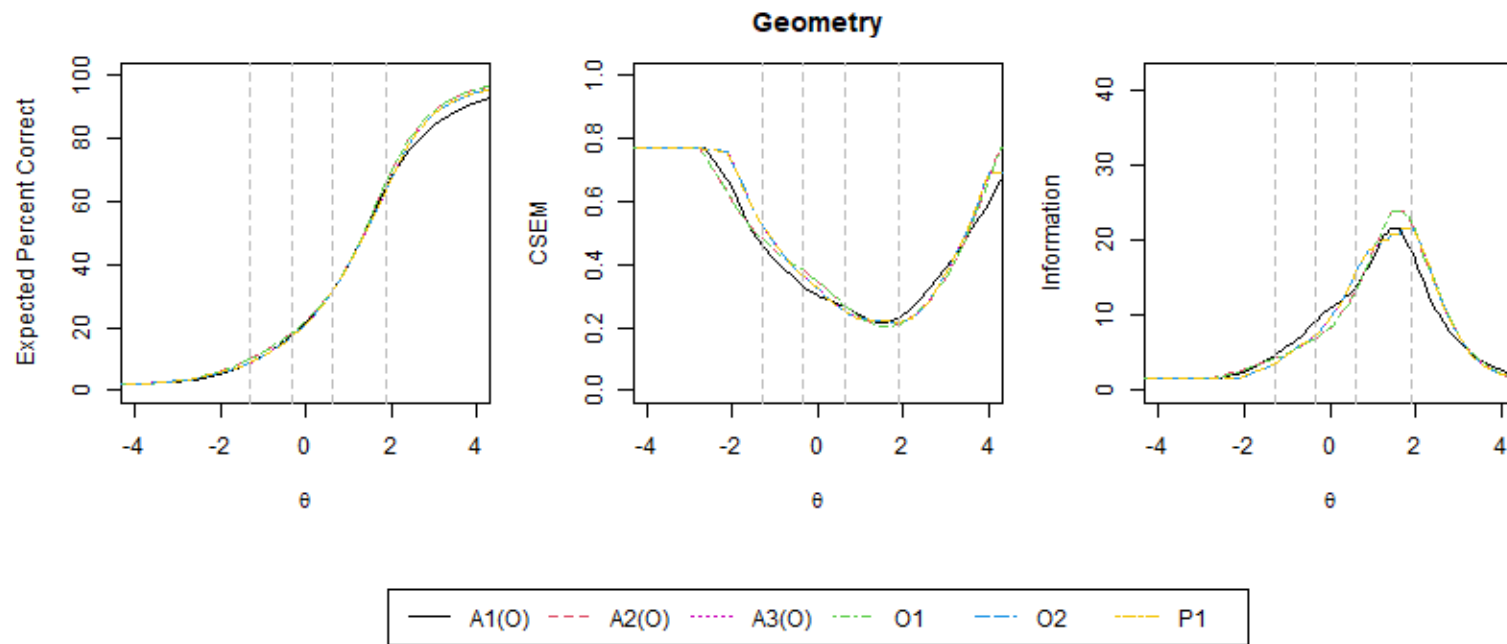


Figure A.12.16 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Geometry

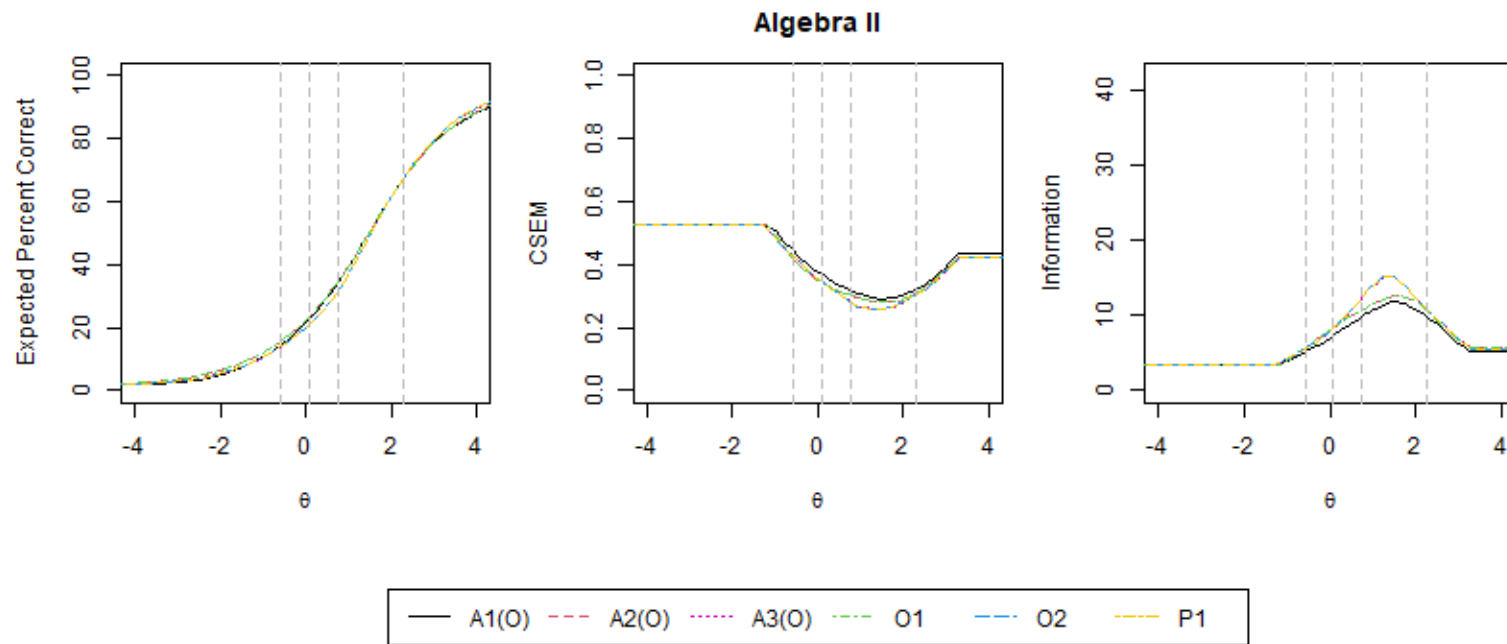


Figure A.12.17 Pre-Equated IRT Test Characteristic Curves, CSEM Curves, and Information Curves Algebra II

Appendix 12.4: Scale Score Cumulative Frequencies

Table A.12.23 Scale Score Cumulative Frequencies: ELA/L Grade 3

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	9,613	4.19	9,613	4.19
655–659	7,122	3.11	16,735	7.30
660–664	99	0.04	16,834	7.35
665–669	8,156	3.56	24,990	10.90
670–674	3,211	1.40	28,201	12.30
675–679	4,915	2.14	33,116	14.45
680–684	7,428	3.24	40,544	17.69
685–689	6,852	2.99	47,396	20.68
690–694	6,310	2.75	53,706	23.43
695–699	6,019	2.63	59,725	26.06
700–704	5,873	2.56	65,598	28.62
705–709	11,265	4.92	76,863	33.54
710–714	5,716	2.49	82,579	36.03
715–719	11,766	5.13	94,345	41.17
720–724	5,983	2.61	100,328	43.78
725–729	12,393	5.41	112,721	49.18
730–734	6,243	2.72	118,964	51.91
735–739	12,773	5.57	131,737	57.48
740–744	6,640	2.90	138,377	60.38
745–749	13,373	5.84	151,750	66.21
750–754	9,509	4.15	161,259	70.36
755–759	6,615	2.89	167,874	73.25
760–764	10,121	4.42	177,995	77.66
765–769	8,565	3.74	186,560	81.40
770–774	5,340	2.33	191,900	83.73
775–779	6,881	3.00	198,781	86.73
780–784	6,259	2.73	205,040	89.46
785–789	3,374	1.47	208,414	90.94
790–794	3,002	1.31	211,416	92.25
795–799	3,907	1.70	215,323	93.95
800–804	2,248	0.98	217,571	94.93
805–809	2,988	1.30	220,559	96.24
810–814	1,602	0.70	222,161	96.94
815–819	1,382	0.60	223,543	97.54
820–824	1,249	0.54	224,792	98.08
825–829	989	0.43	225,781	98.51
830–834	405	0.18	226,186	98.69
835–839	644	0.28	226,830	98.97
840–844	538	0.23	227,368	99.21
845–849	490	0.21	227,858	99.42
850	1,327	0.58	229,185	100.00

Table A.12.24 Scale Score Cumulative Frequencies: ELA/L Grade 4

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	2,573	1.11	2,573	1.11
655–659	2,815	1.21	5,388	2.32
660–664	74	0.03	5,462	2.35
665–669	4,341	1.87	9,803	4.22
670–674	5,553	2.39	15,356	6.61
675–679	2,339	1.01	17,695	7.62
680–684	3,856	1.66	21,551	9.28
685–689	6,389	2.75	27,940	12.03
690–694	8,485	3.65	36,425	15.68
695–699	9,017	3.88	45,442	19.56
700–704	5,266	2.27	50,708	21.83
705–709	9,878	4.25	60,586	26.08
710–714	5,023	2.16	65,609	28.24
715–719	9,940	4.28	75,549	32.52
720–724	10,144	4.37	85,693	36.89
725–729	10,222	4.40	95,915	41.29
730–734	10,703	4.61	106,618	45.89
735–739	10,874	4.68	117,492	50.57
740–744	11,149	4.80	128,641	55.37
745–749	11,166	4.81	139,807	60.18
750–754	13,974	6.02	153,781	66.19
755–759	10,906	4.69	164,687	70.89
760–764	10,266	4.42	174,953	75.31
765–769	9,411	4.05	184,364	79.36
770–774	6,594	2.84	190,958	82.20
775–779	8,181	3.52	199,139	85.72
780–784	7,251	3.12	206,390	88.84
785–789	4,775	2.06	211,165	90.90
790–794	5,616	2.42	216,781	93.31
795–799	3,519	1.51	220,300	94.83
800–804	3,760	1.62	224,060	96.45
805–809	2,025	0.87	226,085	97.32
810–814	1,992	0.86	228,077	98.17
815–819	783	0.34	228,860	98.51
820–824	923	0.40	229,783	98.91
825–829	980	0.42	230,763	99.33
830–834	346	0.15	231,109	99.48
835–839	376	0.16	231,485	99.64
840–844	226	0.10	231,711	99.74
845–849	171	0.07	231,882	99.81
850	435	0.19	232,317	100.00

Table A.12.25 Scale Score Cumulative Frequencies: ELA/L Grade 5

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	1,746	0.73	1,746	0.73
655–659	1,528	0.64	3,274	1.38
660–664	975	0.41	4,249	1.79
665–669	2,504	1.05	6,753	2.84
670–674	1,570	0.66	8,323	3.50
675–679	5,677	2.39	14,000	5.89
680–684	4,255	1.79	18,255	7.68
685–689	7,023	2.95	25,278	10.63
690–694	7,201	3.03	32,479	13.66
695–699	9,311	3.91	41,790	17.57
700–704	6,637	2.79	48,427	20.36
705–709	9,688	4.07	58,115	24.43
710–714	8,223	3.46	66,338	27.89
715–719	8,866	3.73	75,204	31.62
720–724	10,943	4.60	86,147	36.22
725–729	10,788	4.54	96,935	40.76
730–734	10,787	4.54	107,722	45.29
735–739	14,059	5.91	121,781	51.20
740–744	10,818	4.55	132,599	55.75
745–749	13,655	5.74	146,254	61.49
750–754	10,519	4.42	156,773	65.91
755–759	10,223	4.30	166,996	70.21
760–764	12,540	5.27	179,536	75.48
765–769	9,161	3.85	188,697	79.34
770–774	8,262	3.47	196,959	82.81
775–779	7,470	3.14	204,429	85.95
780–784	6,558	2.76	210,987	88.71
785–789	5,843	2.46	216,830	91.16
790–794	3,980	1.67	220,810	92.84
795–799	4,341	1.83	225,151	94.66
800–804	3,594	1.51	228,745	96.17
805–809	2,248	0.95	230,993	97.12
810–814	1,718	0.72	232,711	97.84
815–819	1,820	0.77	234,531	98.61
820–824	966	0.41	235,497	99.01
825–829	540	0.23	236,037	99.24
830–834	573	0.24	236,610	99.48
835–839	305	0.13	236,915	99.61
840–844	357	0.15	237,272	99.76
845–849	173	0.07	237,445	99.83
850	402	0.17	237,847	100.00

Table A.12.26 Scale Score Cumulative Frequencies: ELA/L Grade 6

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	1,711	0.72	1,711	0.72
655–659	—	—	—	—
660–664	2,463	1.03	4,174	1.75
665–669	—	—	—	—
670–674	3,695	1.55	7,869	3.30
675–679	1,585	0.67	9,454	3.97
680–684	5,083	2.13	14,537	6.10
685–689	5,542	2.33	20,079	8.43
690–694	5,610	2.35	25,689	10.78
695–699	9,041	3.79	34,730	14.58
700–704	7,430	3.12	42,160	17.69
705–709	10,447	4.38	52,607	22.08
710–714	10,231	4.29	62,838	26.37
715–719	10,118	4.25	72,956	30.62
720–724	10,063	4.22	83,019	34.84
725–729	15,084	6.33	98,103	41.17
730–734	12,763	5.36	110,866	46.53
735–739	12,445	5.22	123,311	51.75
740–744	12,727	5.34	136,038	57.10
745–749	14,833	6.23	150,871	63.32
750–754	14,327	6.01	165,198	69.34
755–759	11,584	4.86	176,782	74.20
760–764	10,704	4.49	187,486	78.69
765–769	9,989	4.19	197,475	82.88
770–774	8,976	3.77	206,451	86.65
775–779	7,839	3.29	214,290	89.94
780–784	5,464	2.29	219,754	92.23
785–789	4,425	1.86	224,179	94.09
790–794	3,565	1.50	227,744	95.59
795–799	3,000	1.26	230,744	96.85
800–804	1,888	0.79	232,632	97.64
805–809	1,841	0.77	234,473	98.41
810–814	1,059	0.44	235,532	98.86
815–819	860	0.36	236,392	99.22
820–824	624	0.26	237,016	99.48
825–829	363	0.15	237,379	99.63
830–834	379	0.16	237,758	99.79
835–839	102	0.04	237,860	99.83
840–844	94	0.04	237,954	99.87
845–849	127	0.05	238,081	99.92
850	179	0.08	238,260	100.00

Table A.12.27 Scale Score Cumulative Frequencies: ELA/L Grade 7

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	6,029	2.46	6,029	2.46
655–659	1,597	0.65	7,626	3.11
660–664	—	—	—	—
665–669	4,910	2.00	12,536	5.11
670–674	3,085	1.26	15,621	6.37
675–679	2,324	0.95	17,945	7.32
680–684	5,545	2.26	23,490	9.58
685–689	5,700	2.33	29,190	11.91
690–694	5,486	2.24	34,676	14.15
695–699	7,456	3.04	42,132	17.19
700–704	8,131	3.32	50,263	20.51
705–709	7,199	2.94	57,462	23.44
710–714	10,151	4.14	67,613	27.59
715–719	7,530	3.07	75,143	30.66
720–724	10,321	4.21	85,464	34.87
725–729	10,598	4.32	96,062	39.19
730–734	10,569	4.31	106,631	43.50
735–739	13,241	5.40	119,872	48.91
740–744	11,112	4.53	130,984	53.44
745–749	13,407	5.47	144,391	58.91
750–754	10,746	4.38	155,137	63.29
755–759	10,269	4.19	165,406	67.48
760–764	14,488	5.91	179,894	73.39
765–769	8,766	3.58	188,660	76.97
770–774	9,971	4.07	198,631	81.04
775–779	7,278	2.97	205,909	84.01
780–784	6,527	2.66	212,436	86.67
785–789	7,407	3.02	219,843	89.69
790–794	5,047	2.06	224,890	91.75
795–799	4,432	1.81	229,322	93.56
800–804	3,770	1.54	233,092	95.10
805–809	1,752	0.71	234,844	95.81
810–814	2,185	0.89	237,029	96.70
815–819	2,384	0.97	239,413	97.68
820–824	1,039	0.42	240,452	98.10
825–829	1,295	0.53	241,747	98.63
830–834	823	0.34	242,570	98.97
835–839	641	0.26	243,211	99.23
840–844	266	0.11	243,477	99.34
845–849	507	0.21	243,984	99.54
850	1,122	0.46	245,106	100.00

Table A.12.28 Scale Score Cumulative Frequencies: ELA/L Grade 8

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	5,464	2.19	5,464	2.19
655–659	2,371	0.95	7,835	3.14
660–664	2,013	0.81	9,848	3.95
665–669	2,799	1.12	12,647	5.07
670–674	5,665	2.27	18,312	7.34
675–679	3,141	1.26	21,453	8.60
680–684	6,316	2.53	27,769	11.13
685–689	6,393	2.56	34,162	13.70
690–694	6,292	2.52	40,454	16.22
695–699	8,416	3.37	48,870	19.60
700–704	8,298	3.33	57,168	22.92
705–709	5,307	2.13	62,475	25.05
710–714	9,971	4.00	72,446	29.05
715–719	10,120	4.06	82,566	33.11
720–724	9,872	3.96	92,438	37.07
725–729	10,013	4.02	102,451	41.08
730–734	12,787	5.13	115,238	46.21
735–739	10,243	4.11	125,481	50.32
740–744	10,400	4.17	135,881	54.49
745–749	15,440	6.19	151,321	60.68
750–754	10,260	4.11	161,581	64.79
755–759	10,101	4.05	171,682	68.84
760–764	11,961	4.80	183,643	73.64
765–769	9,209	3.69	192,852	77.33
770–774	8,385	3.36	201,237	80.69
775–779	7,776	3.12	209,013	83.81
780–784	7,150	2.87	216,163	86.68
785–789	6,253	2.51	222,416	89.18
790–794	4,282	1.72	226,698	90.90
795–799	4,924	1.97	231,622	92.88
800–804	4,013	1.61	235,635	94.49
805–809	2,731	1.10	238,366	95.58
810–814	1,586	0.64	239,952	96.22
815–819	2,861	1.15	242,813	97.36
820–824	1,188	0.48	244,001	97.84
825–829	1,090	0.44	245,091	98.28
830–834	918	0.37	246,009	98.65
835–839	841	0.34	246,850	98.98
840–844	427	0.17	247,277	99.15
845–849	604	0.24	247,881	99.40
850	1,507	0.60	249,388	100.00

Table A.12.29 Scale Score Cumulative Frequencies: ELA/L Grade 9

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	727	1.51	727	1.51
655–659	3	0.01	730	1.52
660–664	480	1.00	1,210	2.52
665–669	594	1.24	1,804	3.75
670–674	2	0.00	1,806	3.76
675–679	656	1.37	2,462	5.12
680–684	707	1.47	3,169	6.59
685–689	758	1.58	3,927	8.17
690–694	776	1.61	4,703	9.79
695–699	756	1.57	5,459	11.36
700–704	1,527	3.18	6,986	14.54
705–709	823	1.71	7,809	16.25
710–714	1,642	3.42	9,451	19.67
715–719	1,767	3.68	11,218	23.35
720–724	1,837	3.82	13,055	27.17
725–729	1,989	4.14	15,044	31.31
730–734	2,096	4.36	17,140	35.67
735–739	3,172	6.60	20,312	42.27
740–744	2,177	4.53	22,489	46.80
745–749	2,142	4.46	24,631	51.26
750–754	3,236	6.73	27,867	57.99
755–759	2,178	4.53	30,045	62.53
760–764	2,004	4.17	32,049	66.70
765–769	3,053	6.35	35,102	73.05
770–774	1,900	3.95	37,002	77.00
775–779	1,903	3.96	38,905	80.96
780–784	1,655	3.44	40,560	84.41
785–789	1,501	3.12	42,061	87.53
790–794	1,377	2.87	43,438	90.40
795–799	619	1.29	44,057	91.69
800–804	1,121	2.33	45,178	94.02
805–809	465	0.97	45,643	94.99
810–814	815	1.70	46,458	96.68
815–819	343	0.71	46,801	97.40
820–824	288	0.60	47,089	98.00
825–829	230	0.48	47,319	98.47
830–834	210	0.44	47,529	98.91
835–839	174	0.36	47,703	99.27
840–844	128	0.27	47,831	99.54
845–849	85	0.18	47,916	99.72
850	136	0.28	48,052	100.00

Table A.12.30 Scale Score Cumulative Frequencies: ELA/L Grade 10

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	451	6.38	451	6.38
655–659	—	—	—	—
660–664	149	2.11	600	8.49
665–669	84	1.19	684	9.68
670–674	93	1.32	777	10.99
675–679	171	2.42	948	13.41
680–684	153	2.16	1,101	15.58
685–689	155	2.19	1,256	17.77
690–694	169	2.39	1,425	20.16
695–699	144	2.04	1,569	22.20
700–704	131	1.85	1,700	24.05
705–709	327	4.63	2,027	28.67
710–714	148	2.09	2,175	30.77
715–719	189	2.67	2,364	33.44
720–724	201	2.84	2,565	36.29
725–729	294	4.16	2,859	40.44
730–734	154	2.18	3,013	42.62
735–739	318	4.50	3,331	47.12
740–744	314	4.44	3,645	51.56
745–749	317	4.48	3,962	56.05
750–754	143	2.02	4,105	58.07
755–759	301	4.26	4,406	62.33
760–764	319	4.51	4,725	66.84
765–769	238	3.37	4,963	70.21
770–774	301	4.26	5,264	74.47
775–779	238	3.37	5,502	77.83
780–784	177	2.50	5,679	80.34
785–789	179	2.53	5,858	82.87
790–794	178	2.52	6,036	85.39
795–799	192	2.72	6,228	88.10
800–804	164	2.32	6,392	90.42
805–809	115	1.63	6,507	92.05
810–814	61	0.86	6,568	92.91
815–819	113	1.60	6,681	94.51
820–824	69	0.98	6,750	95.49
825–829	45	0.64	6,795	96.12
830–834	86	1.22	6,881	97.34
835–839	27	0.38	6,908	97.72
840–844	43	0.61	6,951	98.33
845–849	26	0.37	6,977	98.70
850	92	1.30	7,069	100.00

Table A.12.31 Scale Score Cumulative Frequencies: Mathematics Grade 3

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	3,819	1.66	3,819	1.66
655–659	—	—	—	—
660–664	3,613	1.57	7,432	3.23
665–669	—	—	—	—
670–674	4,794	2.08	12,226	5.31
675–679	5,553	2.41	17,779	7.73
680–684	2,751	1.20	20,530	8.92
685–689	6,161	2.68	26,691	11.60
690–694	6,598	2.87	33,289	14.47
695–699	6,796	2.95	40,085	17.42
700–704	6,831	2.97	46,916	20.39
705–709	10,652	4.63	57,568	25.02
710–714	9,796	4.26	67,364	29.28
715–719	6,806	2.96	74,170	32.24
720–724	12,921	5.62	87,091	37.86
725–729	9,356	4.07	96,447	41.92
730–734	12,452	5.41	108,899	47.34
735–739	11,851	5.15	120,750	52.49
740–744	8,591	3.73	129,341	56.22
745–749	11,262	4.90	140,603	61.12
750–754	11,087	4.82	151,690	65.94
755–759	10,542	4.58	162,232	70.52
760–764	10,360	4.50	172,592	75.02
765–769	7,347	3.19	179,939	78.22
770–774	9,339	4.06	189,278	82.28
775–779	6,676	2.90	195,954	85.18
780–784	6,323	2.75	202,277	87.93
785–789	5,741	2.50	208,018	90.42
790–794	5,112	2.22	213,130	92.64
795–799	3,115	1.35	216,245	94.00
800–804	2,920	1.27	219,165	95.27
805–809	2,565	1.11	221,730	96.38
810–814	1,092	0.47	222,822	96.86
815–819	2,072	0.90	224,894	97.76
820–824	1,748	0.76	226,642	98.52
825–829	2	0.00	226,644	98.52
830–834	1,287	0.56	227,931	99.08
835–839	—	—	—	—
840–844	992	0.43	228,923	99.51
845–849	—	—	—	—
850	1,131	0.49	230,054	100.00

Table A.12.32 Scale Score Cumulative Frequencies: Mathematics Grade 4

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	2,066	0.89	2,066	0.89
655–659	—	—	—	—
660–664	3,049	1.31	5,115	2.20
665–669	—	—	—	—
670–674	4,795	2.06	9,910	4.25
675–679	2,734	1.17	12,644	5.43
680–684	3,774	1.62	16,418	7.05
685–689	7,679	3.30	24,097	10.34
690–694	8,472	3.64	32,569	13.98
695–699	8,577	3.68	41,146	17.66
700–704	16,951	7.28	58,097	24.94
705–709	7,740	3.32	65,837	28.26
710–714	7,642	3.28	73,479	31.54
715–719	14,089	6.05	87,568	37.59
720–724	13,274	5.70	100,842	43.29
725–729	12,320	5.29	113,162	48.58
730–734	8,843	3.80	122,005	52.37
735–739	11,468	4.92	133,473	57.30
740–744	14,215	6.10	147,688	63.40
745–749	10,587	4.54	158,275	67.94
750–754	10,080	4.33	168,355	72.27
755–759	9,466	4.06	177,821	76.33
760–764	9,252	3.97	187,073	80.30
765–769	8,329	3.58	195,402	83.88
770–774	7,640	3.28	203,042	87.16
775–779	6,944	2.98	209,986	90.14
780–784	4,877	2.09	214,863	92.23
785–789	4,334	1.86	219,197	94.09
790–794	3,800	1.63	222,997	95.73
795–799	2,230	0.96	225,227	96.68
800–804	1,921	0.82	227,148	97.51
805–809	1,586	0.68	228,734	98.19
810–814	1,338	0.57	230,072	98.76
815–819	1,105	0.47	231,177	99.24
820–824	389	0.17	231,566	99.40
825–829	369	0.16	231,935	99.56
830–834	304	0.13	232,239	99.69
835–839	233	0.10	232,472	99.79
840–844	—	—	—	—
845–849	203	0.09	232,675	99.88
850	280	0.12	232,955	100.00

Table A.12.33 Scale Score Cumulative Frequencies: Mathematics Grade 5

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	3,622	1.52	3,622	1.52
655–659	1,384	0.58	5,006	2.10
660–664	—	—	—	—
665–669	5,713	2.39	10,719	4.49
670–674	57	0.02	10,776	4.52
675–679	7,637	3.20	18,413	7.72
680–684	4,979	2.09	23,392	9.81
685–689	3,989	1.67	27,381	11.48
690–694	9,940	4.17	37,321	15.65
695–699	9,877	4.14	47,198	19.79
700–704	10,193	4.27	57,391	24.06
705–709	10,157	4.26	67,548	28.32
710–714	14,131	5.92	81,679	34.24
715–719	13,932	5.84	95,611	40.08
720–724	12,236	5.13	107,847	45.21
725–729	15,082	6.32	122,929	51.53
730–734	13,878	5.82	136,807	57.35
735–739	9,358	3.92	146,165	61.27
740–744	11,652	4.88	157,817	66.16
745–749	10,652	4.47	168,469	70.62
750–754	9,858	4.13	178,327	74.76
755–759	8,885	3.72	187,212	78.48
760–764	12,173	5.10	199,385	83.58
765–769	7,149	3.00	206,534	86.58
770–774	6,565	2.75	213,099	89.33
775–779	5,897	2.47	218,996	91.80
780–784	2,685	1.13	221,681	92.93
785–789	4,955	2.08	226,636	95.01
790–794	2,970	1.25	229,606	96.25
795–799	2,726	1.14	232,332	97.40
800–804	1,456	0.61	233,788	98.01
805–809	1,713	0.72	235,501	98.72
810–814	554	0.23	236,055	98.96
815–819	792	0.33	236,847	99.29
820–824	319	0.13	237,166	99.42
825–829	587	0.25	237,753	99.67
830–834	—	—	—	—
835–839	394	0.17	238,147	99.83
840–844	—	—	—	—
845–849	172	0.07	238,319	99.91
850	226	0.09	238,545	100.00

Table A.12.34 Scale Score Cumulative Frequencies: Mathematics Grade 6

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	3,537	1.48	3,537	1.48
655–659	—	—	—	—
660–664	40	0.02	3,577	1.50
665–669	5,784	2.42	9,361	3.92
670–674	87	0.04	9,448	3.96
675–679	8,960	3.75	18,408	7.71
680–684	118	0.05	18,526	7.76
685–689	11,011	4.61	29,537	12.38
690–694	11,978	5.02	41,515	17.40
695–699	5,288	2.22	46,803	19.61
700–704	11,617	4.87	58,420	24.48
705–709	16,472	6.90	74,892	31.38
710–714	10,252	4.30	85,144	35.68
715–719	14,272	5.98	99,416	41.66
720–724	13,175	5.52	112,591	47.18
725–729	16,210	6.79	128,801	53.98
730–734	11,309	4.74	140,110	58.71
735–739	16,437	6.89	156,547	65.60
740–744	8,590	3.60	165,137	69.20
745–749	13,366	5.60	178,503	74.80
750–754	11,171	4.68	189,674	79.48
755–759	7,975	3.34	197,649	82.83
760–764	8,615	3.61	206,264	86.44
765–769	7,536	3.16	213,800	89.59
770–774	6,311	2.64	220,111	92.24
775–779	4,355	1.83	224,466	94.06
780–784	3,749	1.57	228,215	95.64
785–789	3,156	1.32	231,371	96.96
790–794	2,024	0.85	233,395	97.81
795–799	1,147	0.48	234,542	98.29
800–804	1,406	0.59	235,948	98.88
805–809	740	0.31	236,688	99.19
810–814	614	0.26	237,302	99.44
815–819	260	0.11	237,562	99.55
820–824	225	0.09	237,787	99.65
825–829	379	0.16	238,166	99.81
830–834	—	—	—	—
835–839	262	0.11	238,428	99.92
840–844	1	0.00	238,429	99.92
845–849	—	—	—	—
850	201	0.08	238,630	100.00

Table A.12.35 Scale Score Cumulative Frequencies: Mathematics Grade 7

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	856	0.36	856	0.36
655–659	—	—	—	—
660–664	987	0.42	1,843	0.78
665–669	—	—	—	—
670–674	1,275	0.54	3,118	1.32
675–679	2,133	0.90	5,251	2.23
680–684	2,498	1.06	7,749	3.29
685–689	3,440	1.46	11,189	4.75
690–694	8,501	3.61	19,690	8.35
695–699	9,940	4.22	29,630	12.57
700–704	10,914	4.63	40,544	17.20
705–709	11,088	4.70	51,632	21.91
710–714	10,816	4.59	62,448	26.50
715–719	15,949	6.77	78,397	33.26
720–724	14,088	5.98	92,485	39.24
725–729	13,494	5.73	105,979	44.97
730–734	16,346	6.94	122,325	51.90
735–739	14,640	6.21	136,965	58.11
740–744	19,061	8.09	156,026	66.20
745–749	10,957	4.65	166,983	70.85
750–754	14,632	6.21	181,615	77.06
755–759	8,413	3.57	190,028	80.63
760–764	11,068	4.70	201,096	85.32
765–769	7,951	3.37	209,047	88.69
770–774	7,020	2.98	216,067	91.67
775–779	5,940	2.52	222,007	94.19
780–784	3,916	1.66	225,923	95.86
785–789	2,466	1.05	228,389	96.90
790–794	2,092	0.89	230,481	97.79
795–799	1,746	0.74	232,227	98.53
800–804	987	0.42	233,214	98.95
805–809	774	0.33	233,988	99.28
810–814	636	0.27	234,624	99.55
815–819	3	0.00	234,627	99.55
820–824	469	0.20	235,096	99.75
825–829	4	0.00	235,100	99.75
830–834	309	0.13	235,409	99.88
835–839	—	—	—	—
840–844	62	0.03	235,471	99.91
845–849	104	0.04	235,575	99.95
850	117	0.05	235,692	100.00

Table A.12.36 Scale Score Cumulative Frequencies: Mathematics Grade 8

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	5,867	2.77	5,867	2.77
655–659	2,557	1.21	8,424	3.98
660–664	4,659	2.20	13,083	6.18
665–669	—	—	—	—
670–674	4,406	2.08	17,489	8.27
675–679	6,611	3.12	24,100	11.39
680–684	6,363	3.01	30,463	14.40
685–689	15,191	7.18	45,654	21.58
690–694	8,317	3.93	53,971	25.51
695–699	8,204	3.88	62,175	29.39
700–704	15,515	7.33	77,690	36.72
705–709	14,375	6.79	92,065	43.52
710–714	12,805	6.05	104,870	49.57
715–719	11,003	5.20	115,873	54.77
720–724	9,386	4.44	125,259	59.21
725–729	8,293	3.92	133,552	63.13
730–734	7,394	3.49	140,946	66.62
735–739	12,445	5.88	153,391	72.50
740–744	5,373	2.54	158,764	75.04
745–749	9,280	4.39	168,044	79.43
750–754	5,869	2.77	173,913	82.20
755–759	5,140	2.43	179,053	84.63
760–764	5,961	2.82	185,014	87.45
765–769	5,043	2.38	190,057	89.83
770–774	3,322	1.57	193,379	91.40
775–779	3,737	1.77	197,116	93.17
780–784	2,338	1.11	199,454	94.28
785–789	2,080	0.98	201,534	95.26
790–794	1,862	0.88	203,396	96.14
795–799	1,664	0.79	205,060	96.93
800–804	1,410	0.67	206,470	97.59
805–809	878	0.42	207,348	98.01
810–814	730	0.35	208,078	98.35
815–819	681	0.32	208,759	98.67
820–824	546	0.26	209,305	98.93
825–829	486	0.23	209,791	99.16
830–834	428	0.20	210,219	99.36
835–839	346	0.16	210,565	99.53
840–844	172	0.08	210,737	99.61
845–849	129	0.06	210,866	99.67
850	697	0.33	211,563	100.00

Table A.12.37 Scale Score Cumulative Frequencies: Algebra I

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	851	0.74	851	0.74
655–659	847	0.74	1,698	1.47
660–664	3	0.00	1,701	1.48
665–669	1,320	1.15	3,021	2.62
670–674	1,722	1.50	4,743	4.12
675–679	18	0.02	4,761	4.13
680–684	5,271	4.58	10,032	8.71
685–689	—	—	—	—
690–694	6,727	5.84	16,759	14.55
695–699	3,513	3.05	20,272	17.60
700–704	3,833	3.33	24,105	20.93
705–709	7,218	6.27	31,323	27.20
710–714	3,148	2.73	34,471	29.93
715–719	6,381	5.54	40,852	35.47
720–724	5,829	5.06	46,681	40.53
725–729	5,430	4.71	52,111	45.25
730–734	7,363	6.39	59,474	51.64
735–739	6,530	5.67	66,004	57.31
740–744	3,890	3.38	69,894	60.69
745–749	5,441	4.72	75,335	65.41
750–754	4,960	4.31	80,295	69.72
755–759	5,689	4.94	85,984	74.66
760–764	5,067	4.40	91,051	79.06
765–769	3,300	2.87	94,351	81.92
770–774	3,815	3.31	98,166	85.23
775–779	3,990	3.46	102,156	88.70
780–784	2,759	2.40	104,915	91.09
785–789	2,243	1.95	107,158	93.04
790–794	2,269	1.97	109,427	95.01
795–799	1,387	1.20	110,814	96.21
800–804	1,142	0.99	111,956	97.21
805–809	956	0.83	112,912	98.04
810–814	728	0.63	113,640	98.67
815–819	287	0.25	113,927	98.92
820–824	385	0.33	114,312	99.25
825–829	182	0.16	114,494	99.41
830–834	142	0.12	114,636	99.53
835–839	147	0.13	114,783	99.66
840–844	52	0.05	114,835	99.71
845–849	97	0.08	114,932	99.79
850	242	0.21	115,174	100.00

Table A.12.38 Scale Score Cumulative Frequencies: Geometry

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	56	0.14	56	0.14
655–659	—	—	—	—
660–664	238	0.59	294	0.73
665–669	1	0.00	295	0.73
670–674	—	—	—	—
675–679	558	1.38	853	2.11
680–684	—	—	—	—
685–689	456	1.13	1,309	3.24
690–694	498	1.23	1,807	4.47
695–699	1,445	3.58	3,252	8.05
700–704	935	2.31	4,187	10.36
705–709	1,913	4.73	6,100	15.10
710–714	1,987	4.92	8,087	20.02
715–719	1,906	4.72	9,993	24.73
720–724	1,937	4.79	11,930	29.53
725–729	1,938	4.80	13,868	34.32
730–734	1,817	4.50	15,685	38.82
735–739	2,508	6.21	18,193	45.03
740–744	3,157	7.81	21,350	52.84
745–749	2,859	7.08	24,209	59.92
750–754	2,483	6.15	26,692	66.06
755–759	2,811	6.96	29,503	73.02
760–764	2,340	5.79	31,843	78.81
765–769	2,379	5.89	34,222	84.70
770–774	2,103	5.20	36,325	89.90
775–779	1,371	3.39	37,696	93.30
780–784	995	2.46	38,691	95.76
785–789	737	1.82	39,428	97.58
790–794	327	0.81	39,755	98.39
795–799	350	0.87	40,105	99.26
800–804	114	0.28	40,219	99.54
805–809	81	0.20	40,300	99.74
810–814	39	0.10	40,339	99.84
815–819	25	0.06	40,364	99.90
820–824	19	0.05	40,383	99.95
825–829	6	0.01	40,389	99.96
830–834	7	0.02	40,396	99.98
835–839	5	0.01	40,401	99.99
840–844	—	—	—	—
845–849	—	—	—	—
850	3	0.01	40,404	100.00

Table A.12.39 Scale Score Cumulative Frequencies: Algebra II

Score Band	Count	Percent	Cumulative Count	Cumulative Percent
650–654	263	1.88	263	1.88
655–659	—	—	—	—
660–664	127	0.91	390	2.79
665–669	164	1.17	554	3.97
670–674	225	1.61	779	5.58
675–679	232	1.66	1,011	7.24
680–684	251	1.80	1,262	9.04
685–689	286	2.05	1,548	11.08
690–694	588	4.21	2,136	15.30
695–699	—	—	—	—
700–704	544	3.90	2,680	19.19
705–709	520	3.72	3,200	22.91
710–714	245	1.75	3,445	24.67
715–719	530	3.80	3,975	28.46
720–724	547	3.92	4,522	32.38
725–729	483	3.46	5,005	35.84
730–734	500	3.58	5,505	39.42
735–739	516	3.69	6,021	43.11
740–744	715	5.12	6,736	48.23
745–749	491	3.52	7,227	51.75
750–754	693	4.96	7,920	56.71
755–759	599	4.29	8,519	61.00
760–764	781	5.59	9,300	66.60
765–769	390	2.79	9,690	69.39
770–774	693	4.96	10,383	74.35
775–779	671	4.80	11,054	79.16
780–784	421	3.01	11,475	82.17
785–789	529	3.79	12,004	85.96
790–794	342	2.45	12,346	88.41
795–799	281	2.01	12,627	90.42
800–804	346	2.48	12,973	92.90
805–809	156	1.12	13,129	94.01
810–814	189	1.35	13,318	95.37
815–819	112	0.80	13,430	96.17
820–824	140	1.00	13,570	97.17
825–829	87	0.62	13,657	97.79
830–834	76	0.54	13,733	98.34
835–839	25	0.18	13,758	98.52
840–844	55	0.39	13,813	98.91
845–849	48	0.34	13,861	99.26
850	104	0.74	13,965	100.00

Appendix 12.5: Subgroup Scale Score Performance

Table A.12.40 Subgroup Performance for ELA/L Scale Scores: Grade 3

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		229,184	730.45	44.30	650	850
Gender	Female	112,531	735.44	44.88	650	850
	Male	116,647	725.63	43.19	650	850
Ethnicity	American Indian/Alaska Native	507	728.13	46.08	650	850
	Asian	17,564	762.49	42.61	650	850
	Black or African American	38,247	709.72	40.39	650	850
	Hispanic/Latino	65,085	716.97	41.59	650	850
	Native Hawaiian or Pacific Islander	358	741.97	45.24	650	850
	Two or More Races	9,558	736.35	44.73	650	850
	White	97,143	741.20	40.94	650	850
Economic Status*	Not Economically Disadvantaged	123,355	744.71	42.33	650	850
	Economically Disadvantaged	96,131	711.28	39.65	650	850
English Learner Status	Non-English Learner	183,158	735.44	43.84	650	850
	English Learner	35,250	704.19	37.04	650	850
Disabilities	Students without Disabilities	188,159	735.84	43.24	650	850
	Student with Disability (SWD)	38,991	705.46	40.61	650	850
Reading Summative Score		229,184	42.98	17.62	10	90
Gender	Female	112,531	44.55	17.70	10	90
	Male	116,647	41.46	17.41	10	90
Ethnicity	American Indian/Alaska Native	507	41.67	17.75	10	84
	Asian	17,564	54.70	16.71	10	90
	Black or African American	38,247	35.07	16.16	10	90
	Hispanic/Latino	65,085	37.41	16.36	10	90
	Native Hawaiian or Pacific Islander	358	46.67	17.17	10	90
	Two or More Races	9,558	45.53	17.79	10	90
	White	97,143	47.43	16.53	10	90
Economic Status*	Not Economically Disadvantaged	123,355	48.57	16.89	10	90
	Economically Disadvantaged	96,131	35.40	15.69	10	90
English Learner Status	Non-English Learner	183,158	44.97	17.41	10	90
	English Learner	35,250	32.24	14.47	10	90
Disabilities	Students without Disabilities	188,159	45.06	17.20	10	90
	Student with Disability (SWD)	38,991	33.33	16.37	10	90
Writing Summative Score		229,184	27.46	13.47	10	60
Gender	Female	112,531	29.19	13.45	10	60
	Male	116,647	25.79	13.27	10	60

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Ethnicity	American Indian/Alaska Native	507	27.17	13.81	10	60
	Asian	17,564	36.43	12.14	10	60
	Black or African American	38,247	21.63	12.54	10	60
	Hispanic/Latino	65,085	24.21	13.03	10	60
	Native Hawaiian or Pacific Islander	358	31.01	13.45	10	60
	Two or More Races	9,558	28.69	13.57	10	60
	White	97,143	30.18	12.73	10	60
Economic Status*	Not Economically Disadvantaged	123,355	31.24	12.90	10	60
	Economically Disadvantaged	96,131	22.42	12.59	10	60
English Learner Status	Non-English Learner	183,158	28.70	13.40	10	60
	English Learner	35,250	21.15	12.06	10	60
Disabilities	Students without Disabilities	188,159	28.95	13.22	10	60
	Student with Disability (SWD)	38,991	20.60	12.44	10	60

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.41 Subgroup Performance for ELA/L Scale Scores: Grade 4

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		232,317	737.26	38.96	650	850
Gender	Female	114,052	741.07	38.94	650	850
	Male	118,262	733.59	38.63	650	850
Ethnicity	American Indian/Alaska Native	473	731.41	38.96	650	835
	Asian	18,097	765.46	36.75	650	850
	Black or African American	39,049	717.65	35.84	650	850
	Hispanic/Latino	66,242	725.43	36.51	650	850
	Native Hawaiian or Pacific Islander	401	750.41	35.77	656	850
	Two or More Races	9,301	743.15	39.33	650	850
	White	98,062	747.25	35.58	650	850
Economic Status*	Not Economically Disadvantaged	125,784	750.29	36.57	650	850
	Economically Disadvantaged	97,154	719.52	35.01	650	850
English Learner Status	Non-English Learner	186,611	741.84	38.34	650	850
	English Learner	35,009	712.11	31.99	650	850
Disabilities	Students without Disabilities	187,594	742.97	37.25	650	850
	Student with Disability (SWD)	42,727	713.13	36.93	650	850
Reading Summative Score		232,317	46.01	15.79	10	90
Gender	Female	114,052	46.97	15.59	10	90
	Male	118,262	45.07	15.93	10	90
Ethnicity	American Indian/Alaska Native	473	44.07	16.07	10	90
	Asian	18,097	56.66	15.05	10	90
	Black or African American	39,049	38.50	14.49	10	90
	Hispanic/Latino	66,242	41.09	14.59	10	90
	Native Hawaiian or Pacific Islander	401	50.50	14.30	10	90
	Two or More Races	9,301	48.59	16.03	10	90
	White	98,062	50.08	14.69	10	90
Economic Status*	Not Economically Disadvantaged	125,784	51.12	14.98	10	90
	Economically Disadvantaged	97,154	38.99	14.08	10	90
English Learner Status	Non-English Learner	186,611	47.85	15.57	10	90
	English Learner	35,009	35.76	12.64	10	90
Disabilities	Students without Disabilities	187,594	48.15	15.18	10	90
	Student with Disability (SWD)	42,727	36.90	15.12	10	90
Writing Summative Score		232,317	29.11	12.26	10	60
Gender	Female	114,052	30.71	12.00	10	60
	Male	118,262	27.56	12.31	10	60
Ethnicity	American Indian/Alaska Native	473	27.10	12.40	10	55

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Asian	18,097	37.01	10.21	10	60
	Black or African American	39,049	23.29	11.99	10	60
	Hispanic/Latino	66,242	26.13	12.12	10	60
	Native Hawaiian or Pacific Islander	401	33.32	11.14	10	60
	Two or More Races	9,301	30.51	12.09	10	60
	White	98,062	31.83	11.20	10	60
Economic Status*	Not Economically Disadvantaged	125,784	32.74	11.23	10	60
	Economically Disadvantaged	97,154	24.18	11.91	10	60
English Learner Status	Non-English Learner	186,611	30.31	12.04	10	60
	English Learner	35,009	22.58	11.41	10	60
Disabilities	Students without Disabilities	187,594	30.87	11.64	10	60
	Student with Disability (SWD)	42,727	21.69	12.04	10	60

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.42 Subgroup Performance for ELA/L Scale Scores: Grade 5

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		237,847	738.17	37.48	650	850
Gender	Female	116,510	742.97	37.73	650	850
	Male	121,320	733.54	36.65	650	850
Ethnicity	American Indian/Alaska Native	458	732.74	38.49	650	850
	Asian	18,019	766.41	35.83	650	850
	Black or African American	39,713	720.02	34.19	650	850
	Hispanic/Latino	68,390	727.26	35.24	650	850
	Native Hawaiian or Pacific Islander	416	749.74	35.23	650	845
	Two or More Races	8,968	742.70	36.78	650	850
	White	101,226	747.22	34.59	650	850
Economic Status*	Not Economically Disadvantaged	130,148	749.90	35.83	650	850
	Economically Disadvantaged	98,485	721.76	33.48	650	850
English Learner Status	Non-English Learner	196,956	742.43	36.73	650	850
	English Learner	30,452	709.33	28.50	650	833
Disabilities	Students without Disabilities	191,021	743.87	35.97	650	850
	Student with Disability (SWD)	44,831	714.70	34.52	650	850
Reading Summative Score		237,847	46.03	14.82	10	90
Gender	Female	116,510	47.21	14.85	10	90
	Male	121,320	44.88	14.70	10	90
Ethnicity	American Indian/Alaska Native	458	44.37	15.69	10	90
	Asian	18,019	56.72	14.43	10	90
	Black or African American	39,713	39.25	13.53	10	90
	Hispanic/Latino	68,390	41.47	13.78	10	90
	Native Hawaiian or Pacific Islander	416	49.50	14.04	10	90
	Two or More Races	8,968	48.17	14.64	10	90
	White	101,226	49.66	13.80	10	90
Economic Status*	Not Economically Disadvantaged	130,148	50.57	14.26	10	90
	Economically Disadvantaged	98,485	39.64	13.15	10	90
English Learner Status	Non-English Learner	196,956	47.71	14.52	10	90
	English Learner	30,452	34.51	11.00	10	87
Disabilities	Students without Disabilities	191,021	48.11	14.30	10	90
	Student with Disability (SWD)	44,831	37.42	13.89	10	90
Writing Summative Score		237,847	28.75	13.00	10	60
Gender	Female	116,510	31.06	12.61	10	60
	Male	121,320	26.53	12.99	10	60
Ethnicity	American Indian/Alaska Native	458	26.75	13.17	10	60

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Asian	18,019	37.13	10.90	10	60
	Black or African American	39,713	22.97	12.54	10	60
	Hispanic/Latino	68,390	26.00	12.82	10	60
	Native Hawaiian or Pacific Islander	416	33.36	11.67	10	56
	Two or More Races	8,968	29.64	12.79	10	60
	White	101,226	31.30	12.20	10	60
Economic Status*	Not Economically Disadvantaged	130,148	32.23	12.21	10	60
	Economically Disadvantaged	98,485	23.88	12.57	10	60
English Learner Status	Non-English Learner	196,956	29.93	12.80	10	60
	English Learner	30,452	20.76	11.56	10	60
Disabilities	Students without Disabilities	191,021	30.69	12.41	10	60
	Student with Disability (SWD)	44,831	20.76	12.34	10	60

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.43 Subgroup Performance for ELA/L Scale Scores: Grade 6

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		238,260	737.27	33.89	650	850
Gender	Female	116,268	742.46	33.71	650	850
	Male	121,948	732.31	33.31	650	850
Ethnicity	American Indian/Alaska Native	455	733.10	33.17	650	850
	Asian	18,040	762.47	32.97	650	850
	Black or African American	40,148	721.52	31.25	650	850
	Hispanic/Latino	68,037	727.95	31.94	650	850
	Native Hawaiian or Pacific Islander	376	745.19	32.45	650	830
	Two or More Races	8,758	740.82	33.85	650	850
	White	101,842	744.93	31.41	650	850
Economic Status*	Not Economically Disadvantaged	131,693	747.41	32.40	650	850
	Economically Disadvantaged	98,297	723.06	30.87	650	850
English Learner Status	Non-English Learner	204,279	740.75	33.03	650	850
	English Learner	24,413	706.93	25.38	650	816
Disabilities	Students without Disabilities	192,068	742.77	32.07	650	850
	Student with Disability (SWD)	44,671	714.14	31.65	650	850
Reading Summative Score		238,260	46.11	13.42	10	90
Gender	Female	116,268	47.58	13.25	10	90
	Male	121,948	44.70	13.42	10	90
Ethnicity	American Indian/Alaska Native	455	44.15	12.77	10	90
	Asian	18,040	55.50	13.14	10	90
	Black or African American	40,148	40.16	12.42	10	90
	Hispanic/Latino	68,037	42.28	12.50	10	90
	Native Hawaiian or Pacific Islander	376	48.62	12.92	10	90
	Two or More Races	8,758	47.97	13.59	10	90
	White	101,842	49.19	12.56	10	90
Economic Status*	Not Economically Disadvantaged	131,693	50.04	12.88	10	90
	Economically Disadvantaged	98,297	40.58	12.17	10	90
English Learner Status	Non-English Learner	204,279	47.47	13.07	10	90
	English Learner	24,413	34.11	9.89	10	82
Disabilities	Students without Disabilities	192,068	48.19	12.70	10	90
	Student with Disability (SWD)	44,671	37.36	12.83	10	90
Writing Summative Score		238,260	28.04	12.49	10	60
Gender	Female	116,268	30.39	12.00	10	60
	Male	121,948	25.80	12.53	10	60
Ethnicity	American Indian/Alaska Native	455	26.97	12.70	10	60

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Asian	18,040	36.14	10.44	10	60
	Black or African American	40,148	22.63	12.20	10	60
	Hispanic/Latino	68,037	25.32	12.40	10	60
	Native Hawaiian or Pacific Islander	376	31.39	11.39	10	52
	Two or More Races	8,758	28.73	12.31	10	60
	White	101,842	30.50	11.60	10	60
Economic Status*	Not Economically Disadvantaged	131,693	31.34	11.64	10	60
	Economically Disadvantaged	98,297	23.41	12.22	10	60
English Learner Status	Non-English Learner	204,279	29.14	12.22	10	60
	English Learner	24,413	18.39	10.72	10	60
Disabilities	Students without Disabilities	192,068	29.94	11.88	10	60
	Student with Disability (SWD)	44,671	20.10	11.86	10	60

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.44 Subgroup Performance for ELA/L Scale Scores: Grade 7

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		245,105	739.51	40.53	650	850
Gender	Female	119,492	746.04	39.93	650	850
	Male	125,555	733.29	40.12	650	850
Ethnicity	American Indian/Alaska Native	476	736.54	40.25	650	850
	Asian	18,180	772.06	38.79	650	850
	Black or African American	40,843	721.19	37.47	650	850
	Hispanic/Latino	71,221	728.05	38.58	650	850
	Native Hawaiian or Pacific Islander	437	752.24	36.52	650	850
	Two or More Races	8,555	741.72	40.34	650	850
	White	104,828	748.54	37.14	650	850
Economic Status*	Not Economically Disadvantaged	136,806	751.18	38.71	650	850
	Economically Disadvantaged	100,388	722.81	37.19	650	850
English Learner Status	Non-English Learner	212,988	743.54	39.25	650	850
	English Learner	22,874	699.54	29.75	650	839
Disabilities	Students without Disabilities	198,003	745.94	38.42	650	850
	Student with Disability (SWD)	45,668	712.22	37.97	650	850
Reading Summative Score		245,105	47.11	16.36	10	90
Gender	Female	119,492	48.84	15.98	10	90
	Male	125,555	45.47	16.55	10	90
Ethnicity	American Indian/Alaska Native	476	45.86	16.26	10	90
	Asian	18,180	59.31	15.77	10	90
	Black or African American	40,843	40.07	15.04	10	90
	Hispanic/Latino	71,221	42.21	15.41	10	90
	Native Hawaiian or Pacific Islander	437	51.11	14.73	10	90
	Two or More Races	8,555	48.74	16.53	10	90
	White	104,828	50.91	15.19	10	90
Economic Status*	Not Economically Disadvantaged	136,806	51.72	15.71	10	90
	Economically Disadvantaged	100,388	40.49	14.92	10	90
English Learner Status	Non-English Learner	212,988	48.73	15.84	10	90
	English Learner	22,874	30.99	11.68	10	90
Disabilities	Students without Disabilities	198,003	49.56	15.55	10	90
	Student with Disability (SWD)	45,668	36.70	15.64	10	90
Writing Summative Score		245,105	29.79	12.63	10	60
Gender	Female	119,492	32.34	12.09	10	60
	Male	125,555	27.37	12.66	10	60
Ethnicity	American Indian/Alaska Native	476	29.10	12.53	10	60

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Asian	18,180	38.84	10.80	10	60
	Black or African American	40,843	24.49	12.35	10	60
	Hispanic/Latino	71,221	27.09	12.43	10	60
	Native Hawaiian or Pacific Islander	437	34.34	10.73	10	60
	Two or More Races	8,555	29.75	12.56	10	60
	White	104,828	32.11	11.74	10	60
Economic Status*	Not Economically Disadvantaged	136,806	33.00	11.86	10	60
	Economically Disadvantaged	100,388	25.20	12.34	10	60
English Learner Status	Non-English Learner	212,988	30.85	12.34	10	60
	English Learner	22,874	19.21	10.63	10	60
Disabilities	Students without Disabilities	198,003	31.71	11.94	10	60
	Student with Disability (SWD)	45,668	21.65	12.24	10	60

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.45 Subgroup Performance for ELA/L Scale Scores: Grade 8

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		249,387	737.91	41.92	650	850
Gender	Female	122,108	745.39	41.38	650	850
	Male	127,201	730.72	41.18	650	850
Ethnicity	American Indian/Alaska Native	583	729.82	42.10	650	850
	Asian	18,478	772.51	40.34	650	850
	Black or African American	41,934	719.37	38.38	650	850
	Hispanic/Latino	72,324	727.04	39.44	650	850
	Native Hawaiian or Pacific Islander	392	751.97	42.22	650	850
	Two or More Races	8,334	740.17	41.46	650	850
	White	106,854	746.37	39.20	650	850
Economic Status*	Not Economically Disadvantaged	140,600	749.39	40.61	650	850
	Economically Disadvantaged	101,447	721.44	38.29	650	850
English Learner Status	Non-English Learner	218,320	741.77	40.98	650	850
	English Learner	22,246	698.81	30.03	650	845
Disabilities	Students without Disabilities	202,285	744.13	40.29	650	850
	Student with Disability (SWD)	45,733	710.95	38.12	650	850
Reading Summative Score		249,387	46.83	16.69	10	90
Gender	Female	122,108	48.94	16.44	10	90
	Male	127,201	44.81	16.68	10	90
Ethnicity	American Indian/Alaska Native	583	43.80	16.49	10	90
	Asian	18,478	59.84	16.16	10	90
	Black or African American	41,934	39.94	15.40	10	90
	Hispanic/Latino	72,324	42.33	15.57	10	90
	Native Hawaiian or Pacific Islander	392	51.18	17.00	10	90
	Two or More Races	8,334	48.36	16.66	10	90
	White	106,854	50.20	15.78	10	90
Economic Status*	Not Economically Disadvantaged	140,600	51.26	16.23	10	90
	Economically Disadvantaged	101,447	40.44	15.24	10	90
English Learner Status	Non-English Learner	218,320	48.34	16.31	10	90
	English Learner	22,246	31.43	11.95	10	87
Disabilities	Students without Disabilities	202,285	49.21	16.05	10	90
	Student with Disability (SWD)	45,733	36.48	15.47	10	90
Writing Summative Score		249,387	29.34	12.91	10	60
Gender	Female	122,108	32.08	12.36	10	60
	Male	127,201	26.72	12.87	10	60
Ethnicity	American Indian/Alaska Native	583	27.20	12.99	10	60

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Asian	18,478	38.80	11.29	10	60
	Black or African American	41,934	23.90	12.43	10	60
	Hispanic/Latino	72,324	26.65	12.63	10	60
	Native Hawaiian or Pacific Islander	392	33.95	12.49	10	60
	Two or More Races	8,334	29.40	12.84	10	60
	White	106,854	31.66	12.06	10	60
Economic Status*	Not Economically Disadvantaged	140,600	32.53	12.27	10	60
	Economically Disadvantaged	101,447	24.79	12.48	10	60
English Learner Status	Non-English Learner	218,320	30.39	12.66	10	60
	English Learner	22,246	18.69	10.55	10	60
Disabilities	Students without Disabilities	202,285	31.16	12.37	10	60
	Student with Disability (SWD)	45,733	21.51	12.26	10	60

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.46 Subgroup Performance for ELA/L Scale Scores: Grade 9

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		48,052	746.44	38.52	650	850
Gender	Female	23,441	753.29	37.52	650	850
	Male	24,536	739.85	38.32	650	850
Ethnicity	American Indian/Alaska Native	60	743.88	39.09	650	847
	Asian	4,827	777.27	35.34	650	850
	Black or African American	7,393	730.58	34.77	650	847
	Hispanic/Latino	15,306	735.37	37.25	650	850
	Native Hawaiian or Pacific Islander	97	756.93	38.45	650	847
	Two or More Races	1,125	751.51	37.26	650	850
	White	19,226	753.29	35.33	650	850
Economic Status*	Not Economically Disadvantaged	33,144	752.70	37.66	650	850
	Economically Disadvantaged	14,907	732.53	36.73	650	850
English Learner Status	Non-English Learner	45,847	748.87	37.13	650	850
	English Learner	2,201	695.89	31.36	650	821
Disabilities	Students without Disabilities	38,540	752.77	36.56	650	850
	Student with Disability (SWD)	9,512	720.83	35.50	650	850
Reading Summative Score		48,052	49.33	16.00	10	90
Gender	Female	23,441	50.94	15.57	10	90
	Male	24,536	47.76	16.26	10	90
Ethnicity	American Indian/Alaska Native	60	48.15	16.07	10	82
	Asian	4,827	60.99	15.22	10	90
	Black or African American	7,393	43.25	14.37	10	90
	Hispanic/Latino	15,306	44.30	15.02	10	90
	Native Hawaiian or Pacific Islander	97	52.64	15.66	11	82
	Two or More Races	1,125	52.35	15.64	10	90
	White	19,226	52.55	15.07	10	90
Economic Status*	Not Economically Disadvantaged	33,144	52.07	15.80	10	90
	Economically Disadvantaged	14,907	43.22	14.71	10	90
English Learner Status	Non-English Learner	45,847	50.30	15.52	10	90
	English Learner	2,201	29.12	12.13	10	72
Disabilities	Students without Disabilities	38,540	51.58	15.44	10	90
	Student with Disability (SWD)	9,512	40.21	14.99	10	90
Writing Summative Score		48,052	32.37	11.60	10	60
Gender	Female	23,441	34.92	10.79	10	60
	Male	24,536	29.92	11.82	10	60
Ethnicity	American Indian/Alaska Native	60	31.32	12.96	10	60

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Asian	4,827	40.42	9.23	10	60
	Black or African American	7,393	27.92	11.65	10	60
	Hispanic/Latino	15,306	29.90	11.86	10	60
	Native Hawaiian or Pacific Islander	97	35.57	11.13	10	60
	Two or More Races	1,125	32.92	11.30	10	60
	White	19,226	33.98	10.51	10	60
Economic Status*	Not Economically Disadvantaged	33,144	33.90	11.06	10	60
	Economically Disadvantaged	14,907	28.97	12.04	10	60
English Learner Status	Non-English Learner	45,847	33.03	11.23	10	60
	English Learner	2,201	18.70	10.81	10	60
Disabilities	Students without Disabilities	38,540	34.35	10.62	10	60
	Student with Disability (SWD)	9,512	24.35	11.97	10	60

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.47 Subgroup Performance for ELA/L Scale Scores: Grade 10

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		7,069	740.86	49.34	650	850
Gender	Female	3,446	748.54	47.92	650	850
	Male	3,617	733.48	49.52	650	850
Ethnicity	American Indian/Alaska Native	n/r	n/r	n/r	n/r	n/r
	Asian	283	764.70	46.52	650	850
	Black or African American	3,016	723.24	46.41	650	850
	Hispanic/Latino	1,472	739.63	46.64	650	850
	Native Hawaiian or Pacific Islander	54	756.54	43.47	650	850
	Two or More Races	477	752.41	44.60	650	850
	White	1,590	765.13	45.09	650	850
Economic Status*	Not Economically Disadvantaged	n/r	n/r	n/r	n/r	n/r
	Economically Disadvantaged	2,162	714.49	45.39	650	850
English Learner Status	Non-English Learner	n/r	n/r	n/r	n/r	n/r
	English Learner	468	702.08	39.55	650	819
Disabilities	Students without Disabilities	5,456	749.49	47.48	650	850
	Student with Disability (SWD)	1,324	709.22	44.94	650	850
Reading Summative Score		7,069	48.52	20.41	10	90
Gender	Female	3,446	50.67	19.66	10	90
	Male	3,617	46.44	20.88	10	90
Ethnicity	American Indian/Alaska Native	n/r	n/r	n/r	n/r	n/r
	Asian	283	58.36	19.56	10	90
	Black or African American	3,016	40.89	18.85	10	90
	Hispanic/Latino	1,472	47.50	19.31	10	90
	Native Hawaiian or Pacific Islander	54	53.43	17.58	16	90
	Two or More Races	477	54.47	18.75	10	90
	White	1,590	59.29	18.40	10	90
Economic Status*	Not Economically Disadvantaged	n/r	n/r	n/r	n/r	n/r
	Economically Disadvantaged	2,162	37.16	18.15	10	90
English Learner Status	Non-English Learner	n/r	n/r	n/r	n/r	n/r
	English Learner	468	31.72	15.52	10	78
Disabilities	Students without Disabilities	5,456	51.91	19.62	10	90
	Student with Disability (SWD)	1,324	36.09	19.21	10	90
Writing Summative Score		7,069	30.42	13.09	10	60
Gender	Female	3,446	32.86	12.68	10	60
	Male	3,617	28.07	13.04	10	60
Ethnicity	American Indian/Alaska Native	n/r	n/r	n/r	n/r	n/r

Group Type	Subgroup	N	Mean	SD	Min.	Max.
	Asian	283	36.08	11.42	10	60
	Black or African American	3,016	26.63	12.85	10	60
	Hispanic/Latino	1,472	30.53	12.58	10	60
	Native Hawaiian or Pacific Islander	54	35.06	11.33	10	60
	Two or More Races	477	32.35	12.09	10	60
	White	1,590	35.32	12.34	10	60
Economic Status*	Not Economically Disadvantaged	n/r	n/r	n/r	n/r	n/r
	Economically Disadvantaged	2,162	24.74	12.70	10	60
English Learner Status	Non-English Learner	n/r	n/r	n/r	n/r	n/r
	English Learner	468	22.57	11.57	10	51
Disabilities	Students without Disabilities	5,456	32.53	12.59	10	60
	Student with Disability (SWD)	1,324	22.56	12.40	10	60

Note. SD = standard deviation; n/r = not reported due to n<20.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.48 Subgroup Performance for Mathematics Scale Scores: Grade 3

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		230,053	737.99	39.38	650	850
Gender	Female	112,938	736.26	37.99	650	850
	Male	117,109	739.67	40.61	650	850
Ethnicity	American Indian/Alaska Native	508	738.03	40.49	650	850
	Asian	17,645	771.70	37.65	650	850
	Black or African American	38,055	714.02	34.27	650	850
	Hispanic/Latino	66,053	724.77	34.31	650	850
	Native Hawaiian or Pacific Islander	359	746.36	38.97	650	850
	Two or More Races	9,537	743.19	40.24	650	850
	White	97,179	749.69	35.75	650	850
Economic Status*	Not Economically Disadvantaged	124,021	751.91	37.40	650	850
	Economically Disadvantaged	96,292	719.20	34.00	650	850
English Learner Status	Non-English Learner	182,726	741.86	39.36	650	850
	English Learner	36,588	718.55	33.08	650	850
Disabilities	Students without Disabilities	189,138	742.03	38.44	650	850
	Student with Disability (SWD)	38,894	719.38	38.40	650	850
Language Form	Spanish	5,063	708.73	29.38	650	850

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.49 Subgroup Performance for Mathematics Scale Scores: Grade 4

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		232,955	732.80	35.41	650	850
Gender	Female	114,345	731.52	34.55	650	850
	Male	118,604	734.03	36.18	650	850
Ethnicity	American Indian/Alaska Native	474	729.05	37.12	652	850
	Asian	18,144	764.66	34.52	650	850
	Black or African American	38,800	711.16	29.75	650	850
	Hispanic/Latino	67,133	720.30	30.45	650	850
	Native Hawaiian or Pacific Islander	402	742.33	33.23	662	836
	Two or More Races	9,283	738.17	35.99	650	850
	White	98,031	743.51	32.15	650	850
Economic Status*	Not Economically Disadvantaged	126,394	745.61	33.82	650	850
	Economically Disadvantaged	97,211	715.18	29.68	650	850
English Learner Status	Non-English Learner	185,990	736.55	35.40	650	850
	English Learner	36,354	712.82	28.01	650	836
Disabilities	Students without Disabilities	188,400	736.88	34.68	650	850
	Student with Disability (SWD)	42,583	715.59	33.35	650	850
Language Form	Spanish	4,605	701.37	25.49	650	800

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.50 Subgroup Performance for Mathematics Scale Scores: Grade 5

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		238,545	730.00	35.36	650	850
Gender	Female	116,876	729.42	33.87	650	850
	Male	121,651	730.55	36.72	650	850
Ethnicity	American Indian/Alaska Native	459	726.50	36.25	650	848
	Asian	18,101	763.78	34.60	650	850
	Black or African American	39,488	708.46	29.55	650	850
	Hispanic/Latino	69,287	717.78	30.19	650	850
	Native Hawaiian or Pacific Islander	413	738.15	33.74	650	848
	Two or More Races	8,943	734.82	36.19	650	850
	White	101,199	740.29	32.21	650	850
Economic Status*	Not Economically Disadvantaged	130,744	742.27	34.10	650	850
	Economically Disadvantaged	98,587	712.71	29.53	650	850
English Learner Status	Non-English Learner	196,416	733.48	35.20	650	850
	English Learner	31,776	707.00	26.52	650	848
Disabilities	Students without Disabilities	191,904	734.54	34.37	650	850
	Student with Disability (SWD)	44,664	711.29	33.35	650	850
Language Form	Spanish	4,220	700.54	26.60	650	808

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.51 Subgroup Performance for Mathematics Scale Scores: Grade 6

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		238,630	727.54	32.92	650	850
Gender	Female	116,445	727.17	32.01	650	850
	Male	122,140	727.90	33.76	650	850
Ethnicity	American Indian/Alaska Native	453	724.07	32.63	650	850
	Asian	18,103	758.16	33.36	650	850
	Black or African American	39,760	707.63	28.01	650	850
	Hispanic/Latino	68,908	716.85	28.51	650	850
	Native Hawaiian or Pacific Islander	377	733.89	30.43	650	800
	Two or More Races	8,707	731.56	33.06	650	850
	White	101,730	736.78	29.87	650	850
Economic Status*	Not Economically Disadvantaged	132,235	738.50	31.78	650	850
	Economically Disadvantaged	98,162	712.22	28.11	650	850
English Learner Status	Non-English Learner	203,440	730.72	32.47	650	850
	English Learner	25,760	702.09	24.08	650	839
Disabilities	Students without Disabilities	192,687	731.93	32.06	650	850
	Student with Disability (SWD)	44,443	708.98	29.97	650	850
Language Form	Spanish	3,721	700.92	24.03	650	797

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.52 Subgroup Performance for Mathematics Scale Scores: Grade 7

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		235,692	733.63	29.03	650	850
Gender	Female	115,051	733.05	27.75	650	850
	Male	120,587	734.18	30.20	650	850
Ethnicity	American Indian/Alaska Native	452	731.72	27.93	650	812
	Asian	15,592	759.84	29.99	650	850
	Black or African American	39,850	717.81	24.43	650	850
	Hispanic/Latino	70,316	725.11	25.46	650	850
	Native Hawaiian or Pacific Islander	340	740.81	28.11	663	850
	Two or More Races	7,681	735.08	30.44	650	850
	White	101,026	741.63	27.21	650	850
Economic Status*	Not Economically Disadvantaged	132,518	742.53	28.52	650	850
	Economically Disadvantaged	99,771	721.79	25.20	650	850
English Learner Status	Non-English Learner	206,894	736.33	28.69	650	850
	English Learner	23,753	711.25	21.19	650	850
Disabilities	Students without Disabilities	189,826	737.82	27.84	650	850
	Student with Disability (SWD)	44,452	716.11	27.46	650	850
Language Form	Spanish	3,226	709.76	20.23	650	806

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.53 Subgroup Performance for Mathematics Scale Scores: Grade 8

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		211,563	720.42	37.50	650	850
Gender	Female	102,954	720.51	36.30	650	850
	Male	108,548	720.34	38.61	650	850
Ethnicity	American Indian/Alaska Native	530	715.19	34.84	650	850
	Asian	10,889	753.59	44.23	650	850
	Black or African American	37,791	701.28	29.55	650	850
	Hispanic/Latino	65,025	711.58	32.29	650	850
	Native Hawaiian or Pacific Islander	297	726.94	37.34	650	850
	Two or More Races	7,121	723.12	39.41	650	850
	White	89,514	730.70	36.96	650	850
Economic Status*	Not Economically Disadvantaged	112,639	730.83	38.50	650	850
	Economically Disadvantaged	93,626	707.88	32.22	650	850
English Learner Status	Non-English Learner	182,466	723.51	37.76	650	850
	English Learner	22,612	696.84	25.57	650	850
Disabilities	Students without Disabilities	168,062	725.22	37.26	650	850
	Student with Disability (SWD)	42,161	701.70	32.23	650	850
Language Form	Spanish	2,844	696.27	24.10	650	844

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.54 Subgroup Performance for Mathematics Scale Scores: Algebra I

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		115,174	734.83	36.12	650	850
Gender	Female	56,085	735.00	34.88	650	850
	Male	58,986	734.66	37.27	650	850
Ethnicity	American Indian/Alaska Native	167	734.59	37.89	650	850
	Asian	11,431	767.76	34.96	650	850
	Black or African American	19,311	716.48	30.43	650	850
	Hispanic/Latino	37,136	721.06	31.33	650	850
	Native Hawaiian or Pacific Islander	269	745.09	34.38	650	850
	Two or More Races	3,019	743.77	35.26	650	850
	White	43,602	745.36	32.61	650	850
Economic Status*	Not Economically Disadvantaged	74,152	742.29	36.10	650	850
	Economically Disadvantaged	34,995	718.79	31.10	650	850
English Learner Status	Non-English Learner	99,804	737.72	35.72	650	850
	English Learner	7,893	702.97	25.80	650	850
Disabilities	Students without Disabilities	92,916	739.65	35.44	650	850
	Student with Disability (SWD)	21,857	714.74	31.81	650	850
Language Form	Spanish	2,563	698.87	22.15	650	786

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.55 Subgroup Performance for Mathematics Scale Scores: Geometry

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		40,403	740.79	27.63	650	850
Gender	Female	20,199	740.09	26.49	650	850
	Male	20,166	741.48	28.73	650	838
Ethnicity	American Indian/Alaska Native	60	733.93	30.47	662	807
	Asian	6,628	761.18	23.32	650	850
	Black or African American	6,045	718.89	24.02	650	838
	Hispanic/Latino	8,644	726.09	25.03	650	815
	Native Hawaiian or Pacific Islander	135	738.96	25.56	677	813
	Two or More Races	1,430	745.47	24.78	650	828
	White	17,306	747.63	22.62	650	850
Economic Status*	Not Economically Disadvantaged	27,208	748.10	25.49	650	850
	Economically Disadvantaged	8,194	721.14	25.34	650	832
English Learner Status	Non-English Learner	32,412	744.80	26.55	650	850
	English Learner	1,400	709.77	21.69	650	802
Disabilities	Students without Disabilities	35,482	743.21	26.62	650	850
	Student with Disability (SWD)	4,648	723.71	28.98	650	821
Language Form	Spanish	425	703.60	17.77	650	757

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Table A.12.56 Subgroup Performance for Mathematics Scale Scores: Algebra II

Group Type	Subgroup	N	Mean	SD	Min.	Max.
Full Summative Score		13,965	744.87	42.50	650	850
Gender	Female	6,903	740.74	40.18	650	850
	Male	7,046	748.89	44.30	650	850
Ethnicity	American Indian/Alaska Native	29	743.66	41.51	665	824
	Asian	3,532	775.74	35.85	650	850
	Black or African American	1,049	714.76	35.24	650	836
	Hispanic/Latino	2,997	713.81	36.49	650	847
	Native Hawaiian or Pacific Islander	89	726.35	40.64	650	850
	Two or More Races	671	739.67	36.08	650	850
	White	5,523	748.83	35.26	650	850
Economic Status*	Not Economically Disadvantaged	8,987	758.71	39.05	650	850
	Economically Disadvantaged	2,111	710.17	36.41	650	850
English Learner Status	Non-English Learner	10,785	751.14	42.14	650	850
	English Learner	424	694.75	30.79	650	812
Disabilities	Students without Disabilities	12,531	747.98	41.45	650	850
	Student with Disability (SWD)	1,433	717.69	41.94	650	850
Language Form	Spanish	138	682.98	21.55	650	750

Note. SD = standard deviation.

*Economic status was based on participation in National School Lunch Program (NSLP): receipt of free or reduced-price lunch (FRL).

Appendix 13.1: Reliability by Content and Grade/Subject

Table A.13.1 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 3

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	54	3.69	0.89	603	0.81	117,193	0.90
Gender							
Male	54	3.52	0.88	348	0.79	36,316	0.88
Female	54	3.66	0.89	739	0.79	36,154	0.89
Ethnicity							
Black/African American	54	3.47	0.87	176	0.76	18,749	0.89
Asian/Pacific Islander	54	3.98	0.87	8,053	0.85	9,349	0.87
Hispanic/Latino	54	3.57	0.88	238	0.76	32,478	0.89
American Indian/Alaska Native	54	3.79	0.89	265	0.89	206	0.90
Multiple	54	3.71	0.89	271	0.84	4,879	0.90
White	54	3.78	0.87	153	0.82	51,277	0.88
Special Instruction Needs							
Economically Disadvantaged	54	3.50	0.87	401	0.75	47,431	0.88
Not Economically Disadvantaged	54	3.82	0.87	166	0.81	64,135	0.88
English Learner	54	3.40	0.85	157	0.75	18,143	0.86
Non-English Learner	54	3.74	0.88	412	0.81	94,294	0.89
Students with Disabilities	54	3.31	0.88	540	0.80	15,577	0.90
Students without Disabilities	54	3.77	0.89	82,449	0.88	100,768	0.89
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	54	2.90	0.83	7,524	0.83	7,524	0.83

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.2 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 4

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	70	4.32	0.89	879	0.84	114,919	0.91
Gender							
Male	70	4.15	0.88	1,619	0.80	31,811	0.90
Female	70	4.29	0.89	164	0.79	31,249	0.90
Ethnicity							
Black/African American	70	4.06	0.87	207	0.73	18,548	0.90
Asian/Pacific Islander	70	4.54	0.88	8,361	0.86	9,682	0.89
Hispanic/Latino	70	4.21	0.88	106	0.80	32,630	0.89
American Indian/Alaska Native	70	4.20	0.90	206	0.89	234	0.91
Multiple	70	4.37	0.90	260	0.85	4,661	0.91
White	70	4.41	0.87	347	0.83	48,944	0.89
Special Instruction Needs							
Economically Disadvantaged	70	4.11	0.87	4,391	0.79	45,418	0.89
Not Economically Disadvantaged	70	4.44	0.88	349	0.85	63,652	0.89
English Learner	70	4.01	0.84	1,347	0.75	16,970	0.86
Non-English Learner	70	4.36	0.89	648	0.85	92,935	0.90
Students with Disabilities	70	3.86	0.89	784	0.84	16,588	0.91
Students without Disabilities	70	4.41	0.89	86,821	0.87	97,585	0.90
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	70	4.05	0.91	185	0.91	185	0.91
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	70	3.41	0.84	8,827	0.84	8,827	0.84

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.3 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 5

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	73	4.37	0.90	9,534	0.84	129,428	0.90
Gender							
Male	73	4.12	0.88	1,560	0.76	38,833	0.89
Female	73	4.35	0.89	964	0.79	38,216	0.89
Ethnicity							
Black/African American	73	4.07	0.88	2,024	0.77	20,641	0.89
Asian/Pacific Islander	73	4.60	0.88	253	0.87	10,401	0.89
Hispanic/Latino	73	4.27	0.88	3,358	0.81	36,956	0.89
American Indian/Alaska Native	73	4.21	0.91	177	0.91	249	0.91
Multiple	73	4.42	0.89	258	0.85	4,931	0.90
White	73	4.48	0.88	141	0.83	56,042	0.89
Special Instruction Needs							
Economically Disadvantaged	73	4.14	0.87	4,534	0.77	51,589	0.88
Not Economically Disadvantaged	73	4.51	0.89	358	0.85	72,121	0.89
English Learner	73	3.89	0.81	1,300	0.72	16,655	0.82
Non-English Learner	73	4.43	0.89	7,904	0.84	107,857	0.90
Students with Disabilities	73	3.90	0.88	817	0.83	19,259	0.90
Students without Disabilities	73	4.47	0.89	115	0.86	109,292	0.90
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	73	4.40	0.91	163	0.91	163	0.91
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	73	3.61	0.83	9,371	0.83	9,371	0.83

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.4 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 6

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	73	4.58	0.90	684	0.85	118,609	0.92
Gender							
Male	73	4.22	0.90	171	0.83	36,987	0.91
Female	73	4.54	0.90	139	0.85	108	0.92
Ethnicity							
Black/African American	73	4.27	0.89	172	0.79	19,096	0.90
Asian/Pacific Islander	73	4.84	0.89	8,536	0.87	9,456	0.90
Hispanic/Latino	73	4.43	0.89	217	0.79	33,239	0.91
American Indian/Alaska Native	73	4.52	0.90	199	0.89	225	0.91
Multiple	73	4.65	0.90	3,801	0.89	4,500	0.91
White	73	4.68	0.89	44,795	0.87	136	0.91
Special Instruction Needs							
Economically Disadvantaged	73	4.31	0.89	380	0.81	47,390	0.90
Not Economically Disadvantaged	73	4.74	0.89	269	0.86	66,460	0.91
English Learner	73	3.81	0.83	129	0.73	12,154	0.84
Non-English Learner	73	4.65	0.90	525	0.85	102,314	0.91
Students with Disabilities	73	3.98	0.90	596	0.84	17,573	0.91
Students without Disabilities	73	4.70	0.90	88,019	0.88	100,500	0.91
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	73	4.20	0.93	226	0.93	226	0.93
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	73	3.56	0.86	9,497	0.86	9,497	0.86

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.5 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 7

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	71	4.51	0.91	704	0.82	114,605	0.91
Gender							
Male	71	4.22	0.89	1,638	0.82	35,595	0.90

Table A.13.5 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 7

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Female	71	4.48	0.89	191	0.81	34,776	0.90
Ethnicity							
Black/African American	71	4.27	0.89	2,029	0.81	18,276	0.90
Asian/Pacific Islander	71	4.71	0.90	9,337	0.89	8,865	0.91
Hispanic/Latino	71	4.39	0.90	168	0.80	32,971	0.90
American Indian/Alaska Native	71	4.56	0.91	232	0.91	214	0.91
Multiple	71	4.55	0.90	261	0.81	4,152	0.91
White	71	4.61	0.89	350	0.81	49,979	0.90
Special Instruction Needs							
Economically Disadvantaged	71	4.29	0.89	390	0.80	45,657	0.90
Not Economically Disadvantaged	71	4.64	0.90	287	0.82	64,711	0.91
English Learner	71	3.78	0.81	1,119	0.72	10,689	0.81
Non-English Learner	71	4.57	0.90	593	0.82	100,186	0.91
Students with Disabilities	71	4.00	0.89	561	0.78	17,015	0.91
Students without Disabilities	71	4.62	0.90	141	0.80	97,126	0.91
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	71	3.97	0.93	241	0.93	241	0.93
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	71	3.55	0.85	9,151	0.85	9,151	0.85

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.6 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 8

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	71	4.63	0.90	705	0.84	116,341	0.92
Gender							
Male	71	4.27	0.89	244	0.81	38,003	0.91
Female	71	4.59	0.89	1,007	0.85	37,509	0.91
Ethnicity							
Black/African American	71	4.39	0.89	158	0.73	18,409	0.91
Asian/Pacific Islander	71	4.80	0.89	9,609	0.86	229	0.93
Hispanic/Latino	71	4.52	0.89	196	0.77	33,353	0.91
American Indian/Alaska Native	71	4.47	0.91	263	0.90	282	0.92
Multiple	71	4.70	0.90	3,762	0.88	4,091	0.92
White	71	4.72	0.89	318	0.84	51,276	0.92
Special Instruction Needs							
Economically Disadvantaged	71	4.41	0.89	367	0.80	46,358	0.91
Not Economically Disadvantaged	71	4.75	0.89	293	0.85	66,076	0.92
English Learner	71	3.85	0.81	111	0.51	10,630	0.84
Non-English Learner	71	4.69	0.90	553	0.84	102,188	0.92
Students with Disabilities	71	4.12	0.89	546	0.81	17,213	0.92
Students without Disabilities	71	4.73	0.90	157	0.81	98,682	0.92
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	71	4.02	0.94	287	0.94	287	0.94
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	71	3.71	0.88	8,871	0.88	8,871	0.88

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.7 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 9

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	71	5.08	0.87	4,554	0.81	344	0.91
Gender							
Male	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Female	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Ethnicity							
Black/African American	71	4.96	0.85	966	0.78	243	0.88
Asian/Pacific Islander	71	5.04	0.85	130	0.82	10,372	0.85
Hispanic/Latino	71	4.99	0.86	1,614	0.80	29,275	0.87
American Indian/Alaska Native	71	5.25	0.86	140	0.86	140	0.86
Multiple	71	5.27	0.86	2,155	0.86	2,155	0.86
White	71	5.12	0.85	1,748	0.82	38,037	0.85
Special Instruction Needs							
Economically Disadvantaged	71	4.97	0.86	1,668	0.80	26,216	0.86
Not Economically Disadvantaged	71	5.12	0.87	2,688	0.81	65,553	0.87
English Learner	71	4.10	0.81	154	0.71	4,371	0.82
Non-English Learner	71	5.13	0.86	4,242	0.81	87,813	0.86
Students with Disabilities	71	4.63	0.85	4,554	0.81	103	0.89
Students without Disabilities	71	5.17	0.86	79,936	0.86	239	0.90
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	71	4.17	0.81	4,502	0.81	4,502	0.81

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.8 Summary of Test Reliability Estimates for Subgroups: ELA/L Grade 10

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum N	Reliability Alpha	Maximum N	Reliability Alpha
Total Group	73	4.98	0.88	167	0.79	3,188	0.90
Gender							
Male	73	5.14	0.85	758	0.83	725	0.87
Female	73	5.25	0.84	677	0.82	695	0.86
Ethnicity							
Black/African American	73	4.81	0.86	136	0.77	1,257	0.88
Asian/Pacific Islander	73	5.24	0.86	165	0.84	161	0.88
Hispanic/Latino	73	4.95	0.87	698	0.86	659	0.88
American Indian/Alaska Native	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Multiple	73	5.24	0.85	230	0.84	223	0.85
White	73	5.11	0.87	731	0.84	805	0.90
Special Instruction Needs							
Economically Disadvantaged	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Not Economically Disadvantaged	n/a	n/a	n/a	n/a	n/a	n/a	n/a
English Learner	73	4.46	0.82	206	0.81	181	0.84
Non-English Learner	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Students with Disabilities	73	4.55	0.86	167	0.79	510	0.89
Students without Disabilities	73	5.08	0.87	2,509	0.86	2,533	0.89
Students Taking Accommodated Forms							
ASL	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Closed-Caption	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	73	3.68	0.79	166	0.79	166	0.79

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.9 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 3

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum N	Reliability Alpha	Maximum N	Reliability Alpha
Total Group	52	3.13	0.93	1,043	0.91	78,760	0.93
Gender							
Male	52	3.11	0.93	388	0.90	22,171	0.93
Female	52	3.12	0.92	271	0.91	155	0.93
Ethnicity							
Black/African American	52	2.88	0.91	210	0.89	10,627	0.92
Asian/Pacific Islander	52	3.15	0.91	7,001	0.91	1,764	0.93
Hispanic/Latino	52	3.02	0.91	289	0.89	18,729	0.92
American Indian/Alaska Native	52	3.18	0.93	154	0.93	166	0.93
Multiple	52	3.16	0.93	3,447	0.93	955	0.93
White	52	3.23	0.91	35,969	0.91	9,280	0.93
Special Instruction Needs							
Economically Disadvantaged	52	2.96	0.91	243	0.89	27,047	0.91
Not Economically Disadvantaged	52	3.21	0.92	45,291	0.91	141	0.94
English Learner	52	2.93	0.90	219	0.89	8,067	0.91
Non-English Learner	52	3.16	0.93	719	0.92	21,108	0.93
Students with Disabilities	52	2.94	0.92	9,423	0.90	9,358	0.93
Students without Disabilities	52	3.16	0.92	244	0.91	68,814	0.93
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	52	3.00	0.93	27,682	0.93	30,959	0.93
Students Taking Translated Forms							
Spanish Language	52	2.72	0.87	4,930	0.87	4,930	0.87

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.10 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 4

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	52	3.02	0.93	411	0.90	76,069	0.94
Gender							
Male	52	2.99	0.93	228	0.89	21,118	0.93
Female	52	3.00	0.92	148	0.88	20,956	0.93
Ethnicity							
Black/African American	52	2.73	0.91	103	0.82	10,020	0.92
Asian/Pacific Islander	52	3.15	0.92	7,532	0.91	1,717	0.94
Hispanic/Latino	52	2.89	0.91	249	0.85	18,285	0.92
American Indian/Alaska Native	52	3.10	0.94	150	0.94	120	0.94
Multiple	52	3.08	0.93	3,548	0.93	805	0.94
White	52	3.16	0.92	38,279	0.91	9,165	0.94
Special Instruction Needs							
Economically Disadvantaged	52	2.82	0.91	249	0.88	25,573	0.92
Not Economically Disadvantaged	52	3.15	0.92	338	0.91	12,057	0.94
English Learner	52	2.77	0.89	183	0.87	8,001	0.90
Non-English Learner	52	3.07	0.93	622	0.90	21,941	0.94
Students with Disabilities	52	2.79	0.92	672	0.89	8,778	0.94
Students without Disabilities	52	3.07	0.93	215	0.89	66,730	0.94
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	52	2.84	0.93	32,229	0.93	30,152	0.93
Students Taking Translated Forms							
Spanish Language	52	2.55	0.87	4,454	0.87	4,454	0.87

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.11 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 5

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum N	Reliability Alpha	Maximum N	Reliability Alpha
Total Group	52	3.01	0.92	363	0.87	80,635	0.93
Gender							
Male	52	2.97	0.92	180	0.84	21,786	0.93
Female	52	2.99	0.91	218	0.83	21,519	0.92
Ethnicity							
Black/African American	52	2.70	0.88	197	0.71	10,480	0.90
Asian/Pacific Islander	52	3.19	0.92	6,939	0.92	1,630	0.94
Hispanic/Latino	52	2.87	0.89	339	0.81	19,361	0.91
American Indian/Alaska Native	52	3.17	0.92	133	0.92	134	0.92
Multiple	52	3.06	0.93	3,243	0.93	798	0.93
White	52	3.15	0.91	143	0.89	9,451	0.93
Special Instruction Needs							
Economically Disadvantaged	52	2.79	0.89	206	0.72	26,682	0.90
Not Economically Disadvantaged	52	3.15	0.92	116	0.90	13,602	0.93
English Learner	52	2.68	0.85	224	0.80	6,608	0.87
Non-English Learner	52	3.06	0.92	262	0.87	23,155	0.93
Students with Disabilities	52	2.72	0.90	767	0.84	9,393	0.93
Students without Disabilities	52	3.08	0.92	18,004	0.92	70,683	0.93
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	52	2.81	0.92	29,571	0.91	31,982	0.92
Students Taking Translated Forms							
Spanish Language	52	2.60	0.85	4,053	0.85	4,053	0.85

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.12 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 6

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	52	3.03	0.92	350	0.87	862	0.94
Gender							
Male	52	2.97	0.91	194	0.84	22,763	0.93
Female	52	2.98	0.91	115	0.87	199	0.92
Ethnicity							
Black/African American	52	2.64	0.88	115	0.67	10,965	0.90
Asian/Pacific Islander	52	3.36	0.92	7,392	0.91	1,374	0.94
Hispanic/Latino	52	2.83	0.89	172	0.86	19,987	0.91
American Indian/Alaska Native	52	3.03	0.92	161	0.90	144	0.94
Multiple	52	3.10	0.92	3,335	0.91	3,418	0.93
White	52	3.19	0.91	144	0.87	502	0.93
Special Instruction Needs							
Economically Disadvantaged	52	2.74	0.88	361	0.81	27,279	0.90
Not Economically Disadvantaged	52	3.20	0.92	50,459	0.90	463	0.94
English Learner	52	2.47	0.82	4,665	0.77	5,723	0.86
Non-English Learner	52	3.09	0.92	282	0.87	739	0.94
Students with Disabilities	52	2.65	0.88	11,178	0.85	9,460	0.92
Students without Disabilities	52	3.11	0.92	16,080	0.91	315	0.93
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	52	2.78	0.91	27,292	0.90	29,049	0.92
Students Taking Translated Forms							
Spanish Language	52	2.35	0.84	3,580	0.84	3,580	0.84

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.13 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 7

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	52	3.14	0.91	29,369	0.91	779	0.95
Gender							
Male	52	3.11	0.92	249	0.89	9,551	0.92
Female	52	3.09	0.90	156	0.87	8,331	0.91
Ethnicity							
Black/African American	52	2.80	0.88	143	0.68	10,501	0.88
Asian/Pacific Islander	52	3.38	0.92	6,332	0.91	1,258	0.93
Hispanic/Latino	52	3.00	0.89	12,949	0.88	125	0.92
American Indian/Alaska Native	52	3.12	0.91	132	0.90	155	0.92
Multiple	52	3.16	0.92	2,966	0.92	676	0.93
White	52	3.28	0.91	39,596	0.90	431	0.94
Special Instruction Needs							
Economically Disadvantaged	52	2.92	0.88	317	0.83	15,385	0.89
Not Economically Disadvantaged	52	3.28	0.91	51,156	0.91	432	0.95
English Learner	52	2.64	0.81	4,546	0.79	5,084	0.83
Non-English Learner	52	3.19	0.91	73,554	0.91	696	0.95
Students with Disabilities	52	2.75	0.89	11,232	0.86	9,824	0.92
Students without Disabilities	52	3.22	0.91	15,886	0.91	308	0.95
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	52	2.93	0.91	27,384	0.91	27,154	0.91
Students Taking Translated Forms							
Spanish Language	52	2.63	0.80	2,900	0.80	2,900	0.80

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.14 Summary of Test Reliability Estimates for Subgroups: Mathematics Grade 8

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	52	3.00	0.89	564	0.77	73,879	0.89
Gender							
Male	52	3.03	0.90	209	0.78	25,081	0.91
Female	52	3.04	0.89	182	0.81	24,560	0.90
Ethnicity							
Black/African American	52	2.64	0.80	141	0.52	10,218	0.82
Asian/Pacific Islander	52	3.27	0.92	4,545	0.92	1,053	0.94
Hispanic/Latino	52	2.85	0.84	117	0.75	10,068	0.85
American Indian/Alaska Native	52	3.05	0.84	163	0.84	186	0.85
Multiple	52	3.04	0.90	586	0.90	568	0.91
White	52	3.16	0.88	286	0.80	7,548	0.90
Special Instruction Needs							
Economically Disadvantaged	52	2.79	0.84	280	0.67	15,312	0.85
Not Economically Disadvantaged	52	3.14	0.89	237	0.81	9,954	0.90
English Learner	52	2.51	0.70	4,407	0.68	5,116	0.71
Non-English Learner	52	3.05	0.89	458	0.78	20,998	0.89
Students with Disabilities	52	2.64	0.81	10,596	0.76	9,617	0.87
Students without Disabilities	52	3.08	0.89	158	0.77	15,357	0.89
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	52	2.82	0.88	26,585	0.87	25,908	0.88
Students Taking Translated Forms							
Spanish Language	52	2.53	0.67	2,619	0.67	2,619	0.67

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.15 Summary of Test Reliability Estimates for Subgroups: Algebra I

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	55	2.82	0.91	326	0.90	46,782	0.92
Gender							
Male	55	2.96	0.89	848	0.88	911	0.89
Female	55	2.97	0.87	886	0.86	868	0.87
Ethnicity							
Black/African American	55	2.51	0.86	1,327	0.85	7,532	0.86
Asian/Pacific Islander	55	3.26	0.92	5,211	0.92	470	0.92
Hispanic/Latino	55	2.58	0.86	4,992	0.80	13,826	0.88
American Indian/Alaska Native	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Multiple	55	2.93	0.92	1,366	0.92	125	0.93
White	55	2.96	0.90	165	0.87	2,452	0.91
Special Instruction Needs							
Economically Disadvantaged	55	2.56	0.87	3,750	0.81	12,190	0.88
Not Economically Disadvantaged	55	2.93	0.92	249	0.90	3,851	0.92
English Learner	55	2.20	0.74	2,853	0.62	2,132	0.85
Non-English Learner	55	2.86	0.91	290	0.90	41,572	0.92
Students with Disabilities	55	2.45	0.87	3,709	0.82	7,003	0.90
Students without Disabilities	55	2.90	0.91	180	0.89	5,765	0.92
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	55	2.60	0.91	7,513	0.91	6,831	0.91
Students Taking Translated Forms							
Spanish Language	55	2.10	0.58	2,426	0.58	2,426	0.58

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.16 Summary of Test Reliability Estimates for Subgroups: Geometry

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	55	3.12	0.90	160	0.88	1,472	0.90
Gender							
Male	55	2.93	0.87	756	0.86	717	0.87
Female	55	2.90	0.84	730	0.83	686	0.85
Ethnicity							
Black/African American	55	2.51	0.84	353	0.79	1,940	0.85
Asian/Pacific Islander	55	3.44	0.90	171	0.87	2,726	0.90
Hispanic/Latino	55	2.71	0.85	816	0.75	2,861	0.86
American Indian/Alaska Native	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Multiple	55	3.16	0.89	632	0.89	610	0.90
White	55	3.21	0.87	106	0.80	471	0.88
Special Instruction Needs							
Economically Disadvantaged	55	2.66	0.86	621	0.81	1,942	0.87
Not Economically Disadvantaged	55	3.25	0.89	150	0.86	10,981	0.90
English Learner	55	2.27	0.76	400	0.63	315	0.86
Non-English Learner	55	3.18	0.90	154	0.87	1,159	0.90
Students with Disabilities	55	2.71	0.89	545	0.85	1,542	0.91
Students without Disabilities	55	3.16	0.90	131	0.86	13,876	0.90
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	55	2.83	0.90	1,633	0.89	1,472	0.90
Students Taking Translated Forms							
Spanish Language	55	2.04	0.46	306	0.46	306	0.46

Note. SEM = standard error of measurement; n/a = not applicable.

Table A.13.17 Summary of Test Reliability Estimates for Subgroups: Algebra II

	Avg. Max. Raw Score	Avg. SEM	Avg. Reliability	Minimum Reliability N	Reliability Alpha	Maximum Reliability N	Reliability Alpha
Total Group	55	3.30	0.89	6,087	0.88	697	0.93
Gender							
Male	55	2.98	0.85	650	0.85	673	0.86
Female	55	2.91	0.82	650	0.80	652	0.84
Ethnicity							
Black/African American	55	2.85	0.85	476	0.82	416	0.87
Asian/Pacific Islander	55	3.60	0.88	1,678	0.87	128	0.91
Hispanic/Latino	55	2.80	0.80	553	0.70	956	0.87
American Indian/Alaska Native	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Multiple	55	3.24	0.86	333	0.83	314	0.89
White	55	3.38	0.85	2,609	0.84	2,595	0.87
Special Instruction Needs							
Economically Disadvantaged	55	2.74	0.79	494	0.69	528	0.88
Not Economically Disadvantaged	55	3.50	0.88	4,169	0.87	254	0.93
English Learner	55	2.30	0.66	180	0.57	100	0.81
Non-English Learner	55	3.39	0.89	4,651	0.88	651	0.93
Students with Disabilities	55	2.90	0.87	124	0.74	557	0.91
Students without Disabilities	55	3.35	0.89	5,526	0.88	573	0.94
Students Taking Accommodated Forms							
American Sign Language	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Screen Reader	n/a	n/a	n/a	n/a	n/a	n/a	n/a
Text-to-Speech	55	2.84	0.92	740	0.91	697	0.93
Students Taking Translated Forms							
Spanish Language	55	2.19	0.43	134	0.43	134	0.43

Note. SEM = standard error of measurement; n/a = not applicable.

Appendix 13.2: Reliability of Classification by Content and Grade/Subject

Table A.13.18 Reliability of Classification: Grade 3 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.22	0.03	0.00	0.00	0.00	0.25
	700–724	0.04	0.10	0.05	0.00	0.00	0.19
	725–749	0.00	0.04	0.12	0.04	0.00	0.21
	750–809	0.00	0.00	0.05	0.25	0.03	0.33
	810–850	0.00	0.00	0.00	0.01	0.01	0.02
Decision Consistency	650–699	0.21	0.05	0.01	0.00	0.00	0.27
	700–724	0.04	0.07	0.06	0.01	0.00	0.19
	725–749	0.01	0.04	0.09	0.05	0.00	0.19
	750–809	0.00	0.01	0.07	0.22	0.02	0.32
	810–850	0.00	0.00	0.00	0.02	0.01	0.03

Table A.13.19 Reliability of Classification: Grade 4 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.15	0.02	0.00	0.00	0.00	0.18
	700–724	0.04	0.11	0.04	0.00	0.00	0.19
	725–749	0.00	0.04	0.14	0.05	0.00	0.23
	750–809	0.00	0.00	0.05	0.23	0.03	0.32
	810–850	0.00	0.00	0.00	0.02	0.06	0.08
Decision Consistency	650–699	0.15	0.04	0.01	0.00	0.00	0.19
	700–724	0.04	0.08	0.06	0.01	0.00	0.19
	725–749	0.01	0.05	0.11	0.06	0.00	0.22
	750–809	0.00	0.01	0.06	0.20	0.04	0.30
	810–850	0.00	0.00	0.00	0.04	0.05	0.10

Table A.13.20 Reliability of Classification: Grade 5 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.14	0.02	0.00	0.00	0.00	0.16
	700–724	0.04	0.12	0.05	0.00	0.00	0.21
	725–749	0.00	0.04	0.16	0.05	0.00	0.25
	750–809	0.00	0.00	0.05	0.27	0.02	0.34
	810–850	0.00	0.00	0.00	0.01	0.03	0.04
Decision Consistency	650–699	0.13	0.04	0.00	0.00	0.00	0.17
	700–724	0.04	0.10	0.06	0.01	0.00	0.20
	725–749	0.00	0.05	0.12	0.06	0.00	0.24
	750–809	0.00	0.01	0.06	0.24	0.02	0.33
	810–850	0.00	0.00	0.00	0.02	0.03	0.05

Table A.13.21 Reliability of Classification: Grade 6 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.11	0.02	0.00	0.00	0.00	0.13
	700–724	0.03	0.14	0.05	0.00	0.00	0.22
	725–749	0.00	0.04	0.19	0.05	0.00	0.29
	750–809	0.00	0.00	0.05	0.25	0.02	0.32
	810–850	0.00	0.00	0.00	0.01	0.04	0.05
Decision Consistency	650–699	0.11	0.03	0.00	0.00	0.00	0.14
	700–724	0.04	0.11	0.06	0.00	0.00	0.22
	725–749	0.00	0.05	0.15	0.06	0.00	0.27
	750–809	0.00	0.00	0.07	0.22	0.02	0.31
	810–850	0.00	0.00	0.00	0.03	0.04	0.06

Table A.13.22 Reliability of Classification: Grade 7 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.14	0.02	0.00	0.00	0.00	0.16
	700–724	0.03	0.11	0.04	0.00	0.00	0.19
	725–749	0.00	0.04	0.15	0.05	0.00	0.24
	750–809	0.00	0.00	0.05	0.20	0.03	0.29
	810–850	0.00	0.00	0.00	0.03	0.10	0.13
Decision Consistency	650–699	0.13	0.04	0.01	0.00	0.00	0.18
	700–724	0.03	0.09	0.06	0.01	0.00	0.19
	725–749	0.00	0.05	0.11	0.06	0.00	0.22
	750–809	0.00	0.01	0.06	0.17	0.04	0.27
	810–850	0.00	0.00	0.00	0.04	0.09	0.14

Table A.13.23 Reliability of Classification: Grade 8 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.16	0.03	0.00	0.00	0.00	0.18
	700–724	0.04	0.10	0.05	0.00	0.00	0.19
	725–749	0.00	0.04	0.14	0.05	0.00	0.23
	750–809	0.00	0.00	0.05	0.23	0.03	0.31
	810–850	0.00	0.00	0.00	0.02	0.06	0.08
Decision Consistency	650–699	0.15	0.04	0.01	0.00	0.00	0.20
	700–724	0.04	0.08	0.06	0.01	0.00	0.19
	725–749	0.01	0.05	0.10	0.06	0.00	0.21
	750–809	0.00	0.01	0.07	0.20	0.03	0.30
	810–850	0.00	0.00	0.00	0.04	0.06	0.10

Table A.13.24 Reliability of Classification: Grade 9 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.09	0.02	0.00	0.00	0.00	0.11
	700–724	0.03	0.09	0.04	0.00	0.00	0.17
	725–749	0.00	0.04	0.14	0.05	0.00	0.24
	750–809	0.00	0.00	0.06	0.28	0.04	0.38
	810–850	0.00	0.00	0.00	0.03	0.08	0.11
Decision Consistency	650–699	0.08	0.03	0.01	0.00	0.00	0.12
	700–724	0.03	0.07	0.05	0.01	0.00	0.17
	725–749	0.01	0.05	0.10	0.07	0.00	0.22
	750–809	0.00	0.01	0.07	0.23	0.05	0.36
	810–850	0.00	0.00	0.00	0.05	0.07	0.13

Table A.13.25 Reliability of Classification: Grade 10 ELA/L

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.18	0.03	0.00	0.00	0.00	0.21
	700–724	0.04	0.07	0.05	0.01	0.00	0.16
	725–749	0.01	0.04	0.09	0.05	0.00	0.18
	750–809	0.00	0.01	0.06	0.20	0.04	0.31
	810–850	0.00	0.00	0.00	0.04	0.10	0.14
Decision Consistency	650–699	0.17	0.04	0.02	0.00	0.00	0.23
	700–724	0.04	0.05	0.05	0.02	0.00	0.15
	725–749	0.01	0.04	0.07	0.05	0.00	0.17
	750–809	0.00	0.01	0.06	0.16	0.05	0.29
	810–850	0.00	0.00	0.00	0.06	0.09	0.16

Table A.13.26 Reliability of Classification: Grade 3 Mathematics

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.14	0.02	0.00	0.00	0.00	0.17
	700–724	0.03	0.14	0.04	0.00	0.00	0.21
	725–749	0.00	0.04	0.15	0.04	0.00	0.24
	750–809	0.00	0.00	0.04	0.23	0.02	0.29
	810–850	0.00	0.00	0.00	0.02	0.07	0.09
Decision Consistency	650–699	0.14	0.04	0.00	0.00	0.00	0.18
	700–724	0.03	0.11	0.05	0.00	0.00	0.20
	725–749	0.00	0.05	0.12	0.05	0.00	0.23
	750–809	0.00	0.00	0.05	0.20	0.03	0.29
	810–850	0.00	0.00	0.00	0.03	0.07	0.10

Table A.13.27 Reliability of Classification: Grade 4 Mathematics

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.15	0.03	0.00	0.00	0.00	0.17
	700–724	0.03	0.18	0.04	0.00	0.00	0.26
	725–749	0.00	0.04	0.17	0.04	0.00	0.26
	750–809	0.00	0.00	0.03	0.23	0.01	0.28
	810–850	0.00	0.00	0.00	0.01	0.03	0.04
Decision Consistency	650–699	0.14	0.04	0.00	0.00	0.00	0.18
	700–724	0.04	0.15	0.05	0.00	0.00	0.25
	725–749	0.00	0.06	0.14	0.05	0.00	0.25
	750–809	0.00	0.00	0.05	0.21	0.01	0.28
	810–850	0.00	0.00	0.00	0.01	0.03	0.04

Table A.13.28 Reliability of Classification: Grade 5 Mathematics

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.16	0.03	0.00	0.00	0.00	0.19
	700–724	0.03	0.18	0.04	0.00	0.00	0.26
	725–749	0.00	0.05	0.17	0.04	0.00	0.26
	750–809	0.00	0.00	0.04	0.19	0.02	0.25
	810–850	0.00	0.00	0.00	0.01	0.03	0.04
Decision Consistency	650–699	0.16	0.05	0.00	0.00	0.00	0.21
	700–724	0.04	0.15	0.06	0.00	0.00	0.25
	725–749	0.00	0.06	0.14	0.05	0.00	0.25
	750–809	0.00	0.00	0.05	0.17	0.02	0.25
	810–850	0.00	0.00	0.00	0.02	0.03	0.05

Table A.13.29 Reliability of Classification: Grade 6 Mathematics

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.16	0.03	0.00	0.00	0.00	0.19
	700–724	0.03	0.20	0.05	0.00	0.00	0.28
	725–749	0.00	0.05	0.19	0.04	0.00	0.28
	750–809	0.00	0.00	0.04	0.17	0.01	0.22
	810–850	0.00	0.00	0.00	0.01	0.02	0.03
Decision Consistency	650–699	0.16	0.05	0.00	0.00	0.00	0.21
	700–724	0.04	0.17	0.06	0.00	0.00	0.27
	725–749	0.00	0.06	0.16	0.05	0.00	0.27
	750–809	0.00	0.00	0.05	0.15	0.01	0.22
	810–850	0.00	0.00	0.00	0.01	0.02	0.04

Table A.13.30 Reliability of Classification: Grade 7 Mathematics

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.09	0.02	0.00	0.00	0.00	0.11
	700–724	0.03	0.20	0.05	0.00	0.00	0.28
	725–749	0.00	0.05	0.23	0.04	0.00	0.33
	750–809	0.00	0.00	0.04	0.20	0.01	0.25
	810–850	0.00	0.00	0.00	0.01	0.03	0.03
Decision Consistency	650–699	0.09	0.04	0.00	0.00	0.00	0.13
	700–724	0.04	0.17	0.06	0.00	0.00	0.27
	725–749	0.00	0.06	0.19	0.06	0.00	0.31
	750–809	0.00	0.00	0.06	0.18	0.01	0.25
	810–850	0.00	0.00	0.00	0.02	0.03	0.04

Table A.13.31 Reliability of Classification: Grade 8 Mathematics

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.24	0.06	0.00	0.00	0.00	0.31
	700–724	0.05	0.18	0.05	0.00	0.00	0.29
	725–749	0.00	0.06	0.11	0.04	0.00	0.21
	750–809	0.00	0.00	0.03	0.13	0.01	0.18
	810–850	0.00	0.00	0.00	0.00	0.02	0.02
Decision Consistency	650–699	0.23	0.08	0.01	0.00	0.00	0.32
	700–724	0.06	0.14	0.06	0.01	0.00	0.26
	725–749	0.01	0.07	0.09	0.04	0.00	0.20
	750–809	0.00	0.01	0.05	0.12	0.01	0.19
	810–850	0.00	0.00	0.00	0.01	0.02	0.03

Table A.13.32 Reliability of Classification: Algebra I

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.13	0.03	0.00	0.00	0.00	0.16
	700–724	0.04	0.15	0.05	0.00	0.00	0.24
	725–749	0.00	0.05	0.16	0.05	0.00	0.26
	750–809	0.00	0.00	0.04	0.26	0.01	0.32
	810–850	0.00	0.00	0.00	0.00	0.02	0.02
Decision Consistency	650–699	0.13	0.05	0.01	0.00	0.00	0.18
	700–724	0.04	0.12	0.06	0.01	0.00	0.23
	725–749	0.00	0.06	0.12	0.06	0.00	0.25
	750–809	0.00	0.01	0.06	0.24	0.01	0.32
	810–850	0.00	0.00	0.00	0.01	0.02	0.03

Table A.13.33 Reliability of Classification: Geometry

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.06	0.01	0.00	0.00	0.00	0.07
	700–724	0.02	0.15	0.04	0.00	0.00	0.21
	725–749	0.00	0.05	0.22	0.06	0.00	0.32
	750–809	0.00	0.00	0.05	0.28	0.03	0.36
	810–850	0.00	0.00	0.00	0.01	0.02	0.03
Decision Consistency	650–699	0.06	0.03	0.00	0.00	0.00	0.08
	700–724	0.02	0.13	0.05	0.00	0.00	0.21
	725–749	0.00	0.06	0.18	0.07	0.00	0.31
	750–809	0.00	0.00	0.07	0.24	0.03	0.34
	810–850	0.00	0.00	0.00	0.03	0.02	0.05

Table A.13.34 Reliability of Classification: Algebra II

	Full Summative Scale Score	Level 1	Level 2	Level 3	Level 4	Level 5	Category Total
Decision Accuracy	650–699	0.11	0.03	0.00	0.00	0.00	0.14
	700–724	0.04	0.09	0.04	0.00	0.00	0.17
	725–749	0.00	0.05	0.10	0.05	0.00	0.21
	750–809	0.00	0.00	0.05	0.34	0.03	0.43
	810–850	0.00	0.00	0.00	0.01	0.03	0.05
Decision Consistency	650–699	0.11	0.04	0.01	0.00	0.00	0.16
	700–724	0.04	0.07	0.05	0.01	0.00	0.17
	725–749	0.01	0.05	0.08	0.07	0.00	0.20
	750–809	0.00	0.01	0.06	0.30	0.03	0.41
	810–850	0.00	0.00	0.00	0.04	0.03	0.07

Appendix 14: Quality Testing Standards

Table A.14.1 ELA/L Grade 6 Form 1 Matching Results

ELA/L Grade 6 Form 1	Unmatched		Diff*	Matched		Diff*
	Current Form 1	Original Form 1		Current Form 1	Original Form 1	
Sample Size	119,838	31,031		30,667	30,667	
American Indian/Alaska Native	1.3	0.3	1	0.3	0.3	0
Asian	6.8	6.7	0.1	6.7	6.7	0
Black/African American	14.1	32.8	-18.6	32.2	32.2	0
Hispanic/Latino Ethnicity	31.4	18.9	12.5	19.1	19.1	0
Hawaiian/Pacific Islander	0.2	0.2	0	0.1	0.1	0
White	43.4	36.5	6.9	37	37	0
Two or More Races	2.9	4.7	-1.8	4.7	4.7	0
Female	49.7	49.4	0.3	49.4	49.4	0
Economic Disadvantage	48.3	44.1	4.2	44.5	44.5	0
English Learner	7.2	5.7	1.4	5.6	5.6	0
Students with Disabilities	14.4	13.9	0.5	13.7	13.7	0
Grade 6	100	100	0	100	100	0
Prior Year Scale Score	745	742.3	2.7	742.7	742.7	0
Prior Performance Level 1	10.2	11.7	-1.5	11.4	11.4	0
Prior Performance Level 2	18	19	-1	18.8	18.8	0
Prior Performance Level 3	26.4	26.3	0.1	26.4	26.4	0
Prior Performance Level 4	39.3	38.5	0.9	38.8	38.8	0
Prior Performance Level 5	6.1	4.6	1.5	4.6	4.6	0

*Diff = Current Percent – Original Percent.

Table A.14.2 Mathematics Grade 6 Form 1 Matching Results

Mathematics Grade 6 Form 1	Unmatched		Diff*	Matched		Diff*
	Current Form 1	Original Form 1		Current Form 1	Original Form 1	
Sample Size	95,174	28,514		27,677	27,677	
American Indian/Alaska Native	1.1	0.2	0.9	0.2	0.2	0
Asian	7.6	7	0.6	7.1	7.1	0
Black/African American	11.5	33.4	-21.9	31.6	31.6	0
Hispanic/Latino Ethnicity	28	17.9	10.1	18.5	18.5	0
Hawaiian/Pacific Islander	0.1	0.2	0	0.1	0.1	0
White	48.4	36.5	11.9	37.6	37.6	0
Two or More Races	3.2	4.8	-1.6	4.9	4.9	0
Female	50.2	50	0.2	50.1	50.1	0
Economic Disadvantage	42.6	42.4	0.3	43.2	43.2	0
English Learner	4.6	3.7	0.9	3.5	3.5	0
Students with Disabilities	9.8	11	-1.2	10.6	10.6	0
Grade 6	100	100	0	100	100	0
Prior Year Scale Score	743.9	741.1	2.8	741.7	741.7	0
Prior Performance Level 1	9	12.6	-3.6	12	12	0
Prior Performance Level 2	18.9	20.3	-1.4	20	20	0
Prior Performance Level 3	28.6	25.6	3	25.8	25.8	0
Prior Performance Level 4	35.7	33.8	1.9	34.3	34.3	0
Prior Performance Level 5	7.8	7.8	0	7.8	7.8	0

*Diff = Current Percent – Original Percent.

Table A.14.3 ELA/L Grade 10 Form 1 Matching Results

ELA/L Grade 10 Form 1	Unmatched		Diff*	Matched		Diff*
	Current Form 1	Original Form 1		Current Form 1	Original Form 1	
Sample Size	55,046	27,951		22,970	22,970	
American Indian/Alaska Native	2	0.3	1.7	0.3	0.3	0
Asian	9.3	7.5	1.8	8.6	8.6	0
Black/African American	11.1	33.2	-22	24.1	24.1	0
Hispanic/Latino Ethnicity	32.1	14.9	17.2	17.5	17.5	0
Hawaiian/Pacific Islander	0.2	0.1	0.1	0.1	0.1	0
White	44	39.5	4.5	46.9	46.9	0
Two or More Races	1.3	4.6	-3.3	2.6	2.6	0
Female	50.2	50.5	-0.2	50.5	50.5	0
Economic Disadvantage	35.8	35	0.9	32.6	32.6	0
English Learner	3.2	3.2	0	2.9	2.9	0
Students with Disabilities	15.6	14.7	0.9	14.4	14.4	0
Grade 9	1.3	3.5	-2.2	1.8	1.8	0
Grade 10	98.6	96.5	2.2	98.2	98.2	0
2017 Scale Score	755.5	740	15.5	746.3	746.2	0.1
2017 Performance Level 1	8.8	15.8	-7	11.3	11.3	0
2017 Performance Level 2	13	18.8	-5.8	15.9	15.9	0
2017 Performance Level 3	21.4	23.7	-2.2	24.5	24.5	0
2017 Performance Level 4	39.6	34	5.6	39.1	39.1	0
2017 Performance Level 5	17.3	7.7	9.5	9.3	9.3	0

*Diff = Current Percent – Original Percent.

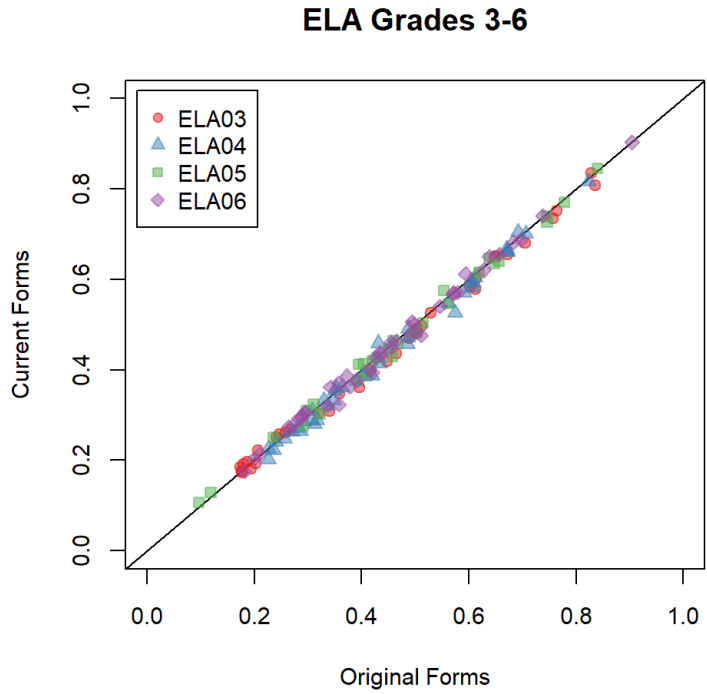


Figure A.14.1 ELA/L Grades 3-6 P-Values

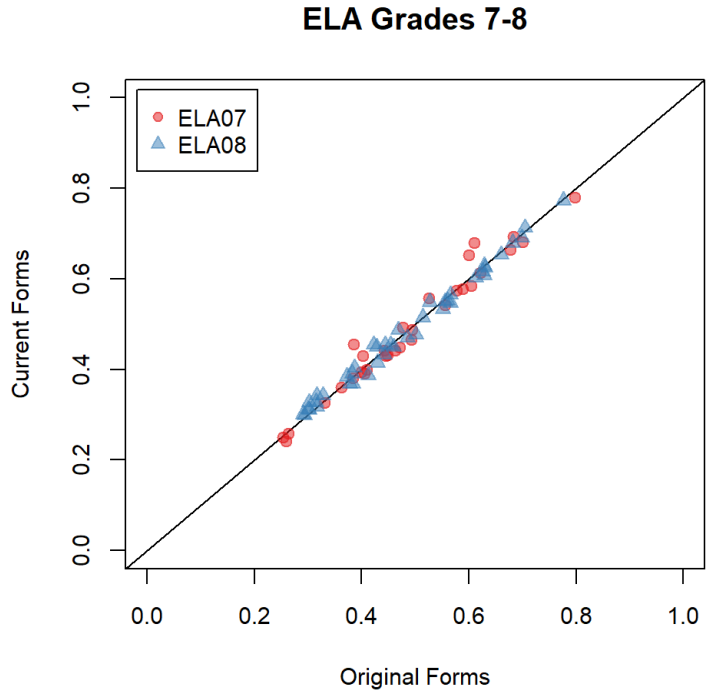


Figure A.14.2 ELA/L Grades 7-8 P-Values

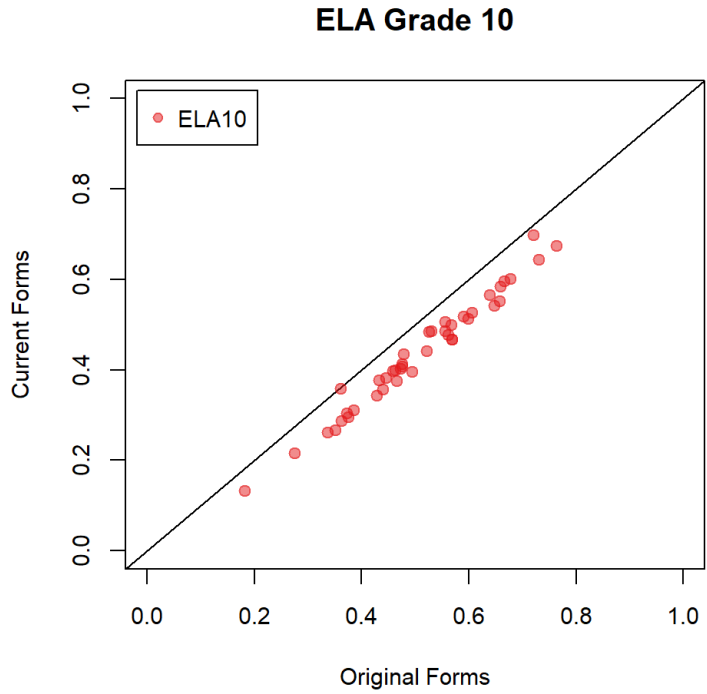


Figure A.14.3 ELA/L Grade 10 P-Values

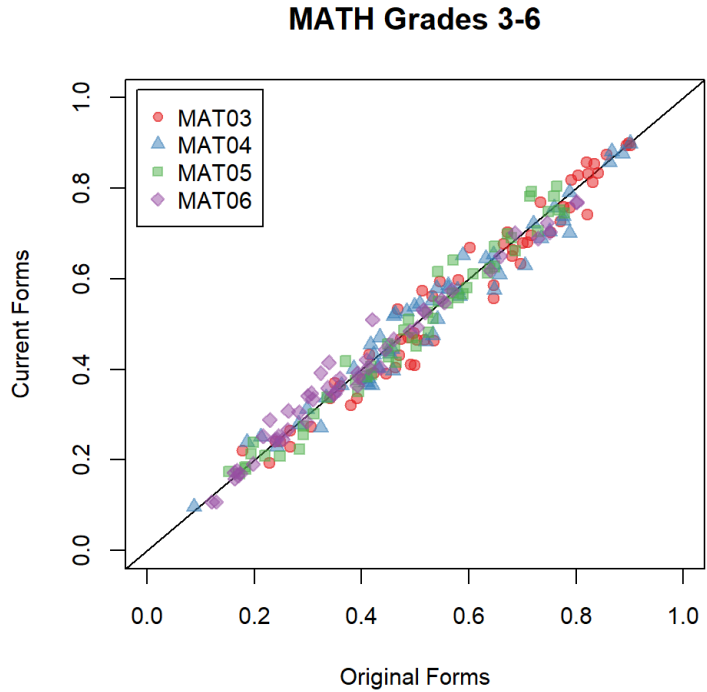


Figure A.14.4 Mathematics Grades 3-6 P-Values

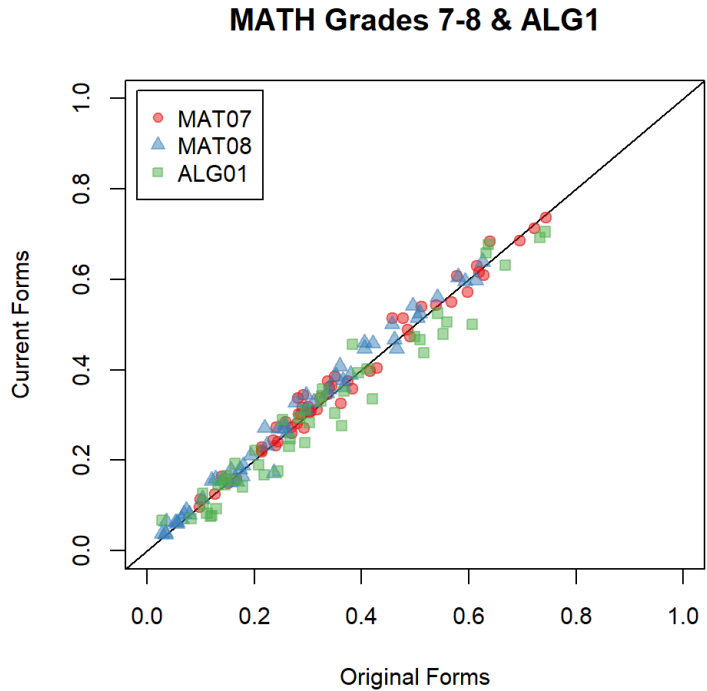


Figure A.14.5 Mathematics Grades 7–8 and Algebra I P-Values

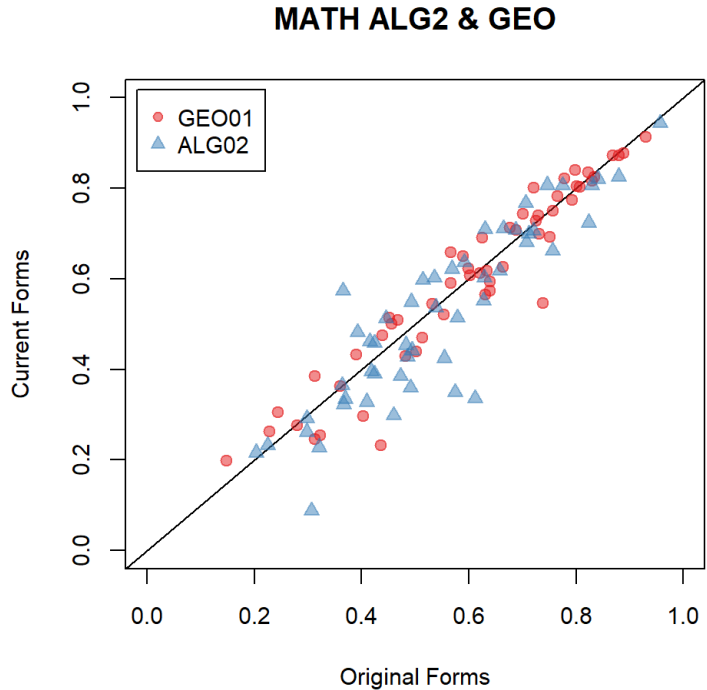


Figure A.14.6 Algebra II and Geometry P-Values

Table A.14.4 Distributions of P-Value Differences* for ELA/L

Grade	N	Min.	25%	Median	75%	Max.
3	34	-0.034	-0.017	-0.01	0.004	0.016
4	42	-0.049	-0.019	-0.01	-0.004	0.028
5	31	-0.029	-0.016	-0.006	0.009	0.021
6	42	-0.035	-0.008	-0.001	0.008	0.02
7	31	-0.026	-0.016	-0.006	0	0.07
8	42	-0.025	-0.01	0	0.011	0.032
10	42	-0.106	-0.085	-0.073	-0.062	-0.003

*Difference = Current P-value – Original P-value.

Table A.14.5 Distributions of P-Value Differences* for Mathematics

Grade/ Course	N	Min.	25%	Median	75%	Max.
3	59	-0.088	-0.038	-0.017	0.018	0.068
4	56	-0.086	-0.036	-0.003	0.016	0.064
5	54	-0.06	-0.023	-0.01	0.011	0.075
6	52	-0.048	-0.009	0	0.015	0.09
7	55	-0.034	-0.006	0.006	0.022	0.057
8	54	-0.065	0.005	0.013	0.025	0.054
Algebra I	48	-0.105	-0.042	-0.019	0.014	0.073
Geometry	55	-0.204	-0.031	0.004	0.04	0.094
Algebra II	51	-0.275	-0.062	-0.022	0.04	0.209

*Difference = Current P-value – Original P-value.

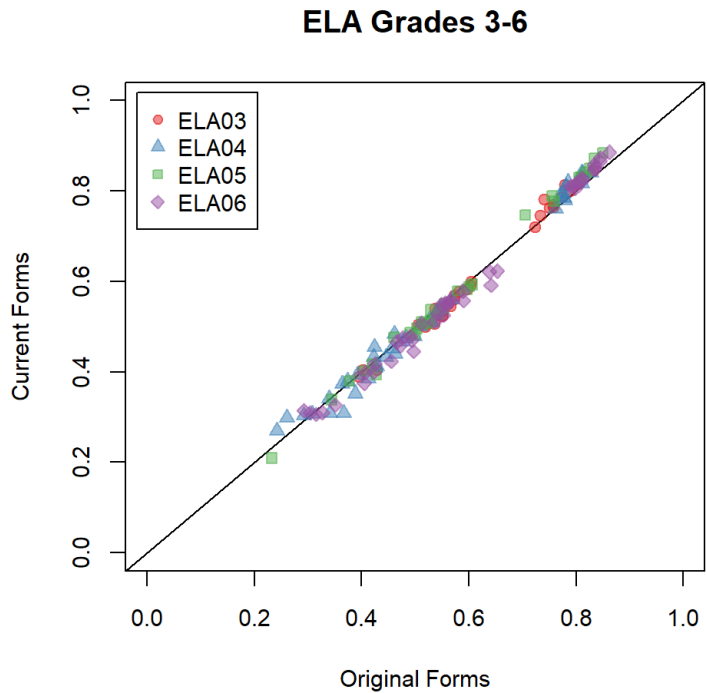


Figure A.14.7 Polyserial Correlations ELA/L Grades 3–6

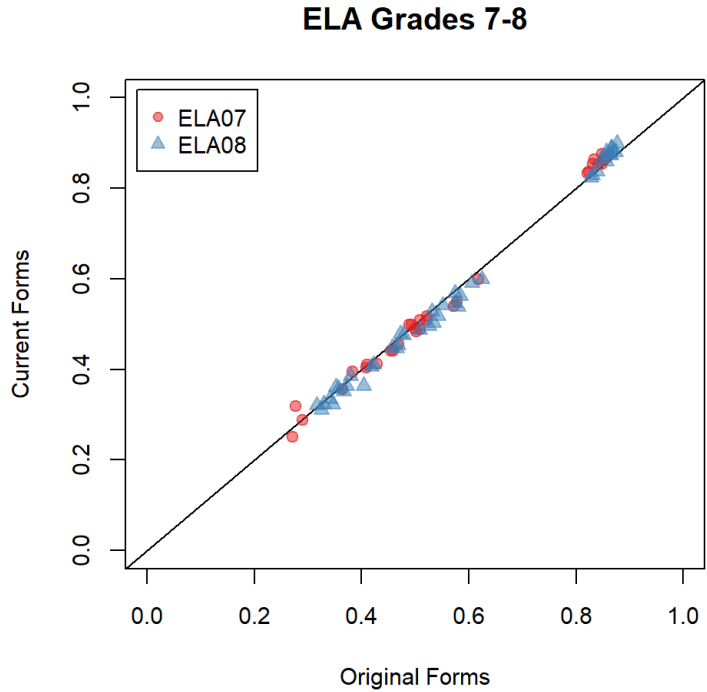


Figure A.14.8 Polyserial Correlations ELA/L Grades 7–8

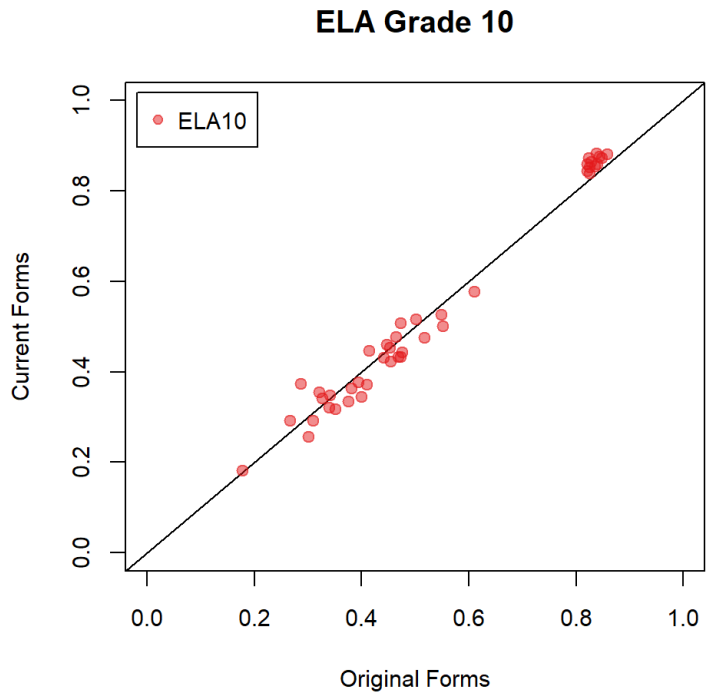


Figure A.14.9 Polyserial Correlations ELA/L Grade 10

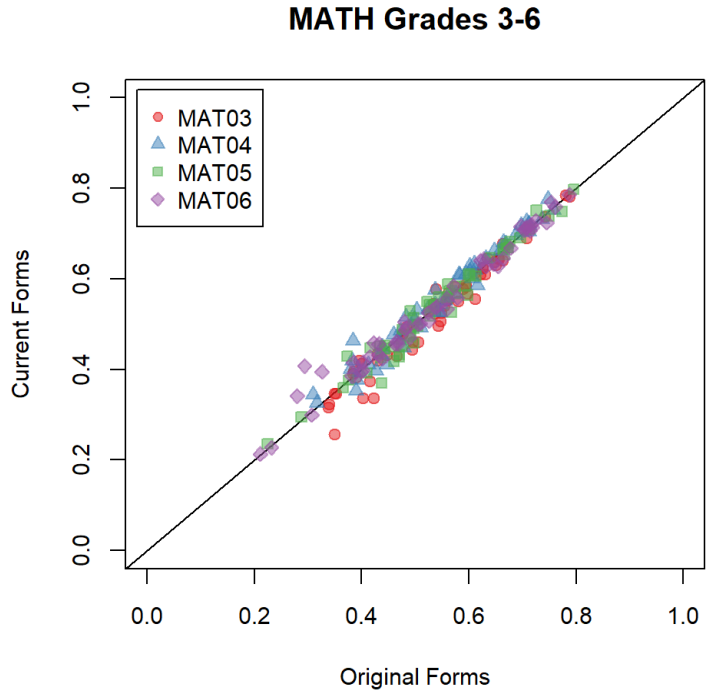


Figure A.14.10 Polyserial Correlations Mathematics Grades 3-6

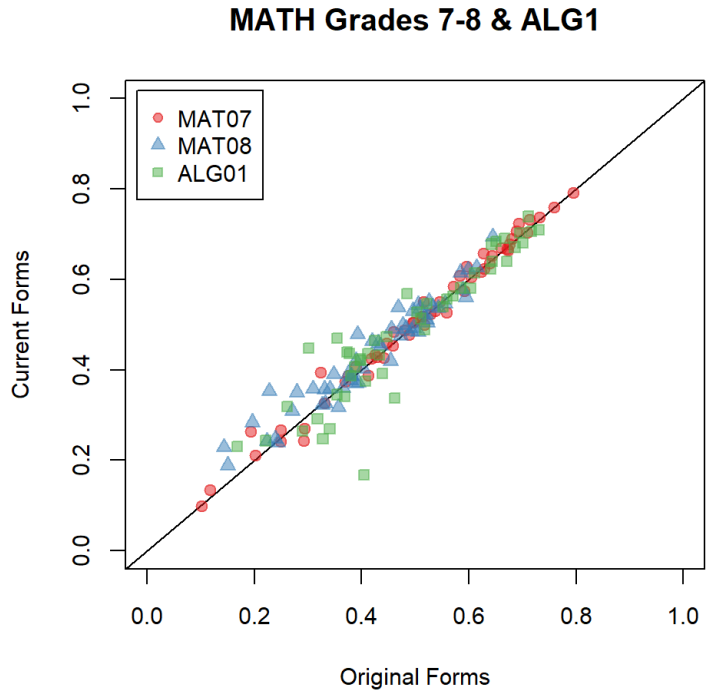


Figure A.14.11 Polyserial Correlations Mathematics Grades 7–8 and Algebra I

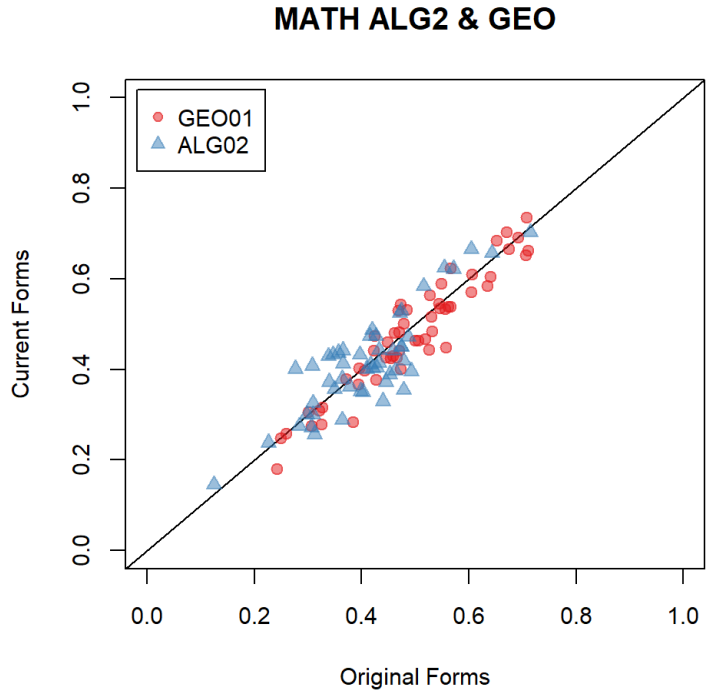


Figure A.14.12 Polyserial Correlations Algebra II and Geometry

Table A.14.6 Distributions of Polyserial Differences* for ELA/L

Grade	N	Min.	25%	Median	75%	Max.
3	34	-0.029	-0.015	-0.004	0.012	0.041
4	42	-0.058	-0.011	0	0.017	0.037
5	31	-0.034	-0.013	-0.003	0.020	0.042
6	42	-0.052	-0.022	-0.008	0.013	0.028
7	31	-0.031	-0.015	0	0.012	0.043
8	42	-0.042	-0.017	-0.007	0.005	0.023
10	42	-0.055	-0.032	0.010	0.026	0.088

*Difference = Current Polyserial – Original Polyserial.

Table A.14.7 Distributions of Polyserial Differences* for Mathematics

Grade/ Course	N	Min.	25%	Median	75%	Max.
3	59	-0.092	-0.022	-0.01	0.004	0.040
4	56	-0.036	-0.004	0.008	0.018	0.079
5	54	-0.067	-0.011	-0.002	0.010	0.056
6	52	-0.026	-0.008	-0.001	0.012	0.113
7	55	-0.050	-0.005	0.005	0.012	0.070
8	54	-0.040	-0.006	0.014	0.034	0.125
Algebra I	48	-0.238	-0.022	0.001	0.025	0.145
Geometry	55	-0.108	-0.037	-0.011	0.012	0.072
Algebra II	51	-0.125	-0.025	0.002	0.052	0.125

*Difference = Current Polyserial – Original Polyserial.

Table A.14.8 DIF Category Cross Tabulations for ELA/L

ELA/L Grades 3–8 & 10	Percent of DIF Calculations		
	None	B DIF (Current)	C DIF (Current)
None	89.9%–96.7%	0%–2.7%	0%–0.4%
B DIF (Original)	0.6%–4.8%	1.2%–2.4%	0%
C DIF (Original)	0%–0.4%	0%–1.8%	0%–1.6%

Table A.14.9 DIF Category Crosstabulations for Mathematics Grades 3–8 and Algebra I

Mathematics Grades 3–8 & Algebra I	Percent of DIF Calculations		
	None	B DIF (Current)	C DIF (Current)
None	94.5%–97.3%	0.2%–2.1%	0%–0.3%
B DIF (Original)	1.4%–2.5%	0.2%–2.2%	0%–0.5%
C DIF (Original)	0%–0.5%	0%–0.5%	0%–0.2%

Table A.14.10 DIF Category Crosstabulations for Algebra II and Geometry

Geometry & Algebra II	Percent of DIF Calculations		
	None	B DIF (Current)	C DIF (Current)
None	73.2%–77.5%	8.6%–12.7%	0%–1.4%
B DIF (Original)	5.9%–7.3%	2%–3.2%	0%–0.5%
C DIF (Original)	1.8%–2.0%	0%–0.9%	0%–3.2%

Table A.14.11 ELA/L Reliability

Grade	Original		Current Form 1				Current Form 2			
	Pts.	Alpha**	Pts.	Alpha	SB	DIFF*	Pts.	Alpha	SB	DIFF*
3	82	0.92	54	0.90	0.89	0.01	55	0.89	0.89	0
4	106	0.92	74	0.89	0.89	0	67	0.88	0.88	0
5	106	0.93	74	0.89	0.89	0	67	0.88	0.89	-0.01
6	109	0.94	74	0.92	0.92	0	70	0.90	0.90	0
7	109	0.94	74	0.91	0.91	0	70	0.90	0.91	-0.01
8	109	0.94	74	0.92	0.92	0	70	0.90	0.91	-0.01
10	109	0.93	74	0.90	0.89	0.01	70	0.88	0.89	-0.01

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

**Alpha = Weighted average of the stratified alphas from Original form 1 and Original form 2.

Table A.14.12 ELA/L Raw Score Standard Error of Measurement

Grade	Original			Current Form 1			Current Form 2		
	RS Points	RS SEM	SEM/ Points	RS Points	RS SEM	SEM/ Points	RS Points	RS SEM	SEM/ Points
3	82	4.42	0.054	54	3.54	0.066	55	3.58	0.065
4	106	5.41	0.051	74	4.46	0.06	67	4.51	0.067
5	106	5.46	0.052	74	4.48	0.061	67	4.48	0.067
6	109	5.53	0.051	74	4.50	0.061	70	4.49	0.064
7	109	5.93	0.054	74	4.71	0.064	70	5.06	0.072
8	109	5.63	0.052	74	4.52	0.061	70	4.69	0.067
10	109	5.95	0.044	74	4.71	0.05	70	5.20	0.06

Note. RS = raw score; SEM = standard error of measurement.

Table A.14.13 ELA/L Scale Score Standard Error of Measurement

Grade	Original Form 1		Original Form 2		Current Form 1		Current Form 2	
	SS Points	SS SEM	SS Points	SS SEM	SS Points	SS SEM	SS Points	SS SEM
3	82	11.6	82	11.8	54	13.8	55	13.9
4	106	10.6	106	10.6	74	12.9	67	13.3
5	106	9.7	106	9.5	74	11.9	67	12.6
6	109	8	109	8.4	74	9.7	70	10.9
7	109	9.7	109	9.7	74	11.9	70	12.9
8	109	9.8	109	9.7	74	11.8	70	12.9
10	109	11.4	109	11.6	74	14.6	70	16.3

Note. SS = scale score; SEM = standard error of measurement.

Table A.14.14 Mathematics Reliability

Grade/ Course	Original		Current Form 1 and Form 2			Diff*
	Points	Alpha**	Points	Alpha**	SB	
3	66	0.94	52	0.92	0.93	-0.01
4	66	0.94	52	0.93	0.93	0
5	66	0.94	52	0.93	0.93	0
6	66	0.95	52	0.93	0.94	-0.01
7	66	0.93	52	0.92	0.91	0.01
8	66	0.87	52	0.86	0.84	0.02
Algebra I	81	0.93	55	0.90	0.90	0
Geometry	81	0.93	55	0.89	0.90	-0.01
Algebra II	81	0.89	55	0.84	0.85	-0.01

Note. SB = Spearman Brown.

**Alpha = Weighted average of the stratified alphas from form 1 and form 2.

Table A.14.15 Mathematics Raw Score Standard Error of Measurement

Grade/Course	Original			Current		
	RS Points	RS SEM	SEM/ Points	RS Points	RS SEM	SEM/ Points
3	66	3.58	0.054	52	3.20	0.062
4	66	3.74	0.057	52	3.32	0.064
5	66	3.69	0.056	52	3.29	0.063
6	66	3.49	0.053	52	3.14	0.060
7	66	3.50	0.053	52	3.10	0.060
8	66	2.96	0.045	52	2.71	0.052
Algebra I	81	3.61	0.045	55	2.88	0.052
Geometry	81	4.21	0.052	55	3.51	0.064
Algebra II	81	4.25	0.052	55	3.50	0.064

Note. RS = scale score; SEM = standard error of measurement.

Table A.14.16 Mathematics Scale Score Standard Error of Measurement

Grade/Course	Original				Current			
	Form 1		Form 2		Form 1		Form 2	
	SS Points	SS SEM	SS Points	SS SEM	SS Points	SS SEM	SS Points	SS SEM
3	66	8.8	66	8.8	52	9.9	52	10.3
4	66	7.9	66	8.4	52	8.9	52	9.2
5	66	8.2	66	7.9	52	9.3	52	9.3
6	66	7.6	66	7.3	52	9.1	52	8.6
7	66	7.5	66	7.3	52	8.3	52	8.1
8	66	11.0	66	11.5	52	12.0	52	13.0
Algebra I	80	8.9	81	8.7	55	10.8	55	10.4
Geometry	81	6.4	81	6.4	55	7.9	55	8.0
Algebra II	81	9.7	81	9.8	55	11.4	55	12.2

Note. RS = scale score; SEM = standard error of measurement.

Table A.14.17 ELA/L Scale Score Descriptive Statistics

Grade	N	Current			Original			Diff*	D
		Mean	Median	SD	Mean	Median	SD		
3	62,753	737.6	739	41.9	739.2	740	42.3	-1.6	-0.04
4	61,139	742.3	742	38.5	744.7	746	37.3	-2.5	-0.06
5	62,463	744.3	743	36.2	744.6	745	35.0	-0.4	-0.01
6	61,173	743.2	744	33.9	742.6	744	32.7	0.6	0.02
7	59,137	746	747	40.8	747.4	749	39.2	-1.4	-0.04
8	58,210	746.6	748	41.5	745.1	746	40.5	1.5	0.04
10	40,163	749	752	46.9	767.1	770	42.7	-18.1	-0.40

Note. SD = standard deviation;

*Diff = Current mean – Original mean.

Table A.14.18 Mathematics Scale Score Descriptive Statistics

Grade	N	Current			Original			Diff*	D
		Mean	Median	SD	Mean	Median	SD		
3	51,957	746.6	747	35.5	748.4	749	36.8	-1.8	-0.05
4	50,277	745.1	747	34.8	746.7	748	34.0	-1.65	-0.05
5	53,131	743.6	743	33.6	744.9	744	33.8	-1.33	-0.04
6	55,342	735.8	736	32.7	736.1	735	32.2	-0.33	-0.01
7	47,340	735.3	735	28.4	735	734	27.7	0.35	0.01
8	28,657	717	715	33.1	713.7	713	31.8	3.27	0.10
Algebra I	35,083	739.7	739	33.4	743.5	742	32.9	-3.82	-0.12
Geometry	3,054	773.4	776.5	24.9	772.6	775	24.7	0.81	0.03
Algebra II	1,576	778.2	779	29.6	782.3	782	28.9	-4.09	-0.14

Note. SD = standard deviation;

*Diff = Current mean – Original mean.

Table A.14.19 ELA/L Writing Claim Score Descriptive Statistics

Grade	N	Current			Original			Diff*	D
		Mean	Median	SD	Mean	Median	SD		
3	62,753	45.3	45	16.8	46.7	47	17.3	-1.4	-0.08
4	61,139	47.2	47	15.5	48.2	48	15.1	-1	-0.07
5	62,463	47.7	47	14.6	48.3	49	14.3	-0.6	-0.04
6	61,173	47.5	47	13.4	47.5	47	13.3	0	0
7	59,137	48.6	49	16.3	49.3	50	16.0	-0.7	-0.04
8	58,210	48.9	48	16.8	48.8	49	16.4	0.1	0.01
10	40,163	49.3	49	18.6	57.2	57	17.8	-7.8	-0.43

Note. SD = standard deviation;

*Diff = Current mean – Original mean.

Table A.14.20 Reading Claim Score Descriptive Statistics

Grade	N	Current			Original			Diff*	D
		Mean	Median	SD	Mean	Median	SD		
3	62,753	29	33	13.5	29.8	32	12.7	-0.8	-0.06
4	61,139	31.6	34	11.7	32.5	34	10.6	-0.9	-0.08
5	62,463	31.0	33	12.6	31.8	33	10.9	-0.8	-0.07
6	61,173	30.5	34	12.4	30.8	33	11.2	-0.3	-0.02
7	59,137	32.4	34	12.4	32.8	35	11.5	-0.4	-0.03
8	58,210	32.0	33	12.9	31.6	34	12.2	0.3	0.03
10	40,163	33.6	35	13.0	37.7	39	11.0	-4.1	-0.34

Note. SD = standard deviation;

*Diff = Current mean – Original mean

Table A.14.21 ELA/L Subclaim Distributions

Form	Level	Percent of Students by Subclaim Performance Level				
		RL	RI	RV	WE	WKL
Current	1	45	42.2	44.9	39.5	38.2
	2	26.3	24.7	23.7	27.3	28.3
	3	28.7	33.1	31.4	33.1	33.4
Original	1	44.5	45.6	44.1	41.9	40
	2	25.2	22.4	24.7	25.4	26.1
	3	30.3	32.1	31.2	32.7	33.9
ES	-	0.02	0.04	0.01	0.03	0.03

Note. RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WE = Writing Written Expression; WKL = Writing Knowledge Language and Conventions; ES = SD = standard deviation;

Table A.14.22 Mathematics Subclaim Distributions

Form	Level	Percent of Students by Subclaim Performance Level			
		A (MC)	C (MR)	D (MP)	B (ASC)
Current	1	33.5	36.7	31	33.5
	2	30.5	27.1	26.4	33.9
	3	36	36.1	42.5	32.6
Original	1	32.6	37.5	32.1	33
	2	29	24.4	25.6	28.3
	3	38.4	38.1	42.2	38.7
ES	-	0.03	0.03	0.01	0.07

Note. MC = Major Content; MR = Mathematical Reasoning; MP = Modeling Practice; ASC = Additional and Supporting Content;

Table A.14.23 ELA/L Subclaim Distribution Comparison: Effect Size

Grade	Subclaim Distribution Effect Size				
	RL	RI	RV	WE	WKL
3	0.01	0.03	0.1	0.14	0.1
4	0.03	0.03	0.08	0.11	0.04
5	0.03	0.03	0.03	0.11	0.08
6	0.02	0.04	0.01	0.03	0.03
7	0.04	0.06	0.05	0.1	0.08
8	0.02	0.05	0.07	0.03	0.04
10	0.19	0.2	0.15	0.15	0.14

Note. RL = Reading Literature; RI = Reading Information; RV = Reading Vocabulary; WE = Writing Written Expression; WKL = Writing Knowledge Language and Conventions.

Table A.14.24 Mathematics Subclaim Distribution Comparison: Effect Size

Grade/ Course	Subclaim Distribution Effect Size			
	A (MC)	C (MR)	D (MP)	B (ASC)
3	0.03	0.01	0.06	0.09
4	0.03	0.02	0.03	0.02
5	0.04	0.11	0.03	0.01
6	0.03	0.03	0.01	0.07
7	0.03	0.19	0.01	0.05
8	0.04	0.13	0.03	0.06
Algebra I	0.05	0.11	0.11	0.06
Geometry	0.03	0.05	0.04	0.02
Algebra II	0.06	0.04	0.16	0.09

Note. MC = Major Content; MR = Mathematical Reasoning; MP = Modeling Practice; ASC = Additional and Supporting Content.

Table A.14.25 ELA/L Longitudinal Scale Score Comparison: Original to Current

Grade	2018 Original SS			2019 Current SS			2019-2018		
	N**	Mean	SD	N**	Mean	SD	Diff*	SD	D
3	265,192	739.7	42.5	257,201	738.5	42.1	-1.2	42.3	-0.03
4	270,283	744.4	37.2	265,584	742.8	38.4	-1.6	37.8	-0.04
5	274,435	743.0	35.3	272,234	744.0	36.5	1.0	35.9	0.03
6	269,341	742.6	33.5	275,880	742.9	34.6	0.3	34.1	0.01
7	266,380	745.5	40.4	270,119	746.7	41.6	1.2	41.0	0.03
8	267,861	744.1	40.5	267,281	746.3	42.2	2.3	41.4	0.05
9	123,153	746.9	39.8	122,200	748.5	40.9	1.6	40.4	0.04
10	118,486	744.2	48.6	118,902	752.3	50.3	8.1	49.5	0.16

Note. SS = scale score; SD = standard deviation;

*Diff = 2019 Current mean – 2018 Original mean.

**All students (not matched samples).

Table A.14.26 ELA/L Longitudinal Scale Score Comparison: Original to Original

Grade	2018 Original			2019 Original			2019-2018		
	N**	Mean	SD	N**	Mean	SD	Diff*	SD	D
3	74,206	735.3	43.4	72,606	737.1	42.5	1.8	43	0.04
4	75,608	741.8	37.9	74,281	741.8	38.2	0	38.1	0
5	74,695	740.4	35.4	75,575	741.8	35.9	1.4	35.7	0.04
6	76,094	739.3	33	79,034	740.6	33.1	1.4	33.1	0.04
7	73,574	742.8	39.8	75,398	745.2	39.6	2.3	39.7	0.06
8	72,661	739.6	40.3	72,976	743	40.8	3.3	40.5	0.08
9	3,449	728.5	39.9	3,468	731.7	40.9	3.2	40.4	0.08
10	72,150	744.2	49.4	74,517	747.8	48.6	3.6	49	0.07

Note. SD = standard deviation;

*Diff = 2019 Current mean – 2018 Original mean.

**All students (not matched samples).

Table A.14.27 Mathematics Longitudinal Scale Score Comparison: Original to Current

Grade	2018 Original			2019 Current			2019-2018		
	N**	Mean	SD	N**	Mean	SD	Diff*	SD	D
3	267,990	742.6	36.7	259,115	743.1	36.5	0.5	36.6	0.01
4	272,625	738.1	33.6	267,191	739.3	34.9	1.2	34.3	0.03
5	275,716	738.2	33.6	273,312	737.8	33.1	-0.4	33.4	-0.01
6	270,735	734.7	31.9	276,652	732.6	32.7	-2.1	32.3	-0.07
7	262,841	736.6	29.5	265,978	737.2	30.6	0.6	30.1	0.02
8	224,120	727.5	37.3	226,912	728.0	38.5	0.6	37.9	0.02
Algebra I	136,154	742.5	37.1	134,975	740.0	36.7	-2.6	36.9	-0.07
Geometry	112,873	732.6	27.4	105,676	731.9	29.5	-0.7	28.4	-0.02
Algebra II	20,658	714.8	33.2	21,414	712.4	34.8	-2.4	34.0	-0.07

Note. SD = standard deviation;

*Diff = 2019 Current mean – 2018 Original mean.

**All students (not matched samples).

Table A.14.28 Mathematics Longitudinal Scale Score Comparison: Original to Original

Grade	2018 Original			2019 Original			2019-2018		
	N**	Mean	SD	N**	Mean	SD	Diff*	SD	D
3	80,700	741.9	39.1	79,361	741.7	38.2	-0.2	38.7	0
4	82,028	737.9	34.8	80,844	739.5	35.8	1.6	35.3	0.05
5	80,953	738	34.9	81,733	738.7	34.4	0.7	34.6	0.02
6	76,153	732.9	32.4	79,141	731.6	32.8	-1.4	32.7	-0.04
7	62,141	731.5	28.9	63,242	731.3	28.7	-0.1	28.8	0
8	41,129	714.6	34.4	40,263	710.2	32.8	-4.3	33.6	-0.13
Algebra I	82,923	736.5	36.3	86,205	734.3	35	-2.1	35.7	-0.06
Geometry	7,110	726.1	24.6	6,967	727.5	27.2	1.5	25.9	0.06
Algebra II	2,841	727.6	33.6	2,943	725.5	34.1	-2.2	33.9	-0.06

Note. SD = standard deviation;

*Diff = 2019 Current mean – 2018 Original mean.

**All students (not matched samples).

Table A.14.29 ELA/L Longitudinal Regression

Grade (Prior Grade)	Sample Size			R2		
	Original- Current	Original- Original	All	Full	Reduced	Change
4 (3)	251,957	70,459	322,416	0.6486	0.648	0.0007
5 (4)	258,568	71,980	330,548	0.6948	0.6948	0
6 (5)	261,213	69,545	330,758	0.6967	0.6966	0.0001
7 (6)	255,849	70,466	326,315	0.7093	0.709	0.0004
8 (7)	253,432	68,542	321,974	0.7263	0.7261	0.0002
9 (8)	109,156	3,015	112,171	0.7306	0.7306	0.0001
10 (8)	103,001	53,963	156,964	0.6598	0.6338	0.026

Table A.14.30 Mathematics Longitudinal Regression

Grade (Prior Grade)	Sample Size			R2		
	Original- Current	Original- Original	All	Full	Reduced	Change
4 (3)	254,114	75,024	329,138	0.7335	0.7332	0.0003
5 (4)	260,243	76,369	336,612	0.7286	0.7283	0.0003
6 (5)	261,817	73,544	335,361	0.7121	0.712	0.0001
7 (6)	251,850	59,342	311,192	0.7391	0.7388	0.0003
8 (7)	213,821	37,357	251,178	0.6821	0.6795	0.0026
A1 (7, 8)	105,010	50,900	155,910	0.6443	0.642	0.0023
GE (A1)	92,531	11,117	103,648	0.6769	0.6707	0.0062
A2 (A1, GE)	60,547	4,136	64,683	0.6793	0.6766	0.0027

Note. A1 = Algebra I; GE = Geometry; A2 = Algebra II.

Table A.14.31 ELA/L Grade 3 Performance Level Comparison

Level	N Count		Percent		Diff
	Current	Original	Current	Original	
1	12,869	12,533	20.5	20	0.5
2	11,212	10,901	17.9	17.4	0.5
3	13,896	12,699	22.1	20.2	1.9
4	21,847	23,625	34.8	37.6	-2.8
5	2,929	2,995	4.7	4.8	-0.1

Note. Cramer's V Effect Size = .03.

Table A.14.32 Mathematics Grade 3 Performance Level Comparison

Level	N Count		Percent		Diff
	Current	Original	Current	Original	
1	5,315	5,430	10.2	10.5	-0.2
2	8,385	7,462	16.1	14.4	1.8
3	12,854	13,100	24.7	25.2	-0.5
4	19,894	19,503	38.3	37.5	0.8
5	5,509	6,462	10.6	12.4	-1.8

Note. Cramer's V Effect Size = .04.

Table A.14.33 Performance Level Comparison Summary: Effect Sizes

ELA/L		Mathematics	
Grade	Cramer's V Effect Size	Grade/ Course	Cramer's V Effect Size
3	0.03	3	0.04
4	0.04	4	0.03
5	0.04	5	0.03
6	0.02	6	0.02
7	0.02	7	0.02
8	0.04	8	0.06
10	0.20	Algebra I	0.09
		Geometry	0.04
		Algebra II	0.07

Table A.14.34 College and Career Readiness Comparison Summary: Effect Sizes

Proportion of Students at or Above the CCR Cut							
ELA/L				Mathematics			
Grade	Current	Original	Cohen's h^*	Grade/Course	Current	Original	Cohen's h^*
3	0.39	0.42	-0.06	3	0.49	0.50	-0.02
4	0.43	0.46	-0.05	4	0.46	0.48	-0.03
5	0.45	0.46	-0.03	5	0.43	0.44	-0.02
6	0.43	0.43	-0.01	6	0.34	0.34	0
7	0.48	0.50	-0.04	7	0.30	0.30	0
8	0.48	0.47	0.01	8	0.18	0.14	0.09
10	0.51	0.68	-0.35	Algebra I	0.38	0.42	-0.09
				Geometry	0.87	0.86	0.03
				Algebra II	0.86	0.89	-0.09

*Computed as Current proportion – Original proportion.

Table A.14.35 ELA/L Classification Accuracy

Grade	Performance Level Classification			College- and Career-Readiness* Classification		
	Current	Original	Cohen's h	Current	Original	Cohen's h
3	0.71	0.75	-0.10	0.90	0.92	-0.05
4	0.68	0.74	-0.13	0.89	0.91	-0.06
5	0.72	0.78	-0.15	0.90	0.92	-0.08
6	0.74	0.79	-0.13	0.91	0.92	-0.06
7	0.71	0.77	-0.13	0.91	0.93	-0.06
8	0.71	0.77	-0.13	0.91	0.93	-0.07
10	0.67	0.77	-0.23	0.90	0.93	-0.10

Table A.14.36 ELA/L Classification Consistency

Grade	Performance Level Classification			College and Career Readiness* Classification		
	Current	Original	Cohen's h	Current	Original	Cohen's h
3	0.61	0.66	-0.10	0.86	0.88	-0.06
4	0.57	0.64	-0.15	0.85	0.88	-0.07
5	0.62	0.70	-0.17	0.86	0.89	-0.09
6	0.64	0.71	-0.15	0.87	0.89	-0.08
7	0.60	0.67	-0.15	0.87	0.90	-0.07
8	0.62	0.69	-0.15	0.87	0.90	-0.08
10	0.57	0.69	-0.25	0.86	0.90	-0.12

Table A.14.37 Mathematics Classification Accuracy

Grade/ Course	Performance Level Classification			College- and Career-Readiness* Classification		
	Current	Original	Cohen's <i>h</i>	Current	Original	Cohen's <i>h</i>
3	0.75	0.78	-0.06	0.91	0.93	-0.05
4	0.78	0.80	-0.05	0.92	0.92	-0.02
5	0.77	0.79	-0.04	0.92	0.93	-0.02
6	0.77	0.81	-0.10	0.92	0.94	-0.05
7	0.77	0.79	-0.04	0.92	0.93	-0.03
8	0.71	0.73	-0.04	0.92	0.93	-0.06
Algebra I	0.74	0.79	-0.11	0.91	0.92	-0.06
Geometry	0.81	0.85	-0.11	0.96	0.96	-0.03
Algebra II	0.82	0.86	-0.1	0.92	0.95	-0.10

Table A.14.38 Mathematics Classification Consistency

Grade/ Course	Performance Level Classification			College- and Career- Readiness* Classification		
	Current	Original	<i>h</i>	Current	Original	<i>h</i>
3	0.66	0.69	-0.07	0.88	0.90	-0.06
4	0.69	0.72	-0.06	0.89	0.89	-0.03
5	0.68	0.70	-0.05	0.89	0.90	-0.02
6	0.68	0.73	-0.12	0.89	0.91	-0.06
7	0.68	0.70	-0.05	0.89	0.90	-0.04
8	0.61	0.63	-0.05	0.88	0.90	-0.07
Algebra I	0.65	0.70	-0.13	0.87	0.89	-0.07
Geometry	0.73	0.78	-0.13	0.94	0.94	-0.04
Algebra II	0.74	0.79	-0.12	0.89	0.92	-0.12

Table A.14.39 ELA/L Grade 6 Performance Level Comparison

Level	Original to Current			Original to Original		
	Current States 2018	Current States 2019	Diff	Original States 2018	Original States 2019	Diff
1	10.2	11.3	1.1	12.4	12.6	0.2
2	20.1	17.9	-2.2	21.3	18.8	-2.5
3	28	28.5	0.5	27.7	27.5	-0.2
4	33.3	33.8	0.5	32.1	34.3	2.2
5	8.3	8.4	0.1	6.6	6.8	0.2
Cramer's V Effect Size = .03				Cramer's V Effect Size = .03		

Table A.14.40 Mathematics Grade 6 Performance Level Comparison

Level	Original to Current			Original to Original		
	Current States 2018	Current States 2019	Diff	Original States 2018	Original States 2019	Diff
1	13.4	14.4	1	15.7	17.5	1.8
2	25.9	28.0	2.1	26.1	25.9	-0.2
3	28.4	27.4	-0.9	26.8	26.8	0
4	27.4	25.5	-1.9	26.9	25.4	-1.5
5	5	4.7	-0.3	4.5	4.3	-0.2
Cramer's V Effect Size = .03				Cramer's V Effect Size = .03		

Table A.14.41 Performance Level Comparison Summary: Effect Sizes

ELA/L		Mathematics			
Grade	Original to Current	Original to Original	Grade/Course	Original to Current	Original to Original
3	0.02	0.03	3	0.04	0.05
4	0.03	0.02	4	0.05	0.02
5	0.02	0.03	5	0.06	0.05
6	0.03	0.03	6	0.03	0.03
7	0.02	0.03	7	0.03	0.06
8	0.04	0.05	8	0.04	0.08
9	0.04	0.05	Algebra I	0.10	0.05
10	0.09	0.04	Geometry	0.07	0.06
			Algebra II	0.05	0.05

Table A.14.42 ELA/L Reading Claim Reliability

Grade	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	46	0.9	30	0.86	0.85	0.01
4	64	0.88	42	0.83	0.83	0
5	64	0.9	42	0.85	0.86	-0.01
6	64	0.91	42	0.87	0.87	0
7	64	0.91	42	0.86	0.87	-0.01
8	64	0.9	42	0.85	0.86	-0.01
10	64	0.89	42	0.82	0.84	-0.02

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

Table A.14.43 ELA/L Writing Claim Reliability

Grade	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	36	0.85	24	0.79	0.79	0
4	42	0.86	28	0.8	0.8	0
5	42	0.86	29	0.8	0.81	-0.01
6	45	0.87	30	0.82	0.82	0
7	45	0.88	30	0.83	0.83	0
8	45	0.89	30	0.85	0.84	0.01
10	45	0.88	30	0.84	0.83	0.01

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

Table A.14.44 ELA/L Reading Information (RI) Subclaim Reliability

Grade	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	17	0.74	11	0.68	0.65	0.03
4	26	0.76	16	0.62	0.66	-0.04
5	23	0.75	14	0.56	0.65	-0.09
6	24	0.76	16	0.67	0.68	-0.01
7	24	0.81	14	0.66	0.71	-0.05
8	21	0.78	15	0.71	0.72	-0.01
10	30	0.8	19	0.68	0.72	-0.04

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

Table A.14.45 ELA/L Reading Literature (RL) Subclaim Reliability

Grade	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	19	0.8	11	0.71	0.7	0.01
4	26	0.73	17	0.66	0.64	0.02
5	26	0.79	17	0.74	0.71	0.03
6	26	0.84	18	0.76	0.78	-0.02
7	25	0.79	17	0.7	0.72	-0.02
8	26	0.79	16	0.69	0.7	-0.01
10	20	0.7	14	0.61	0.62	-0.01

Note. SB = Spearman Brown.

*Diff: Current Alpha – Spearman Brown (SB) Prophecy

Table A.14.46 ELA/L Reading Vocabulary (RV) Subclaim Reliability

Grade	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	10	0.68	8	0.61	0.63	-0.02
4	12	0.61	9	0.56	0.54	0.02
5	15	0.75	11	0.67	0.69	-0.02
6	14	0.72	8	0.58	0.56	-0.02
7	15	0.66	11	0.62	0.59	0.03
8	17	0.69	11	0.53	0.59	-0.06
10	14	0.6	10	0.47	0.52	-0.05

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

Table A.14.47 ELA/L Writing Knowledge and Conventions (WKL) Subclaim Reliability

Grade	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	9	0.87	6	0.82	0.82	0
4	9	0.88	6	0.84	0.83	0.01
5	9	0.88	6	0.84	0.83	0.01
6	9	0.89	6	0.85	0.84	0.01
7	9	0.89	6	0.86	0.84	0.02
8	9	0.91	6	0.87	0.87	0
10	9	0.89	6	0.86	0.84	0.02

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy

Table A.14.48 ELA/L Written Expression (WE) Subclaim Reliability

Grade	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	27	0.81	18	0.74	0.74	0
4	33	0.83	22	0.77	0.76	0.01
5	33	0.81	23	0.72	0.75	-0.03
6	36	0.86	24	0.81	0.8	0.01
7	36	0.88	24	0.85	0.83	0.02
8	36	0.9	24	0.86	0.86	0
10	36	0.88	24	0.85	0.83	0.02

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

Table A.14.49 Mathematics Subclaim A Reliability

Grade/Course	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	28	0.91	20	0.86	0.88	-0.02
4	31	0.9	21	0.86	0.86	0
5	30	0.9	20	0.86	0.86	0
6	26	0.88	20	0.83	0.85	-0.02
7	29	0.87	20	0.84	0.82	0.02
8	27	0.77	20	0.74	0.71	0.03
Algebra I	26	0.79	17	0.72	0.71	0.01
Geometry	30	0.84	18	0.79	0.76	0.03
Algebra II	25	0.74	16	0.66	0.65	0.01

Note. SB = Spearman Brown.

*Diff: Current Alpha – Spearman Brown (SB) Prophecy.

Table A.14.50 Mathematics Subclaim B Reliability

Grade/Course	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	12	0.76	10	0.69	0.73	-0.04
4	9	0.72	9	0.72	0.72	0
5	10	0.71	10	0.7	0.71	-0.01
6	14	0.77	10	0.67	0.71	-0.04
7	11	0.67	10	0.64	0.65	-0.01
8	13	0.53	10	0.49	0.46	0.03
Algebra I	17	0.73	9	0.64	0.59	0.05
Geometry	19	0.79	12	0.65	0.7	-0.05
Algebra II	20	0.7	12	0.55	0.58	-0.03

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

Table A.14.51 Mathematics Subclaim C Reliability

Grade/Course	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	14	0.62	10	0.48	0.54	-0.06
4	14	0.79	10	0.76	0.73	0.03
5	14	0.71	10	0.62	0.64	-0.02
6	14	0.78	10	0.71	0.72	-0.01
7	14	0.64	10	0.52	0.56	-0.04
8	14	0.59	10	0.54	0.51	0.03
Algebra I	14	0.75	10	0.7	0.68	0.02
Geometry	14	0.64	10	0.6	0.56	0.04
Algebra II	14	0.55	10	0.44	0.47	-0.03

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

Table A.14.52 Mathematics Subclaim D Reliability

Grade/Course	Original		Current		SB	Diff*
	Points	Alpha	Points	Alpha		
3	12	0.76	12	0.75	-	-
4	12	0.66	12	0.66	-	-
5	12	0.74	12	0.73	-	-
6	12	0.71	12	0.69	-	-
7	12	0.73	12	0.74	-	-
8	12	0.5	12	0.52	-	-
Algebra I	18	0.75	15	0.69	0.71	-0.02
Geometry	18	0.7	15	0.64	0.66	-0.02
Algebra II	18	0.59	15	0.56	0.55	0.01

Note. SB = Spearman Brown.

*Diff = Current Alpha – Spearman Brown (SB) Prophecy.

Appendix 15: Growth

Appendix 15 provides the summary growth results for subgroups for grades 4 through 11 ELA/L and mathematics 4 through 8 and high school. Grade 9 ELA/L, Algebra II, and Integrated Mathematics I and II do not have sufficient sample sizes for subgroup summary analysis.

Table A.15.1 Summary of Student Growth Percentile Estimates for Subgroups: Grade 4 ELA/L

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	49,196	49.67	14.11	50
Female	47,458	50.31	14.06	50
Ethnicity				
White	51,770	51.44	13.99	52
African American	12,578	42.71	14.41	40
Asian/Pacific Islander	5,241	59.02	13.78	63
American Indian/Alaska Native	207	45.45	14.43	44
Hispanic	21,886	48.48	14.24	48
Multiple	4,703	50.78	13.95	51
Special Instruction Needs				
Economically Disadvantaged	40,641	45.49	14.24	44
Not Economically Disadvantaged	56,014	53.25	13.98	55
English Learner (EL)	14,992	47.00	14.37	46
Non-English Learner	81,663	50.54	14.03	51
Students with Disabilities (SWD)	16,915	41.49	14.39	38
Students without Disabilities	79,740	51.79	14.02	53

Table A.15.2 Summary of Student Growth Percentile Estimates for Subgroups: Grade 5 ELA/L

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	50,517	48.42	14.33	48
Female	48,422	51.64	13.91	52
Ethnicity				
White	53,370	50.53	13.85	51
African American	12,790	45.64	14.89	44
Asian/Pacific Islander	5,293	60.02	13.40	64
American Indian/Alaska Native	186	49.89	14.28	50.5
Hispanic	22,422	48.79	14.55	48
Multiple	4,596	50.66	14.12	51
Special Instruction Needs				
Economically Disadvantaged	41,252	46.51	14.59	45
Not Economically Disadvantaged	57,692	52.49	13.80	53
English Learner (EL)	12,853	46.21	15.07	44
Non-English Learner	86,091	50.56	13.99	51
Students with Disabilities (SWD)	17,328	41.53	15.22	38
Students without Disabilities	81,616	51.79	13.90	53

Table A.15.3 Summary of Student Growth Percentile Estimates for Subgroups: Grade 6 ELA/L

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	51,180	47.85	13.79	47
Female	48,204	52.28	13.52	53
Ethnicity				
White	53,631	51.49	13.43	52
African American	12,982	45.27	14.30	43
Asian/Pacific Islander	5,245	56.51	13.44	59
American Indian/Alaska Native	173	47.55	13.67	46
Hispanic	22,641	47.79	13.91	47
Multiple	4,473	50.00	13.61	50
Special Instruction Needs				
Economically Disadvantaged	41,534	46.56	14.05	45
Not Economically Disadvantaged	57,858	52.47	13.38	54
English Learner (EL)	10,517	43.85	14.59	41
Non-English Learner	88,875	50.73	13.55	51
Students with Disabilities (SWD)	17,314	41.29	14.66	38
Students without Disabilities	82,078	51.84	13.45	53

Table A.15.4 Summary of Student Growth Percentile Estimates for Subgroups: Grade 7 ELA/L

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	50,999	47.61	13.98	47
Female	48,315	52.51	13.80	54
Ethnicity				
White	53,630	51.74	13.84	52
African American	12,925	44.79	13.95	43
Asian/Pacific Islander	5,264	58.51	13.88	62
American Indian/Alaska Native	213	51.16	13.55	49
Hispanic	22,702	46.93	13.97	46
Multiple	4,322	49.70	13.96	50
Special Instruction Needs				
Economically Disadvantaged	40,845	46.27	13.91	45
Not Economically Disadvantaged	58,481	52.60	13.88	54
English Learner (EL)	9,526	42.89	14.16	40
Non-English Learner	89,800	50.75	13.86	51
Students with Disabilities (SWD)	17,135	42.45	14.14	40
Students without Disabilities	82,191	51.57	13.84	52

Table A.15.5 Summary of Student Growth Percentile Estimates for Subgroups: Grade 8 ELA/L

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	50,489	46.95	13.94	46
Female	47,703	53.18	13.84	55
Ethnicity				
White	53,427	50.84	13.84	51
African American	13,197	46.73	14.10	45
Asian/Pacific Islander	5,096	57.57	13.93	61
American Indian/Alaska Native	218	48.89	13.58	50
Hispanic	21,949	48.16	13.88	47
Multiple	4,089	49.28	13.97	49
Special Instruction Needs				
Economically Disadvantaged	40,330	47.40	13.91	46
Not Economically Disadvantaged	57,871	51.77	13.88	52
English Learner (EL)	8,635	44.34	14.22	42
Non-English Learner	89,566	50.52	13.86	51
Students with Disabilities (SWD)	16,560	43.94	14.37	41
Students without Disabilities	81,641	51.20	13.79	52

Table A.15.6 Summary of Student Growth Percentile Estimates for Subgroups: Grade 4 Mathematics

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	48,093	50.44	13.48	51
Female	46,365	49.71	13.49	50
Ethnicity				
White	51,416	50.45	13.05	51
African American	12,328	44.39	14.79	42
Asian/Pacific Islander	5,187	59.13	13.14	63
American Indian/Alaska Native	197	49.16	13.55	48
Hispanic	20,401	50.15	13.88	50
Multiple	4,661	51.37	13.49	51
Special Instruction Needs				
Economically Disadvantaged	39,044	46.91	14.05	46
Not Economically Disadvantaged	55,415	52.32	13.08	53
English Learner (EL)	13,590	49.62	14.13	50
Non-English Learner	80,869	50.16	13.37	50
Students with Disabilities (SWD)	16,526	43.36	14.15	41
Students without Disabilities	77,933	51.51	13.34	52
Spanish Language Form	1,268	42.63	14.57	41

Table A.15.7 Summary of Student Growth Percentile Estimates for Subgroups: Grade 5 Mathematics

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	49,440	49.34	14.11	49
Female	47,335	50.72	14.18	51
Ethnicity				
White	53,050	50.56	13.71	51
African American	12,483	44.07	15.47	41
Asian/Pacific Islander	5,260	59.32	13.52	63
American Indian/Alaska Native	175	53.42	13.85	54
Hispanic	20,963	49.89	14.64	50
Multiple	4,557	49.74	14.06	50
Special Instruction Needs				
Economically Disadvantaged	39,735	46.91	14.83	46
Not Economically Disadvantaged	57,045	52.17	13.67	53
English Learner (EL)	11,639	48.20	15.20	47
Non-English Learner	85,141	50.26	14.00	50
Students with Disabilities (SWD)	16,962	41.68	14.93	38
Students without Disabilities	79,818	51.78	13.98	53
Spanish Language Form	1,212	43.79	15.07	42

Table A.15.8 Summary of Student Growth Percentile Estimates for Subgroups: Grade 6 Mathematics

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	50,175	49.59	14.66	49
Female	47,242	50.45	14.61	51
Ethnicity				
White	53,275	51.66	14.29	52
African American	12,619	42.78	15.75	40
Asian/Pacific Islander	5,205	58.76	13.93	63
American Indian/Alaska Native	165	49.52	14.29	46
Hispanic	21,490	48.28	15.02	47
Multiple	4,428	48.98	14.56	49
Special Instruction Needs				
Economically Disadvantaged	40,200	46.47	15.26	45
Not Economically Disadvantaged	57,224	52.49	14.19	53
English Learner (EL)	9,684	42.27	15.93	39
Non-English Learner	87,740	50.86	14.49	51
Students with Disabilities (SWD)	16,941	40.97	15.40	37
Students without Disabilities	80,483	51.91	14.47	53
Spanish Language Form	858	44.62	15.83	42

Table A.15.9 Summary of Student Growth Percentile Estimates for Subgroups: Grade 7 Mathematics

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	48,644	50.81	15.24	51
Female	46,050	49.07	15.45	49
Ethnicity				
White	52,045	50.48	15.14	51
African American	12,299	45.45	15.96	44
Asian/Pacific Islander	4,958	56.11	15.30	58
American Indian/Alaska Native	196	52.79	15.55	52.5
Hispanic	21,179	50.26	15.47	50
Multiple	3,836	48.44	15.42	48
Special Instruction Needs				
Economically Disadvantaged	39,749	47.63	15.60	47
Not Economically Disadvantaged	54,955	51.65	15.15	52
English Learner (EL)	8,686	45.73	16.11	44
Non-English Learner	86,018	50.39	15.26	51
Students with Disabilities (SWD)	16,357	41.86	15.85	38
Students without Disabilities	78,347	51.65	15.23	52
Spanish Language Form	444	44.72	16.10	42

Table A.15.10 Summary of Student Growth Percentile Estimates for Subgroups: Grade 8 Mathematics

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	48,284	48.89	16.65	48
Female	45,576	50.97	16.85	51
Ethnicity				
White	51,836	51.18	16.10	52
African American	12,607	44.40	18.68	42
Asian/Pacific Islander	4,805	56.59	14.72	59
American Indian/Alaska Native	197	48.42	16.98	47
Hispanic	20,646	48.70	17.64	48
Multiple	3,614	49.08	16.94	48
Special Instruction Needs				
Economically Disadvantaged	39,280	47.12	17.92	46
Not Economically Disadvantaged	54,589	51.90	15.90	53
English Learner (EL)	7,980	46.14	19.03	44
Non-English Learner	85,889	50.25	16.53	50
Students with Disabilities (SWD)	15,864	44.92	18.51	43
Students without Disabilities	78,005	50.91	16.39	51
Spanish Language Form	273	47.37	19.97	44

Table A.15.11 Summary of Student Growth Percentile Estimates for Subgroups: Algebra I

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	742	50.16	15.82	49.5
Female	713	48.89	15.66	49
Ethnicity				
White	531	52.68	15.49	54
African American	142	46.05	15.84	46
Asian/Pacific Islander	120	51.17	15.92	54.5
American Indian/Alaska Native	—	—	—	—
Hispanic	384	46.43	15.91	45
Multiple	243	48.91	15.75	47
Special Instruction Needs				
Economically Disadvantaged	—	—	—	—
Not Economically Disadvantaged	1,455	49.54	15.74	49
English Learner (EL)	135	44.88	16.22	40
Non-English Learner	1,320	50.02	15.69	51
Students with Disabilities (SWD)	291	44.67	16.43	42
Students without Disabilities	1,164	50.76	15.57	51
Spanish Language Form				
	—	—	—	—

Note. “—” indicates insufficient sample for student growth percentile calculation for these tests.

Table A.15.12 Summary of Student Growth Percentile Estimates for Subgroups: Geometry

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	969	49.60	15.30	50
Female	892	49.66	15.18	49
Ethnicity				
White	721	53.28	14.99	56
African American	164	46.25	16.17	46
Asian/Pacific Islander	189	45.51	15.31	43
American Indian/Alaska Native	—	—	—	—
Hispanic	398	46.05	15.61	43
Multiple	343	50.03	14.88	49
Special Instruction Needs				
Economically Disadvantaged	—	—	—	—
Not Economically Disadvantaged	1,861	49.63	15.24	50
English Learner (EL)	97	39.79	16.69	34
Non-English Learner	1,764	50.17	15.16	50
Students with Disabilities (SWD)	254	43.50	16.93	40
Students without Disabilities	1,607	50.60	14.97	51
Spanish Language Form				
	—	—	—	—

Note. “—” indicates insufficient sample for student growth percentile calculation for these tests.

Table A.15.13 Summary of Student Growth Percentile Estimates for Subgroups: Algebra II

	Total Sample Size	Average SGP	Average Standard Error	Median SGP
Gender				
Male	802	51.01	15.56	50
Female	756	48.07	16.25	48
Ethnicity				
White	627	52.15	15.56	52
African American	138	45.99	17.19	45
Asian/Pacific Islander	165	47.56	15.42	49
American Indian/Alaska Native	—	—	—	—
Hispanic	349	46.01	16.51	42
Multiple	240	51.92	15.48	55
Special Instruction Needs				
Economically Disadvantaged	—	—	—	—
Not Economically Disadvantaged	1,558	49.58	15.89	49
English Learner (EL)	69	43.22	17.97	36
Non-English Learner	1,489	49.88	15.80	50
Students with Disabilities (SWD)	202	44.92	17.11	42.5
Students without Disabilities	1,356	50.28	15.71	50
Spanish Language Form				
	—	—	—	—

Note. “—” indicates insufficient sample for student growth percentile calculation for these tests.