



Wisconsin Forward Exam Spring 2023 Technical Report

**Submitted to
Wisconsin Department of Public Instruction
October 2023**

DATA RECOGNITION
DRC
CORPORATION

Copyright

Developed and published under contract with the Wisconsin Department of Public Instruction by Data Recognition Corporation, 13490 Bass Lake Road, Maple Grove, MN 55311.

Copyright © 2023 by the Wisconsin Department of Public Instruction. All rights reserved. Only State of Wisconsin educators and citizens may copy, download and/or print the document, located online at <http://dpi.wi.gov>. Any other use or reproduction of this document, in whole or in part, requires written permission of the Wisconsin Department of Public Instruction.

Foreword

The technical information herein is intended for use by those who evaluate tests, interpret scores, or use test results in making educational decisions. It is assumed that the reader has technical knowledge of test construction and measurement procedures as stated in the *Standards for Educational and Psychological Testing* (American Educational Research Association, American Psychological Association, & National Council on Measurement in Education, 2014).

Table of Contents

Copyright	i
Foreword	ii
Appendices	iv
List of Tables	v
List of Figures	ix
Executive Summary	1
E.1 Overview of the Wisconsin Forward Exam.....	1
E.2 Administration	2
E.3 Student Performance.....	2
E.4 Validity of Intended Interpretation of Test Scores	4
Part 1: Overview	7
1.1 Historical Background.....	7
1.2 Uses of Test Scores	9
1.3 Technical Report Structure.....	11
Part 2: Validity Framework	15
2.1 Sources of Validity Evidence	15
2.2 Summary of Validity Evidence for Wisconsin Forward Exam.....	16
Part 3: Test Content and Test Development	20
3.1 Test Blueprints	21
3.2 Test Design.....	21
3.3 Universal Design.....	22
3.4 Item Development Process.....	24
3.5 Field-Testing	26
3.6 Form Development.....	27
3.7 DPI Approvals.....	31
3.8 Summary	31
Part 4: Test Administration	38
4.1 Student Participation	39
4.2 Standardized Test Administration.....	40
4.3 Accessibility Resources.....	40
4.4 Test Security.....	44
4.5 Test Administration Training	46
4.6 Summary	49
Part 5: Scoring	59
5.1 Multiple-Choice and Multi-Select Item Scoring Process.....	59
5.2 Technology-Enhanced, Short-Answer, and Evidence-Based Selected Response Item Scoring Process....	59
5.3 Scoring of Text-Dependent Analysis Items	60
5.4 Inter-rater Reliability.....	64
5.5 Summary	65
Part 6: Psychometric Analyses	69
6.1 Overview of the Operational Test Data Analysis.....	69
6.2 Classical Item Analysis: Item Level Statistics	69
6.3 Test-Level Statistics	73
6.4 Item Response Theory Methodology	75
6.5 Summary	93

Part 7: Standard Setting	172
7.1 Background Information	172
7.2 Standard Setting Methodology and Process.....	173
7.3 Performance Level Descriptors	174
7.4 Cut Scores	175
7.5 Summary	175
Part 8: Studies of Reliability	177
8.1 Measures of Internal Consistency and Standard Error of Measurement.....	178
8.2 Classification Consistency and Accuracy	182
8.3 Inter-rater Reliability for TDA Items	186
8.4 Summary	188
Part 9: Studies of Construct-Related Validity	203
9.1 Differential Item Functioning.....	204
9.2 Validity Evidence Based on Internal Test Structure	207
9.3 Validity Evidence Based on Relationship with Other Variables	209
9.4 Test Integrity: Data Forensic Analyses	213
9.5 Summary	213
Part 10: Test Results	227
10.1 Types of Reports	227
10.2 Scale Scores Summary Statistics.....	231
10.3 Performance Level Classifications.....	235
10.4 Standard Performance Index for Content Standards	238
10.5 Longitudinal Comparisons of Test Scores	241
10.6 Summary	243
Part 11: Summary and Recommendations	278
References	279

Appendices

Appendix A: Item Review Training Slides
Appendix B: Data Review Training Slides
Appendix C: Spring 2023 English Language Arts Operational Test Maps
Appendix D: Spring 2023 Mathematics Operational Test Maps
Appendix E: Spring 2023 Science Operational Test Maps
Appendix F: Spring 2023 Social Studies Operational Test Maps
Appendix G: Test Participation Rates by Subgroup
Appendix H: Classical Item Analysis Results
Appendix I: Conditional Standard Error of Measurement with Cut Scores
Appendix J: Classification Consistency and Accuracy Indices by Subgroup
Appendix K: Wisconsin Standard Performance Index Score Computation

List of Tables

Table E-1 Test Participation Rates in Spring 2023.....	5
Table E-2 Percentage of Students Classified as <i>Proficient</i> or <i>Advanced</i> from 2016 through 2023, English Language Arts.....	6
Table E-3 Percentage of Students Classified as <i>Proficient</i> or <i>Advanced</i> from 2016 through 2023, Mathematics.....	6
Table E-4 Percentage of Students Classified as <i>Proficient</i> or <i>Advanced</i> from 2019 through 2023, Science.....	6
Table E-5 Percentage of Students Classified as <i>Proficient</i> or <i>Advanced</i> in 2022 and 2023, Social Studies.....	6
Table 3-1 Test Blueprints for English Language Arts Grades 3–8.....	32
Table 3-2 Test Blueprints for Mathematics Grades 3–8.....	32
Table 3-3 Test Blueprints for Science Grades 4 and 8.....	33
Table 3-4 Test Blueprints for Social Studies Grades 4, 8, and 10.....	33
Table 3-5 Item Type Descriptions for Items on the Wisconsin Forward Exam.....	34
Table 3-6 Test Design for English Language Arts.....	35
Table 3-7 Test Design for Mathematics.....	35
Table 3-8 Test Design for Science.....	35
Table 3-9 Test Design for Social Studies.....	36
Table 3-10 Elements of Universal Design.....	36
Table 3-11 College- and Career-Ready Item Bank Development Activities.....	36
Table 3-12 Items Reviewed during Summer 2021 Item Review.....	37
Table 3-13 Items Reviewed during Summer 2022 Item Data Review.....	37
Table 4-1 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 3.....	49
Table 4-2 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 4.....	50
Table 4-3 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 5.....	51
Table 4-4 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 6.....	52
Table 4-5 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 7.....	53
Table 4-6 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 8.....	54
Table 4-7 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 10.....	55
Table 4-8 Summary Table of Manual Materials.....	56
Table 5-1 TDA Item Scoring Guidelines, Grades 3–8.....	65
Table 5-2 TDA Item Non-scorable Codes, Grades 3–8.....	67
Table 5-3 TDA Item Score Distribution.....	67
Table 5-4 TDA Item Score Distribution: AI Engine vs. Human Scorer.....	68
Table 5-5 TDA Item Percentage Score Distribution: AI Engine vs. Human Scorer.....	68
Table 6-A Example of Item Parameters for a Test.....	88
Table 6-B Example of Item Response Pattern.....	88
Table 6-1 Summary of Flagged Operational Items on the Wisconsin Forward Exam.....	93
Table 6-2 Items Flagged for Classical Item Analysis Statistics, English Language Arts.....	94

Table 6-3 Items Flagged for Classical Item Analysis Statistics, Mathematics.....	95
Table 6-4 Items Flagged for Classical Item Analysis Statistics, Science and Social Studies	96
Table 6-5 Percentage of Students Attempting Last Operational Item in Test	96
Table 6-6 Item Analysis, English Language Arts Grade 3	97
Table 6-7 Item Analysis, English Language Arts Grade 4	98
Table 6-8 Item Analysis, English Language Arts Grade 5	99
Table 6-9 Item Analysis, English Language Arts Grade 6	100
Table 6-10 Item Analysis, English Language Arts Grade 7	101
Table 6-11 Item Analysis, English Language Arts Grade 8	102
Table 6-12 Item Analysis, Mathematics Grade 3	103
Table 6-13 Item Analysis, Mathematics Grade 4	104
Table 6-14 Item Analysis, Mathematics Grade 5	106
Table 6-15 Item Analysis, Mathematics Grade 6	108
Table 6-16 Item Analysis, Mathematics Grade 7	110
Table 6-17 Item Analysis, Mathematics Grade 8	112
Table 6-18 Item Analysis, Science Grade 4.....	114
Table 6-19 Item Analysis, Science Grade 8.....	115
Table 6-20 Item Analysis, Social Studies Grade 4	116
Table 6-21 Item Analysis, Social Studies Grade 8	117
Table 6-22 Item Analysis, Social Studies Grade 10	118
Table 6-23 Test-Level Descriptive Statistics.....	119
Table 6-24 Calibration Sample Demographics Compared to Population, English Language Arts.....	120
Table 6-25 Calibration Sample Demographics Compared to Population, Mathematics	126
Table 6-26 Calibration Sample Demographics Compared to Population, Science	132
Table 6-27 Calibration Sample Demographics Compared to Population, Social Studies.....	134
Table 6-28 Items Flagged Based on Yen’s Q1	137
Table 6-29 Equating Evaluation Results, Stocking and Lord Method	138
Table 6-30 Scale Transformation Constants.....	138
Table 6-31 Scoring Table for English Language Arts Grade 3	139
Table 6-32 Scoring Table for English Language Arts Grade 4	140
Table 6-33 Scoring Table for English Language Arts Grade 5	141
Table 6-34 Scoring Table for English Language Arts Grade 6	142
Table 6-35 Scoring Table for English Language Arts Grade 7	143
Table 6-36 Scoring Table for English Language Arts Grade 8	144
Table 6-37 Scoring Table for Mathematics Grade 3	145
Table 6-38 Scoring Table for Mathematics Grade 4	146
Table 6-39 Scoring Table for Mathematics Grade 5	147
Table 6-40 Scoring Table for Mathematics Grade 6	148
Table 6-41 Scoring Table for Mathematics Grade 7	149

Table 6-42 Scoring Table for Mathematics Grade 8	150
Table 6-43 Scoring Table for Science Grade 4.....	151
Table 6-44 Scoring Table for Science Grade 8.....	152
Table 6-45 Scoring Table for Social Studies Grade 4	153
Table 6-46 Scoring Table for Social Studies Grade 8	154
Table 6-47 Scoring Table for Social Studies Grade 10	155
Table 6-48 Numbers and Percentages of Students at LOSS and HOSS.....	155
Table 7-1 Policy Performance Level Descriptors for the Wisconsin Forward Exam.....	176
Table 7-2 Wisconsin Forward Exam Cut Scores.....	176
Table 8-A Example Contingency Table with Three Cut Scores.....	182
Table 8-B Example Classification Table for One Cut Point (C1)	184
Table 8-C Data Structure 1: Enumeration by Response	186
Table 8-D Data Structure 2: Cross-Tabulation of Score 1 and Score 2.....	187
Table 8-1 Cronbach’s Alpha Reliability Coefficients for Total Group and Subgroups	189
Table 8-2 Standard Error of Measurement for Total Group and Subgroups	190
Table 8-3 Cronbach’s Alpha Reliability Coefficients for Content Standards and Domains	191
Table 8-4 Standard Error of Measurement per Content Standards and Domains.....	192
Table 8-5 Classification Consistency and Classification Accuracy for English Language Arts Grade 3	193
Table 8-6 Classification Consistency and Classification Accuracy for English Language Arts Grade 4	193
Table 8-7 Classification Consistency and Classification Accuracy for English Language Arts Grade 5	194
Table 8-8 Classification Consistency and Classification Accuracy for English Language Arts Grade 6	194
Table 8-9 Classification Consistency and Classification Accuracy for English Language Arts Grade 7	195
Table 8-10 Classification Consistency and Classification Accuracy for English Language Arts Grade 8	195
Table 8-11 Classification Consistency and Classification Accuracy for Mathematics Grade 3	196
Table 8-12 Classification Consistency and Classification Accuracy for Mathematics Grade 4	196
Table 8-13 Classification Consistency and Classification Accuracy for Mathematics Grade 5	197
Table 8-14 Classification Consistency and Classification Accuracy for Mathematics Grade 6	197
Table 8-15 Classification Consistency and Classification Accuracy for Mathematics Grade 7	198
Table 8-16 Classification Consistency and Classification Accuracy for Mathematics Grade 8	198
Table 8-17 Classification Consistency and Classification Accuracy for Science Grade 4.....	199
Table 8-18 Classification Consistency and Classification Accuracy for Science Grade 8.....	199
Table 8-19 Classification Consistency and Classification Accuracy for Social Studies Grade 4	200
Table 8-20 Classification Consistency and Classification Accuracy for Social Studies Grade 8	200
Table 8-21 Classification Consistency and Classification Accuracy for Social Studies Grade 10	201
Table 8-22 Inter-Rater Reliability, English Language Arts.....	202
Table 9-1 Items Flagged for DIF in English Language Arts	215
Table 9-2 Items Flagged for DIF in Mathematics	216
Table 9-3 Items Flagged for DIF in Science.....	217
Table 9-4 Items Flagged for DIF in Social Studies	217

Table 9-5 Correlations between English Language Arts Test Domains..... 218

Table 9-6 Correlations between Content Standards, English Language Arts..... 219

Table 9-7 Correlations between Content Standards, Mathematics..... 220

Table 9-8 Correlations between Content Standards, Science..... 220

Table 9-9 Correlations between Content Standards, Social Studies..... 221

Table 9-10 Principal Components Analysis..... 221

Table 9-11 Correlations between Content Area Scale Scores..... 222

Table 9-12 Correlations between Content Area Scale Scores by Gender..... 222

Table 9-13 Correlations between Content Area Scale Scores by Ethnicity/Race..... 223

Table 9-14 Correlations between Content Area Scale Scores by English Proficiency Status..... 224

Table 9-15 Correlations between Content Area Scale Scores by Economic Status..... 224

Table 9-16 Correlations between Content Area Scale Scores by Disability Status..... 225

Table 9-17 Partial Correlations between Content Area Scale Scores..... 225

Table 9-18 Comparison of Most Recent Wisconsin NAEP and Spring 2023 Wisconsin Forward Exam
Impact Data..... 226

Table 10-1 Scale Score Descriptive Statistics for Total Population..... 244

Table 10-2 Scale Score Descriptive Statistics by Subgroup, English Language Arts..... 245

Table 10-3 Scale Score Descriptive Statistics by Subgroup, Mathematics..... 248

Table 10-4 Scale Score Descriptive Statistics by Subgroup, Science..... 251

Table 10-5 Scale Score Descriptive Statistics by Subgroup, Social Studies..... 252

Table 10-6 Score Ranges and Associated Impact Data, English Language Arts..... 254

Table 10-7 Score Ranges and Associated Impact Data, Mathematics..... 254

Table 10-8 Score Ranges and Associated Impact Data, Science..... 254

Table 10-9 Score Ranges and Associated Impact Data, Social Studies..... 254

Table 10-10 Percentage of Students in Each Performance Level by Subgroup, English Language Arts..... 255

Table 10-11 Percentage of Students in Each Performance Level by Subgroup, Mathematics..... 257

Table 10-12 Percentage of Students in Each Performance Level by Subgroup, Science..... 259

Table 10-13 Percentage of Students in Each Performance Level by Subgroup, Social Studies..... 260

Table 10-14 Summary Statistics for Content Standards Raw and SPI Scores, English Language Arts..... 261

Table 10-15 Summary Statistics for Domain Raw and SPI Scores, English Language Arts..... 263

Table 10-16 Summary Statistics for Content Standards Raw and SPI Scores, Mathematics..... 264

Table 10-17 Summary Statistics for Content Standards Raw and SPI Scores, Science..... 265

Table 10-18 Summary Statistics for Content Standards Raw and SPI Scores, Social Studies..... 265

Table 10-19 SPI Cut Scores, English Language Arts..... 266

Table 10-20 SPI Cut Scores, Mathematics..... 268

Table 10-21 SPI Cut Scores, Science..... 270

Table 10-22 SPI Cut Scores, Social Studies..... 271

Table 10-23 Longitudinal Comparison of State-Level Participation Rates and Scale Score Means,
English Language Arts..... 272

Table 10-24 Longitudinal Comparison of State-Level Participation Rates and Scale Score Means, Mathematics.....	273
Table 10-25 Longitudinal Comparison of State-Level Participation Rates and Scale Score Means, Science ...	274
Table 10-26 Longitudinal Comparison of State-Level Participation Rates and Scale Score Means, Social Studies.....	274
Table 10-27 Longitudinal Comparison of State-Level Impact Data, English Language Arts.....	275
Table 10-28 Longitudinal Comparison of State-Level Impact Data, Mathematics.....	276
Table 10-29 Longitudinal Comparison of State-Level Impact Data, Science	277
Table 10-30 Longitudinal Comparison of State-Level Impact Data, Social Studies.....	277

List of Figures

Figure 6-A Examples of Likelihood Functions or the Probability of Each Ability Level Estimate (or Scale Score).....	91
Figure 6-1 Anchor Set Test Characteristic Curves, English Language Arts Grades 3 through 8.....	156
Figure 6-2 Anchor Set Test Characteristic Curves, Mathematics Grades 3 through 8.....	157
Figure 6-3 Anchor Set Test Characteristic Curves, Science Grades 4 and 8.....	158
Figure 6-4 Anchor Set Test Characteristic Curves, Social Studies Grades 4, 8, and 10	159
Figure 6-5 Test Characteristic Curves, English Language Arts	160
Figure 6-6 Standard Error Curves, English Language Arts.....	161
Figure 6-7 Scale Scores and Growth at Quartiles, English Language Arts	162
Figure 6-8 Test Characteristic Curves, Mathematics.....	163
Figure 6-9 Standard Error Curves, Mathematics	164
Figure 6-10 Scale Scores and Growth at Quartiles, Mathematics	165
Figure 6-11 Test Characteristic Curves, Science.....	166
Figure 6-12 Standard Error Curves, Science	167
Figure 6-13 Scale Scores at Quartiles, Science	168
Figure 6-14 Test Characteristic Curves, Social Studies.....	169
Figure 6-15 Standard Error Curves, Social Studies	170
Figure 6-16 Scale Scores at Quartiles, Social Studies	171

Executive Summary

This report is a technical summary of the 2023 administration of the Wisconsin Forward Exam in English Language Arts (ELA) and Mathematics (administered in grades 3 through 8), Science (administered in grades 4 and 8), and Social Studies (administered in grades 4, 8, and 10).

The Wisconsin Forward Exam assessments are designed to measure students' knowledge of ELA, Mathematics, Science, and Social Studies, and they are aligned with Wisconsin Academic Standards. The ELA, Mathematics, and Science test forms administered in Spring 2023 were developed by Data Recognition Corporation (DRC) using DRC's college- and career-ready item bank. The Spring 2023 Social Studies assessments contained Wisconsin-owned items. One new operational test form was developed for the Spring 2023 administration for each grade and content area. All assessments except for Braille and accommodated paper-based forms were administered online.

E.1 Overview of the Wisconsin Forward Exam

The Wisconsin Forward Exam is designed to measure Wisconsin Academic Standards, which define the knowledge and skills students need in each grade level to succeed in college, other postsecondary training, and careers.

The Wisconsin ELA and Mathematics grade-level tests have undergone multiple alignment changes since their first administration in the 2005–06 school year, with the latest changes in the 2015–16 administration, which was also the first administration year of the tests under the Wisconsin Forward Exam program. The current ELA and Mathematics assessments are aligned to Wisconsin Academic Standards adopted in 2010. The reporting scales for the ELA and Mathematics tests were established after the Spring 2016 test administration, and the performance level cut scores were set in Summer 2016. The ELA and Mathematics 2015–16 results are considered the baseline for year-to-year student performance comparisons. The 2022–23 ELA and Mathematics assessments are statistically linked to the established scales, allowing for test score comparability from Spring 2016 to Spring 2023 within each content area.

The Science assessments (grades 4 and 8) have been on a different trajectory. A change to the Science test blueprint and design was made for the Spring 2019 operational test administration. New Science tests, aligned to the new Wisconsin Standards for Science (WSS) adopted in 2017 and the Next Generation Science Standards (NGSS), were developed and administered to Wisconsin students for the first time in Spring 2019. Due to the change of standards, new scales were developed for the new Science tests and new performance level cut scores were set after the Spring 2019 test administration. The 2022–23 Science assessments are statistically linked to the scales established in Spring 2019, allowing for test score comparability across the last four administrations.

A change to the Social Studies test blueprint and design was made for the Spring 2022 operational test administration. The new Social Studies assessments for grades 4, 8, and 10 are aligned to the new Wisconsin Standards for Social Studies that were adopted in 2018. The new reporting scales for the Social Studies tests were established after the Spring 2022 test administration, and new performance level cut scores were set for

these assessments in Spring 2022. The Social Studies 2021–22 results are considered a new baseline for year-to-year student performance comparisons. The 2022–23 Social Studies assessments are statistically linked to the scales established in Spring 2022, allowing for test score comparability across the last two administrations.

All Wisconsin assessments are administered online and contain various item types, including multiple-choice (MC), multi-select (MS), technology-enhanced (TE), evidence-based selected response (EBSR), and short-answer (SA). Braille, print-on-demand, and Spanish translation forms that contain the same items as regular online operational test forms are also available to students who need them.

E.2 Administration

In Spring 2023, Wisconsin administered summative assessments in ELA and Mathematics to students in grades 3 through 8. Science assessments were administered to students in grades 4 and 8, and Social Studies assessments were administered in grades 4, 8, and 10. The Wisconsin Forward Exam was administered from March 20 to April 28, 2023. Test administration is discussed in Part 4 of this report.

A total of 454 public school districts, 400 choice schools, and 5 private schools had students who participated in at least one Wisconsin Forward Exam test in grades 3 through 8 or in grade 10 (Social Studies only). Table E-1 shows test participation rates in Spring 2023. For the purposes of this report, participation rate is defined as the percentage of students who received a valid scale score compared to the total number of students expected to take the test. The “Enrolled” column shows the total number of students expected to take the test in Spring 2023. The “Number Tested” and “Percent Tested” columns show the number and percentage of students who participated in the test and received a valid scale score. The test participation rates for grades 3 through 7 ranged from approximately 95% to 96% across all content areas. The test participation rates for grade 8 were approximately 94% across all content areas. The participation rate for Social Studies grade 10 was approximately 88%. The Spring 2023 participation rates for all grades were similar to the Spring 2022 participation rates. Further analysis of the Spring 2023 participation rates is provided in Part 4 of this report.

E.3 Student Performance

This is the seventh year of the ELA and Mathematics scores being reported on the scales established in Spring 2016. Spring 2023 also marks the fourth year of the Science assessments measuring the new Wisconsin Standards for Science and the second year of the Social Studies assessments measuring the new Wisconsin Standards for Social Studies. Tables E-2 and E-3 present the percentages of students classified as *Proficient* or *Advanced* from 2016 through 2023 in ELA and Mathematics, respectively. Table E-4 shows the percentages of students classified as *Proficient* or *Advanced* from 2019 to 2023 in Science. Due to setting new scales and performance cut scores for Science after the Spring 2019 test administration, student results in Science are not directly comparable between the Spring 2019 administration and previous administrations and the previous data are not reported in this table. Student results are comparable between the Spring 2019 and Spring 2023 administrations for Science. Table E-5 shows the percentages of students classified as *Proficient* or *Advanced* in 2022 and 2023 in Social Studies. New performance level cut scores were established for Social Studies after the

2021–22 test administration. Therefore, student performance in Social Studies in the Spring 2022 administration is directly comparable only with the Spring 2023 student performance.

Caution should be used when making statewide data comparisons over time for ELA, Mathematics, and Science. Due to the COVID-19 pandemic and disruptions to student learning in the 2020–21 school year, the participation rates for Spring 2021 were considerably lower than in previous years or the current year. In addition, the makeup of the Spring 2021 tested population was not representative of the Spring 2021 enrolled population. Therefore, the longitudinal data trend might have not accurately represented changes in state-level student performance between 2019 and 2021 and again between 2021 and subsequent administrations.

The percentages of students classified as *Proficient* or *Advanced* in ELA in 2023 ranged from approximately 37% for grade 3 to approximately 45% for grade 4. The percentages of students classified as *Proficient* or *Advanced* in Mathematics in 2023 ranged from 31% for grade 8 to 48% for grade 3. Approximately 51% of students were classified as *Proficient* or *Advanced* in Science grade 4, and approximately 49% of students were classified as *Proficient* or *Advanced* in Science grade 8. The percentages of students classified as *Proficient* or *Advanced* in Social Studies in 2023 were about 59% in grades 4 and 8 and approximately 47% in grade 10. More details on student performance are provided in Part 10 of this report.

When year-to-year performance trends were considered, more students were classified in the *Proficient* or *Advanced* performance categories in Spring 2023 compared to Spring 2022 for ELA grades 3, 4, 6, 7, and 8. These increases ranged from about 2% for grade 3 and 6 to over 4% for grade 8. A small decrease of less than 2% of students classified as *Proficient* or *Advanced* between Spring 2022 and 2023 was observed for ELA grade 5.

For Mathematics, more students were classified in the *Proficient* or *Advanced* categories in Spring 2023 compared to Spring 2022 in all grades. These differences ranged from less than 1% in grades 3 and 8 to approximately 3% in grade 6.

Negligible changes of less than half a percent of students in the *Proficient* or *Advanced* categories between Spring 2022 and Spring 2023 were found for Science grades 4 and 8 and Social Studies grades 4, 8, and 10.

It was also observed that the percentages of students in the *Proficient* or *Advanced* performance categories continued to be lower in Spring 2023 than in the pre-pandemic administration of Spring 2019 for most ELA, Mathematics, and Science grades. The exception was ELA grade 4, where an increase of about 2% of students classified in the two highest levels between Spring 2019 and 2023 was found. Minor increases of less than 1% of students in the *Proficient* or *Advanced* categories between Spring 2019 and Spring 2023 were found for ELA grade 8 and Mathematics grades 4 and 5. This observed change in performance between Spring 2019 and Spring 2023 should be interpreted in the context of circumstances related to the long-lasting effects of the COVID-19 pandemic (including school closures, nonstandard instruction delivery modes in the 2020–21 school year, and potential diminished opportunity to learn for students) and slow educational recovery.

E.4 Validity of Intended Interpretation of Test Scores

Most sections of this report are designed to provide validity evidence to support the use and intended interpretation of the Wisconsin Forward Exam ELA, Mathematics, Science, and Social Studies test scores. Test scores are used to identify strengths and areas for improvement in Wisconsin's student performance, to inform stakeholders (teachers, school administrators, district administrators, DPI staff members, parents, and the public) about the state's status with respect to its progress toward meeting the academic performance standards of the state, and to meet the requirements of the state's accountability program. Part 2 of this report provides the validity framework and a summary of the validity evidence for the Wisconsin Forward Exam.

Evidence of validity based on test content was supported by the test specifications, including the test design and test blueprint. Wisconsin assessments were developed in alignment with Wisconsin Academic Standards. A rigorous item review and test form development process were implemented to select items from DRC's college- and career-ready item pool for the ELA, Mathematics, and Science assessments and from Wisconsin-owned pool of items, written by DRC content specialists and reviewed by Wisconsin educators, for Social Studies. More details on test content and test development are provided in Part 3 of this report.

With the exceptions of Braille and a limited number of paper-based test forms, Wisconsin Forward Exam assessments were administered online in a standardized manner, further supporting validity of the intended score interpretation. Universal tools were available for all students to use. Designated supports and accommodations were available to students for whom such aids were deemed appropriate and/or indicated in their Individualized Education Programs. More details on test administration and use of universal tools, designated supports, and accommodations are provided in Part 4 of this report.

Scoring of technology-enhanced, short-answer, multi-select, evidence-based selected-response, and text-dependent analysis items followed predefined scoring criteria. The multiple-choice, multi-select, technology-enhanced, short-answer, and evidence-based selected-response items were autoscored. ELA text-dependent analysis items were scored via artificial intelligence engine supported by human reader score verification. The inter-rater reliability statistics demonstrated that the text-dependent analysis items were scored reliably (refer to Parts 5 and 8 for details).

The test scaling and equating was conducted using item response theory (IRT) methodology. Students' scale scores were derived using item parameters estimated after the Spring 2023 test administration. The IRT models used for Wisconsin Forward Exam scaling were appropriate for the test data supporting the operational data analysis and ensuring that the test items, as well as the overall tests, were functioning appropriately. For details on test scaling and equating, refer to Part 6. The cut scores used to classify students into different performance levels and associated performance level descriptors were established during the Summer 2016 standard setting for ELA and Mathematics, the Spring 2019 standard setting for Science, and the Spring 2022 standard setting for Social Studies in a collaborative and participatory process, further supporting the validity and interpretation of the Wisconsin Forward Exam scores (refer to Part 7 for details).

Evidence of construct-related validity—supporting the intended interpretation of test scores and their use—was provided through studies of test reliability, evaluation of test fairness, evaluation of internal test structure, and evaluation of the relationship of test scores with external variables. The reliability analysis results indicated that the Wisconsin Forward Exam tests produce scores that would be relatively stable if the tests were administered repeatedly under similar conditions (refer to Part 8 of this report for details).

Test and item fairness were evaluated through differential item functioning analysis (refer to Part 9 of this report for details). The assumption that the content area Wisconsin Forward Exam tests were unidimensional (i.e., each grade-level test measured one primary dimension) was confirmed through principal component analysis. The evidence of the validity of the intended interpretation of the Wisconsin Forward Exam test scores based on the relationships with other variables was evaluated through the correlations computed between the ELA, Mathematics, Science, and Social Studies scale scores. The student scores were found to be highly, but not perfectly, related to each other, suggesting that while different constructs are being measured, the two assessments may also be tapping into a similar knowledge base or general underlying ability. When considering the Wisconsin Academic Standards and the percentages of students classified as *Proficient* or *Advanced* (based on the Wisconsin Forward Exam cut scores for ELA, Mathematics, and Science), the Wisconsin Forward Exam impact data are in alignment with the National Assessment of Educational Progress (NAEP) impact data. This provides evidence of the relationship between the state assessments and the national assessments in these content areas (see Part 9 of this report for details).

Finally, in Part 10 of this report, test results are presented in the context of score reports that aid the user in understanding the meaning of the test results. The current administration test results are presented for the total population and subgroups of students. The longitudinal test results are also presented for all content areas. Monitoring group performance is possible if the test content and the construct measured by the test are comparable from year to year and if the scores are reported on the same scale used in previous years.

Table E-1 Test Participation Rates in Spring 2023

Grade	Enrolled	English Language Arts		Mathematics		Science		Social Studies	
		Number Tested	Percent Tested	Number Tested	Percent Tested	Number Tested	Percent Tested	Number Tested	Percent Tested
3	61228	58497	95.54	58722	95.91				
4	61689	58996	95.63	59165	95.91	59141	95.87	59131	95.85
5	62058	59386	95.69	59577	96.00				
6	62310	59412	95.35	59570	95.60				
7	63623	60413	94.95	60559	95.18				
8	66060	62249	94.23	62360	94.40	62289	94.29	62261	94.25
10	69876							61819	88.47

Table E-2 Percentage of Students Classified as *Proficient* or *Advanced* from 2016 through 2023, English Language Arts

Grade	English Language Arts						
	2016	2017	2018	2019	2021	2022	2023
3	43.13	41.83	39.75	38.69	34.56	34.93	37.18
4	43.30	46.72	43.91	42.98	40.12	41.60	45.01
5	42.47	46.42	44.17	40.06	37.52	40.28	38.44
6	42.58	45.26	42.86	40.96	38.45	37.98	40.05
7	41.98	43.63	45.15	44.87	42.92	38.26	40.90
8	41.56	41.12	37.33	37.03	35.66	33.46	37.62

Note: Caution should be exercised when interpreting the Spring 2021 statewide data due to participation rates below 90%.

Table E-3 Percentage of Students Classified as *Proficient* or *Advanced* from 2016 through 2023, Mathematics

Grade	Mathematics						
	2016	2017	2018	2019	2021	2022	2023
3	48.00	48.03	49.83	49.44	44.99	47.20	47.81
4	44.20	43.50	44.46	45.05	41.07	43.73	45.30
5	44.08	44.46	45.95	46.58	41.59	44.80	47.18
6	42.84	43.61	43.96	42.49	35.57	38.76	41.75
7	39.26	39.29	38.97	38.83	34.84	33.73	35.18
8	33.86	34.62	36.61	35.85	30.00	30.31	31.19

Note: Caution should be exercised when interpreting the Spring 2021 statewide data due to participation rates below 90%.

Table E-4 Percentage of Students Classified as *Proficient* or *Advanced* from 2019 through 2023, Science

Grade	Science			
	2019	2021	2022	2023
4	52.78	51.16	50.85	51.25
8	53.95	51.47	48.90	48.86

Note: Caution should be exercised when interpreting the Spring 2021 statewide data due to participation rates below 90%.

Table E-5 Percentage of Students Classified as *Proficient* or *Advanced* in 2022 and 2023, Social Studies

Grade	Social Studies	
	2022	2023
4	58.89	59.12
8	58.73	59.06
10	47.67	47.24

Note: Caution should be exercised when interpreting the Spring 2022 statewide data for grade 10 due to participation rates below 90%.

Part 1: Overview

The *Wisconsin Forward Exam Spring 2023 Technical Report* documents the processes and procedures applied in test development, administration, and scoring, as well as the assessment results. This report also provides evidence in support of the validity and reliability of the testing program in adherence to the *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014). This report demonstrates that the Spring 2023 Wisconsin Forward Exam adhered to the appropriate standards and practices of educational assessment. Ultimately, this report provides evidence that valid inferences about Wisconsin student performance can be derived from this assessment.

1.1 Historical Background

The Improving America's Schools Act of 1994 required that states establish challenging academic standards as well as aligned annual assessments. The Goals 2000: Educate America Act and the Elementary and Secondary Education Act (ESEA) spelled out additional requirements to ensure that citizens receive coherent information about whether and to what degree students are meeting rigorous academic standards. This Technical Report is an important part of meeting those requirements.

Wisconsin students in grades 4, 8, and 10 began taking the Wisconsin Knowledge and Concepts Examination (WKCE) norm-referenced assessments in the 1997 school year. At that time and in the following years, *TerraNova*TM tests developed by CTB/McGraw-Hill (1997, 2000, 2009) were used. The selection of those tests was partly predicated on an awareness of the academic standards being developed. In January 1998, the Wisconsin Model Academic Standards (WMAS) were adopted. These new standards were the work of the Governor's Commission on Wisconsin Model Academic Standards, chaired by then Lieutenant Governor Scott McCallum and the Wisconsin Department of Public Instruction (DPI). The assessments aligned to WMAS would measure student performance in the same subjects as the *TerraNova* tests.

Beginning in the 2005–06 school year, the federal No Child Left Behind Act (NCLB) required all states to test all students in Reading and Mathematics in grades 3 through 8 and once in high school (in grade 10 under Wisconsin law § 118.30). Based on the NCLB legislation, student performance, reported in terms of proficiency categories, was used to determine the Adequate Yearly Progress (AYP) of students at the school, district, and state levels. Beginning with the 2007–08 school year, states were also required to administer Science assessments at least once in grades 3–5, once in grades 6–9, and once in grades 10–12.

It was within this policy context that the WKCE was constructed, as a criterion referenced test, for the Fall 2005 administration, replacing the previously existing norm referenced WKCE in Reading and Mathematics. The criterion-referenced WKCE was designed specifically for Wisconsin students to measure their performance on the WMAS. These assessments were designed to evaluate students' knowledge and to measure achievement in the basic skills taught in schools at grades 3–8 and 10. The Fall 2013 WKCE was the ninth administration of these assessments and the last administration of Reading, ELA, and Mathematics. The assessments in Science and Social Studies under the existing WKCE model continued to be administered until Fall 2014.

A major change in the Wisconsin assessments occurred for the 2014–15 test administration. First, the ELA and Mathematics assessments were moved from the Fall testing window to the Spring testing window. Second, the new ELA and Mathematics tests for grades 3–8 developed for the Spring 2015 administration consisted of new Smarter Balanced Assessment Consortium (SBAC) items aligned to the Common Core State Standards (CCSS). Thus, the 2014–15 ELA and Mathematics assessments were not comparable content- and construct-wise to the assessments administered in prior years. Third, while the prior years' assessments included CTB's *TerraNova* items that yielded norm-referenced scores, the 2014–15 assessments did not include such items. Fourth, the regular versions of the 2014–15 assessments were administered as fixed forms in the online mode, in contrast to the previous assessments, which were all administered in the paper-and-pencil mode. Fifth, TE item types were introduced in the 2014–15 online test administration. Last, the student test scores for ELA and Mathematics were reported on SBAC scales and the students were classified into performance levels based on SBAC cut scores. Further details on the structure and reporting of the Spring 2015 ELA and Mathematics assessments (called the Wisconsin Badger Exam) can be found at <https://dpi.wi.gov/assessment/historical/smarter>.

The ELA and Mathematics assessments underwent yet another change in the 2015–16 administration year. The Wisconsin DPI partnered with DRC to develop new ELA and Mathematics assessments for grades 3–8 for the Spring 2016 administration. The items contained in these assessments were drawn from DRC's nationally field-tested college- and career-ready (CCR) item bank and aligned with Wisconsin Academic Standards for ELA and Mathematics. The new assessment program is called the Wisconsin Forward Exam, and the new ELA and Mathematics tests were administered online in Spring 2016. Since the new assessments did not contain any items from the 2014–15 Wisconsin Badger Exam tests, the new scales were not statistically linked to the previous scales. The new reporting scales for the ELA and Mathematics tests were developed after the Spring 2016 test administration, and the new performance level cut scores were set for these assessments in Summer 2016.

Science (grades 4 and 8) and Social Studies (grades 4, 8, and 10) assessments have been on a different trajectory, as they continued to be aligned with the WMAS. However, the test administration for these assessments was moved from the Fall window to the Spring window for the 2015–16 administration year. The items contained in the Science and Social Studies tests were mainly drawn from the pool of previously administered items, but new items were also included. Several of the previously administered items were edited to improve item quality and reflect test content changes over time. Despite the fact that many Science and Social Studies items in the Spring 2016 administration came from the previous item pool, statistically linking the Spring 2016 forms to the previous forms was not recommended due to the change of the testing window and the numerous changes to the items themselves. Instead, similar to what was done for the ELA and Mathematics assessments, new scales were developed for the Science and Social Studies tests under the new Wisconsin Forward Exam program. Following the new scale development, the new performance level cut scores were set for Science and Social Studies in Summer 2016.

Details regarding development, scaling, reporting, and standard setting for all Spring 2016 assessments are included in the *Wisconsin Forward Exam Spring 2016 Technical Report* available at <https://dpi.wi.gov/assessment/forward/resources>.

Spring 2023 was the seventh administration year for the Wisconsin Forward Exam in ELA and Mathematics. Spring 2023 was also the fourth administration year for the new Wisconsin Forward Exam in Science, aligned to the new WSS and the NGSS. The new Science assessments focus on understanding content linked to work with science and engineering practices and crosscutting concepts as detailed in the *National Research Council Framework for K–12 Science Education* (National Research Council, 2012). The Spring 2023 ELA and Mathematics assessments were statistically linked to their respective Spring 2016 scales, and the Science assessments were statistically linked to the Spring 2019 Science scales, allowing for student score comparisons across administrations for these content areas.

Spring 2023 was the second administration year for the Wisconsin Forward Exam in Social Studies aligned to the new Wisconsin Standards for Social Studies. New scales were developed, and the new performance level cut scores were set for Social Studies in Spring 2022, allowing for student score comparisons between Spring 2022 and Spring 2023 for Social Studies.

This Technical Report documents all aspects of the 2022–23 testing cycle. The structure of this report mirrors the testing cycle. A brief content summary of the report is provided later in this part of the report.

1.2 Uses of Test Scores

Validity is the overarching component of the Wisconsin Forward Exam program. The following excerpt is from the *Standards for Educational and Psychological Testing* (hereafter the Standards) (AERA, APA, & NCME, 2014):

Ultimately, the validity of an intended interpretation of test scores relies on all the available evidence relevant to the technical quality of a testing system. Different components of validity evidence . . . include evidence of careful test construction; adequate score reliability; appropriate test administration and scoring; accurate score scaling, equating, and standard setting; and careful attention to fairness for all test takers, as appropriate to the test interpretation in question. (p. 22)

As stated by the *Standards*, the validity of a testing program hinges on the use of the test scores. Validity evidence that supports the uses of the Wisconsin Forward Exam scores is provided in this Technical Report. This section examines some possible uses of the Wisconsin Forward Exam scores.

Part 2 of this Technical Report provides a summary of the evidence of the validity of intended score interpretation of the Wisconsin Forward Exam. Parts 3 through 10 of the report contain details of the validity evidence as well as technical support for some of the interpretations and uses of test scores. The information in Parts 3 through 10 also provides a firm foundation of evidence that the Wisconsin Forward Exam measures what it is intended to measure. However, this Technical Report cannot anticipate all possible interpretations and uses of the Wisconsin Forward Exam scores. It is recommended that policy and program evaluation studies, in accordance with the *Standards*, be conducted to support some of the uses of the Wisconsin Forward Exam scores.

The validity of a test score ultimately rests on how that test score is used. To understand whether a test score is being used properly, one must first understand the purpose of the test. The intended uses of the Wisconsin Forward Exam scores include the following:

- Identifying students’ strengths and areas in need of improvement
- Communicating expectations for all students
- Evaluating school-, district-, and state-level programs
- Informing stakeholders (i.e., teachers, school administrators, district administrators, DPI staff members, parents, and the public) about the status of the progress toward meeting academic achievement standards of the state
- Meeting the requirements of the state’s accountability program

This Technical Report refers to the use of the test-level scores (scale scores and performance levels) and standard-level (reporting category) scores (standard performance index [SPI] scores and performance levels).

1.2.1 Test-Level Scores

At the test level, an overall scale score that is based on student performance on the entire test is reported. In addition, an associated level of performance is reported. These scores indicate, in varying ways, a student’s achievement in ELA, Mathematics, Science, or Social Studies. Test level scores are reported at four levels: state, school district, school, and student.

Two types of test-level scores are reported to indicate a student’s achievement on the Wisconsin Forward Exam: (1) the scale score and (2) its associated level of performance.

Scale Scores

A scale score indicating a student’s performance is determined for each content area. The overall scale score for a content area quantifies the achievement being measured by the ELA, Mathematics, Science, or Social Studies test. In other words, the scale score represents the student’s level of performance, where higher scale scores indicate higher levels of performance on the test and lower scale scores indicate lower levels of performance.

Levels of Performance

A student’s performance on the ELA, Mathematics, Science, or Social Studies Wisconsin Forward Exam is reported in one of four levels of performance: *Below Basic*, *Basic*, *Proficient*, or *Advanced*. The cut scores for the levels of performance for ELA and Mathematics were recommended by Wisconsin educators at the standard setting workshop in June 2016. The cut scores for Science were established during the standard setting workshop in May 2019 and the cut scores for Social Studies were established during the standard setting workshop in May 2022. The cut scores reflect the expectations of Wisconsin educators of what Wisconsin students should know and be able to do in ELA, Mathematics, Science, and Social Studies (see Part 7 of this report for a brief description of the Wisconsin Forward Exam standard setting).

Use of Test-Level Scores

The Wisconsin Forward Exam scale scores and performance levels provide summary evidence of student achievement in ELA, Mathematics, Science, and Social Studies. Classroom teachers may use these scores as evidence of student achievement in these content areas. At the aggregate level, district and school administrators may use this information for activities such as curriculum planning. The results presented in this Technical Report provide evidence that the scale scores are valid and reliable indicators of student performance in ELA, Mathematics, Science, and Social Studies.

1.2.2 Standard-Level Subscores and Performance Levels

The standard-level subscores (i.e., the SPI scores) indicate student performance on a content standard and can be interpreted as an estimate of the number of items a student would be expected to answer correctly if there had been 100 similar items for a given reporting category. The SPI scores are criterion-referenced scores, in that they estimate how much a student knows in a clearly defined skill domain (i.e., the criterion). The SPI scores are computed for content standards measured by at least four items.

Based on their SPI scores, students are classified in one of the four content category performance levels: *Below Basic*, *Basic*, *Proficient*, or *Advanced*. The SPI cut scores separating these performance levels are derived as expected percentages of possible score points for a given standard (content category) for students whose total test score is at the corresponding total test cut score (*Basic*, *Proficient*, or *Advanced*).

Use of the Standard-Level Subscores

The purpose of reporting SPI scores on the Wisconsin Forward Exam is to show the relationship between the overall achievement being measured (represented by the test score) and the skills within each of the content standards associated with the content area. Teachers may use the SPI scores for individual students as indicators of strengths and needs, but the SPI scores are best corroborated by other evidence, such as homework, class participation, diagnostic test scores, or observation. Part 3 of this Technical Report provides evidence of content validity that supports the use of the standard-level subscores. Part 9 of this Technical Report provides evidence of construct validity that further supports the use of these subscores.

District and school administrators may compare their results by content standard and grade level with the state results to better understand students' strengths and needs within a particular content area and grade level. Caution should be exercised when comparing standard-level subscores across years because different items will contribute to these subscores and these items may vary in difficulty between test forms or test administrations.

1.3 Technical Report Structure

This Technical Report documents, in the subsequent parts, the major activities of the testing cycle. It provides comprehensive details that confirm that the processes and procedures applied in the Wisconsin Forward Exam adhere to appropriate professional standards and practices of educational assessment. Ultimately, this report provides evidence that valid inferences about Wisconsin student performance can be derived from the Wisconsin Forward Exam. An overview of the subsequent parts within this report is provided below.

Part 2: Validity Framework

Part 2 of the Technical Report discusses the concept of validity evidence. This Technical Report is composed of evidence that supports the use of the Wisconsin Forward Exam ELA, Mathematics, Science, and Social Studies scores. In Part 2, some of the uses of the Wisconsin Forward Exam scores are discussed.

Part 3: Test Content and Test Development

Part 3 of this report describes the test blueprint, test design, the item development and test form development process, and some aspects of the content-related validity of the Wisconsin Forward Exam. More specifically, it describes how DRC and DPI collaborated to ensure that the appropriate content was included in the Wisconsin Forward Exam and to ensure that the test items adequately sampled the domain of content knowledge necessary to make legitimate inferences about student performance. The Wisconsin Academic Standards were the basis of the test blueprints and item specifications for their respective content areas. Wisconsin educators were involved in reviewing the items in all content areas to ensure the appropriateness of the test to the standards. The first item review for grades 3–8 in ELA and Mathematics occurred in December 2015. The first item review for new assessments for grades 4 and 8 in Science occurred in August 2017. In addition, the first item review for Social Studies assessments in grades 4, 8, and 10, measuring new Wisconsin Standards for Social Studies occurred in August 2019. Each year after that, with the exception of year 2020, new items were reviewed and added to the Wisconsin pool of items for future field-testing. The item reviews served to establish the accessibility of the items and reading passages. Simultaneously, DRC created the test specifications documents that were later approved by DPI and will continue to serve as a foundation for item and test development. Additional item reviews, supported by the item data, occurred after each field test administration and were conducted by DPI content experts. The purpose of these reviews was to refine the pool of items from which the subsequent operational test forms were selected.

Part 3 also presents the Wisconsin Forward Exam design and discusses features of the Spring 2023 Wisconsin Forward Exam forms. The Spring 2023 test forms adhered to the approved test blueprints, test designs, and psychometric specifications.

Part 4: Test Administration

Part 4 describes test administration and accommodations. It also provides information on student participation in the ELA, Mathematics, Science, and Social Studies assessments in Spring 2023. In the 2022–23 school year, the Wisconsin Forward Exam was administered to Wisconsin students for the seventh time.

The Spring 2023 Wisconsin Forward Exam was an online assessment with a corresponding print-on-demand form at each grade level. Student responses to the print-on-demand form were transcribed by a proctor into the online assessment system. Other variations of the forms included stacked Spanish translation forms, video sign language, and closed-captioning. These were provided in an online format at each grade level.

Test administration was conducted during a six-week window from March 20 to April 28, 2023. All testing was conducted online, administered via DRC’s INSIGHT platform.

Part 4 of the Technical Report serves to describe the processes and activities implemented and information disseminated to help ensure standardized test administration procedures and, thus, uniform test administration conditions for students.

Part 5: Scoring

Part 5 documents the scoring process for different item types: scanning of multiple choice (MC) items and multi-select (MS) items; autoscoring of technology-enhanced (TE) items, short-answer (SA) items, and evidence-based selected response (EBSR) items; and artificial intelligence (AI) scoring and handscoring of text-dependent analysis (TDA) items. The description of the handscoring process includes the development and review of the scoring rubrics, anchor (sample) paper selection, training of scoring personnel, ongoing quality assurance, and a systematic review of the resulting score distributions supporting reliable and valid reported test scores. The scoring rubric used in the handscoring of the TDA writing items is presented in detail.

Part 6: Psychometric Analyses

The Spring 2023 administration year is the seventh administration year for the Wisconsin Forward Exam in all grades and content areas. Part 6 discusses characteristics of the sample of student data used for data analysis and describes the classical and item response model (IRT) procedure implemented to analyze the Wisconsin Forward Exam test data. The results of item analysis and test statistical properties are presented in this Part. The Wisconsin Forward Exam data were calibrated using two different item response theory (IRT) models, one for MC items and one for non-MC items. Evaluation of the sufficiency of the IRT model results included model-to-data fit and the standard error of measurement (SEM). The equating of Spring 2023 assessments to their respective reporting scales was performed using the Stocking and Lord procedure. Item-pattern scoring was applied to the Spring 2023 Wisconsin Forward Exam. As discussed in Part 6, item-pattern scoring is generally recommended over number-correct scoring because it produces more accurate scores for individual students. Part 6 also explains how a student's scale score is derived from the raw score using item-pattern scoring.

Part 7: Standard Setting

Part 7 provides a brief overview of the standard setting process, during which the performance level cut scores were set for the ELA and Mathematics tests in Summer 2016, for the Science tests in Spring 2019, and for the Social Studies tests in Spring 2022. The standard setting methodology and results, including short performance level descriptors and cut scores, are presented.

Part 8: Studies of Reliability

Part 8 elaborates on the reliability of the test based on results presented in previous parts of the report. Standard error of measurement (SEM) was assessed for raw scores and scale scores. Internal consistency was evaluated for all tests for the total student population and for subgroups identified by gender, race/ethnicity, economic status, disability status, accommodation use, and English language proficiency. Classification consistency and accuracy were estimated for performance classification. In addition, inter-rater reliability was computed for TDA items on ELA tests that were scored using the AI scoring engine with human scorer verification.

Part 9: Studies of Construct-Related Validity

Part 9 provides additional construct-related validity evidence supporting the Wisconsin Forward Exam. An analysis of differential item functioning is presented. Principal component analysis, correlations among content standards (reporting category scores), and a relationship between the Wisconsin Forward Exam scores and external variables are presented in the context of construct validity. Forensic analysis procedures, implemented to detect possible aberrant testing behavior, are also discussed.

Part 10: Test Results

Part 10 includes short descriptions of reports provided to end users, including individual student reports and aggregate reports. It also contains information on the results of the Spring 2023 Wisconsin Forward Exam administration. Detailed summary statistics of the total scale scores and performance levels and the SPI scores are provided for the total population and for subgroups identified by gender, race/ethnicity, economic status, disability status, accommodation use, and English language proficiency. Longitudinal results are also presented for ELA, Mathematics, and Science.

Part 11: Summary and Recommendations

Key findings of the Spring 2023 Wisconsin Forward Exam administration are presented in the body of the report. However, some issues of a more technical nature that stand out as key recommendations and summary statements that should be considered in subsequent administrations are presented in Part 11. Recommendations based on the Spring 2023 Wisconsin Forward Exam administration may cover different phases of the testing cycle: item development; scoring; and psychometric, or measurement-based, research and evaluation.

Part 2: Validity Framework

Validity is the overarching component of the Wisconsin Forward Exam program. The following excerpt is from the *Standards* (AERA, APA, & NCME, 2014):

Ultimately, the validity of an intended interpretation of test scores relies on all the available evidence relevant to the technical quality of a testing system. Different components of validity evidence . . . include evidence of careful test construction; adequate score reliability; appropriate test administration and scoring; accurate score scaling, equating, and standard setting; and careful attention to fairness for all test takers, as appropriate to the test interpretation in question. (p. 22)

As stated by the *Standards*, the validity of a testing program hinges on the interpretation of the test scores. Validity evidence that supports the uses of the Wisconsin Forward Exam test scores is provided in this Technical Report.

The purpose of test score validation is not to validate the test itself but to validate interpretations of the test scores for particular purposes or actions. Test score validation is not a quantifiable property but an ongoing process, beginning at initial conceptualization and continuing throughout the entire assessment process. Every aspect of an assessment provides evidence in support of (or in challenge to) the validity of an intended interpretation of test scores, including design, content specifications, item development, psychometric quality, and inferences made from the results.

2.1 Sources of Validity Evidence

The sources of validity evidence described in the *Standards* (AERA et al. 2014, pp. 26–31) include evidence based on test content, evidence based on response processes, evidence based on internal test structure, evidence based on relationships with other variables, and evidence based on consequences of testing. These sources of validity evidence are briefly described below.

Validity evidence based on test content can be obtained from an analysis of the relationship between the content of a test and the construct it is intended to measure (AERA et al. 2014, p.14). It refers to traditional forms of content validity evidence and is supported by a correspondence between test content and a specification of the content domain. This type of evidence can be demonstrated through consistent adherence to test blueprints, through a high-quality test development process that includes review of items for accessibility to English language learners and students using testing accommodations, and through alignment studies.

Validity evidence based on response processes relies in large degree on the evaluation of the cognitive processes of examinees responding to various types of items and the relationship between these processes and the construct being measured. Direct evidence based on response processes typically comes from analyzing test takers' individual responses or from questioning test takers from various groups that make up the intended test-taking population about their performance or response strategies on specific items (AERA et al. 2014, p.15). Such evidence can be gathered through cognitive labs conducted as part of the field test data analysis. Validity

evidence based on response process is also supported by a relationship between the item type, format, and content and the construct being measured. For example, if a test is intended to measure a certain set of skills, it is important to determine whether the items included in the test are, in fact, designed to measure these skills or knowledge. In addition, evaluation of student written responses (e.g., text-dependent analysis) further contributes to the validity evidence based on response processes. In such cases, validity evidence includes the extent to which the processes of item response scoring, whether by a human reader or by an artificial intelligence engine, are consistent with the intended interpretation of scores. For example, scorers are expected to apply particular criteria in scoring students' responses and not be influenced by factors that are irrelevant to the intended interpretation of the scores (AERA et al., 2014, pp. 15–16). Recruitment and training of human scorers as well as monitoring the artificial intelligence scoring processes and results, contribute to the validity evidence based on response processes.

Validity evidence based on internal test structure refers to the fact that “analyses of the internal structure of a test can indicate the degree to which the relationships among test items and test components conform to the construct on which the proposed test score interpretations are based” (AERA et al., 2014, p. 16). Such analyses may include statistical analyses of items and subscores conducted to investigate the dimensionality of an assessment. Procedures for gathering such evidence may include factor analysis for single assessments and evaluation of the continuity of the construct across grades for vertically scaled assessments. Internal test structure can also be evaluated using indices of measurement precision such as test reliability, decision accuracy and consistency, generalizability coefficients, and standard errors of measurement. Evaluation of the correlation coefficients that measure the relationship between the content standard (domain) scores and studies of whether test items may function differently for different subgroups of students are additional sources of validity evidence based on internal test structure.

Validity evidence based on relationships to other variables refers to “evidence about the degree to which these relationships are consistent with the construct underlying the proposed test score interpretations” (AERA et al., 2014, p. 16). In educational testing, such evidence is often gathered through studies of correlations between the test scores and measures of different or similar constructs. As stated in the *Standards*, relationships between test scores and other measures intended to assess the same or similar constructs provide convergent evidence, whereas relationships between test scores and measures of different constructs provide discriminant evidence (AERA et al., 2014, pp. 16–17).

Validity evidence based on the consequences of testing is ultimately determined by the stakeholders. Stakeholders decide the purpose and interpretation of scores within their system of reporting and accountability. DRC provides information about test content and technical quality but does not decide the use of test scores. As such, the validity evidence based on consequences of testing is not addressed in this report.

2.2 Summary of Validity Evidence for Wisconsin Forward Exam

In this Technical Report, validity evidence is presented in relation to test content, response processes, internal test structure, and relationship with other variables. Gathering validity evidence related to test consequences is beyond the scope of this Technical Report.

Parts 3 through 10 of this Technical Report provide evidence for the uses as well as technical support for some of the interpretations and uses of test scores. As the Technical Report progresses part by part, it moves through the phases of the testing cycle. Each part of the Technical Report details the procedures and processes applied in the Wisconsin Forward Exam program as well as the test results. Each part highlights the meaning and significance of the procedures, processes, and results in terms of validity evidence or a relationship to the *Standards*. A summary of Wisconsin Forward Exam validity evidence as documented in Parts 3 through 10 is presented here.

Part 3 of the Technical Report documents evidence of the content-related validity demonstrated through each Wisconsin Forward Exam assessment's consistent adherence to the assessment blueprints, which were constructed by DPI based on the Wisconsin Academic Standards. This part of the report also presents the test design and describes the key development tasks related to creating the Spring 2023 Wisconsin ELA, Mathematics, Science, and Social Studies operational test forms. This part documents the involvement of Wisconsin educators, DPI, and DRC in the item review and test development process. The test development process and the involvement of Wisconsin educators in that process forms an important part of the validity of the entire Wisconsin Forward Exam program. The knowledge, expertise, and professional judgment offered by Wisconsin educators ultimately ensures that the content of the Wisconsin Forward Exam forms an adequate and representative sample of appropriate content and that the content forms a legitimate basis upon which to derive valid conclusions about student achievement. The blueprint and design as well as the item and test development activities described in Part 3 explain how specific development processes provide evidence in support of the validity of an intended interpretation of test scores, primarily based on the test content and through the use of expert professional judgment from Wisconsin educators and from DRC test development specialists. The foundational documents—test blueprints and test designs—developed and approved during the initial phases of test development served as critical guides throughout the development and field-testing of items. These documents contribute to ensuring that each form of the test accurately measures the content in consistent and stable ways, thus providing evidence supporting using test scores as an indicator of student achievement of Wisconsin standards.

Part 3 provides evidence to support the validity of an intended interpretation of test scores based on test content of the Wisconsin Forward Exam and address AERA, APA, & NCME (2014) Standards 3.1, 3.2, 3.9, 4.0, 4.1, 4.7, and 4.12.

Part 4 of the Technical Report discusses the processes, procedures, and policies that guide the administration of the Wisconsin Forward Exam, including accommodations, security, and procedures provided to test administrators and school personnel. The following AERA, APA, & NCME (2014) Standards are addressed: 3.4, 3.5, 4.15, 4.16, 6.1, 6.2, 6.3, 6.4, 6.6, and 6.7. The process, procedures, and policies detailed in this section contribute to the validity of an intended interpretation of test scores by reducing the impact of construct-irrelevant variables (e.g., nonstandard administration methods, limitations associated with student disabilities, security breaches) on test performance.

Part 5 of the Technical Report demonstrates adherence to AERA, APA, & NCME (2014) Standards 4.18, 4.20, 6.8, and 6.9. It describes how MC, MS, EBSR, SA, and TE items were autoscored and how text-dependent

analysis (TDA) items were scored by the artificial intelligence (AI) engine followed by a human reader verification. Training of the AI engine as well as the handscoring process, including training and selection of readers, and validation of scoring accuracy, are discussed and support the validity of an intended interpretation of test scores based on the response processes. The procedures described in this section contribute to the evidence of the validity of an intended interpretation of test scores of the Wisconsin Forward Exam by preventing hardware- or software-related errors in machine scoring and reducing construct-irrelevant score variance associated with variations in readers' interpretation and application of scoring rubrics for TDA items.

Part 6 describes the classical and IRT item and test analysis, including item calibration, test equating, and test scaling. The calibration, equating, and scaling methods as well as the processes and procedures for deriving scale scores from response patterns are also described in this part of the Technical Report. Some references to introductory and advanced discussions of IRT are provided. Several axes upon which to evaluate the calibration, equating, and scaling procedures, such as the models and data used, the software applied, the vertical relationship across grades, the estimation of parameters, the fit, the SEM, and the IRT scoring method, are discussed. Part 6 of this report addresses AERA, APA, & NCME (2014) Standards 1.8, 2.13, 4.14, 5.2, 5.13, 5.15, and 7.2. These processes and procedures contribute to the validity of an intended interpretation of test scores of the ELA, Mathematics, Science, and Social Studies assessments by providing the opportunity to evaluate items contributing to the accurate and reliable measurement of the intended constructs and by ensuring the stability of the Wisconsin Forward assessments. The results of the psychometric analyses contribute to the validity evidence based on the internal test structure.

Part 7 of the Technical Report provides a summary of the Wisconsin Forward Exam standard setting for ELA and Mathematics, conducted in June 2016, for Science, conducted in May 2019, and for Social Studies, conducted in May 2022, during which the cut scores were set for the four content areas. The process of the standard setting adhered to AERA, APA, & NCME (2014) Standards 5.21 and 5.22, providing evidence of the procedural validity of the standard setting process, methodology, and outcomes.

Part 8 demonstrates adherence to the *Standards* (AERA, APA, & NCME, 2014) through analyses of the reliability of the Spring 2023 ELA, Mathematics, Science, and Social Studies assessments. It presents a reliability analysis using Cronbach's alpha, SEM, and CSEM results and a detailed analysis of classification consistency and classification accuracy for the total student population and by subgroup. The results of the inter-rater reliability for the ELA text dependent items are also discussed in this part of the report. These analyses address AERA, APA, & NCME (2014) Standards 2.0, 2.3, 2.7, 2.11, 2.13, 2.14, and 2.16. The results of the reliability studies indicate that the Wisconsin Forward Exam tests produce scores that would be stable if the test were administered repeatedly under similar conditions. Reliability is a prerequisite to score validity, and the analyses in this part contribute to the evidence of the validity of an intended interpretation of test scores based on the internal test structure by establishing the reliability of the ELA, Mathematics, Science, and Social Studies scores and proficiency classifications.

As presented in Part 9, additional metrics with which the validity of an intended interpretation of test scores of the ELA, Mathematics, Science, and Social Studies assessments was examined included evaluation of the performance of subgroups of students on the individual test items. As described in Part 9, the issue of item and

test fairness is considered during the item development, item review, and test form construction processes and is formally assessed through an analysis of DIF. It is possible for items to function differently across different population groups, and it is also possible that results for an item do not reflect student ability but instead reflect irrelevant information influenced by demographic factors. The DIF analysis serves to determine whether that possibility occurred and, if so, to what degree, item by item, for each of the categories of gender, race/ethnicity, and other demographic subgroups. The evaluation of item and test fairness addresses AERA, APA, & NCME (2014) Standards 3.1, 3.2, 3.3, and 3.6.

Also included in Part 9 is additional evidence of the construct-related validity based on the internal test structure, gathered through the analysis of the relationships among test items and test components that conform to the test construct, which in turn provides a basis for test score interpretation. The assumption that the content area Wisconsin Forward Exam tests were unidimensional (that is, each grade-level test measured one primary dimension) was confirmed through principal component analysis. In addition, the relationship between the content area reporting category subscores was explored and validated through the measures of correlations between the reporting category scores within a content area. These analyses addressed AERA, APA, & NCME (2014) Standards 1.13 and 1.21.

The relationship between the Wisconsin Forward Exam scale scores and other variables was examined to provide evidence of the construct validity based on the relationships with other variables. These analyses included measures of cross-content correlations of the ELA, Mathematics, Science, and Social Studies scores for the total population and by subgroups and comparisons of student performance on the Wisconsin Forward Exam with student performance on the National Assessment of Educational Progress (NAEP). These analyses are in alignment with multiple best practices of the testing industry (AERA et al., 2014) and are also presented in Part 9 of the report.

Part 10 of the Technical Report contains descriptions of the score reports available to end users. It also provides information on the results of the Spring 2023 administration and on the longitudinal data trends for ELA, Mathematics, Science, and Social Studies. AERA, APA, & NCME (2014) Standards 5.1, 6.10, 7.0, 7.1, and 12.18 are addressed in Part 10.

While the information in Parts 3 through 10 provides a firm foundation of evidence that the Wisconsin Forward Exam tests measure what they are intended to measure, this Technical Report cannot anticipate all possible interpretations and uses of the Wisconsin Forward Exam ELA, Mathematics, Science, and Social Studies scores. It is recommended that policy and program evaluation studies, in accordance with the *Standards* (AERA et al., 2014), be conducted to support some of the uses of the ELA, Mathematics, Science, and Social Studies scores.

Part 3: Test Content and Test Development

The purpose of this section is to describe how DRC, DPI, and Wisconsin educators collaborated through a series of test development processes to ensure that appropriate content was included in the Wisconsin Forward Exam and to ensure that test items adequately sampled the domain of content knowledge necessary to make accurate inferences about student performance. Part 3 documents the test blueprints, test designs, item development process, review and field-testing of new items, and the test form development process for the Spring 2023 administration.

This part of the Technical Report is particularly relevant to AERA, APA, & NCME (2014) Standards 3.1, 3.2, 3.9, 4.0, 4.1, 4.7, and 4.12. Each of these Standards and the way each Standard is addressed will be presented in this section of the report. AERA, APA, & NCME (2014) Standard 4.0 states the following:

Tests and testing programs should be designed and developed in a way that supports the validity of interpretations of the test scores for their intended uses. Test developers and publishers should document steps taken during the design and development process to provide evidence of fairness, reliability, and validity for intended uses for individuals in the intended examinee population. (p. 85)

The test blueprint and item development activities described in this part explain how specific development processes provided evidence to support test validity, primarily content validity, through the use of expert professional judgment from Wisconsin DPI and from DRC test development specialists. The foundational documents—test blueprints and test designs—developed and approved during the initial phases of the project served as critical guides throughout development of the test forms. These documents contributed to ensuring that each test form accurately measured the content in consistent and stable ways, thus providing evidence supporting using the test as an indicator of student achievement of Wisconsin standards.

The Wisconsin Forward Exam ELA, Mathematics, Science, and Social Studies domains are generally defined as the knowledge and skills that are identified within the Wisconsin Academic Standards for these content areas. The framework of Wisconsin Academic Standards, in turn, is based on prior consensus among DPI, Wisconsin educators, and experienced subject-matter experts that the framework represents what is important for teachers to teach and students to learn.

Evidence of validity based on test content includes information about the test specifications, including the test design and test blueprint. Test development involves creating a design framework from the statement of the construct to be measured. The primary consideration in the development of the Wisconsin Forward Exam test specifications was the assessment alignment with the Wisconsin Academic Standards. Constraints of the assessment program and state policy decisions were also taken into consideration in development of the test specifications.

The Wisconsin Forward Exam test specifications consist of a test blueprint and a test design for each grade level and content area. In partnership with DRC, DPI created test blueprints and test designs. DRC and DPI content experts scrutinized each blueprint to ensure optimal content coverage and efficient use of time and resources.

3.1 Test Blueprints

AERA, APA, & NCME (2014) Standard 4.1 states the following:

Test specifications should describe the purpose(s) of the test, the definition of the construct or domain measured, the intended examinee population, and interpretations for intended uses. The specifications should include a rationale supporting the interpretations and uses of test results for the intended purpose(s). (p. 85)

The key structural aspect of the Wisconsin Forward Exam for ELA, Mathematics, Science, and Social Studies is the assessment blueprint that specifies the target score points for each grade and content strand or domain. These assessment blueprints were developed by Wisconsin DPI who made recommendations for the test content for each grade and content area, seeking to ensure optimal content coverage of the Wisconsin Forward Exam assessments. In general, each blueprint represents content sampling proportions that reflect the intended emphasis in instruction and mastery for each content area and at each grade level. Specifications for a range of items organized by standard and item type demonstrate the desired proportions within the summative assessment. In summary, the Wisconsin Forward Exam assessment blueprint for a given grade and content area provides guidance on how the standards are measured.

The test blueprints specify the number of item points for each reporting category and subskill as well as the allowable depth of knowledge (DOK) levels for the respective reporting categories. The process used for developing the blueprints for the Wisconsin Forward Exam was a collaborative effort between DRC and DPI. The DPI-approved blueprints can be found in Tables 3-1 through 3-4.

3.2 Test Design

The test design for the 2023 operational assessments included the use of items reviewed and approved by Wisconsin educators and DPI. Information concerning the item development process can be found in Section 3.4. Various item types were included in the Wisconsin Forward Exam in order to best assess students' understanding of the standards. A description of item types included in the Wisconsin Forward Exam is presented in Table 3-5. The following sections provide detailed information about the test design of the content areas assessed on the Spring 2023 Wisconsin Forward Exams.

3.2.1 English Language Arts

Table 3-6 shows the ELA test design, including the number of passages, items, and points at each grade level. This table also identifies the various item types that appeared on the ELA forms and the points for item scoring. There was one set of operational test items administered in each grade. No new items were field-tested in Spring 2023. Detailed descriptions of the item types are provided in Table 3-5 of this report.

The ELA section of the Forward Exam is divided into four sessions: text-dependent writing prompt, writing/language, listening, and reading. Students were able to take the sessions in any order. Recommended

testing times for all sessions were included in the test design document as well as in the test administration manual.

3.2.2 Mathematics

Table 3-7 shows the Mathematics test design, including the number of items and points at each grade level. There was one set of operational test items administered in each grade. No new items were field-tested in Spring 2023.

The Mathematics section of the exam was divided into two testing sessions, with students able to take the sessions in either order. In grades 3–5, no calculator was allowed for any of the Mathematics items. In grades 6–8, no calculator was allowed for the first session and the second session allowed students to use an embedded calculator. Recommended testing times for both sessions were included in the test design document as well as in the test administration manual.

3.2.3 Science

Table 3-8 shows the Science test design, including the number of items and points at each grade level that were used in the core and embedded field test positions. There was one common set of core operational items in each of the twenty test forms at each grade level. The unique items in each test form were field test items.

The Science section of the exam was divided into three testing sessions, with students being allowed to take the sessions in any order. Recommended testing times for all sessions were included in the test design document as well as in the test administration manual.

3.2.4 Social Studies

Table 3-9 shows the Social Studies test design, including the number of items and points at each grade level that were used in the operational test. No field test items were administered in 2023. The Social Studies exam included two test sessions that could be administered in either order. The Social Studies exams at grades 4, 8, and 10 included custom items developed specifically for the Wisconsin Forward Exam. Recommended testing times for both sessions were included in the test design document as well as in the test administration manual.

3.3 Universal Design

Assessments that are universally designed allow for the participation of the widest possible range of students, resulting in more valid inferences about student performance. Universally designed grade-level assessments may reduce the need for accommodations by reducing or eliminating access barriers associated with the tests themselves. Table 3-10 presents the elements of universal design that were implemented on the Wisconsin Forward Exam (Thompson & Thurlow, 2002).

These elements of universal design are relevant to both item development and form construction. This section addresses how the elements of universal design were addressed in the construction of the Spring 2022 test forms in compliance with AERA, APA, & NCME (2014) Standard 3.1, which states the following:

Those responsible for test development, revision, and administration should design all steps of the testing process to promote valid score interpretations for intended score uses for the widest possible range of individuals and relevant subgroups in the intended population. (p. 63)

A goal of universal design is to measure the performance of students with a wide range of abilities and skills, ensuring that students with diverse learning needs receive opportunities to demonstrate competence on the same content. To accommodate the greatest number of students for the Wisconsin Forward Exam, the assessments include simple, clear, and intuitive instructions and procedures; maximum readability and comprehensibility; and maximum legibility. These design components are addressed primarily through the physical layout and formatting of the online test forms as well as the paper-based test forms used for accommodations. The page specifications define how directions and test items are placed on the pages, the location and appearance of headers and footers, the spacing between an item stem and the answer choices, and other page elements to ensure a consistent, legible appearance of online forms and paper-based test forms. Written instructions at the beginning of each test session are clearly and simply stated, and the wording of such instructions is standardized as much as possible across content areas and grade levels to ensure clarity and consistency.

AERA, APA, & NCME (2014) Standard 3.9 states the following:

Test developers and/or test users are responsible for developing and providing test accommodations, when appropriate and feasible, to remove construct-irrelevant barriers that otherwise would interfere with examinees' ability to demonstrate their standing on the target constructs. (p. 67)

Students with disabilities or students who are English Learners may be provided with test administration accommodations based on their Individualized Education Programs (IEPs). Accommodation code definitions can be found in the Accessibility Guide available on the “Wisconsin Forward Exam Accommodations and Supports” page on DPI’s website: <https://dpi.wi.gov/assessment/forward/accommodations>.

Braille test version was available for each grade and content area to enable students who are blind to participate in the Wisconsin Forward Exam testing. Braille forms for all grades and content areas were created by DRC test developers and consisted of the same items as those included in the regular operational online test forms. Specific recommendations on how to transcribe items into Braille were provided by an independent Braille expert who collaborated with the Braille publisher to produce the Braille version of the Wisconsin Forward Exam assessment and teacher’s notes that accompany the Braille forms.

3.4 Item Development Process

New operational test forms were developed for each grade and content area for the Spring 2023 administration.

ELA, Mathematics, and Science test forms were selected from DRC’s College- and Career-Ready (CCR) item bank. DRC’s CCR item bank contains nationally field-tested CCR items that support the next generation of standards and assessments. Items from the CCR bank are aligned to the College and Career Readiness standards in ELA and Mathematics grades 3–8. Science items are aligned to Wisconsin’s Standards for Science and enhanced by the Next Generation Science Standards (NGSS) based on the National Research Council’s Framework for K–12 Science Education. The item bank is designed to support states like Wisconsin that have adopted more rigorous content standards, curricula, and assessments that better prepare students for college and careers.

Alignment to standards, grade-level appropriateness, depth of knowledge (DOK), item/task level of complexity, estimated difficulty level, relevancy of context, rationale for distractors, style, accuracy, and correct terminology were major considerations in the item development process. DRC’s item development process for the CCR item bank followed the *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014). DRC’s item development work was and continues to be designed to produce reliable and instructionally valid tests that reflect the complete range of performance articulated in the AERA, APA, and NCME standards.

Furthermore, DRC’s item development work adheres to the Principles of Universal Design (Thompson, Johnstone, & Thurlow, 2002) and reflects how items and tests must lend themselves to accessibility by diverse groups of students. Members of DRC’s item development team have received direct training from the National Center on Educational Outcomes (NCEO). Therefore, DRC employs the Principles of Universal Design throughout all stages of both the item development process and the test development process.

All DRC’s ELA, Mathematics, and Science items that appear on the Wisconsin Forward Exam were reviewed for content and for fairness not only by DRC’s content experts but also by a panel of external experts and Wisconsin educators. The external reviewers had a broad range of experience in the educational field. All the reviewers had bachelor’s, master’s, or doctoral degrees and teaching experience in their specific areas of expertise. Table 3-11 provides a high-level sequence of the activities that occurred in the development of the DRC CCR item bank.

Wisconsin-owned Social Studies items were developed by DRC internal item writers and contracted external consultant item writers. All item writers were required to hold a bachelor’s degree or higher in the content area related to the subject for which they would be writing items, in curriculum and instruction, and/or in a related field. They also needed to have three or more years of teaching experience in the content area for which they would be writing items and two or more years developing items to adhere to client specifications. DRC external item writers are trained annually at a series of content-focused writing training sessions. DRC item and test development staff prepared all materials for the item writing training session, including but not limited to the Wisconsin Standards for Social Studies; guidelines for adhering to the Principles of Universal Design and other

accessibility guidelines, including guidelines for English language learners; guidelines for freedom from issues of bias, fairness, and sensitivity; item specifications (including guidelines for writing items to cover a range of difficulty, a range of subject matter, and items focused on specific performance level descriptor alignments); item writing templates by item type, including scoring guidelines for those items scored using automatic scoring; and sample items. Social Studies item development was coordinated by DRC.

Social Studies items underwent reviews by DRC content experts as well as DRC bias and sensitivity experts. All Social Studies items were also reviewed and approved by committees of Wisconsin educators. The efforts by DRC in developing items are in alignment with multiple best practices of the testing industry and, in particular, support the following AERA, APA, & NCME (2014) Standards:

Standard 3.1 Those responsible for test development, revision, and administration should design all steps of the testing process to promote valid score interpretations for intended score uses for the widest possible range of individuals and relevant subgroups in the intended population. (p. 63)

Standard 3.2 Test developers are responsible for developing tests that measure the intended construct and for minimizing the potential for tests' being affected by construct irrelevant characteristics, such as linguistic, communicative, cognitive, cultural, physical, or other characteristics. (p. 64)

As stated earlier, Wisconsin licensed ELA, Mathematics, and Science items from DRC's CCR item bank. Due to the state-specific nature of the Social Studies standards, DPI owns the items for that content area.

3.4.1 Reading Passage and Item Reviews

The test items typically begin their life cycle two years prior to their operational administration. New ELA, Mathematics, Science, and Social Studies passages and/or items were first reviewed and approved for placement on the Wisconsin Forward Exam by both DPI and Wisconsin educators. Educators from across the state reviewed items in an online format so that items could be evaluated in the same testing engine and style in which they are presented to students during the actual administration. The item reviews were held virtually in August 2021 for all content areas and facilitated by DRC. An example of the training PowerPoint presentation used at the reviews can be found in Appendix A of this report.

Table 3-12 shows the number of items taken to the 2021 item review by grade and content area. Using the approved test blueprints as a guide, DRC content specialists determined the focus of the items that would be taken to item review. Using an electronic tally sheet, Wisconsin educators made determinations about standard alignment, depth-of-knowledge levels, and key(s). They noted any bias and sensitivity concerns and had the opportunity to determine whether items were accepted as is or accepted with revisions. They also had the opportunity to register a "dissenting view," which indicated that the committee preferred the item not be selected to appear on the Wisconsin Forward Exam in a field test position.

Items that were approved by the Wisconsin educators were then included in the next field test administration in Spring 2022 for all content areas. The purpose of the field tests was to expand the pool of items eligible for inclusion in the subsequent operational test forms of the Forward Exam.

3.5 Field-Testing

ELA, Mathematics, Science, and Social Studies items approved for the field test administration during the Summer 2021 item review were field-tested in Spring 2022 during the operational test administration. Field test items were fully embedded in the operational forms, and students were not able to distinguish between the operational and field test items. The field test items were embedded in several test forms administered in each grade and content area. Each test form contained the same operational test items and unique field test items. The test forms were spiraled at the student level within a grade and a content area. A total of 811 new items were field-tested for ELA. A total of 755 items were field-tested for Mathematics. A total of 187 items were field-tested for Science, and a total of 152 items were field-tested for Social Studies in Spring 2022.

The most recent field test of new Science items occurred in Spring 2023. A total of 100 new items were field-tested in each Science grade. No new items were field-tested for ELA, Mathematics, and Social Studies in Spring 2023.

3.5.1 Statistical Analysis of Field Test Data

Following the field test data acquisition, the field test data analyses were conducted. The analyses included classical item analysis, differential item functioning (DIF) analysis, and item response theory (IRT). The classical item analysis included computation and evaluation of the following statistics: item p -values (difficulty), item-total test correlation, percentage of students selecting incorrect responses, point-biserial correlation for incorrect responses for the multiple choice (MC) items, score point distribution for items worth more than 1 point, and omit rates for all items. More details on classical item analysis methodology are provided in Part 6 of this report.

DIF was conducted for all field test items to examine potential item bias and to determine whether item performance differences between identifiable subgroups were due to factors other than student ability, making the items unfairly difficult for a particular subgroup in the student population. DIF analyses were conducted based on gender, race/ethnicity, socioeconomic status, disability status, and English language proficiency (ELP) groups. More details on the DIF methodology are provided in Part 9 of this report.

As the last step of the field test data analysis, the field test items were calibrated and equated to operational test scales using the IRT methodology (explained in detail in Part 6 of this report). The field test items were equated to their respective operational test scales in Spring 2022. All operational test items contained in the Spring 2022 operational test forms served as anchor items to place the field test items on the operational test scales using the Stocking and Lord (1983) equating procedure.

The field test item statistics are used as a means of detecting items that deserve closer scrutiny, rather than as a mechanism for automatic retention or rejection. Toward this end, a set of criteria was used as a screening tool to identify items that needed a closer review. For an item to be flagged for an additional review, the criteria included any of the following:

- p -value <0.20 or >0.90

- item-total test correlation (point biserial for MC items) <0.15
- positive point biserial on a distractor for an MC item
- omit rate >5%
- fewer than 5 students at any score point for multi-point items (EBSR and TDA item types in ELA)
- large differential item functioning (DIF) status

Items flagged for any of the above reasons were reviewed by the content area specialists prior to their review by DPI.

3.5.2 Item Data Review

In the preceding section, it was stated that test development content area specialists used certain statistics from item and DIF analyses of the 2022 field test to identify items for further review. The flagging criteria for this purpose were specified in the previous section. Items without statistical flags were regarded as statistically acceptable and were not included in the data review. Likewise, items of extremely poor statistical quality were regarded as unacceptable and needed no further review. Such items were excluded from the Wisconsin item pool prior to the data review with DPI. The remaining flagged items were regarded by DRC content area test development specialists and a DRC psychometrician as needing further review. The intent was to capture all items that needed an additional review based on their statistical properties; thus, the criteria employed for item flagging tended to intentionally overidentify potential item issues.

The review of the new field test items with data was conducted in August 2022. This review was conducted online. As in past item review meetings, reviewers were first trained by a representative from DRC's staff with regard to the statistical indices used in item evaluation. This was followed by a discussion with examples concerning reasons that an item might be retained regardless of the statistics. The review process involved a brief exploration of possible reasons for the statistical profile of an item (e.g., possible bias, grade appropriateness, instructional issues) and a decision regarding acceptance. DRC content area test development specialists facilitated the review of the items. DPI reviewed the pool of field test items and made recommendations on each item and/or scenario/passage. The training presentation used at the data review meeting may be found in Appendix B. A summary of the data review results, including the number of items that were field-tested, the number and percentage of items with statistical flags, and the number and percentage of items rejected by DPI during the data review, is presented in Table 3-13. Items accepted for subsequent use in the Wisconsin Forward Exam were included in the pool of items for the Spring 2023 operational test form selection.

3.6 Form Development

The creation of test forms in a typical test development cycle involves the expertise of multiple DRC departments and DPI. The Wisconsin Forward Exam test development process complied with the following AERA, APA, & NCME (2014) Standards:

Standard 4.1 Test specifications should describe the purpose(s) of the test, the definition of the construct or domain measured, the intended examinee population, and interpretations for intended uses. The specifications should include a rationale supporting the interpretations and uses of test results for the intended purpose(s). (p. 85)

Standard 4.7 The procedures used to develop, review, and try out items and to select items from the item pool should be documented. (p. 87)

Standard 4.12 Test developers should document the extent to which the content domain of a test represents the domain defined in the test specifications. (p. 89)

The DRC team works cooperatively with DPI content and assessment specialists to select passages and prompts with associated content-specific items for the online assessments. The DRC team constructs forms that comply with the approved test blueprints and form construction guidelines. DRC uses an integrated team approach to test development, which includes content area specialists, psychometricians, and scoring specialists working as a unit in collaboration with DPI content experts.

3.6.1 Wisconsin Forward Test Form Creation

As stated earlier in the report, new operational test forms were developed for each grade and content area, for the Spring 2023 administration. As a first step in building the online assessments, the DRC team prepared all eligible items in DRC's item banking system, which is called IDEAS. The form, format, extent, and organization of items in their respective test sessions were determined in consultation with DPI.

Following the preparation of all necessary materials and resources, forms construction began. The construction of the test forms themselves was a collaborative effort within DRC's integrated development team of assessment specialists, psychometric services specialists, and scoring specialists.

Before test forms were created, passages, item/performance tasks, and artwork were carefully selected. The following process was used for item selection:

- Using the pool of vendor-owned items for ELA, Mathematics, and Science and Wisconsin owned items for Social Studies, DRC test development specialists first selected items to match the approved test blueprints.
- DRC test development specialists checked to see that each item clearly aligned with the standards where applicable and that each item with available item statistics met psychometric guidelines for inclusion in the test.
- DRC test development specialists verified that each item met technical quality requirements for well-crafted items, including that each item
 - had one clearly correct answer (or answers if the item was multi-select);
 - used clear and concise wording;

- was grammatically correct;
- had an appropriate range of difficulty;
- was free of any offensive, inappropriate, or biased content; and
- met the Principles of Universal Design and maximum accessibility.

In addition to content requirements, the following statistical criteria were used in item selection:

- Test length and item types match the DPI-approved test design.
- Content coverage matches the DPI-approved test blueprint.
- Items had acceptable statistics which included:
 - p -value between 0.20 and 0.90
 - Item-total test correlation >0.15
 - Omit rates $<3\%$
 - Acceptable fit statistics (no misfit flag)
 - No large DIF—If an item with large DIF had to be included in the test to maintain blueprint coverage, the item was examined to determine whether any content reason exists for the DIF flag (sometimes items demonstrate statistical bias but no content reason can be determined for the bias).

The statistical properties of the Spring 2022 test forms were used as targets for selection of the Spring 2023 ELA, Mathematics, Science, and Social Studies test forms. The item selection was conducted in two phases.

In the first phase, the anchor (linking) items were selected. The anchor items are used for the statistical linking of the new forms to previous test forms on already established test scales. All anchor items on the Spring 2023 test forms were selected from the Spring 2022 operational item pool. The anchor set was selected as a “mini” version of the full operational test for each grade level and content area in regard to its length, content coverage, and psychometric properties.

The length of the anchor sets was at least one-third of the length of the total test. The items included in the anchor sets met the same blueprint specifications as the full test in regard to the percentage of score points measuring each content standard. In addition, the psychometric properties of the anchor sets matched the corresponding properties of the target forms as closely as possible. Anchor selections were reviewed and approved by a DRC psychometrician.

In the second phase of the item selection process, non-anchor operational items were selected. The non-anchor operational items came from the operational and field test pools of items previously administered in Wisconsin. The non-anchor operational items were selected using the item selection guidelines presented earlier in this section. Full form selections were reviewed and approved by a DRC psychometrician.

In constructing the final forms, the DRC content area test development specialists followed the guidelines provided below:

- Forms included adequate standards coverage as required by test blueprints.
- No item in a form “clued” another item on that same form.
- Forms were diverse in terms of artwork and graphics.
- Forms included a wide range of topics and a variety of questions.
- Correct answer distributions were reasonable across MC items on the form.
- Forms did not contain any items that had been released to the public.
- DPI reviewed and gave final approval of all online test forms.

The test maps in Appendices C, D, E, and F provide details on the operational items on the Spring 2023 Wisconsin Forward Exam per grade and content area. The test maps include the session number, item sequence, item type, item usage, item maximum score, depth-of-knowledge level, standard code, and domain name. The ELA test maps are included in Appendix C, the Mathematics test maps are contained in Appendix D, the Science test maps are provided in Appendix E, and the Social Studies test maps are given in Appendix F.

3.6.2 Item and Form Quality Reviews

In all phases of the item and form development process, content area test development specialists and editorial specialists reviewed items and passages for technical quality; alignment with the standards; issues of bias, fairness, and sensitivity; depth of knowledge; estimated difficulty; and adherence to the Principles of Universal Design in all steps of the forms creation and forms review processes. The aim for this team approach was to conduct a multitiered internal review of all passages and items prior to submission for review by DPI and then, with approval by DPI prior to submission, for review by Wisconsin educators to ensure that all items align with Wisconsin’s standards and adhere to DPI’s standards for high-quality items.

DRC content and editorial teams reviewed all passages and items to ensure that they possessed the following traits:

- content alignment or congruence with the knowledge and skills specified in the standards;
- a range of estimated difficulty levels;
- appropriate grade-level vocabulary, subject matter, and assumed student knowledge;
- freedom from issues or concerns regarding bias, sensitivity, or fairness;
- accessibility, following the Principles of Universal Design; and
- correct grammar, usage, and structure/format.

As a part of DRC’s internal review of the items and test forms, the test development team members and graphic specialists ensured that item art could be reproduced clearly and accurately when electronically displayed and when used in the print-on-demand forms.

Test specifications were reviewed to identify any potential display requirements that may present challenges in an electronic display environment. Display tolerances are impacted by line thickness, percentage of screening for shading, specialized fonts and symbols, photographs, and color. These are defined in the early stages of the item and test development process to help guide the delineation of style requirements and specifications.

Item art was produced using transparent vector graphics that allow for adjustments without the breakdown of image clarity, which is common with lower-quality formats, and that provide for the online accommodation of alternate background colors. The DRC multitiered quality assurance process made certain that converted item art was carefully compared to the original format throughout the test development and production process.

In reviewing forms in the online environment, multiple reviewers checked passages and items on the multiple electronic platforms on which students took the test to ensure a smooth testing experience.

3.7 DPI Approvals

DPI had the opportunity to review passages and items placed on the Spring 2023 Wisconsin Forward Exam during the following phases:

- prior to item content review in Summer 2021
- at item content review in Summer 2021
- during data review of flagged field test items in Summer 2022
- during the form construction process in Fall 2022
- before the Spring 2023 test administration

Prior to the opening of the testing window, all online forms were made accessible to DPI for review in DRC’s secure INSIGHT testing engine.

3.8 Summary

In summary, the Spring 2023 Wisconsin Forward Exam assessment adhered to the Wisconsin test blueprints and test designs for each grade level and content area. The items included in the Spring 2023 Wisconsin Forward Exam were reviewed by DRC, DPI, and Wisconsin educators for issues regarding accessibility, bias, sensitivity, and content. During the reviews, experts identified (1) issues that could negatively affect a student’s ability to access stimuli and items, (2) content in stimuli and items that could unfairly affect a student’s response because of the student’s background, (3) developmental appropriateness, and (4) the alignment of stimuli and items to the content specifications. Item content was checked for the accuracy of the content, answer keys, and scoring rules. Following Spring 2022 field-testing, items flagged for accessibility, bias and sensitivity, and/or other content concerns were further reviewed by DRC and DPI to determine whether these flagged items

should be removed from the Wisconsin item pool prior to the form construction of the Wisconsin Forward Exam. The efforts and procedures used in the development of the Spring 2023 Wisconsin Forward Exam forms balanced the content and psychometric requirements for the form development. The psychometric properties of the ELA, Mathematics, Science, and Social Studies test forms were comparable to the psychometric properties of the Spring 2022 (target) forms. Overall, the process implemented in the Spring 2023 operational form development was in alignment with multiple best practices of the testing industry.

Table 3-1 Test Blueprints for English Language Arts Grades 3–8

Domain (Reporting Category)	Depth of Knowledge	Total Points by Grade					
		3	4	5	6	7	8
Reading		22	24	24	24	24	24
Key Ideas and Details	grade 3: 1–3 grades 4–8: 2–3	6–12	6–12	6–12	6–12	6–12	6–12
Craft and Structure/Integration of Knowledge and Ideas	all grades: 2–3	4–10	4–10	4–10	4–10	4–10	4–10
Vocabulary Use—Includes Language Standards 4 and 5	grades 3–5: 1–3 grades 6–8: 2–3	4–6	4–6	4–6	4–6	4–6	4–6
Literature		about 60%	about 60%	about 60%	about 50%	about 50%	about 50%
Informational Text		about 40%	about 40%	about 40%	about 50%	about 50%	about 50%
Writing/Language		24	24	24	24	24	24
Text Types and Purposes/Text-Dependent Analysis	all grades: 1–3	10–14	10–14	10–14	10–14	10–14	10–14
Research	all grades: 1–3	6–8	6–8	6–8	6–8	6–8	6–8
Language Conventions	all grades: 1–3	6–8	6–8	6–8	6–8	6–8	6–8
Listening	all grades: 2–3	7	8	8	8	8	8
ELA Points Total		53	56	56	56	56	56

Table 3-2 Test Blueprints for Mathematics Grades 3–8

Reporting Category	Depth of Knowledge	Total Points by Grade					
		3	4	5	6	7	8
Operations and Algebraic Thinking	grades 3–5: 1–3	8–10	9–11	8–10			
Number and Operations in Base Ten	grades 3–5: 1–3	7–9	8–10	8–10			
Number and Operations—Fractions	grades 3–5: 1–3	7–9	9–11	8–10			
Measurement and Data	grades 3–5: 1–3	9–11	9–11	9–11			
Geometry	grades 3–8: 1–2	6–8	6–8	8–10	6–8	9–11	9–11
Ratios and Proportional Relationships	grades 6–7: 1–3				6–8	7–9	
The Number System	grades 6–8: 1–3				10–12	6–8	7–9
Expressions and Equations	grades 6, 8: 1–3 grade 7: 1–2				10–12	9–11	9–11
Statistics and Probability	grades 6–8: 1–3				9–11	10–12	7–9
Functions	grade 8: 1–3						9–11
Mathematics Points Total		42	46	46	46	46	46

Table 3-3 Test Blueprints for Science Grades 4 and 8

Reporting Category	Depth of Knowledge	Total Points by Grade	
		4	8
Practices and Crosscutting Concepts in Life Science	grades 4, 8: 2-3	8-12	8-12
Practices and Crosscutting Concepts in Physical Science	grades 4, 8: 2-3	8-12	8-12
Practices and Crosscutting Concepts in Earth and Space Science	grades 4, 8: 2-3	8-12	8-12
Practices and Crosscutting Concepts in Engineering	grades 4, 8: 2-3	8-12	8-12
Science Total Points		40	40

Table 3-4 Test Blueprints for Social Studies Grades 4, 8, and 10

Reporting Category	Depth of Knowledge	Total Points by Grade		
		4	8	10
Geography (Inquiry Practices and Processes)	grade 4: 1-3 grade 8: 2-3 grade 10: 2	8-12	8-12	8-10
History (Inquiry Practices and Processes)	all grades: 2-3	8-12	8-12	8-10
Political Science (Inquiry Practices and Processes)	grade 4: 2-3 grade 8: 1-3 grade 10: 2	6-8	6-8	8-10
Economics (Inquiry Practices and Processes)	all grades: 2-3	6-8	6-8	6-8
Behavioral Sciences (Inquiry Practices and Processes)	all grades: 2-3	6-8	6-8	6-8
Social Studies Total Points		40	40	40

Table 3-5 Item Type Descriptions for Items on the Wisconsin Forward Exam

Item Type	Name	Description
EBSR	Evidence-Based Selected Response	Each evidence-based selected response item has two parts, and each two-part item is designed to elicit an evidence-based response from a student who has read a literature text passage, an informational text passage, or a writing concept. In part one, which is similar to a multiple-choice item, the student analyzes a passage or writing concept and chooses the best answer from four response options. In part two, the student uses evidence from the passage or writing concept to select one or more answers based on the response to part one. EBSR items can worth one or two points.
MC	Multiple Choice	Each multiple-choice item has four response options, only one of which is correct. Multiple-choice items are used to assess a variety of skill levels, from short-term recall of information to inference and problem-solving. Each of these items is worth one point.
MS	Multiple Select	Each multiple-select item requires a student to evaluate information presented and respond by choosing two or more correct responses. Multiple-select items can be used to assess multiple skills and concepts in a given content area. MS items can worth one or two points.
SA	Short Answer	Each short-answer item requires a student to enter a short numeric or algebraic response. These items are designed to assess a student’s ability to formulate a solution to a pure or applied math problem without the assistance of response options. The short-answer items are scored on a 0–1-point scale using item-specific autoscoring rules.
TDA	Text-Dependent Analysis	Each text-dependent analysis item is a text-based analysis based on a passage or a multiple-passage set that each student has read during the assessment. Both literary and informational texts are addressed through this item type. Students must draw on basic writing skills while inferring and synthesizing information from the passage in order to develop a comprehensive, holistic essay response. The demands required of a student’s reading and writing skills in response to a TDA item coincide with the similar demands required for a student to be college and career ready. The TDA prompts are scored using a holistic scoring guideline on a 1–4-point scale. A weight of 2 is applied to the item scores in the computation of the student total test raw scores and scale scores. That is, the TDA prompts contribute up to 8 raw score points toward the student total test raw score. This item type is supported by all Wisconsin ELA standards across all grades for both Reading Literature and Reading Informational Texts and by Writing standards 1, 2, 3, 4, and 9 across all grades. The TDA items are scored using artificial intelligence (AI) scoring, with an appropriate level of human scoring to validate the AI algorithms for all TDA items used in the Wisconsin ELA grades 3–8 assessments.
TE	Technology Enhanced	Each technology-enhanced item is designed to elicit evidence of a broad range of student understanding. A student interacts with the enhanced features of these computer-delivered, autoscorable test items to show understanding of skills and concepts. Item types such as drag-and-drop, hot-spot, number line and coordinate graphing, data displays, matching interaction, and drop-down menus are just some of the technology-enhanced items presented to a student. TE items can worth one or two points.

Table 3-6 Test Design for English Language Arts

Test Design		Grade					
		3	4	5	6	7	8
Number of Passage Sets	Literature	2-3	2-3	2-4	2-4	2-4	2-4
	Informational	1-2	1-2	1-3	2-4	2-4	2-4
	Listening	2-3	2-3	2-3	2-3	2-3	2-3
Number of Core (OP) Items	Item Types: MC/TE (1 pt)	20-30	22-34	22-34	22-34	22-34	22-34
	Item Types: MS/TE/EBSR (2 pts)	7-12	7-13	7-13	7-13	7-13	7-13
	Item Type: TDA (4 pts x 2)	1	1	1	1	1	1
	Total Core Items	32-38	32-40	32-40	32-40	32-40	32-40
Total Core Points		53	56	56	56	56	56
Total Estimated Testing Time (minutes)		130	130	130	130	130	130

Note: TDA items are scored using a 1–4-point scoring rubric. A weight of 2 is applied to item scores in the computation of the student total test raw scores and scale scores.

Table 3-7 Test Design for Mathematics

Test Design		Grade					
		3	4	5	6	7	8
Number of Core (OP) Items	Item Types: MC/SA(1 pt)	32-36	36-40	36-40	36-40	36-40	36-40
	Item Type: TE (1 pt)	6-10	6-10	6-10	6-10	6-10	6-10
	Total Core Items	42	46	46	46	46	46
Total Core Points		42	46	46	46	46	46
Total Estimated Testing Time (minutes)		90	90	90	105	105	115

Table 3-8 Test Design for Science

Test Design		Grade	
		4	8
Number of Core (OP) Items	Item Types: MC/MS/TE/EBSR (1 pt)	40	40
Total Core Points		40	40
Embedded Field Test (FT)	Number of Forms	20	20
	FT Items per Form	5	5
	Total Field Test Items	100	100
Total Items (Core + FT) per Form		45	45
Total Estimated Testing Time (minutes)		105	105

Table 3-9 Test Design for Social Studies

Test Design		Grade		
		4	8	10
Number of Core (OP) Items	Item Types: MC/TE/MS/EBSR (1 pt)	40	40	40
Total Core Points		40	40	40
Total Estimated Testing Time (minutes)		70	70	70

Table 3-10 Elements of Universal Design

Element	Explanation
Inclusive Assessment Population	Tests designed for state, district, or school accountability must include every student except those in the alternate assessment, and this is reflected in assessment design and field-testing procedures.
Precisely Defined Constructs	The specific constructs tested must be clearly defined so that all construct-irrelevant cognitive, sensory, emotional, and physical barriers can be removed.
Accessible, Unbiased Items	Accessibility is built into items from the beginning, and bias review procedures ensure that quality is retained in all items.
Amenable to Accommodations	The test design facilitates the use of needed accommodations.
Simple, Clear, and Intuitive Instructions and Procedures	All instructions and procedures are simple, clear, and presented in understandable language.
Maximum Readability and Comprehensibility	Readability and plain language guidelines are followed (e.g., sentence length and number of difficult words are kept to a minimum) to produce readable and comprehensible text.
Maximum Legibility	Characteristics that ensure easy decipherability are applied to text, tables, figures, illustrations, and response formats.

Table 3-11 College- and Career-Ready Item Bank Development Activities

Steps in Item Development	DRC College- and Career-Ready Item Bank Development Activity
1	Establish item/passage development specifications and style guides and prepare item writing training manuals.
2	Determine item development plans.
3	Train item writers and/or passage developers in the project requirements and specifications.
4	Develop passages and write items.
5	Review, edit, code, and track items and produce graphics.
6	Produce review forms for content and bias/fairness/sensitivity reviews by external reviewers.
7	Modify items based on external reviewers' recommendations.
8	Review and approve field test-ready items and passages.
9	Develop field test forms and administer field test.
10	Internally review field test item data.
11	Approve items to be included in the item bank.

Table 3-12 Items Reviewed during Summer 2021 Item Review

Grade	Number of Items			
	English Language Arts	Mathematics	Science	Social Studies
3	72	80		
4	72	80	103	42
5	72	80		
6	81	80		
7	70	80		
8	70	80	98	66
10				62
TOTAL	437	480	201	170

Table 3-13 Items Reviewed during Summer 2022 Item Data Review

Content Area	Grade	Number of Items in 2022 Field Test	Field Test Items Flagged for Poor Statistics or DIF		Field Test Items Rejected at Data Review for Statistical or Content-Related Reasons	
			Number of Items	Percentage of All Field Test Items	Number of Items	Percentage of All Field Test Items
English Language Arts	3	144	27	18.75	8	5.56
	4	144	33	22.92	5	3.47
	5	144	37	25.69	5	3.47
	6	128	27	21.09	4	3.13
	7	127	25	19.69	6	4.72
	8	124	23	18.55	1	0.81
Mathematics	3	128	26	20.31	14	10.94
	4	128	31	24.22	17	13.28
	5	127	30	23.62	23	18.11
	6	116	37	31.90	16	13.79
	7	128	58	45.31	50	39.06
	8	128	48	37.50	39	30.47
Science	4	94	21	22.34	15	15.96
	8	93	36	38.71	28	30.11
Social Studies	4	37	5	13.51	0	0.00
	8	62	18	29.03	12	19.35
	10	53	13	24.53	11	20.75

Part 4: Test Administration

In the Spring of 2023, Wisconsin administered assessments in ELA and Mathematics for grades 3–8. Science was administered in grades 4 and 8, and Social Studies was administered in grades 4, 8, and 10. The test administration window was March 20–April 28, 2023. Part 4 of the Technical Report provides information on student participation rates in the Spring 2023 assessments and describes a set of standardized procedures and policies applied to administer the Wisconsin Forward Exam. The issue of test security in test administration, which has important implications for the integrity of the results and, thus, the validity of Wisconsin Forward Exam scores, is also discussed. Documentation citing the written procedures provided to test administrators and school personnel to standardize the administration of the test is provided in this part as well. The following AERA, APA, & NCME (2014) Standards are addressed in Part 4: 3.4, 3.5, 4.15, 4.16, 6.1, 6.2, 6.3, 6.4, 6.6, and 6.7. Each standard will be explained within the relevant section of this part of the report.

DPI is committed to the proposition that all schools and all students will be held accountable to a common set of high academic content standards, the Wisconsin Academic Standards. As an alternate assessment for students being instructed using alternate academic achievement standards, the Wisconsin Essential Elements, the Dynamic Learning Maps assessment measures the academic progress of students with the most significant cognitive disabilities in the subject areas of ELA and Mathematics at grades 3–11 and Science at grades 4 and 8–11. A teacher rater form is used to assess these students in Social Studies at grades 4, 8, and 10.

All other students are accountable to the grade-level knowledge and skills outlined in the Wisconsin Academic Standards. Those students who have an IEP, a 504 plan (under Section 504 of the Rehabilitation Act of 1973) or are identified as limited English proficient (LEP) or formerly LEP may be eligible to receive testing accommodations or supports. Accommodations and supports are practices and procedures that provide equitable access to grade-level content. They are intended to reduce or eliminate the effects of a student’s disability or level of language acquisition; they do not reduce learning expectations. DPI guidance makes it clear that the accommodations or supports provided to a student must be consistent with classroom instruction, classroom assessments, and district and state assessments. It is important to note that while some accommodations or supports may be appropriate for instructional use, they may not be appropriate for use on a standardized assessment. AERA, APA, & NCME (2014) Standard 6.2 states the following:

When formal procedures have been established for requesting and receiving accommodations, test takers should be informed of these procedures in advance of testing. (p. 115)

An overview of the types of accommodations and supports available to students and the guidelines for test administration conditions are described below. Additionally, IEP teams were directed to the Wisconsin Forward Exam Accommodations and Supports page at <http://dpi.wi.gov/assessment/forward/accommodations> for guidance regarding all available accommodations and supports intended to provide equitable access to grade-level content and assessments.

District Assessment Coordinators (DACs) indicated which accommodations and supports were to be available for use by each student within the student learning profile in the DRC INSIGHT Portal. All student accommodations and supports are managed and can be monitored through the DRC INSIGHT Portal. This system is the interface to the administrative functions of the DRC INSIGHT Online Learning System, where students interface with their online assessments. As a function of this roles-based system, the primary users of the DRC INSIGHT Portal were DACs and School Assessment Coordinators (SACs) who were assigned permissions accordingly for security purposes. The major functions are those of managing users and managing students. As such, the DRC INSIGHT Portal was used to manage and update student information, including demographic and accommodations/accessibilities information. All DRC INSIGHT Portal user roles and permission levels were approved by DPI.

4.1 Student Participation

For the purposes of this report, the test participation rate is defined as the percentage of students who received a valid scale score compared to the total number of students who were scheduled to take the test. The test participation rates for the students in grades 3 through 7 ranged from approximately 95% to 96% across content areas. The test participation rates for the students in grade 8 were approximately 94% across all content areas. The participation rate in Social Studies grade 10 was approximately 88%. These participation rates were comparable to the Spring 2022 participation rates.

The participation rates were comparable for male and female students in all grade levels, for students in different ethnic groups in grades 3 through 5 and for economically disadvantaged and not economically disadvantaged students in grades 3 through 5, in all content areas. Lower participation rates were observed for African American and American Indian students in grades 6 through 8 compared to their White peers. These differences ranged from approximately 3% to 4%, depending on the grade level and content area. The participation rates of Hispanic and White students in grades 6 through 8 were comparable.

The participation rates were comparable for fully English proficient and limited English proficiency students for Mathematics. The limited English proficiency students participated at a lower rate (difference of approximately 3% to 4%) compared to the fully English proficient students in all grades of ELA assessments. Larger differences in test participation rates were found for groups of students with and without disabilities. Between 86% and 91% of students with disabilities participated in the assessments compared to over 95% of their peers without disabilities in grades 3 through 8.

The discrepancies in participation rates between subgroups were larger in Social Studies grade 10 compared to other grades. The participation rates for African American, Hispanic, and American Indian students were approximately 67%, 86%, and 74%, respectively, compared to the participation rate of 92% for White students. Students with disabilities and economically disadvantaged students participated at approximately 76% and 83%, respectively, compared to over 90% of their peers without disabilities or students not considered to be economically disadvantaged. Detailed information on the test participation rates in Spring 2023 for all students and disaggregated by demographic characteristics is provided in Appendix G.

4.2 Standardized Test Administration

Unstandardized testing conditions can pose a serious threat to test validity by adding construct-irrelevant variance to the test scores. McCallin (2006) described a number of such threats to validity, including alterations in test administration requirements (e.g., changing time limits, modifying test instructions, giving hints to examinees), variability across test sites.

(e.g., differences in facilities/equipment, inadvertent posting of instructional aids in classrooms), interruptions during test sessions (e.g., power outages, relocation of students during testing, disturbances, other distractions), test administrator practices that may exacerbate test anxiety in particular students, practices that elicit test wiseness, and security breaches that may result in the exposure of test forms or items. Construct-irrelevant variance may exert a systematic effect on the scores of individual students or groups of students, resulting in an overestimation or underestimation of their true abilities.

Standardized test administration, extensive training of the test scorers and artificial intelligence (AI) engine, and rigorous scoring rules for autoscored items for the Wisconsin Forward Exam comply with AERA, APA, & NCME (2014) Standards 3.4 and 3.5.

Standard 3.4 Test takers should receive comparable treatment during the test administration and scoring process. (p. 65)

Standard 3.5 Test developers should specify and document provisions that have been made to test administration and scoring procedures to remove construct-irrelevant barriers for all relevant subgroups in the test-taker population. (p. 65)

The standardized Wisconsin Forward Exam test administration procedures described in this part of this report were designed to address these potential threats to validity through the use of comprehensive security measures and the provision of detailed Test Administration Manuals and other training materials for DACs, SACs, and TAs.

4.3 Accessibility Resources

Accommodations were allowed for eligible individual students participating in the Wisconsin Forward Exam. Accommodations provided to a student must be documented in a current IEP and used during routine instruction. IEP teams were directed to refer to the Wisconsin Forward Exam accommodations policy and guidance at <https://dpi.wi.gov/assessment/forward/accommodations>.

It is important to note that students were provided access to a range of supports that included universal tools (available to all students), designated supports, and accommodations, including the Braille version of the Wisconsin Forward Exam, based on students' needs. Those supports are defined as follows.

4.3.1 Universal Tools

Universal tools are accessibility features that are available to all students based on student preference and selection. These accessibility features of the assessment are either provided as digitally delivered components of the test administration system (embedded) or separate from it (non-embedded).

Embedded Universal Tools (Online)

- Calculators
- Click to Enlarge
- Cross-Off Tools
- Flag/Mark for Review
- Help/What's This?
- Highlighter
- Go to Question
- Keyboard Navigation
- Line Guide
- Magnifier Tool (Zoom)
- Measuring Tools
- Pause (Breaks)
- Review Page
- Sticky Notes (Digital Notepad)
- Test Directions
- Tool Tips
- Writer's Checklist (ELA TDA Session only)

Non-embedded Universal Tools (Standard)

- Graph Paper
- Scratch Paper
- Writer's Checklist (ELA TDA Session only)

4.3.2 Designated Supports

Designated supports are those features that are available for use by any student for whom the need has been indicated by an educator or team of educators (with parent/guardian and student input as appropriate) and are part of the student's classroom instruction. They are either provided as part of the online test administration

system or separate from it (i.e., embedded or non-embedded). All embedded and non-embedded designated supports must be entered into the DRC INSIGHT Portal prior to test administration. Embedded and non-embedded supports will appear on student test tickets.

Embedded Designated Supports (Online)

- Color Choices (CC)
- Contrasting Color (CTC)
- Reverse Contrast (RC)
- Masking (MSK)
- Stacked Translations (Spanish)
- Text-to-Speech (TTS)

Non-embedded Designated Supports (Standard)

- Amplification Device
- Word-to-Word Bilingual Dictionary
- Color Overlay
- Magnification
- Noise Buffers
- Read Aloud in English
- Read Aloud in Spanish
- Scribe
- Separate Setting
- Small Group Translation
- Translator/Interpreter

4.3.3 Accommodations

Accommodations are features that increase equitable access but do not compromise the grade-level standard or intended outcome of the assessment. They are available for students for whom there is a documented need in the IEP or 504 accommodation plan and who use similar accommodation as part of their classroom instruction. Accommodations are either provided as part of the online test administration system or separate from it (i.e., embedded or non-embedded). All embedded and non-embedded accommodations must be entered into the DRC INSIGHT Portal prior to test administration. Embedded and non-embedded accommodations will appear on student test tickets.

Embedded Accommodations (Online)

- Video Sign Language (VSL)
- Closed-Captioning (C CAP)

Non-embedded Accommodations (Standard)

- Abacus
- Alternate Response Options
- Braille (Unified English Braille) (BRL)
- Calculator
- Listening Scripts (LS)
- Multiplication Table
- Print-on-Demand (POD)
- Read Aloud (Reading Passages)

4.3.4 Translation

For the Spring 2023 Wisconsin Forward Exam administration, the State of Wisconsin used an embedded stacked Spanish translation for Mathematics, Science, and Social Studies items. For ELA assessments, only the test directions are available in stacked translation. The stacked Spanish translation is a designated support for students who are native Spanish speakers and are identified as limited English proficient to demonstrate their knowledge on the Wisconsin Forward Exam. In addition to the embedded stacked translation, bilingual word lists and a translation of the test directions are allowable designated supports.

DPI recognizes that approximately five percent of the Wisconsin limited English proficient population speaks a language other than Spanish, and specific guidelines are provided for these students. Districts that serve students who speak languages other than Spanish may have used qualified translators to provide oral translation support to students. However, the use of translation support was restricted to Mathematics, Science, and Social Studies tests, given that the test constructs are not specific to the English language. DPI recommended that educators consult the list of allowable accommodations and supports (referenced above) to create the most appropriate testing situation for their students.

4.3.5 Additional Accessibility Resources

Additional accessibility resources and guidance included the following:

- **Multiplication Table:** This resource is a non-embedded accommodation available for students who have it in their IEP or 504 plan for grades 4–8 Mathematics.
- **Read Aloud Guidelines:** This document outlines the qualifications, guidelines, and procedures required for a test reader. The test reader must sign the Read Aloud Agreement to Maintain Security

and Confidentiality prior to test administration. Completed agreement forms should be retained by the Site Assessment Coordinator.

- Scribing Guidelines: This document outlines the qualifications, guidelines, and procedures required when using a scribe.
- Interpreter Guidelines: This document outlines the qualifications, guidelines, and procedures required when using an interpreter.

Tables 4-1 through 4-7 provide the list of accommodations or designated supports made available for the Spring 2023 Wisconsin Forward Exam along with the number and percentage of students provided these accommodations or supports. The counts are based on the accommodations and designated supports selected via the DRC INSIGHT Portal. Scores of assessments taken with accommodations were included with the results for students who took these tests under standard conditions and presented at the school, district, and state levels.

4.4 Test Security

Maintaining the security of all test materials is crucial to preventing the possibility of random or systematic errors, such as an unauthorized exposure of test items that would affect the valid interpretation of test scores. Several test security measures have been implemented for the Wisconsin Forward Exam with compliance to the following AERA, APA, & NCME (2014) Standards:

Standard 6.6 Reasonable efforts should be made to ensure the integrity of test scores by eliminating opportunities for test takers to attain scores by fraudulent or deceptive means. (p. 116)

Standard 6.7 Test users have the responsibility of protecting the security of test materials at all times. (p. 117)

The primary goal of test security is to protect the integrity of the assessments and ensure that scores retain their interpretability. To ensure that trends in achievement results can be calculated across years and to provide longitudinal data, a certain number of test questions must be repeated from year to year. If any of these questions are made public, the validity of the test may be compromised. Because the Wisconsin Forward Exam is administered virtually 100 percent online, printed test materials are limited to the very few cases where a student requires a printed version of the test as provided in the IEP (i.e., Braille and Print-on-Demand), so the assessment exposure is limited to those educators who require access for those purposes. DPI and DRC ensured that all who had access to any materials associated with the Wisconsin Forward Exam understood the critical need for test security. They presented security requirements during the pre-test workshops and outlined the acceptable and unacceptable test preparation and administration practices. The Wisconsin Forward Exam was administered under secure testing conditions established by DPI.

Other security measures for Wisconsin Forward Exam test administrations are described below:

- The use of any unauthorized electronic device is prohibited during testing.

- Password-protected, role-based administrator access to all test setup, management, and reporting functions is required.
- Student Test Login Tickets provide secure student access to the test using a unique username and password.
- Test content is securely transferred using leading encryption technologies; content is decrypted when the student login is validated.
- Decrypted test content is purged from the system’s memory upon completion of the test session.
- Device lockdown during testing prevents students from copying, pasting, printing, and accessing other applications.
- If the test is paused, content is removed from the screen to ensure security of test content. The system will time out and close the test after a defined period of inactivity.
- Extensive software quality assurance tests ensure that all data are scanned, captured, and accurately scored in the secure database and all associated reports contain accurate data.

The online systems provided by DRC that are associated with the administration of the Wisconsin Forward Exam have all been designed to provide the level of security required by DPI and described in the DPI Test Security Manual for its assessment programs. Student testing environments are designed to ensure the protection of responses as well as student data (as required under the federal Family Educational Rights and Privacy Act). DRC’s information security policies and procedures are based on the National Institute of Standards and Technology (NIST) criteria (NIST Standard 800-53). This is a nationally recognized standard for information security practices.

4.4.1 Secure Student Access

Students are required to provide a valid username and password to access the online testing system. The Test Administrator (TA) provides each student with a Student Test Login Ticket, which contains the student’s username and a unique, pre-generated password. A separate, unique password is generated for each assessment, ensuring that students can only access the content designated for that particular test. Passwords are generated randomly for each student to use. Test tickets are generated from within the DRC INSIGHT Portal secure administrative system, which is pre populated with student records. As an additional security measure, after a student logs in, a Student Verification Page prompts the student to verify their profile information, including any assigned accommodations, prior to initiating the test. The student’s name is also displayed on the screen during the test, providing an additional verification check for the student and the TA.

Test tickets and rosters are considered secure materials. Therefore, it is recommended that test tickets be printed as close to the date of testing as possible, and sites are instructed to keep test tickets and rosters in a secure location until the session is scheduled to begin. Test tickets are distributed just prior to students logging in and are collected after all students have logged in and begun testing; directions also include a request to count the number of tickets that are distributed and collected after students log in to make sure the numbers of tickets are

the same. After a testing session is complete, all test tickets are returned to the Site Assessment Coordinator for secure destruction or secure storage.

4.4.2 Test Security during Breaks

Test security must be maintained during all breaks within a testing session. To lessen the risk of a security breach occurring during these breaks, students requiring the use of restroom facilities must be escorted by either a proctor or a test examiner. In addition, students must not be allowed to use any form of wireless communication during these breaks.

4.5 Test Administration Training

DRC provided pre-recorded training for DACs, SACs, and TAs for the Spring 2023 administration of the Wisconsin Forward Exam. The webinars were recorded by DRC. The purpose of the webinars was to keep districts and schools informed about policies and procedures related to the Wisconsin Forward Exam administration. The information covered in the webinars included standardizing the administration of the Wisconsin Forward Exam, maintaining the security of the assessments, allowing access to the assessments for special populations by providing appropriate designated supports or accommodations, and providing guidance on appropriate interpretations of the test results. These communication efforts by DPI and the ancillary information developed by DRC are in alignment with multiple best practices of the testing industry and, in particular, support the following AERA, APA, & NCME (2014) Standards:

Standard 4.15 The directions for test administration should be presented with sufficient clarity so that it is possible for others to replicate the administration conditions under which the data on reliability, validity, and (where appropriate) norms were obtained. Allowable variations in administration procedures should be clearly described. The process for reviewing requests for additional testing variations should also be documented. (p. 90)

Standard 4.16 The instructions presented to test takers should contain sufficient detail so that test takers can respond to a task in the manner that the test developer intended. When appropriate, sample materials, practice or sample questions, criteria for scoring, and a representative item identified with each item format or major area in the test's classification or domain should be provided to the test takers prior to the administration of the test, or should be included in the testing material as part of the standard administration instructions. (p. 90)

Standard 6.1 Test administrators should follow carefully the standardized procedures for administration and scoring specified by the test developer and any instructions from the test user. (p. 114)

Standard 6.2 When formal procedures have been established for requesting and receiving accommodations, test takers should be informed of these procedures in advance of testing. (p. 115)

Standard 6.3 Changes or disruptions to standardized test administration procedures or scoring should be documented and reported to the test user. (p. 115)

Standard 6.4 The testing environment should furnish reasonable comfort with minimal distractions to avoid construct-irrelevant variance. (p. 116)

In order to ensure standardized testing administration for all students, a DAC/SAC Administration Management section was included in the Test Administration Manual and made available to all assessment coordinators. The guide included the following topics:

- Testing Roles and Responsibilities
- Test Security
- Resources and Training Materials
- Test Schedules
- DRC INSIGHT Portal
- Accessibility
- Student Transfers
- Prior to the Close of the Testing Window
- Data and Reporting

Test Administration Manuals, made available to all TAs, included the following topics:

- Key Dates
- TA Responsibilities
- Test Times
- Test Security
- Accessibility Information
- Prior to Testing Instructions
- Test Tickets
- During Testing Information
- Test Administration Script

These topics were also addressed in the recorded trainings that were posted for online access.

Student Preparation for Online Testing

Prior to testing, schools and districts were encouraged to provide students with time to complete both a tutorial video series and an online tools training. Sample test items were also provided for each grade and content area.

Student and Administrator Tutorial Videos

Student and administrator tutorial videos were available for students and TAs to become familiar with the online testing environment. Tutorials could be viewed as a class or at an individual student machine by launching INSIGHT and clicking on DRC INSIGHT Online Assessment Tutorials.

Online Tools Training

The Online Tools Training (OTT) was provided for students to have a hands on opportunity to practice the types of items and tools available in the online testing system. The OTTs were available publicly for practice using a Chrome browser. Users (at home or school) could visit <https://dpi.wi.gov/assessment/forward/sample-items> to access the public OTTs. The OTTs could also be accessed on student testing devices once INSIGHT was installed. General OTTs were made available for each content area and grade level. Separate OTTs were available for students to practice using VSL closed-captioning, TTS, Spanish translation, masking, and color choice tools. VSL and Spanish OTTs were available by grade band (3–5, 6–8, and 10). The OTTs were not scored and were not intended for content practice.

Item Samplers

Item samplers were developed to be used by both educators and students to gain familiarity with the various item types and their varying functionalities. The format appears as a “guided practice test” in the online, PDF, and Braille versions of the tests.

Accommodated versions of the item samplers, reflecting the Wisconsin Forward Exam, were produced, including TTS, stacked Spanish translation (in Mathematics, Science, and Social Studies), and VSL. All tools and supports available in the test engine were applied to this student online experience.

Access to the item samplers was granted through the OTT menu page. A username and password were displayed on the login screen. The “click to enlarge” item displayed the answer key and scoring guide for each item online. In addition, a paper answer key and scoring guide were provided as a document for posting.

Administration Supports before and after Testing

With a few exceptions (i.e., accommodated student versions), the Wisconsin Forward Exam was administered fully online. Because DRC produced a variety of Wisconsin-specific manuals with process reviews by DRC program management staff, DRC editorial staff, and DPI staff, substantial consideration was given to the information required for successful online testing to occur. DPI provided final approval for each document prior to delivery and public posting.

Table 4-8 displays a list of electronic manual materials that DRC developed in conjunction with DPI. A final PDF of each deliverable was provided to DPI to post to the DPI informational website to allow districts to review and/or print.

For additional and more specific information related to test administration, refer to the Test Administration Manual that is available online at <https://dpi.wi.gov/assessment/forward/resources#manuals>.

4.6 Summary

This part of the report provides information on student participation rates and summarizes the processes and activities implemented and the information disseminated to help ensure standardized test administration procedures and, thus, uniform test administration conditions for students. It describes how the test administration procedures implemented for the Wisconsin Forward Exam were in alignment with best practices of the testing industry.

Table 4-1 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 3

Grade 3 Accommodation or Support	English Language Arts		Mathematics	
	N Count	Percent	N Count	Percent
Used Braille [BRL]	1	0.00	1	0.00
Used Print on Demand [POD]	1	0.00	1	0.00
Used Bilingual Dictionary			76	0.13
Used Magnification	69	0.12	67	0.11
Used Noise Buffers	697	1.19	670	1.14
Used Read Aloud	350	0.60	373	0.64
Used Scribe	495	0.85	405	0.69
Used Separate Setting	8333	14.25	8405	14.31
Used Alternate Response Options	12	0.02	17	0.03
Used Read Aloud (Reading Passages)	4	0.01		
Provided Color Choices [CC]	52	0.09	51	0.09
Used Contrasting Color [CTC]	39	0.07	38	0.06
Used Reverse Contrast [RC]	22	0.04	31	0.05
Used Masking [MSK]	646	1.10	633	1.08
Used Text-to-Speech [TTS]	13817	23.62	13978	23.80
Used Spanish Translation [ST]	370	0.63	711	1.21
Used Video Sign Language [VSL (ASL)]	13	0.02	14	0.02
Used Color Overlay	19	0.03	17	0.03
Amplification Device	81	0.14	80	0.14
Small Group Translation	90	0.15	149	0.25
Translator/Interpreter	24	0.04	63	0.11
Read Aloud in Spanish	76	0.13	135	0.23
Used Closed Captioning [C CAP] ELA	57	0.10		
Used Listening Scripts [LS] ELA	9	0.02		
Used Abacus Math			37	0.06
Used Non-embedded Calculator Math			185	0.32

Table 4-2 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 4

Grade 4 Accommodation or Support	English Language Arts		Mathematics		Science		Social Studies	
	N Count	Percent	N Count	Percent	N Count	Percent	N Count	Percent
Used Braille [BRL]	1	0.00	1	0.00	1	0.00	1	0.00
Used Print on Demand [POD]	4	0.01	4	0.01	4	0.01	4	0.01
Used Bilingual Dictionary			150	0.25	165	0.28	164	0.28
Used Magnification	55	0.09	55	0.09	55	0.09	54	0.09
Used Noise Buffers	763	1.29	762	1.29	760	1.29	759	1.28
Used Read Aloud	297	0.50	306	0.52	298	0.50	304	0.51
Used Scribe	489	0.83	376	0.64	372	0.63	373	0.63
Used Separate Setting	8759	14.85	8826	14.92	8673	14.66	8654	14.64
Used Alternate Response Options	11	0.02	10	0.02	10	0.02	11	0.02
Used Read Aloud (Reading Passages)	2	0.00						
Provided Color Choices [CC]	73	0.12	73	0.12	72	0.12	72	0.12
Used Contrasting Color [CTC]	49	0.08	49	0.08	49	0.08	49	0.08
Used Reverse Contrast [RC]	33	0.06	33	0.06	33	0.06	34	0.06
Used Masking [MSK]	654	1.11	659	1.11	656	1.11	654	1.11
Used Text-to-Speech [TTS]	13259	22.47	13352	22.57	13276	22.45	13289	22.47
Used Spanish Translation [ST]	417	0.71	660	1.12	619	1.05	607	1.03
Used Video Sign Language [VSL (ASL)]	23	0.04	24	0.04	24	0.04	23	0.04
Used Color Overlay	22	0.04	22	0.04	22	0.04	22	0.04
Amplification Device	64	0.11	64	0.11	64	0.11	64	0.11
Small Group Translation	124	0.21	192	0.32	183	0.31	182	0.31
Translator/Interpreter	29	0.05	60	0.10	58	0.10	56	0.09
Read Aloud in Spanish	105	0.18	160	0.27	161	0.27	162	0.27
Used Closed Captioning [C CAP] ELA	71	0.12						
Used Listening Scripts [LS] ELA	7	0.01						
Used Abacus Math			26	0.04				
Used Non-embedded Calculator Math			218	0.37				
Used Multiplication Table Math			1936	3.27				

Table 4-3 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 5

Grade 5 Accommodation or Support	English Language Arts		Mathematics	
	N Count	Percent	N Count	Percent
Used Braille [BRL]	2	0.00	2	0.00
Used Print on Demand [POD]	1	0.00	1	0.00
Used Bilingual Dictionary			109	0.18
Used Magnification	85	0.14	85	0.14
Used Noise Buffers	610	1.03	610	1.02
Used Read Aloud	294	0.50	308	0.52
Used Scribe	403	0.68	298	0.50
Used Separate Setting	8425	14.19	8522	14.30
Used Alternate Response Options	10	0.02	10	0.02
Used Read Aloud (Reading Passages)	0	0.00		
Provided Color Choices [CC]	63	0.11	60	0.10
Used Contrasting Color [CTC]	61	0.10	60	0.10
Used Reverse Contrast [RC]	34	0.06	33	0.06
Used Masking [MSK]	566	0.95	567	0.95
Used Text-to-Speech [TTS]	11590	19.52	11723	19.68
Used Spanish Translation [ST]	412	0.69	679	1.14
Used Video Sign Language [VSL (ASL)]	13	0.02	13	0.02
Used Color Overlay	25	0.04	25	0.04
Amplification Device	82	0.14	82	0.14
Small Group Translation	96	0.16	126	0.21
Translator/Interpreter	22	0.04	52	0.09
Read Aloud in Spanish	54	0.09	92	0.15
Used Closed Captioning [C CAP] ELA	75	0.13		
Used Listening Scripts [LS] ELA	10	0.02		
Used Abacus Math			17	0.03
Used Non-embedded Calculator Math			315	0.53
Used Multiplication Table Math			2312	3.88

Table 4-4 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 6

Grade 6 Accommodation or Support	English Language Arts		Mathematics	
	N Count	Percent	N Count	Percent
Used Braille [BRL]	1	0.00	1	0.00
Used Print on Demand [POD]	1	0.00	1	0.00
Used Bilingual Dictionary			196	0.33
Used Magnification	52	0.09	52	0.09
Used Noise Buffers	372	0.63	372	0.62
Used Read Aloud	248	0.42	263	0.44
Used Scribe	274	0.46	203	0.34
Used Separate Setting	7722	13.00	7771	13.05
Used Alternate Response Options	6	0.01	5	0.01
Used Read Aloud (Reading Passages)	1	0.00		
Provided Color Choices [CC]	45	0.08	45	0.08
Used Contrasting Color [CTC]	89	0.15	89	0.15
Used Reverse Contrast [RC]	28	0.05	28	0.05
Used Masking [MSK]	500	0.84	502	0.84
Used Text-to-Speech [TTS]	9568	16.10	9619	16.15
Used Spanish Translation [ST]	335	0.56	488	0.82
Used Video Sign Language [VSL (ASL)]	12	0.02	12	0.02
Used Color Overlay	16	0.03	16	0.03
Amplification Device	54	0.09	54	0.09
Small Group Translation	76	0.13	120	0.20
Translator/Interpreter	18	0.03	35	0.06
Read Aloud in Spanish	29	0.05	68	0.11
Used Closed Captioning [C CAP] ELA	74	0.12		
Used Listening Scripts [LS] ELA	14	0.02		
Used Abacus Math			22	0.04
Used Non-embedded Calculator Math			537	0.90
Used Multiplication Table Math			2799	4.70

Table 4-5 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 7

Grade 7 Accommodation or Support	English Language Arts		Mathematics	
	N Count	Percent	N Count	Percent
Used Braille [BRL]	3	0.00	3	0.00
Used Print on Demand [POD]	0	0.00	0	0.00
Used Bilingual Dictionary			232	0.38
Used Magnification	76	0.13	76	0.13
Used Noise Buffers	302	0.50	303	0.50
Used Read Aloud	256	0.42	283	0.47
Used Scribe	255	0.42	187	0.31
Used Separate Setting	7674	12.70	7729	12.76
Used Alternate Response Options	11	0.02	10	0.02
Used Read Aloud (Reading Passages)	5	0.01		
Provided Color Choices [CC]	80	0.13	61	0.10
Used Contrasting Color [CTC]	150	0.25	114	0.19
Used Reverse Contrast [RC]	104	0.17	81	0.13
Used Masking [MSK]	496	0.82	455	0.75
Used Text-to-Speech [TTS]	9093	15.05	9216	15.22
Used Spanish Translation [ST]	354	0.59	481	0.79
Used Video Sign Language [VSL (ASL)]	16	0.03	16	0.03
Used Color Overlay	18	0.03	18	0.03
Amplification Device	61	0.10	60	0.10
Small Group Translation	107	0.18	140	0.23
Translator/Interpreter	19	0.03	35	0.06
Read Aloud in Spanish	28	0.05	54	0.09
Used Closed Captioning [C CAP] ELA	92	0.15		
Used Listening Scripts [LS] ELA	23	0.04		
Used Abacus Math			16	0.03
Used Non-embedded Calculator Math			670	1.11
Used Multiplication Table Math			2763	4.56

Table 4-6 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 8

Grade 8 Accommodation or Support	English Language Arts		Mathematics		Science		Social Studies	
	N Count	Percent	N Count	Percent	N Count	Percent	N Count	Percent
Used Braille [BRL]	6	0.01	6	0.01	6	0.01	6	0.01
Used Print on Demand [POD]	2	0.00	2	0.00	2	0.00	2	0.00
Used Bilingual Dictionary			225	0.36	222	0.36	220	0.35
Used Magnification	70	0.11	69	0.11	68	0.11	68	0.11
Used Noise Buffers	244	0.39	243	0.39	241	0.39	240	0.39
Used Read Aloud	223	0.36	236	0.38	230	0.37	244	0.39
Used Scribe	214	0.34	167	0.27	168	0.27	169	0.27
Used Separate Setting	7842	12.60	7891	12.65	7752	12.45	7772	12.48
Used Alternate Response Options	14	0.02	12	0.02	12	0.02	12	0.02
Used Read Aloud (Reading Passages)	6	0.01						
Provided Color Choices [CC]	133	0.21	77	0.12	74	0.12	75	0.12
Used Contrasting Color [CTC]	150	0.24	107	0.17	107	0.17	108	0.17
Used Reverse Contrast [RC]	97	0.16	61	0.10	63	0.10	62	0.10
Used Masking [MSK]	431	0.69	380	0.61	380	0.61	379	0.61
Used Text-to-Speech [TTS]	8788	14.12	8888	14.25	8763	14.07	8702	13.98
Used Spanish Translation [ST]	361	0.58	471	0.76	493	0.79	533	0.86
Used Video Sign Language [VSL (ASL)]	10	0.02	10	0.02	10	0.02	10	0.02
Used Color Overlay	14	0.02	14	0.02	14	0.02	14	0.02
Amplification Device	49	0.08	46	0.07	46	0.07	46	0.07
Small Group Translation	92	0.15	127	0.20	125	0.20	128	0.21
Translator/Interpreter	14	0.02	31	0.05	29	0.05	29	0.05
Read Aloud in Spanish	32	0.05	59	0.09	66	0.11	70	0.11
Used Closed Captioning [C CAP] ELA	87	0.14						
Used Listening Scripts [LS] ELA	18	0.03						
Used Abacus Math			10	0.02				
Used Non-embedded Calculator Math			810	1.30				
Used Multiplication Table Math			2591	4.15				

Table 4-7 Number and Percentage of Students Using Accommodations or Designated Supports, Grade 10

Grade 10	Social Studies	
Accommodation or Support	N Count	Percent
Used Braille [BRL]	2	0.00
Used Print on Demand [POD]	3	0.00
Used Bilingual Dictionary	192	0.31
Used Magnification	40	0.06
Used Noise Buffers	67	0.11
Used Read Aloud	155	0.25
Used Scribe	64	0.10
Used Separate Setting	5016	8.11
Used Alternate Response Options	2	0.00
Provided Color Choices [CC]	27	0.04
Used Contrasting Color [CTC]	20	0.03
Used Reverse Contrast [RC]	9	0.01
Used Masking [MSK]	121	0.20
Used Text-to-Speech [TTS]	3563	5.76
Used Spanish Translation [ST]	488	0.79
Used Video Sign Language [VSL (ASL)]	14	0.02
Used Color Overlay	11	0.02
Amplification Device	32	0.05
Small Group Translation	59	0.10
Translator/Interpreter	8	0.01
Read Aloud in Spanish	24	0.04

Table 4-8 Summary Table of Manual Materials

Material	Configuration
DRC INSIGHT Portal Guide: Managing Users, Students, and Testing	<p>The DRC INSIGHT Portal Guide includes the following information:</p> <ul style="list-style-type: none"> • Managing user’s own DRC INSIGHT Portal account • Managing other DRC INSIGHT Portal users • Adding and editing students and student demographics, accommodations, and testing codes • Viewing, adding, and editing student test session information • Printing and managing student test tickets • Transferring students between schools and districts • Entering Not-Tested or Invalidation Codes • Unlocking or purging a student test • Managing test sessions • Monitoring testing status
Accessibility Guide	<p>The Accessibility Guide outlines the various accessibility options available to students taking the Wisconsin Forward Exam. Guidelines for using the various accessibility features are also included.</p>
Student/Administrator Tutorials	<p>The Student Tutorial includes 11 videos intended for students in grades 4–10 and 7 videos for students in grade 3. It is designed to show students the interface of the online testing system and familiarize them with the tools and features available. It is intended to accompany the Online Tools Training (OTT).</p> <p>The 2023 tutorials also include 11 videos for test administrators to familiarize them with the administrative features and functionality of the DRC INSIGHT Portal as well as the accessibility features of the Wisconsin Forward Exam.</p>
Item Samplers	<p>The item samplers can be used by both educators and students to gain familiarity with the types of items and their functionalities. The format appears as a “guided practice test” in the online, PDF, and Braille versions. Accommodations, universal tools, and supports are available in the test engine for the item samplers.</p> <p>Item samplers are accessible through the OTT menu page. The PDF versions include the answer key and scoring guide for each item.</p>
Online Tools Training (OTT)	<p>The OTT is a hands-on opportunity for students to become familiar with logging in, navigating, using tools, using accessibility features, reviewing, and submitting the test prior to signing in to an actual test. It is designed to be a second step after viewing the student tutorials.</p>

Material	Configuration
TAM (Test Administration Manual)	<p>The TAM is a document intended for test administrators (TAs) and proctors. It includes the following information:</p> <ul style="list-style-type: none"> • Key dates • TA and proctor responsibilities • Test times • Test security • Accessibility information • Procedures for before and during testing • Test ticket management • Test administration scripts <p>The TAM also includes a DAC/SAC Administration Management section, which contains the following:</p> <ul style="list-style-type: none"> • Roles and responsibilities • Test security • Resources and training materials • Test schedules • DRC INSIGHT Portal and DRC INSIGHT secure browser • Accessibility • Student transfers • Procedures to be completed prior to the close of the testing window • Data and reporting
Technology User Guide (TUG)	<p>The TUG is a document intended for Technology Coordinators. It is split into four volumes and includes detailed instructions on the installation and configuration of INSIGHT and the Central Office Services (COS) for all supported platforms.</p>
Interpretive Guide	<p>The Interpretive Guide is a document that includes the following information:</p> <ul style="list-style-type: none"> • Interpreting Wisconsin Forward Exam scores • Accessing Individual Student Reports (ISRs) and interactive summary reports via the DRC INSIGHT Portal
Technology Readiness Package	<p>The Technology Readiness Package is a suite of documents and tools for Technology Coordinators to prepare for the Wisconsin Forward Exam that includes the following:</p> <ul style="list-style-type: none"> • What is new and changing • Assessing online testing readiness • Capacity estimator • COS decision guide • Extended retries • Headset guidance • Keyboard settings • Installation of COS and INSIGHT • Installation of INSIGHT App • Evaluation and troubleshooting • System requirements • Technology overview presentation • Technology Coordinator Checklist • Technology FAQ

Material	Configuration
<p>Technical Report</p>	<p>The Technical Report is a manual that covers all grades and all psychometric details associated with administering the Wisconsin Forward Exam. The Technical Report provided by DRC presents thorough documentation to demonstrate the assessment validity. The document contains the following information:</p> <ul style="list-style-type: none"> • Description of the item pool used in the Wisconsin form-development process • Description of the test administration process and test security • Scoring of various types of items • Summary information of student performance (including means and standard deviations of scale scores, percentage of examinees within each performance level for each content area and grade level, and scale score distribution tables) • Item- and test-level analysis information for each content area and grade level, test scaling procedure, and student scoring process • Evidence of test validity
<p>Data Forensic Report</p>	<p>A separate Data Forensic Report includes analyses of the following:</p> <ul style="list-style-type: none"> • Evaluation of response changes • Evaluation of student response time to items

Part 5: Scoring

The purpose of Part 5 is to demonstrate adherence to AERA, APA, & NCME (2014) Standard 4.18, which states the following:

Procedures for scoring and, if relevant, scoring criteria, should be presented by the test developer with sufficient detail and clarity to maximize the accuracy of scoring. Instructions for using rating scales or for deriving scores obtained by coding, scaling, or classifying constructed responses should be clear. This is especially critical for extended response items such as performance tasks, portfolios, and essays. (p. 91)

Part 5 describes

- the scoring process of multiple-choice (MC) and multi-select (MS) items;
- the autoscoring process of technology-enhanced (TE), short-answer (SA), and evidence-based selected response (EBSR) items; and
- the scoring of text-dependent analysis (TDA) items, including
 - scoring rubrics,
 - artificial intelligence (AI) scoring process,
 - handscoring process,
 - scoring personnel selection,
 - anchor papers selection, and
 - TDA item scores distribution.

5.1 Multiple-Choice and Multi-Select Item Scoring Process

Responses to MC and MS items were captured during the online test administration. In the case of the Braille or paper-and-pencil form administrations, student responses to these items were transcribed into the online system by a TA. All MC and MS items had one and only one correct item response or a combination of responses.

5.2 Technology-Enhanced, Short-Answer, and Evidence-Based Selected Response Item Scoring Process

All TE, SA, and EBSR items were processed through DRC's autoscoring engine and scored according to the assigned scoring rules. DRC ensured that all rubrics and scoring rules were verified for accuracy before scoring any of these items. DRC established an adjudication process for these items and any gridded responses to verify that correct answers were identified. The quality control process for DRC's TE, SA, and EBSR item scoring included the following:

- A scoring rubric was created for each TE, SA, and EBSR item. It was similar to describing the one correct answer for dichotomously scored items (scored as either right or wrong). For ELA EBSR items worth 2 points, the rubric described in detail the type of response that could receive partial credit for 1 score point.
- The information from the scoring rubric was entered into the scoring system within the item banking system so that all information about the item resided in one place, along with the item image and other metadata. This scoring information designated specific information that varied by item type. For example, for a drag-and-drop item, the information included which objects are to be placed into which drop region to receive credit.
- The information was then verified by another autoscoring expert.
- After testing started, reports were generated that showed every response, how many students gave each response, and the score the scoring system provided.
- The scoring was then checked against the scoring rubric using two levels of verification.
- If any discrepancies were found, the scoring information was modified and verified again. Scoring was then rerun. This checking and modification process continued until no other issues were found.
- As a final check, a final report was run that showed all student responses, along with frequencies and received scores.

In the case of the Braille or paper-and-pencil form administrations, student responses to paper-and-pencil TE, SA, EBSR, or TE-equivalent items were transcribed (entered) into the online system by a TA.

5.3 Scoring of Text-Dependent Analysis Items

Sections 5.3 and 5.4 document the scoring processes used for TDA items. This documentation forms part of the validity evidence supporting the scoring process used for these items. Sections 5.3 and 5.4 describe the scoring rubrics, the scoring process, the selection of sample (anchor) papers used to train scoring personnel, the process of selecting personnel, and the distributions of scores for TDA items.

5.3.1 Description of Scoring Rubrics and Non-score Codes

In the 2023 test administration, the ELA forms in grades 3–8 contained one TDA item at each grade level. As stated in Part 3, Table 3-1, of this report, the TDA prompts are scored using a holistic scoring guideline on a 1–4 point scale. A weight of 2 is later applied to the item scores in order to compute the student total test raw scores and scale scores. That is, the TDA prompts will contribute up to 8 raw score points toward the student total test raw score.

The TDA responses were scored using an AI engine, and then validation scoring was performed by human scorers on approximately 10 percent of the AI scored responses. Table 5-1 presents the scoring rubric. In cases where student responses could not be scored, a non-score code was used. The non-score codes are presented in Table 5-2. All non-score codes were converted to a score of “0” in the derivation of student total test scores.

5.3.2 Artificial Intelligence Scoring

DRC partnered with Measurement Incorporated (MI) to score the TDA tasks. MI employed its essay scoring engine (called Project Essay Grade or PEG) to score all student responses. The AI model for scoring the Wisconsin student responses was built by first having DRC expert scorers score a representative sample of Wisconsin responses twice, independently, and resolve any scores that did not agree. While the engine only requires one score per response to build a model, the second score provides necessary information about how well two humans are able to agree on a score, which is then used as a benchmark for how well the engine's predictions should agree with the human scores. Approximately 3,000 student responses per grade from the Spring 2023 test administration were selected, handscored independently twice, and used in the AI model building.

The engine training sets consisting of scored sample responses and corresponding scores were delivered to the AI team at MI for model development. MI's linguistics experts, software developers, psychometricians, and human computer interaction specialists created task-specific algorithms that were then used to predict how humans would score these responses. The PEG team applied a standard stratified random sampling to all training sets, which is designed to produce two subsets of approximately 1,700 "training responses" and approximately 300 "validation responses," which approximated the score point distribution of the full training sets. The training responses were used to build the scoring model. The validation responses were used to verify the accuracy of AI scoring.

To build a scoring model, the engine analyzes the training set and calculates features that pertain to the content in question. The engine then sends the features to dozens of different algorithms that compete to see which ones can best associate the text features with the human assigned scores. These algorithms draw on many of the latest advances in the field of machine learning to generate both linear and nonlinear models. The strongest models are then automatically blended together to create a final model that retains the best elements from the various algorithms.

When the engine builds a model, it selects the model elements that maximize scoring accuracy for the data in question. Therefore, it is important to choose an agreement statistic on which the engine can optimize its models in such a way that the final model will exhibit reliable, accurate scoring. The inter-rater reliability of two human scorers is often measured via exact and adjacent agreement or the Pearson product-moment correlation coefficient (Pearson's r). It has been found that using quadratic weighted kappa, which has become the industry standard for AI scoring as the optimization and evaluation metric, leads to the most reliable and accurate scoring. Quadratic weighted kappa as a metric can detect changes in mean difference and variance between scorers and is, therefore, well suited for comparing the accuracy of AI scoring with that of human scoring and measuring the agreement of two independent human scorers.

MI's AI scoring software flagged student responses that could not be AI scored. The software has various triggers for identifying alert responses and responses in which it has low confidence. These responses lack proper development, lack enough content to be scored, are written in an unsupported language, contain inappropriate language, or represent a bad faith effort to complete the test (e.g., repeated text, off-topic text).

These responses that could not be scored by AI were routed to DRC for human scoring with a condition code indicating why the response could not be AI scored.

5.3.3 Handscoring Process

Human scoring of TDA items is referred to as “handscoring.” The scoring personnel who score TDA items are referred to as scorers. The scorers were trained using customized training materials, such as the anchor papers described in Section 5.3.5. Once qualified, scorers were required to maintain accuracy standards throughout the project. These requirements were assessed primarily through each scorer’s daily agreement rates with the AI scores (described below) and through targeted read-behinds with team leaders (described below). Reports were generated daily and monitored by the scoring director, team leaders, and project manager. Any scorers falling below the established quality standards for any item were retrained with the supervisors, who monitored scoring trends (such as difficulty with any particular score point). These scorers also received additional reviews and read-behinds. Failure to recalibrate resulted in dismissal from the scoring assignment. This process was in place throughout the entire handscoring window.

5.3.4 Handscoring System

Scoreboard, DRC’s handscoring system, was used to score TDA items as a validation method and to resolve cases where the AI engine returned a non-scorable condition code. Scoreboard presented images of rendered online responses to trained scorers who assigned scores for the TDA items. Images of each student’s responses were automatically routed to designated groups of scorers who were trained and qualified to score these items.

5.3.5 Anchor Papers and Training Papers

DRC’s project managers and scoring directors used the scoring guidelines, or rubrics, and Wisconsin student responses to select a representative sampling of student responses for each score point to create the necessary training materials for operational scoring. The sample reflects the various common response types produced for a specific item. The responses were then assembled into training/qualifying sets and shared with at least one other content expert for review. The scoring director for the specific grade took detailed notes, capturing scores and specific rationales for each score. Each grade and TDA item progressed in the same manner, using the same process. Once all sets were reviewed and scored, each grade-level scoring director had fully annotated a set of anchor papers, training papers, and qualifying papers. These anchor, training, and qualifying papers were then used to train a select group of scorers who scored the student responses that were used to train the AI engine in a process called model building. For this model-building activity, each student response was independently scored by two separate scorers. If there was any disagreement between the two readers, the scores were adjudicated to 100 percent agreement. Approximately 3,000 responses per grade were then delivered to the AI vendor to build the AI engine model. Once the model was built, the AI engine scored the remaining Wisconsin student responses. Upon completion of the AI scoring, a random sample consisting of approximately 10 percent of the student responses scored by the AI engine was sent to DRC for a human read. DRC then scored the 10 percent read-behind sample using a group of scorers trained to qualification standards to ensure consistency. The 10 percent read-behind with human scorers served as a validation check of the AI engine scoring data.

5.3.6 Scoring Personnel and Qualifications

AERA, APA, & NCME (2014) Standard 4.20 specifies the following:

The process for selecting, training, qualifying, and monitoring scorers should be specified by the test developer. The training materials, such as the scoring rubrics and examples of test takers' responses that illustrate the levels on the rubric score scale, and the procedures for training scorers should result in a degree of accuracy and agreement among scorers that allows the scores to be interpreted as originally intended by the test developer. Specifications should also describe processes for assessing scorer consistency and potential drift over time in raters' scoring. (p. 92)

DRC recruited, trained, and managed personnel to complete all the handscoring operations within the timelines of the contract. The recruitment process and requirements of the scorers, team leaders, and scoring supervisors are described in the following sections.

Scorers—The DRC scorer pool included many retired and current educators, engineers, editors, published authors, and individuals with advanced degrees. The minimum qualification for all scorers was a bachelor's degree. Scorers were required to participate in training and successfully pass a qualification round. Once qualified, scorers could start scoring, but throughout the scoring process, scorer performance was assessed by a scoring director, a team leader, and the project manager through read-behinds and reviews of inter-rater reliability statistics as described in Sections 5.3.8 and 5.4.

Team Leaders—Team leaders were selected on the basis of their ability to maintain a high degree of scoring accuracy and consistency, often across multiple content areas and grades. Team leaders were also required to possess good interpersonal and leadership skills in order to be effective when training and counseling scorers. Each team leader was responsible for a small team of scorers. In addition to performing read-behinds on scorers, team leaders also coached scorers when needs were identified through data review or otherwise by supervisory staff.

Scoring Directors—Scoring directors comprised the core group at DRC who directed and organized the scoring process and trained team leaders and scorers. Scoring directors had extensive experience as team leaders prior to their qualification and selection, and all had previous scoring director experience. Scoring directors were content area experts. They oversaw all team leaders and scorers.

5.3.7 Scorer Training

AERA, APA, & NCME (2014) Standard 6.9 specifies the following:

Those responsible for test scoring should establish and document quality control processes and criteria. Adequate training should be provided. The quality of scoring should be monitored and documented. Any systematic source of scoring errors should be documented and corrected. (p. 118)

Qualification was a critical task in the training process and the final determinant of scorer readiness. All scorers, including team leaders, were required to achieve a certain level of scoring accuracy in the qualifying round that

followed training. The standard to which they were held was the industry standard for TDA items: at least 70% exact agreement. Only those who were successfully validated were qualified as scorers to score tests.

5.3.8 Monitoring the Scoring Process

AERA, APA, & NCME (2014) Standard 6.8 states the following:

Those responsible for test scoring should establish scoring protocols. Test scoring that involves human judgment should include rubrics, procedures, and criteria for scoring. When scoring of complex responses is done by computer, the accuracy of the algorithm and processes should be documented. (p. 118)

The read-behind was used as a valuable monitoring technique. Each team leader was able to read a random selection of a scorer's scored responses. This reading could be targeted at the item and score-point level. The team leader scored the same item the scorer had scored. The scores (the scorer score and the team leader score) were compared, and if they agreed, the team leader was able to offer feedback, which enhanced the scorer's confidence and ability to score quickly and accurately. However, if a scorer strayed from the standards established in the training samples, the aberrant scoring was detected, and the team leader was able to offer guidance necessary to refocus the scorer's effort. Read-behinds by team leaders were more frequent for the scorers who had inconsistent scores, thus correcting any scoring variations. For aberrant or inconsistent scoring, DRC has the capability to wholesale drop scores and have the responses rescored if deemed necessary.

5.3.9 Final Scores

All TDA responses were sent to the AI engine for scoring. The AI scores were the final scores (i.e., scores of record). In all cases where the AI engine returned a non-scorable condition code, the student responses were reviewed and scored by humans and a resolution was reached. If a human scorer was able to assign a score for a response that the AI engine was not able to score, then a score from a human scorer became the score of record.

5.4 Inter-rater Reliability

A random 10 percent of the AI-scored responses were sent to human scorers for second reads to validate the AI scores. The statistics for the inter-rater reliability were calculated for all TDA items. To determine the reliability of scoring, the score distribution and percentage of agreement of the two readers were examined. In this section, the distribution of TDA item scores is presented. Additional inter-rater reliability measures, including intra-class correlation and weighted kappa statistics, are presented in Part 8 of the Technical Report.

5.4.1 Distribution of TDA Item Scores

Table 5-3 shows the score and non-scorable code distributions for TDA items. The presented scores, on a 1–4-point scale, are from the AI engine supplemented by non-scorable responses resolved by human scorers. Students who did not attempt the TDA item session of the ELA assessments are not included in this table.

Table 5-4 shows the score and non-scorable code distributions for TDA items for responses selected for the second read (handscoring). Table 5-5 shows the associated percentages of scores and non-scorable codes for TDA items for responses selected for the second read. In both tables, Scorer 1 is the AI engine and Scorer 2 is a human scorer. It should be noted that all non-scorable responses returned by the AI engine were reviewed by the scoring directors and assigned either a specific condition code or a score. The data in the non-scorable code columns in Tables 5-4 and 5-5 show the numbers and percentages of the non-scorable responses from the AI engine and detailed condition codes for these responses assigned by the human scorers (scoring directors).

As shown in Tables 5-4 and 5-5, there was a generally acceptable degree of agreement between the AI engine and the human scorers, with the differences being approximately 2 percent or less. The exceptions were differences at score point 1 in grades 3 and 8, score point 2 in grades 3 and 7, and score point 3 in grade 8. The differences between the AI engine and the human scorers were between 2.5 percent and over 3 percent in these cases. Greater differences between the AI engine and the human scorers were generally found at score points 1 and 2 compared to the differences at points 3 and 4 at all grade levels.

5.5 Summary

Taken together, the information presented in this part of the Technical Report summarizes the scoring procedures for different types of items and the steps taken by DRC to ensure accuracy in the TE item scoring, AI scoring, and handscoring processes. The score distribution statistics from the AI engine and the human scorer presented in Section 5.4 demonstrate that the items were scored reliably during the scoring process. These efforts by DRC follow multiple best practices of the testing industry and support AERA, APA, & NCME (2014) Standard 4.18, as presented in Part 5.

Table 5-1 TDA Item Scoring Guidelines, Grades 3–8

Score Value	Score Description	Scoring Rubrics
4	Demonstrates effective analysis of text and skillful writing	<ul style="list-style-type: none"> • Effective addressing of all parts of the task to demonstrate an in-depth understanding of the text(s) • Strong organizational structure and focus on the task with logically grouped and related ideas, including an effective introduction, development, and conclusion • Thorough analysis based on explicit and implicit meanings from the text(s) to support claims, opinions, and ideas • Substantial, accurate, and direct reference to the text(s) using an effective combination of details, examples, quotes, and/or facts • Substantial reference to the main ideas and relevant key details of the text(s) • Skillful use of transitions to link ideas within categories of textual and supporting information • Effective use of precise language and domain-specific vocabulary drawn from the text(s) • Few errors, if any, in sentence formation, grammar, usage, spelling, capitalization, and punctuation that do not interfere with meaning

Score Value	Score Description	Scoring Rubrics
3	Demonstrates adequate analysis of text and appropriate writing	<ul style="list-style-type: none"> • Adequate addressing of all parts of the task to demonstrate a sufficient understanding of the text(s) • Appropriate organizational structure and focus on the task with logically grouped and related ideas, including a clear introduction, development, and conclusion • Clear analysis based on explicit and implicit meanings from the text(s) to support claims, opinions, and ideas • Sufficient, accurate, and direct reference to the text(s) using an appropriate combination of details, examples, quotes, and/or facts • Sufficient reference to the main ideas and relevant key details of the text(s) • Appropriate use of transitions to link ideas within categories of textual and supporting information • Appropriate use of precise language and domain-specific vocabulary drawn from the text(s) • Some errors in sentence formation, grammar, usage, spelling, capitalization, and punctuation that seldom interfere with meaning
2	Demonstrates limited analysis of text and inconsistent writing	<ul style="list-style-type: none"> • Inconsistent addressing of some parts of the task to demonstrate a partial understanding of the text(s) • Weak organizational structure and focus on the task with ineffectively grouped ideas, including a weak introduction, development, and/or conclusion • Inconsistent analysis based on explicit and/or implicit meanings from the text(s) that ineffectively supports claims, opinions, and ideas • Limited and/or vague reference to the text(s) using some details, examples, quotes, and/or facts • Limited reference to the main ideas and relevant details of the text(s) • Limited use of transitions to link ideas within categories of textual and supporting information • Inconsistent use of precise language and domain-specific vocabulary drawn from the text(s) • Errors in sentence formation, grammar, usage, spelling, capitalization, and punctuation that may interfere with meaning
1	Demonstrates minimal analysis of text and inadequate writing	<ul style="list-style-type: none"> • Minimal addressing of part(s) of the task to demonstrate an inadequate understanding of the text(s) • Minimal evidence of an organizational structure and focus on the task with arbitrarily grouped ideas that may or may not include an introduction, development, and/or conclusion • Minimal analysis based on the text(s) that may or may not support claims, opinions, and ideas • Insufficient reference to the text(s) using few details, examples, quotes, and/or facts • Minimal reference to the main ideas and relevant details of the text(s) • Few, if any, transitions to link ideas • Little or no use of precise language or domain-specific vocabulary drawn from the text(s) • Many errors in sentence formation, grammar, usage, spelling, capitalization, and punctuation that often interfere with meaning

Table 5-2 TDA Item Non-scorable Codes, Grades 3–8

Non-scorable Code	Definition/Example/Notes
B—Blank	A response that is completely blank. This includes responses that <ul style="list-style-type: none"> are completely erased (so that words are unreadable). are completely crossed out (so that words are unreadable). are online and consist solely of “white space” (e.g., spaces, tabs, returns).
R—Refusal	A response that indicates a refusal to attempt the task. This includes the following examples: <ul style="list-style-type: none"> “I don’t care”; “I’m not taking this test”; “This is stupid”; “I won’t do it”; “you can’t make me answer this question” “I don’t know”; “IDK”; “we never learned this”; “X”; “NA” Unrelated song lyrics/rap lyrics/poetry (e.g., the lyrics to “Hotel California” in answer to a writing prompt asking whether backpacks should be allowed in class) Intentionally off-task response (e.g., a detailed description of what the student ate for breakfast that morning in answer to a question about Mozart’s childhood) This also includes responses that consist solely of scribbles, random keystrokes (“yyyyyyy”; “av:aeoiahvb”; “e, hrrttuuvv”), indecipherable writing/keystrokes (“swensts mengetstets arawnstets”) emoticons, stray marks, doodles, drawings, circles, underlines, a couple of random letters (not a word), or other evidence that no attempt was made to address the task.
N—Non-scorable	This category includes <ul style="list-style-type: none"> responses written entirely in a language other than English. responses that are completely illegible due to poor handwriting.* online or typed responses that are incoherent due to consisting of incomprehensible strings of words that are not clearly a Refusal or Off Topic (e.g., “best day school teacher inspired so I car”) responses too insufficient to be assessed by the criteria on the rubric. (for TDAs only) responses that address some part of the question but do not contain any logical/accurate/relevant reference to the passage(s) or any ideas contained in the passage(s). (for TDAs only) responses that consist solely, or almost solely, of text copied directly from the passage(s). *If a response is difficult to read, every effort is made to read the response. Multiple people, including a team leader and/or a scoring director, will attempt to decipher the response, and the original answer document will be reviewed if necessary. If, ultimately, only a portion of the response is legible, that verbiage will be scored on its own merits.
T—Off Topic	A response makes no reference to the item or (if applicable) the passage provided but does not seem to constitute an intentional refusal. If any part of the response relates to the item in any way, score the response.
C—Copied Item/Directions	A response consists of text copied from the item and/or test directions.

Note: Crossed out but legible/partially legible responses are scored according to the rubric based on whatever verbiage is legible.

Table 5-3 TDA Item Score Distribution

Grade	Item Number	Total Count	Item Score				Non-scorable Code				
			1	2	3	4	B	C	N	R	T
3	1	58451	35423	14381	1511	25	316	122	5779	812	82
4	1	58944	28829	15674	7379	79	172	120	6161	493	37
5	1	59333	32281	19362	4895	52	107	49	2224	307	56
6	1	59359	27322	21283	7197	796	159	170	1307	959	166
7	1	60345	20596	23092	10379	2266	140	240	3204	402	26
8	1	62166	15307	21006	16749	4791	182	325	2964	809	33

Table 5-4 TDA Item Score Distribution: AI Engine vs. Human Scorer

Grade	Scorer	Total Count	Score Count				Non-scorable Code Count				
			1	2	3	4	B	C	N	R	T
3	Scorer 1 (AI Engine)	15709	6869	1775	206	22			6837		
	Scorer 2 (Human)	15709	7436	1260	154	22	27	122	5790	816	82
4	Scorer 1 (AI Engine)	15001	5148	2036	946	22			6849		
	Scorer 2 (Human)	15001	5375	2058	615	104	16	120	6179	497	37
5	Scorer 1 (AI Engine)	9759	4442	2111	496	51			2659		
	Scorer 2 (Human)	9759	4649	1997	401	53	16	50	2228	309	56
6	Scorer 1 (AI Engine)	11368	5389	2412	801	120			2646		
	Scorer 2 (Human)	11368	5231	2524	817	150	23	172	1316	969	166
7	Scorer 1 (AI Engine)	15813	5057	4610	1862	370			3914		
	Scorer 2 (Human)	15813	4739	5234	1714	212	13	240	3226	409	26
8	Scorer 1 (AI Engine)	13383	4240	2732	1771	456			4184		
	Scorer 2 (Human)	13383	4582	3024	1333	260	12	328	2996	815	33

Note: TDA items are weighted x 2 in computation of student scores.

Table 5-5 TDA Item Percentage Score Distribution: AI Engine vs. Human Scorer

Grade	Scorer	Total Count	Score Percentage				Non-scorable Code Percentage				
			1	2	3	4	B	C	N	R	T
3	Scorer 1 (AI Engine)	15709	43.73	11.30	1.31	0.14			43.52		
	Scorer 2 (Human)	15709	47.34	8.02	0.98	0.14	0.17	0.78	36.86	5.19	0.52
4	Scorer 1 (AI Engine)	15001	34.32	13.57	6.31	0.15			45.66		
	Scorer 2 (Human)	15001	35.83	13.72	4.10	0.69	0.11	0.80	41.19	3.31	0.25
5	Scorer 1 (AI Engine)	9759	45.52	21.63	5.08	0.52			27.25		
	Scorer 2 (Human)	9759	47.64	20.46	4.11	0.54	0.16	0.51	22.83	3.17	0.57
6	Scorer 1 (AI Engine)	11368	47.40	21.22	7.05	1.06			23.28		
	Scorer 2 (Human)	11368	46.02	22.20	7.19	1.32	0.20	1.51	11.58	8.52	1.46
7	Scorer 1 (AI Engine)	15813	31.98	29.15	11.78	2.34			24.75		
	Scorer 2 (Human)	15813	29.97	33.10	10.84	1.34	0.08	1.52	20.40	2.59	0.16
8	Scorer 1 (AI Engine)	13383	31.68	20.41	13.23	3.41			31.26		
	Scorer 2 (Human)	13383	34.24	22.60	9.96	1.94	0.09	2.45	22.39	6.09	0.25

Note: TDA items are weighted x 2 in computation of student scores.

Part 6: Psychometric Analyses

This part of the Technical Report describes the analyses that were conducted with the ELA, Mathematics, Science, and Social Studies operational test data. These analyses included a classical item analysis and examination of the raw scores and an item response theory (IRT) analysis involving test calibration, scaling, and equating. These analyses were conducted using the calibration samples.

6.1 Overview of the Operational Test Data Analysis

This part of the Technical Report, the classical item statistics, including aggregate raw score statistics and individual item-level statistics are presented first. Next, the analyses involving test calibrating, equating, scaling, and student scoring that occurred for the Wisconsin Forward Exam after the 2023 test administration are described. The calibration samples are presented, followed by the data calibration results, including the model-data fit for the Wisconsin Forward Exam data. If the IRT models fit the empirical item response distributions for the population (i.e., Wisconsin students) for which generalizations are made, then the claim is strengthened that the scores are valid indicators of an underlying ability. The lowest obtainable scale score (LOSS) and highest obtainable scale score (HOSS) for the Wisconsin Forward tests are also presented.

Part 6 demonstrates adherence in the Wisconsin Forward Exam program data analysis to AERA, APA, & NCME (2014) Standards 1.8, 4.14, 5.2, 5.13, 5.15, and 7.2. Each standard will be explicated within the appropriate section of this part. Standard 7.2 provides general guidance that is relevant to this part:

The population for whom a test is intended and specifications for the test should be documented. If normative data are provided, the procedures used to gather the data should be explained; the norming population should be described in terms of relevant demographic variables; and the year(s) in which the data were collected should be reported. (p. 126)

6.2 Classical Item Analysis: Item Level Statistics

Three statistics are frequently used in item analysis: the proportion correct (p -value), the item-total correlation coefficient, and the omit rate for the item.

The p -value is an indication of the difficulty of an item. The p -value for an MC item or any item with a maximum score of 1 represents the proportion of students who answered the item correctly. If all students answered a given item correctly, its p -value would be 1.0. If only 30% of students answered the question correctly, the p -value would be 0.30. The lower the p -value is, the more difficult the item is. Item p -value is a good indication of difficulty, as it takes student performance into account and it makes comparing items in terms of a common statistic very simple. A test made up of items well distributed across the range of item difficulty levels is desirable because it supports the assessment of students at all ability levels.

The p -value for an item worth more than 1 point (e.g., EBSR item type) represents the mean proportion of possible raw score points that students actually obtained for the item. A p -value of 0.33 for an item with a

maximum item score greater than 1 would indicate that, on average, students obtained one-third of the possible points for the item. If a p -value were 0.75, this would indicate a much easier item where, on average, students scored 75% of the maximum possible points for the item. Therefore, the p -value indicates difficulty for such items as well, with lower p -values indicating more difficult items.

The item-total correlation indicates the extent to which individual test items provide reliable measurement of the construct being measured by the total test, and it is an index of the item's ability to discriminate between high-ability and low-ability students. For dichotomously scored items, the item-total correlations are computed as point-biserial correlations between the score on the item and the score on the remaining items in the test. For multi-point items, the item total correlations are computed as Pearson product-moment correlations between the score on the item and the score on the remaining items in the test.¹ The item-total correlation coefficients can range from -1.0 to +1.0. A large positive value (such as 0.40) indicates a strong relationship between a score on an individual item and the total score, with students who earn high scores on the total test tending to score higher on the item than students with low scores on the total test. A low positive value (such as 0.10) indicates a weak relationship between scores on the item and the total score, while a negative value indicates that students who do well on the total test tend to score lower on the item than students who do poorly on the total test.

For MC items, the point-biserial correlation between each distractor and the total score was also calculated. In most cases, items will have negative correlations for each distractor and the total score. However, a weak positive correlation for a distractor does not necessarily mean that the item is defective, provided that the distractor correlation is substantially smaller than the item-total correlation for the correct response. In some cases, it may simply mean that the particular distractor is attractive to moderate-ability students and unattractive to low-ability students.

The omit rate is also computed for each item, reflecting the percentage of students who did not respond to the item. A high omit rate can indicate an especially difficult item or, if located near the end of the test, it can indicate what is referred to as a "speeded" test, where students have insufficient time to respond to all items.

The examination of omit rates complies with AERA, APA, & NCME (2014) Standard 4.14. This standard is concerned with the speededness of a test:

For a test that has a time limit, test development research should examine the degree to which scores include a speed component and should evaluate the appropriateness of that component, given the domain the test is designed to measure. (p. 90)

¹ For both the point-biserial and the Pearson correlations, the studied item is excluded from the computation of the total score so as to not artificially inflate the correlation statistic. This effect would be most noticeable for items worth several points.

For the Spring 2023 Wisconsin Forward Exam, items were flagged for further investigation in the following situations:

- The p -value was less than 0.20. Such a p -value indicates a difficult item, where fewer than 20% of students obtained the correct answer.
- The item-total correlation was less than 0.15 for the correct answer. A low value may indicate that the item is not providing a high degree of discrimination between high-ability and low-ability students, and, in addition, it may be an indication that the correct answer is in question.
- A distractor had a positive correlation with the total test score.
- The omit rate was greater than 3%.

Flagging an item for investigation is just one aspect of a complete evaluation of an item, and flagged items are not necessarily defective. It is desirable to include a small number of items with very high p -values (easy items) or very low p -values (difficult items) in order to provide more reliable measurement at the extreme high and low levels of ability and to fully represent the range of difficulty for particular content standards. In this case, the flagging of p -values is a useful way of verifying that the number of extremely easy or difficult items is relatively small and consistent with the purposes of the test. Thus, flagged items do not necessarily indicate a challenge to test validity, because items have been found to be appropriate during item reviews.

Omit rates may reflect a number of different properties, and an item that is omitted by more than 3 percent of students (the Wisconsin Forward Exam flagging criterion) is not necessarily problematic. Omit rates are often higher for non-MC items than for MC items because students who are fairly certain they do not know the answer may be inclined to simply skip the item altogether rather than taking the time to form a response. Items with high omit rates are referred to content specialists for further review to ensure there is no unintended ambiguity in the items. If these flagged items are judged to be clear and provide a valid measurement of the intended knowledge, skill, or ability, then they are retained on the test.

Items flagged for a low item-total correlation or for a positive distractor-total test correlation are more troublesome because these statistics show the relationship of each option to the construct being measured. In determining whether these items should be retained or removed from scoring, it is important to consider the relative magnitude of the correlation between the correct response and the total score and between the distractor and the total score. In most cases, removing an item with a modest item-total correlation and negative correlations for all of the distractors will actually lower the reliability of the total test, so it is generally preferable to retain these items. The same is true of an item with a small positive correlation for one of the distractors and a much larger positive correlation for the correct response. However, an item that exhibits a low correlation for the correct response in combination with a positive correlation for one or more distractors is likely to degrade the accuracy of the measurement and lower the reliability of the test. Such items should be removed from scoring.

Overall, 53 operational items across all Wisconsin assessments were flagged on the Spring 2023 operational tests as meeting one or more of the investigational criteria bulleted above. More flagged items were found in Mathematics than in ELA, Science, or Social Studies.

Table 6-1 shows the number of scored items in the Spring 2023 Wisconsin Forward Exam operational tests flagged for these conditions by grade and content area. Because some items were flagged for more than one condition, the number of flags may be greater than the number of flagged items.

The flagged items were referred to DRC’s content specialists for further review to ensure that the items were unambiguous and the answer keys were correct. As part of this review, DRC’s content experts also evaluated each flagged item against the Wisconsin Forward Exam depth-of-knowledge criteria to ensure that the cognitive demands of the item reflected the skills and knowledge that the item was designed to measure. Tables 6-2, 6-3, and 6-4 provide more information about the flagged items.

6.2.1 Flagging for a Positive Distractor Correlation

In Tables 6-2 through 6-4, the distractor correlation coefficients are provided for items that were flagged because of positive distractor correlations. The distractor correlations tend to be small and are generally much smaller than the item-total correlations for the correct answer. The majority of items flagged for a positive correlation between a distractor and the total test had a correlation close to 0 for the distractor and an acceptable correlation for the correct answer. More items were flagged for positive distractor correlation in Mathematics than in other content areas. All flagged items were judged to be acceptable based on their content and other statistics and were retained in order to meet the Wisconsin Forward Exam test blueprints.

6.2.2 Flagging for the Item-Total Correlation

One item per grade was flagged for item-total test correlation <0.15 in ELA grade 3 and Mathematics grades 5 and 7. Two items were flagged in Mathematics grade 6. The item-total test correlations for the flagged items ranged from 0.11 to 0.14.

6.2.3 Flagging for p -Value

One item in Mathematics grade 7, three items in Mathematics grade 8, and one item in Science grade 4 were flagged for p -values <0.20 . The flagged items had p -values between 0.14 and 0.19. While these statistics indicate items that were difficult, the number of items flagged for difficulty was small. No operational items were flagged for difficulty in ELA or Social Studies.

6.2.4 Flagging for Omit Rate

No operational items on the Wisconsin Forward Exam were flagged for an omit rate of higher than 3%. Most of the items had omit rates of less than 1%.

6.2.5 Speededness

The degree to which a test is speeded can be evaluated by examining the percentage of students who fail to respond to the final items on a test or the last items in a timed section. One criterion of test speededness currently in use in the testing industry is a rule introduced by Educational Testing Services, which stipulates that at least 80% of test takers should be able to answer all of the items and all test takers should be able to answer at least 75% of the items (Swineford, 1956). However, a more stringent requirement is often applied, considering tests to be non-speeded only if at least 95% of examinees attempt the final item. As shown in Table 6-5, the Wisconsin Forward Exam satisfies this more stringent requirement, with over 99% of the examinees attempting the final item in each of the four content areas.

6.2.6 Supplemental Tables on Classical Item Analysis

Tables 6-6 through 6-22 present more comprehensive results from the classical item analysis for all the items retained in each grade and content area. In those tables, the item type (e.g., MC, EBSR), item p -value, item-total test correlation, and omit rates are presented. Tables 6-6 through 6-22 also show the item numbers, which can be used to understand the locations of test items as students actually encountered them on the test.

The numbers of flagged items across grade and content areas are summarized in Table 6-1. As indicated above, relatively few items were flagged. The item analysis indicated that the p -values of the items in the operational tests were well distributed throughout the range of difficulty levels, with reasonably high point-biserial correlations for most items. Detailed item analysis results including distractor statistics for MC items and score point distributions for non-multiple-choice items are included in Appendix H.

6.3 Test-Level Statistics

Test-level statistics, including test reliability, were computed for the Spring 2023 Wisconsin Forward Exam data for students with complete operational test data. These statistics are presented in Table 6-23. To facilitate interpretation of the test-level statistics, Table 6-23 provides the maximum possible score, the number of students, a measure of test difficulty, the standard deviation (SD) of raw scores, the skewness of the raw score distribution, the kurtosis, the minimum obtained score, the maximum obtained score, the reliability (Cronbach's alpha), and the standard error of measurement (SEM) for raw scores. These measurements are further explained below. Readers can refer to Tables 3-6 through 3-9 for a count of the number of items in the test and the number of score points corresponding to each test.

The mean raw score varies by grade and content area and, specifically, in the context of the maximum possible score points. In ELA, for example, the maximum possible raw score is 53 in grade 3 and 56 in grades 4 through 8. In Mathematics, the maximum possible raw score is 42 in grade 3 and 46 in grades 4 through 8. The maximum possible raw score is 40 in both Science grades and in all Social Studies grades.

Test difficulty is computed as the mean raw score divided by the maximum possible score points. Test difficulty ranges from 0 to 1.0. A larger test difficulty value indicates a mean raw score that is closer to the maximum possible score and, therefore, indicates an easier test. A smaller test difficulty value indicates a mean raw score

that is further from the maximum possible score and, therefore, indicates a more difficult test. Consider an example: A test difficulty statistic would be 0.90 if a mean score of 45 were obtained on a test with a maximum possible score of 50. This would be considered an easier test. On the other hand, test difficulty would be 0.50 if a mean raw score of 25 were obtained on the same test. This would then be considered a more difficult test. For example, the Mathematics grade 3 test mean raw score is 22.44 and the maximum possible score is 42, resulting in the test mean p -value of approximately 0.54. Evaluation of the mean p -values indicates that with the exception of Mathematics grade 3, Mathematics tests were more difficult for students than ELA, Science, and Social Studies tests.

Table 6-23 also shows the skewness and kurtosis statistics for each distribution of raw scores. Skewness and kurtosis describe the shape of a distribution. When a distribution is perfectly normal, skewness is zero. A negative skew has a long tail on the left side of the distribution because of the presence of some low scores, and, because the mean is sensitive to extreme scores, it indicates that most student scores are clustered on the high end of the scale. A positive skew indicates a distribution with some very high scores and a larger number of scores below the mean. Kurtosis describes a distribution in terms of its shape relative to a perfectly normal distribution. When a distribution is perfectly normal, kurtosis is zero. A negative kurtosis statistic indicates a distribution that is flatter than a perfectly normal curve, and a positive kurtosis statistic indicates a distribution that has more scores in the center of the score distribution (making it peaked) than a perfectly normal curve. Table 6-23 reveals that, in most cases, Wisconsin Forward Exam students are not normally distributed along the test scale in each grade and content area. Although this has implications for practitioners who wish to use the Wisconsin Forward Exam raw scores in statistical analyses (since normality of the data cannot be assumed), from a criterion-referenced testing standpoint, it indicates that students on the whole are mastering the Wisconsin Standards for ELA, Science, and Social Studies. The Mathematics assessments in grades 4 through 8 and, to some degree, the Social Studies grade 10 assessment tended to be more difficult, however, showing most of the scores clustered below the mean (as indicated by positively skewed score distributions).

In addition, Table 6-23 shows that the minimum obtained score in ELA grade 4, all Mathematics and Science grades, and Social Studies grades 8 and 10 was zero, meaning that at least one student failed all items for each of those tests. In ELA grades 3 and 5 through 8 in and Social Studies grade 4, the minimum test score was 1, indicating that no student failed to respond correctly to all items. With the exception of ELA grades 3 through 6, the maximum obtained scores were equal to the maximum number of points possible on the test in all grades, meaning that at least one student obtained the full score for all items on each of those tests. For example, as displayed in Table 6-23, in Mathematics grade 3, there was at least one student who failed all items and at least one student who obtained the maximum raw score of 42.

A reliable test is one with high reliability as represented by statistics such as Cronbach's alpha and a low SEM. When interpreting reliability statistics, readers should note that test length (number of items and score points) is one of the important factors that influences reliability statistics and SEM. These concepts are described further in Part 8. For present purposes, the reader should note that measurement error is associated with every test score. A student's true score is the hypothetical average score that would result if the test could be administered repeatedly without the effects of practice or fatigue. Obtained scores should not be regarded as absolute but as one point within a range that, with a certain degree of probability, includes a student's true score.

The test-level statistics for each content area are summarized and discussed below using the measurements described above.

English Language Arts

- Test difficulty ranged from 0.54 to 0.56.
- Reliability coefficient Alpha was relatively high in every grade (0.88 to 0.89).
- SEM ranged from 3.34 to 3.55.

Mathematics

- Test difficulty ranged from 0.44 to 0.54, with generally lower difficulty in lower grades and higher difficulty in higher grades.
- Reliability coefficient Alpha was high in every grade (0.91 to 0.93).
- SEM ranged from 2.67 to 2.92.

Science

- Test difficulty was 0.53 in grade 4 and 0.52 in grade 8.
- Reliability coefficient Alpha was 0.91 in grade 4 and 0.89 in grade 8.
- SEM was 2.68 and 2.78 for grades 4 and 8, respectively.

Social Studies

- Test difficulty ranged from 0.52 to 0.58.
- Reliability coefficient Alpha was 0.91 in grades 4 and 8 and 0.90 in grade 10.
- SEM ranged from 2.69 to 2.79.

6.4 Item Response Theory Methodology

This section of the report outlines the item response theory (IRT) methodology including item calibration, test equating and test scaling, as well as the methodology of computation of the scale scores based on Wisconsin Forward Exam test data. Readers should note that calibration, equating, and scoring using IRT are mathematically complex and computationally intensive processes. A full understanding of these topics requires a background in psychometrics. However, in order to make these processes more accessible and transparent to a wider range of audiences, a brief, nontechnical explanation of the IRT process and how scale scores are derived from student raw responses is provided. Additional references are also suggested to interested readers.

6.4.1 Item Calibration

This section of the report outlines the calibration procedures and results for the Spring 2023 Wisconsin Forward Exam. Student responses on the Wisconsin Forward Exam are inputted into complex mathematical algorithms

designed to model the relationship between a student’s ability in a content area and a test item. The group of algorithms is collectively known as item response theory (IRT). Wisconsin Forward Exam scores are established through the processes of calibration, scaling, and item-pattern scoring.

Calibration is the mathematical process of estimating characteristics of individual items. These characteristics are termed “item parameters.” Section 6.4.1 serves to explain this process, beginning with a description of the calibration methods that were applied to the Spring 2022 Wisconsin Forward Exam, followed by a presentation of a calibration sample, and a discussion of the calibration models and the software used. The results of the calibration process, using model-to-data fit statistics, and the outcomes of test scaling are also discussed in this section.

6.4.1.1 Calibration Models

The three-parameter logistic (3PL) model and the two-parameter partial credit (2PPC) IRT model (Bock & Aitkin, 1981; Thissen, 1982) were used to estimate parameters for multiple choice (MC) items and constructed-response (CR) items, respectively. All non MC items, including technology-enhanced (TE) items, evidence-based selected response (EBSR) items, short-answer (SA) items, and text-dependent analysis (TDA) items, were treated as CR items in calibrations. Item parameters for items contained in all Wisconsin assessments were estimated using a marginal maximum-likelihood procedure.

Under the 3PL model, the probability that a student with a trait or scale score θ will respond correctly to MC item j is

$$P_j(\theta) = c_j + (1 - c_j) / [1 + \exp(-1.7a_j(\theta - b_j))].$$

In the equation, a_j is the item discrimination, b_j is the item difficulty, and c_j is the probability of a correct response by a very low-ability student. Under the 2PPC model, the probability that a student with a trait or scale score θ will respond in category k to partial-credit item j is

$$P_{jk}(\theta) = \exp(z_{jk}) / \sum_{i=1}^{m_j} \exp(z_{ji}),$$

Where $z_{jk} = (k - 1)f_j - \sum_{i=0}^{k-1} g_{ji}$ and $g_{j0} = 0$ for all j .

The summary output of the 3PL and 2PPC models is in two different metrics. The discrimination and location parameters for the MC items are in the traditional 3PL metric and are labeled a and b , respectively. In the 2PPC model, f (alpha) and g (gamma) are analogous to a and b , where alpha is the discrimination parameter and gamma over alpha (g/f) is the location in which adjacent trace lines cross on the ability scale. Because of the different metrics used, the 3PL parameters a and b are not directly comparable to the 2PPC parameters g and f ; however, they can be converted to a common metric. The two metrics are related by $a = f/1.7$ and $b = g/f$ (Burket, 2002). As a result of this procedure, the MC and CR items are placed on the same scale. Note that for

the 2PPC model, there are $m_j - 1$ (where m_j is a score level j) independent g 's and one f , for a total of m_j independent parameters estimated for each item, while there is one a and one b per item in the 3PL model.

Using the 3PL/2PPC models for estimation of ELA, Mathematics, Science, and Social Studies, item parameters were consistent with the past methodology (except for the 2014–15 administration for ELA and Mathematics) implemented for Wisconsin assessments. Item parameters estimated after the 2023 test administration were used to score the responses of Wisconsin students who took these tests.

6.4.1.2 Calibration Sample

The calibration of the Wisconsin Forward Exam occurred after the Spring 2023 test administration and was based on the student data acquired during the entire testing window. This section provides information on the comparability of the calibration sample to the census data in terms of demographic characteristics in adherence to Standard 1.8 of the AERA, APA, & NCME (2014) *Standards*:

The composition of any sample of test takers from which validity evidence is obtained should be described in as much detail as is practical and permissible, including major relevant socio-demographic and developmental characteristics. (p. 25)

The calibration samples consisted of the student data acquired during the entire testing window and included records with complete operational test data. That is, to be included in the sample, a student had to respond to all test questions and not have any missing responses. Students from public, choice, and private schools were included in the calibration data. Students who had non-tested codes indicated in the data were excluded from the calibration samples. The characteristics of the calibration samples are presented in Tables 6-24 through 6-27 for ELA, Mathematics, Science, and Social Studies, respectively. The grade 3 through 8 calibration samples consisted of more than 99 percent of students who later received valid test scores in these grades across all content areas. The grade 10 calibration sample consisted of over 98% of students who later received valid test scores. In addition, calibration samples for grades 3 through 8 consisted of approximately 94% to 95% of all students scheduled to participate in the assessment. The grade 10 calibration sample for Social Studies included approximately 87% of all students scheduled to participate in the assessment. The differences in the percentages of students who were included in the calibration samples and students who received valid scale scores compared to students who were scheduled to test are reported in Tables 6-24 through 6-27 in the top row (“All Students”) for each grade.

When the calibration sample characteristics were compared with the characteristics of students who later received valid test scores, it was found that the differences across subgroups were less than a quarter of a percent for any subgroup in grades 3 through 8 in all content areas. These differences were about half a percent or less for Social Studies grade 10.

When the calibration sample characteristics were compared to the characteristics of all students scheduled to test, it was found that, with few exceptions, the differences across subgroups were less than one percent across all content areas in grades 3 through 8. The exceptions were differences in the percentages of students with disabilities included in the calibration samples compared to the percentages of students with disabilities

scheduled to test. These differences were just over 1% in grades 7 and 8 across all content areas. Larger differences between the calibration sample and the population of students scheduled to test were found for grade 10 Social Studies. White students were overrepresented by more than 3% in the calibration sample compared to the population of students scheduled to test. African American and economically disadvantaged students were underrepresented by over 2% in the calibration sample, and students with disabilities were underrepresented by close to 2% in the calibration sample compared to the population of students scheduled to test. The difference between the calibration sample and the population of students scheduled to test was less than half a percent for all other subgroups for Social Studies grade 10. No adjustment to the calibration sample was made for grade 10.

6.4.1.3 Calibration Procedure

The calibrations were conducted separately for each grade level and content area using the marginal maximum-likelihood procedures implemented with the expected maximum algorithm (Bock & Aitkin, 1981; Thissen, 1982). In a process of item calibration, the number of estimation cycles was set to 99 with the convergence criterion of 0.001 for all content areas. The maximum value of a -parameter was set to 5.0, and the range for b -parameter was set between -7.5 and 7.5. For all items, the estimated a - and b -parameters were within the prescribed parameter ranges. The c -parameters for anchor items were fixed to their Spring 2022 values. It should be noted that there was a small number of items with the default value for the c -parameter on all tests. When the algorithm, which is implemented to calibrate the items, encounters difficulty estimating the c -parameter, it assigns a default c -parameter value of 0.20.

6.4.1.4 Calibration Software

Calibration of the Wisconsin Forward Exam data was performed using PARDUX software (Burket, 2002; Shu 2020). PARDUX is designed to produce a single scale by jointly analyzing data resulting from students' responses to both MC items and CR items for assessments that include both item types. In PARDUX, items are calibrated based on IRT, using the 3PL model (Lord & Novick, 1968) for MC items and the 2PPC model (Yen, 1993) for CR items.

PARSCALE, MULTILOG, and BIGSTEPS are among the most widely known and used IRT programs. Extensive simulation studies and comparisons between PARDUX and MULTILOG (Thissen, 1990), PARSCALE (Muraki & Bock, 1991), and BIGSTEPS (Wright & Linacre, 1992) have shown that PARDUX provides precise parameter and ability estimates and performs as well or more efficiently than these programs (Fitzpatrick, 1991; Fitzpatrick and Julian, 1996). Extensive research with simulation data has also shown that the IRT procedures used for calibration and scaling of Wisconsin assessments produce accurate vertical scaling (Yen & Burket, 1997). PARDUX software has been regularly updated by adding advanced estimation algorithms and other features, keeping it up-to-date to face the challenges of large-scale assessments (Shu, 2020).

6.4.1.5 Calibration Results

This section describes the calibration results in terms of the estimation of item parameters and model-to-data fit for all content areas and grades.

IRT Item Parameters

During calibration, items may not converge, meaning the characteristics of the items will not be determined. When this occurs, items may be suppressed from student scoring and future assessments. In Spring 2023, no non-convergence issues occurred for any item on the operational tests.

IRT Item Fit

The calibration process produces ability and item parameter estimates that can be used to predict student response patterns to each item. For example, based on the item parameter estimates for item difficulty and item discrimination, low-ability students are expected to be less likely to answer a difficult and highly discriminating item correctly than higher-ability students. After parameters are produced, the predicted scoring patterns can be compared to the observed scoring patterns in what are referred to as item-to-model fit comparisons. Where there is little difference between the predicted scoring patterns and the observed scoring patterns, the model can be said to “fit” the data.

A procedure developed by Yen (1981) was used to assess model-to-data fit for all test items. In this procedure, students are rank ordered on the basis of their $\hat{\theta}$ values and sorted into ten cells, with 10 percent of the sample in each cell. Each item j in each decile i has a response from N_{ij} examinees. The fitted IRT models are used to calculate an expected proportion E_{ijk} of examinees who respond to item j in category k . The observed proportion O_{ijk} is also tabulated for each decile. The fit index for item i is

$$Q_{1j} = \sum_{i=1}^{10} \sum_{k=1}^{m_j} \frac{N_{ij}(O_{ijk} - E_{ijk})^2}{E_{ijk}}.$$

Q_{1j} should be approximately chi-square distributed with degrees of freedom (DF) equal to the number of “independent” cells, $10(m_j - 1)$, minus the number of estimated parameters. For the 3PL model, $m_j = 2$, so $DF = 10(2 - 1) - 3 = 7$. For the 2PPC model, $DF = 10(m_j - 1) - m_j = 9m_j - 10$.

DRC evaluated item-to-model fit in a two-step process. First, item-to-model fit information was obtained for each item using a Z -statistic. The Z -statistic is an index of the degree to which obtained proportions of students with each item score match the proportions predicted by the estimated student ability and item parameters. When the difference between the obtained proportions of students with each item score and the proportions predicted by the estimated student ability and item parameters reached a certain threshold, the item was flagged for “misfit.”

The Z -statistic is a transformation of the chi-square (Q_1) statistic that takes into account differing numbers of score levels as well as sample size using the equation

$$Z_j = \frac{(Q_{1j} - DF_j)}{\sqrt{2DF_j}},$$

where Q_{1j} is the item chi-square statistic, j is an item, and DF is the degrees of freedom for a given item j .

Because the value of Z increases as the sample size increases, the critical values for Z were established using the following equation (Yen & Candell, 1991):

$$Z_{crit,j} = \frac{4N_j}{1500},$$

Where $Z_{crit,j}$ is the critical value of Z for item j and N_j is the number of students who responded to item j . These values and the associated chi-squares (Q_1) are computed for ten intervals corresponding to deciles of the ability distribution (Yen, 1984).

Table 6-28 presents items that were flagged for less-than-optimal fit when the obtained Z -statistic exceeded the critical Z -statistic value. This table specifies the content area, grade level, item number in the calibration, calibration model (3PL for MC items and 2PPC for non-MC items), N size (i.e., the number of students who took this item), Z , and critical Z as described previously. Sixteen items were flagged for poor fit for ELA, thirteen items were flagged for Mathematics, five items were flagged for Science, and four items were flagged for Social Studies. Most of the flagged items were non-MC items. For example, item #30 for ELA grade 3 was flagged because the observed Z of 294.09 was larger than the critical Z value of 155.65 based on a sample size of 58,368. For many of the flagged items, the observed Z and the critical Z were not very far apart, indicating small or moderate misfit; however, the misfit was larger for some items (e.g., item #39 in ELA grade 5, item #15 in ELA grade 6, item #23 in ELA grade 7, and item #34 in ELA grade 8).

In order to evaluate item-to-model fit further, DRC inspected the observed-to-predicted item characteristic curve (ICC) for each flagged item. These ICCs simultaneously plot the characteristics of an item (e.g., item difficulty, item discrimination, level of guessing) using IRT model predictions and the observed student responses. The ICCs show exactly where along the ability continuum the misfit occurs and the extent of the misfit.

All cases of MC items flagged for misfit had empirical (observed) information that differed from the model in the lower-ability range, where there are fewer students to provide information at the tail end of the distribution. Similarly, for CR items, there were, in general, fewer students at the lower score levels, which provides less information at the tail ends of the student distribution. Items that only show misfit at the tail ends of the distribution provide stable information about the majority of students—those in the middle range of the distribution. However, if the misfit happens around the middle of the ability range, where there are many students, this may be a concern and may lead to the item being dropped from the item pool.

In a large-scale assessment such as the Wisconsin Forward Exam, with 17 combinations of grades and content areas, it is expected that some items will be flagged for misfit. As noted, the difference between the obtained Z -statistic and the critical Z -statistic was often small or moderate. Items flagged for misfit were reported to the DRC Test Development team for additional review. Such items are flagged in the Wisconsin Forward Exam

item bank and are avoided during the form selection process unless there is a compelling reason that they should be included, such as meeting the test blueprint.

6.4.2 Test Equating

Test equating is the statistical process of placing scores from two or more parallel assessments onto a common scale, resulting in direct comparability of scores from two different test forms. A common-item design was used to link the assessments from 2023 to the established ELA, Mathematics, Science, and Social Studies scales for the Wisconsin Forward Exam. Sets of items that were administered to Wisconsin students in the Spring 2022 operational test administrations and that were repeated in the Spring 2023 assessments served as the anchor sets in each assessment. The anchor sets constituted at least one-third of the Spring 2023 assessments and were representative of the Spring 2023 test content. After the item calibration, item parameters were linked to the Wisconsin Forward Exam scales using the Stocking & Lord (1983) equating procedure.

Standard 5.13 of the AERA, APA, & NCME (2014) *Standards* states the following:

When claims of form-to-form score equivalence are based on equating procedures, detailed technical information should be provided on the method by which equating functions were established and on the accuracy of the equating functions. (p. 105)

The Stocking & Lord procedure minimizes the mean squared difference between the two test characteristic curves (TCCs), one based on estimates from the previous calibration and the other based on transformed estimates from the current calibration. Let $\hat{\Psi}_j$ be the TCC based on estimates from a previous calibration and $\hat{\Psi}_j^*$ be the TCC based on transformed estimates from the current calibration:

$$\hat{\Psi}_j = \hat{\Psi}(\theta_j) = \sum_{i=1}^n P_i(\theta_j; a_i, b_i, c_i)$$
$$\hat{\Psi}_j^* = \hat{\Psi}(\theta_j) = \sum_{i=1}^n P_i\left(\theta_j; \frac{a_i}{A}, Ab_i + B, c_i\right)$$

The TCC method determines the equating constants (A and B) by minimizing the following quadratic loss function (F):

$$F = \frac{1}{N} \sum_{a=1}^N \left(\hat{\Psi}_j - \hat{\Psi}_j^* \right)^2.$$

The Stocking & Lord equating procedure is commonly used in large-scale assessments. The standard error of the equating (SEE) is difficult and cumbersome to estimate for IRT equating procedures like the Stocking & Lord procedure (Kolen & Brennan, 1995; Michaelides & Haertel, 2004). The estimation of the SEE is beyond the scope of this report.

6.4.2.1 Evaluation of Anchor Items

AERA, APA, & NCME (2014) Standard 5.15 requires information about the anchors, stating the following:

In equating studies that employ an anchor test design, the characteristics of the anchor test and its similarity to the forms being equated should be presented, including both content specifications and empirically determined relationships among test scores. If anchor items are used in the equating study, the representativeness and psychometric characteristics of the anchor items should be presented. (p. 105)

Two statistical methods were used to evaluate anchor items: (1) iterative linking (Candell & Drasgow, 1988) using Stocking & Lord's (1983) TCC method and (2) differences between the item-ability regression curves.

Test Characteristic Curve Method

The Stocking & Lord (1983) procedure, also called the TCC method, for which the mathematical equation was provided in a previous section of this document, minimizes the mean squared difference between the two TCCs, one based on estimates from the previous calibration and the other based on transformed estimates from the current calibration.

Differential item functioning was evaluated by examining previous (input) and transformed (estimated) item parameters. Items with an absolute difference of parameters greater than two times the root mean square deviation were flagged for review. These differences were monitored by plotting input and estimated item parameters.

Item Response Theory Item-Ability Regression Curves

Differences between the item-ability regression curves of the anchor items in the Spring 2023 Wisconsin Forward Exam administration were also compared to previous calibrations from Spring 2022. The differences between the item curves were evaluated using the following statistics:

- UnWtd Mean = Average signed difference in estimated probability
- UnWtd Mean Abs = Average absolute (unsigned) difference in estimated probability
- UnWtd RMSD = Root mean squared difference
- Wtd Mean = Weighted average signed difference in estimated probability
- Wtd Mean Abs = Weighted average absolute (unsigned) difference in estimated probability
- Wtd RMSD = Weighted root mean squared difference

Both unweighted and weighted versions of these statistics were calculated. Unweighted differences give equal weight to differences across the ability spectrum. Weighted differences assign weights according to the number of test takers that are impacted (that is, the frequency distribution of estimated student abilities during the calibration).

For the six statistics listed above, differences greater than ± 0.10 are considered large and differences between ± 0.07 and ± 0.10 are considered moderate.

Additionally, the maximum absolute difference (Max Abs) was identified. For Max Abs, large differences are those greater than ± 0.15 and moderate differences are all differences between ± 0.125 and ± 0.15 .

6.4.2.2 Removal of Anchor Items

One of the key requirements of anchor items in deriving valid and reliable linking results is that the anchor items form a miniature of the test in terms of content coverage, or test blueprint. While dropping a flagged anchor item based solely on statistical criteria has its simplicity, this option may change the content coverage and invalidate results. Before an anchor item is dropped from an anchor set, the item characteristics, adequacy of the content coverage, and impact on the size of the anchor set must be evaluated.

An item may be removed from the anchor set only if it adversely affects the quality of scaling, not the desirability of the results. Therefore, DRC does not consider how the removal of an item affects the overall mean scale score or the impact data (i.e., percentage of students in each achievement level) when recommending items for removal.

Items removed from the anchor set are still scored as part of the whole test. DRC recommends that the anchor items be considered for exclusion from the Wisconsin Forward Exam equating sets under the following conditions:

1. An item may be a candidate for removal if it is flagged for moderate or large differences on at least four of the seven statistics (listed in Section 6.4.2.1) considered when examining the differences between the IRT item-ability regression curves.
2. Removal of the item will only be considered after alternative explanations have been considered that may explain shifts in performance. For example, performance on the anchor item may improve because of a statewide initiative emphasizing instruction on a particular set of skills. In this case, improved performance on the item represents true growth in that area. Removing the anchor item may artificially lower test scores.
3. Removal of the item may not significantly alter the content distribution of the anchor set. The distribution of the anchor items across the content standards should remain within 10 percent of the Wisconsin Forward Exam test blueprint.
4. The number of remaining items will remain at an acceptable level of anchor set reliability. Operationally, this means the anchor set will still be representative of the total test blueprint and the anchor set may not be less than 20 percent of the total test length.

Flagged items are reviewed by DRC test development experts to verify that no changes to item content or format occurred between the administration in which the anchor items were used and the current administration. In addition, for the flagged non-MC anchor items, verification that no changes to scoring rubrics occurred between the two administrations is performed.

6.4.2.3 Evaluation of Equating Results

Table 6-29 provides equating results for the TCC method for all content areas. This table summarizes the following information for each grade and content area: number of anchors, number of iterations, quadratic loss function (F), correlation between the a -parameter input and estimates, correlation between the b -parameter input and estimates, number of a - and b -parameter outliers as indicated by the root mean square deviation method, and equating constants (A and B).

The overall alignment of the anchor TCCs was very good for all grades in all content areas. Figures 6-1 through 6-4 show the TCC alignment of the anchor set before and after equating for all grades of ELA, Mathematics, Science, and Social Studies. In these figures, the input anchor set TCC (before equating) is indicated by the dashed red line and the new anchor estimate TCC (after equating) is indicated by the solid blue line. The correlations between the a -parameter input and estimates were at least 0.94 or higher, and the correlations between the b -parameter input and estimates were 0.96 or higher for all grades and content areas. One anchor item was flagged as an a -parameter outlier in each of the following grades: ELA grades 3, 7 and 8; Mathematics grades 5 and 6; Science grade 4; and Social Studies grades 4 and 8. Two anchor items were flagged as a -parameter outliers in ELA grade 4. One anchor item was flagged as a b -parameter outlier in each of the following grades: ELA grades 4 through 7, Mathematics grades 4 through 8, and Social Studies grades 8 and 10. Overall, the number of anchor items flagged using the TCC method was small.

No anchor items were flagged for potential removal from the anchor sets using the IRT item-ability regression method. Consequently, no anchor items were removed from equating.

6.4.3 Test Scaling and Scale Evaluation

The purpose of scaling a test is to enhance its validity by increasing the comparability of test takers' scores. This section explicates the way in which the Wisconsin Forward Exam scales are produced to comply with Standard 5.2 of the AERA, APA, & NCME (2014) *Standards*, which states the following:

The procedures for constructing scales used for reporting scores and the rationale for these procedures should be described clearly. (p. 102)

The Wisconsin Forward Exam scales were established for ELA and Mathematics after the Spring 2016 test administration. New scales were established for Science assessments after the Spring 2019 test administration. New scales were established for Social Studies assessments after the Spring 2022 test administration. In this section, the results of the test scaling of the Wisconsin Forward Exam are described and evaluated.

Following the test equating, the equated item parameter estimates in the theta metric were transformed into the scale score metric for the purpose of the evaluation of the scale properties. The scale evaluation included

- evaluation of the TCCs,
- evaluation of the standard error (SE) curves, and
- examination of the growth at quartiles.

The scaling constants, $M1$ and $M2$, used to transform equated item parameters and ability estimates in the theta metric into the scale score metric are the same as the scaling constants used in the Spring 2016 scale development for ELA and Mathematics. The $M1$ and $M2$ scaling constants for Science were established after the Spring 2019 test administration. The $M1$ and $M2$ scaling constants for Social Studies were established after the Spring 2022 test administration. These scaling constants are presented in Table 6-30.

ELA Scale

Test Characteristic Curves—Figure 6-5 shows the TCCs for ELA tests. As shown in Figure 6-5, the ELA TCCs are ordinal at most ability levels, indicating that the difficulty of these assessments increases as the grade level increases. The exception is the grade 7 and grade 8 TCCs crossing at the lower ability level, indicating that the grade 7 assessments may be more difficult for lower ability students than the grade 8 assessments.

It should be noted that while TCC ordinality is a desirable property of a vertical scale, the lack of it does not necessarily affect student scores or grade-to-grade growth interpretation. As demonstrated by the pattern of scale scores at quartiles (see the “Growth at Quartiles” paragraph below) for grades 3–8, student ability on ELA assessments increases as grade level increases at all grade levels, indicating grade-to-grade growth.

Standard Error Curves—The SE curves for ELA presented in Figure 6-6 are generally U-shaped, indicating smaller errors around ability estimates that are roughly in the middle of the scale score distribution. The SE is expected to be higher at the top and bottom ends of the ability scale, where fewer items measuring very high- and very low-achieving students are found. Overall, the SEs around the scale score were found to be reasonable for ELA assessments (for more details, see Section 8.1.1 of this report).

Scale Scores and Growth at Quartiles—The estimated scale scores for the ELA calibration sample at the 25th, 50th, and 75th percentiles for all grade levels are presented in Figure 6-7. It can be observed that the scale scores increase as the percentile increases within each grade. Consistent with the properties of a vertical scale, the scale scores also increase at the same percentile across grade levels, indicating growth, albeit nonuniform, on the ELA ability scale as students move from one grade to the next.

Mathematics Scale

Test Characteristic Curves—Figure 6-8 shows the TCCs for Mathematics assessments, which are on a vertical scale. As observed in Figure 6-8, the TCCs for Mathematics are ordinal, indicating that the difficulty of the assessments increases as the grade level increases. The exception is the grade 5 and grade 6 TCCs crossing at the lower ability level, indicating that the grade 6 assessments may be more difficult for lower ability students than the grade 5 assessments.

Standard Error Curves—The SE curves for Mathematics presented in Figure 6-9 are U shaped (as expected), indicating smaller errors around ability estimates that are roughly in the middle of the scale score distribution. The SE is expected to be higher at the top and bottom ends of the ability scale, where fewer items measuring these students are found. Overall, the SEs around the scale score were found to be reasonable for Mathematics assessments (for more details, see Section 8.1.1 of this report).

Scale Scores and Growth at Quartiles—The estimated scale scores for the calibration sample at the 25th, 50th, and 75th percentiles for all grade levels are presented in Figure 6-10. It can be observed that the scale scores increase as the percentile increases within each grade level. The scale scores also increase at the same percentile across grade levels, indicating growth on the Mathematics ability scale as students move from one grade to the next.

Science Scale

Test Characteristic Curves—Although the Science assessments are not vertically scaled, the TCCs for grades 4 and 8 are presented together in Figure 6-11 for comparison purposes. The TCCs are S-shaped, indicating increasing probability of a higher test score as a student’s ability increases. The grade 4 and grade 8 TCCs are parallel to each other, indicating similar overall test discrimination of the two assessments.

Standard Error Curves—Figure 6-12 shows the SE curves for Science grades 4 and 8. The SE curves are U-shaped, indicating smaller errors around ability estimates that are approximately in the middle of the scale score distribution. The SE is expected to be higher at the top and bottom ends of the ability scale, where fewer items measuring these students are found. Overall, the SEs around the scale score were found to be reasonable for Science assessments (for more details, see Section 8.1.1 of this report).

Scale Scores at Quartiles—The estimated scale scores for the Science calibration sample at the 25th, 50th, and 75th percentiles for both grade levels are presented in Figure 6-13. The data pattern presented in this figure indicates that the scale scores increase as the percentile increases within each grade level. Because the Science assessments are not on a vertical scale, it is not appropriate to compare scale scores between grades.

Social Studies Scale

Test Characteristic Curves—Although the Social Studies assessments are not vertically scaled, the TCCs for grades 4, 8, and 10 are presented together in Figure 6-14 for comparison purposes. The TCCs are S-shaped, indicating increasing probability of a higher test score as a student’s ability increases. The grade 4 and grade 8 TCCs are parallel to each other, indicating similar overall test discrimination of the two assessments.

Standard Error Curves—Figure 6-15 shows Social Studies SE curves for grades 4, 8, and 10. The SE curves are U-shaped, indicating smaller errors around ability estimates that are approximately in the middle of the scale score distribution. The SE is expected to be higher at the top and bottom ends of the ability scale, where fewer items measuring these students are found. Overall, the SEs around the scale score were found to be reasonable for Social Studies assessments (for more details, see Section 8.1.1 of this report).

Scale Scores at Quartiles—The estimated scale scores for the Social Studies calibration sample at the 25th, 50th, and 75th percentiles for all grade levels are presented in Figure 6-16. The data pattern presented in this figure indicates that the scale scores increase as the percentile increases within each grade level. Because the Social Studies assessments are not on a vertical scale, it is not appropriate to compare scale scores between grades.

6.4.4 Derivation of Scale Scores

A scale score can be interpreted as a highly probable estimate of a student’s ability in a given content area. Scale scores are based on the student’s responses to all items on a given test and account for the characteristics of the items that are on the test (such as item difficulty). Item parameters estimated after the Spring 2023 test administration were used to derive student scale scores.

Scale scores in the Wisconsin Forward Exam are based on the theoretical models of the item response process described above and elaborated upon below. The essential idea behind these models is that the probability of a correct response to a given item is a function of examinee ability and the characteristics of the item, such as the difficulty of the item. It is expected that as examinee ability increases, the probability of a correct response to a given item also increases, given certain conditions and assumptions. This description applies specifically to MC items; non-MC items are treated as CR items and are handled slightly differently, but they follow a logic that is essentially the same.

Whether looking at an individual item or at a group of items that make up a complete test, IRT uses probability models to describe the relationship between a student’s ability and that student’s observed scores. As described above, the 3PL model is used to estimate the probability of a correct response for each of the MC items. The model is provided here because its components are reviewed in the following paragraphs.

$$P(u_i = 1 | \theta) = c_i + \frac{1 - c_i}{1 + e^{-1.7a_i(\theta - b_i)}} \quad (1)$$

In this model, θ denotes a measured ability (e.g., ELA ability) and u_i represents an observed score on a particular item. For MC items, the observed score u_i is either 0 or 1, indicating either an incorrect or a correct response, respectively. For an MC item, the probability model can be denoted as $P(u_i = 1 | \theta)$. That is, P is an estimation of the probability that a student with an ability value θ would answer item i correctly.

The terms on the right side of the equation above (a_i , b_i , c_i) represent the parameters in the model: discrimination, difficulty (or location), and a pseudo-guessing factor. Discrimination refers to how well an item sorts students by ability level, difficulty represents the difficulty of the item or its location on an ability continuum, and the pseudo-guessing factor represents the probability of a low-ability student guessing the correct response.

Given any particular response pattern ($u_1 u_2 \cdots u_n$) on a test with some number of items (n items), the “likelihood function,” or the probability that a student with a given ability value (θ) would produce this particular response pattern, is given by

$$P(u_1 u_2 \cdots u_n | \theta) = \prod_{i=1}^n P(u_i | \theta). \quad (2)$$

The formula indicates that the “estimated maximum likelihood” IRT item-pattern scoring method searches for the ability estimate (θ_0) that maximizes the probability function in the equation shown above (2) and assigns an ability estimate (θ_0) as the test score for the student with the response pattern ($u_1 u_2 \cdots u_n$). In other words, the

scale score is the most likely, or most probable, estimate of student ability, produced in a context in which item parameters are known and based on all the items in a given test.

As indicated, the item-pattern scoring method takes into account not only a student’s total raw score but also the psychometric characteristics of all items the student responded to, including the items the student responded to incorrectly.

Consider the following example. Suppose six examinees in grade 4 take an ELA test with 30 MC items. Suppose further that the properties, or parameters, of the items on that test are as follows (see Table 6-A).

Table 6-A Example of Item Parameters for a Test

Item	Discrimination (<i>a</i>)	Location (<i>b</i>)	Guessing (<i>c</i>)	Item	Discrimination (<i>a</i>)	Location (<i>b</i>)	Guessing (<i>c</i>)
1	0.0341	318.75	0.16	16	0.0398	286.13	0.13
2	0.0342	244.62	0.20	17	0.0523	290.65	0.26
3	0.0234	257.56	0.20	18	0.0387	280.23	0.14
4	0.0306	235.00	0.20	19	0.0329	315.71	0.21
5	0.0125	342.39	0.17	20	0.0370	287.88	0.25
6	0.0305	261.51	0.16	21	0.0387	280.25	0.18
7	0.0316	296.93	0.19	22	0.0321	285.86	0.17
8	0.0228	252.70	0.20	23	0.0219	302.52	0.13
9	0.0383	266.28	0.20	24	0.0551	301.11	0.26
10	0.0229	308.84	0.11	25	0.0165	324.24	0.19
11	0.0536	259.00	0.21	26	0.0279	297.19	0.11
12	0.0478	245.19	0.20	27	0.0423	296.06	0.28
13	0.0418	276.25	0.28	28	0.0658	324.76	0.21
14	0.0377	287.60	0.23	29	0.0488	281.56	0.32
15	0.0177	316.08	0.24	30	0.0237	345.32	0.37

Now suppose that the student response patterns for these six examinees are as follows, where 0 represents an incorrect response and 1 represents a correct response (see Table 6-B).

Table 6-B Example of Item Response Pattern

Student	Response Pattern ($u_1 u_2 \dots u_n$)	Raw Score	Item-Pattern Score
Pam	1000011001010000000000000101	7	140
Craig	101010101010101010101010101010	15	246
Vicki	010101010101010101010101010101	15	266
Tom	001100110011001100110011001101	15	259
Evan	110011001100110011001100110010	15	265
Dan	1111111111111111111111111011111	29	379

The first student, Pam, answered 7 of the items correctly and obtained a scale score of 140, which is equal to the lowest point on the scale score range, called the lowest obtainable scale score, or LOSS. The next four students each answered 15 out of 30 items correctly, but the response pattern of each of these students is different. The raw score of each of these students is 15. However, the maximum likelihood item-pattern scoring method produced a different scale score for each examinee. Scale scores were 246 for Craig, 266 for Vicki, 259 for

Tom, and 265 for Evan. These scores can be accounted for by considering the pattern of the student responses on the test in conjunction with the properties (or parameters) of the items as shown in Table 6-A. By referring to Table 6-A, the reader can observe that Vicki and Evan answered some difficult and highly discriminating items correctly, whereas Craig and Tom did not. The remaining student, Dan, scored 29 out of the 30 items correctly and obtained a scale score of 379, which is near the upper limit of the scale score range, called the highest obtainable scale score, or HOSS.

Figure 6-A shows the probability of each ability estimate (or scale score) for the six examinees. The total scale score range for the test is plotted on the horizontal axis. As indicated by the two vertical lines in the plot, the lower and upper limits of the scale score range are 140 and 420, respectively. The likelihood, or probability, of all possible ability estimates for each examinee is plotted on the vertical axis and ranges from 0 to 1.0. The higher the likelihood, the more probable it is that the ability estimate accurately reflects the examinee's ability level.

As indicated above, scale scores are the most likely, or the maximum likelihood, estimates of examinee ability. As can be observed for Vicki, Tom, and Evan, scores that are plus or minus only a few scale score points are markedly less likely estimates of the students' abilities. The same is true for Craig and Dan, though to a slightly lesser extent. In the case of Pam, a few scores were almost as likely as the maximum likelihood estimate reported. Those scores that appear to be more likely than the reported score are outside of the scale score range of the test (below the LOSS).

There are two IRT-based scoring methods generally used for large-scale assessments: number-correct scoring and item-pattern scoring. Item-pattern scoring may be recommended over number-correct scoring for several reasons. Two reasons, accuracy and reliability, are pertinent for the present purposes.

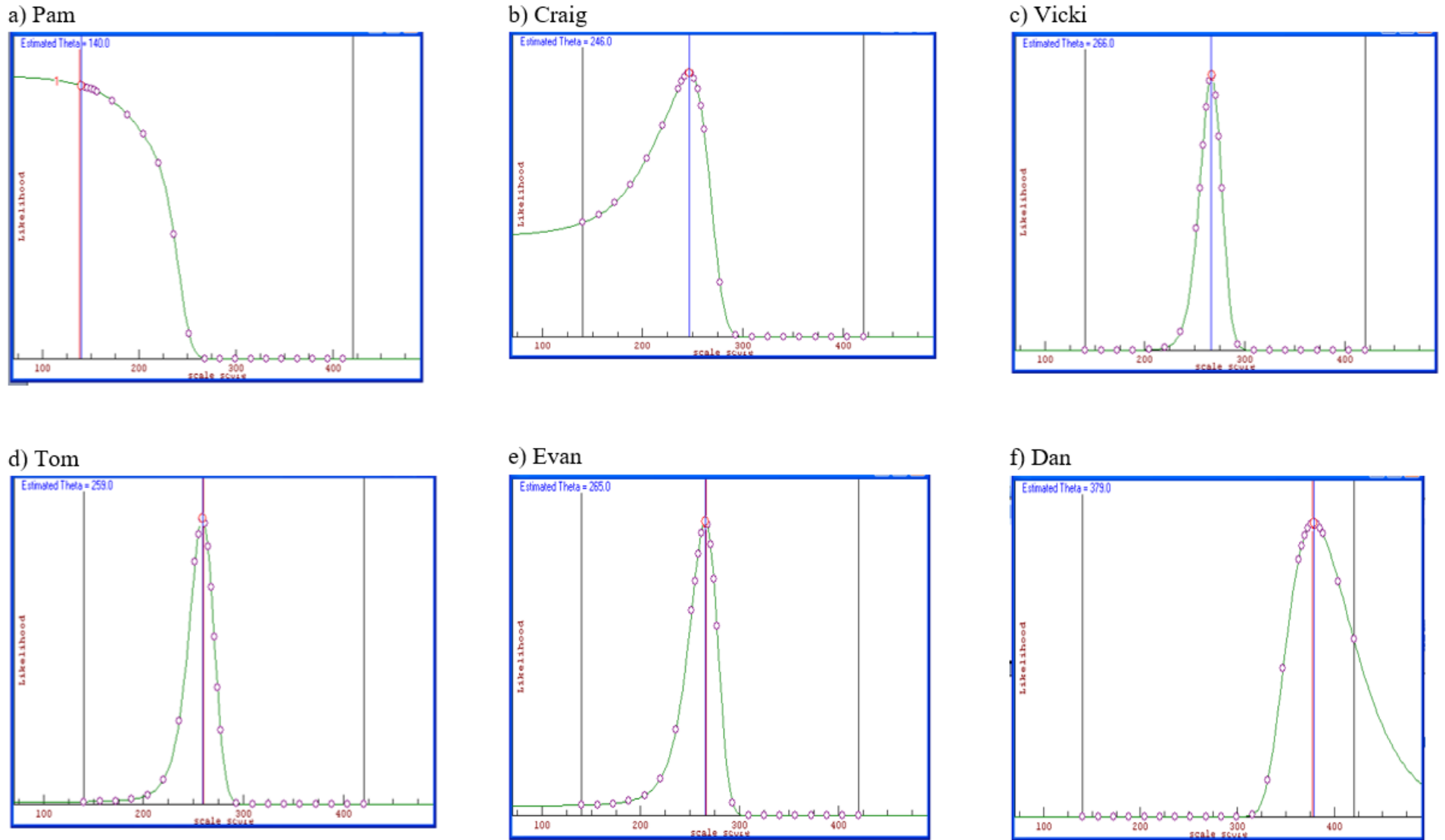
First, item-pattern scoring generally produces more accurate scores for individual students. Specifically, it produces a smaller conditional standard error of measurement (CSEM) across the scale score range for a given test compared to number-correct scoring. The smaller the CSEM, the more confident one can be in the accuracy of the test results. The increase in accuracy provided by item-pattern scoring is equivalent, on average, to an increase by approximately 15 to 20 percent in test length (Yen, 1984; Yen & Candell, 1991).

Second, reliability tends to be higher using item-pattern scoring, which means (a) fewer items are needed to achieve a given level of reliability and (b) a given test with a given number of items will have higher reliability than it would when using number-correct scoring. Yen (1984) has demonstrated that an equivalent level of reliability for a 20-item test scored by the number correct scoring method could be obtained with a 16- or 17-item test scored by the item pattern scoring method.

Because of the nature of item-pattern scoring, a scoring table showing a simple, direct conversion of raw score to scale score cannot be generated for the Spring 2023 Wisconsin Forward Exam. However, scoring tables showing an approximate raw score-to-scale score relationship and the associated CSEM can be produced, and these are provided in Tables 6-31 through 6-47. These tables are provided to illustrate the approximate raw score-to-scale score relationship for each unique raw score and do not include all combinations of raw score-to-scale score associations.

Several supplements to this simplified outline of IRT are available. Introductory discussions of IRT can be found in *Educational Measurement* (Linn, 1989) or Chapter 11 in *Introduction to Measurement Theory* (Allen & Yen, 1979). More advanced discussions of partial-credit models may be found in Muraki (1990, 1992), Yen (1993), and van der Linden & Hambleton (1997). For additional information on the technical details of item-pattern scoring, readers can also refer to Yen & Candell (1991).

Figure 6-A Examples of Likelihood Functions or the Probability of Each Ability Level Estimate (or Scale Score)



Note: The circular dots in the likelihood functions indicate that the software program used is searching for a maximum likelihood estimate (scale score) for the student.

6.4.5 Lowest and Highest Obtainable Scale Scores

As previously established, a scale score is a maximum likelihood ability estimate. The maximum-likelihood procedure cannot produce scale score estimates for students with perfect scores or scores below the scoring level expected by guessing. Although maximum likelihood estimates are available for students with extreme scores other than zero or a perfect score, these estimates generally have large SEMs. Therefore, scores are established for these extreme highs and lows based on a rational, but necessarily non-maximum, likelihood procedure. These values are set separately by grade and are called the LOSS and the HOSS. The LOSS and HOSS values for ELA and Mathematics were established after the Spring 2016 test administration and remained unchanged through the Spring 2023 test administration. New LOSS and HOSS values were established after the Spring 2019 test administration for Science and after the Spring 2022 test administration for Social Studies.

Table 6-48 shows the number and percentage of students at the LOSS and the HOSS. In general, there should not be many students clustered at the LOSS or HOSS. A high proportion of students at the LOSS or HOSS may indicate a floor or ceiling effect.

It should be noted that for ELA and Mathematics, the LOSS and HOSS values were set in such a way during the Spring 2016 scale development that they increase as the grade level increases. Setting increasing LOSS values as the grade level increases is an important property of a vertical scale and constrains student ability in each grade in such a way that the lowest-ability students in a given grade will always have a higher scale score than the lowest-ability students in a grade below and a lower scale score than the lowest-ability students in a grade above. Conversely, setting increasing HOSS values as the grade level increases constrains student ability in each grade in such a way that the highest-ability students in a given grade will always have a higher scale score than the highest-ability students in a grade below and a lower scale score than the highest-ability students in a grade above.

Approximately one-tenth of one percent of students or less received the lowest obtainable scale scores in all ELA and Science grades. In addition, less than one percent of students scored at the LOSS in Social Studies grades 4 and 8. The percentages of students scoring at the LOSS was higher for all Mathematics grades and ranged from just over 1% to less than 2% across all grades. Approximately 3% of students scored at the LOSS in Social Studies grade 10. These percentages of students at the LOSS in Mathematics and Social Studies grade 10 were investigated. It was found that students who scored this way typically correctly answered very few MC items and no non-MC items, which resulted in their LOSS values. For these students to receive a scale score above the LOSS, they would need to correctly answer more items, including some non-MC items. Non-MC items do not assume guessing, so the correct responses tend to represent student ability more accurately.

No students scored at the HOSS in ELA grades 4, 5, and 6. One student received the highest obtainable score in ELA grade 3, and one student received it in ELA grade 7. Six students scored at the HOSS in ELA grade 8. Less than half of one percent of students in all grades of Mathematics, Science, and Social Studies scored at the HOSS.

6.5 Summary

In summary, the overall purpose of the test psychometric data analysis, including scaling and equating, is to ensure that the test items, as well as the overall test, are functioning appropriately. It also helps maintain the test scale so that test results may be appropriately compared across years. The data analyses undertaken by DRC are in alignment with multiple best practices of the testing industry and, in particular, support the following AERA, APA, & NCME (2014) Standards: 1.8, 4.14, 5.2, 5.13, 5.15, and 7.2.

Table 6-1 Summary of Flagged Operational Items on the Wisconsin Forward Exam

Content	Grade	# of Items Flagged	Number of Flags			
			Correlation <0.15	Distractor Correlation >0	Omit >3%	p-Value <0.20
ELA	3	3	1	2	0	0
	4	0	0	0	0	0
	5	3	0	3	0	0
	6	1	0	1	0	0
	7	3	0	3	0	0
	8	2	0	2	0	0
Mathematics	3	2	0	2	0	0
	4	6	0	6	0	0
	5	6	1	6	0	0
	6	5	2	5	0	0
	7	4	1	3	0	1
	8	8	0	5	0	3
Science	4	1	0	0	0	1
	8	3	0	3	0	0
Social Studies	4	2	0	2	0	0
	8	0	0	0	0	0
	10	4	0	4	0	0
Total		53	5	47	0	5

Note: The number of flags may be greater than the number of flagged items.

Table 6-2 Items Flagged for Classical Item Analysis Statistics, English Language Arts

Grade	Content	Item	Item Type	<i>p</i> -Value	Corr	Percent Omit	Flags				
							Corr	Distractor	Distractor Statistic	Omit	<i>p</i> -Value
3	ELA	5	MS	0.30	0.12	0.22	+				
3	ELA	6	MC	0.67	0.23	0.18		+	0.03		
3	ELA	20	MC	0.45	0.28	0.21		+	0.03		
5	ELA	3	MC	0.28	0.21	0.08		+	0.02		
5	ELA	26	MC	0.45	0.21	0.50		+	0.04		
5	ELA	32	MC	0.31	0.20	0.18		+	0.05		
6	ELA	26	MC	0.44	0.29	0.25		+	0.06		
7	ELA	22	MC	0.41	0.21	0.42		+	0.06		
7	ELA	26	MC	0.40	0.18	0.28		+	0.01		
7	ELA	36	MC	0.34	0.27	0.24		+	0.06		
8	ELA	14	MC	0.39	0.16	0.14		+	0.02		
8	ELA	31	MC	0.43	0.21	0.24		+	0.06		

Table 6-3 Items Flagged for Classical Item Analysis Statistics, Mathematics

Grade	Content	Item	Item Type	p-Value	Corr	Percent Omit	Flags				
							Corr	Distractor	Distractor Statistic	Omit	p-Value
3	Math	23	MC	0.34	0.27	0.11		+	0.05		
3	Math	28	MC	0.34	0.24	0.15		+	0.03		
4	Math	12	MC	0.45	0.23	0.12		+	0.01		
4	Math	14	MC	0.42	0.49	0.14		+	0.01		
4	Math	23	MC	0.27	0.31	0.13		+	0.02		
4	Math	26	MC	0.27	0.38	0.09		+	0.02		
4	Math	30	MC	0.36	0.21	0.11		+	0.08		
4	Math	42	MC	0.50	0.16	0.18		+	0.05		
5	Math	2	MC	0.54	0.23	0.09		+	0.03		
5	Math	11	MC	0.43	0.31	0.13		+	0.03		
5	Math	23	MC	0.35	0.43	0.12		+	0.01		
5	Math	31	MC	0.34	0.17	0.30		+	0.12		
5	Math	38	MC	0.28	0.14	0.13	+	+	0.16		
5	Math	44	MC	0.31	0.48	0.14		+	0.00		
6	Math	5	MC	0.33	0.11	0.13	+	+	0.03		
6	Math	11	MC	0.31	0.29	0.15		+	0.02		
6	Math	24	MC	0.41	0.19	0.40		+	0.17		
6	Math	39	MC	0.34	0.11	0.23	+	+	0.11		
6	Math	43	MC	0.34	0.21	0.22		+	0.02		
7	Math	9	SA	0.14	0.44	0.50					+
7	Math	31	MC	0.31	0.15	0.43		+	0.10		
7	Math	35	MC	0.38	0.13	0.49	+	+	0.09		
7	Math	40	MC	0.50	0.18	0.37		+	0.12		
8	Math	1	SA	0.19	0.47	0.15					+
8	Math	5	MC	0.30	0.47	0.11		+	0.02		
8	Math	8	MC	0.34	0.23	0.14		+	0.18		
8	Math	20	TE	0.18	0.38	0.18					+
8	Math	24	MC	0.40	0.25	0.31		+	0.02		
8	Math	35	TE	0.19	0.62	0.43					+
8	Math	38	MC	0.42	0.29	0.35		+	0.04		
8	Math	44	MC	0.37	0.22	0.42		+	0.02		

Table 6-4 Items Flagged for Classical Item Analysis Statistics, Science and Social Studies

Grade	Content	Item	Item Type	p-Value	Corr	Percent Omit	Flags				
							Corr	Distractor	Distractor Statistic	Omit	p-Value
4	Science	18	TE	0.19	0.32	0.16					+
8	Science	8	MC	0.40	0.23	0.22		+	0.01		
8	Science	12	MC	0.46	0.22	0.26		+	0.00		
8	Science	16	MC	0.48	0.15	0.11		+	0.01		
4	Socials	10	MC	0.51	0.29	0.13		+	0.04		
4	Socials	17	MC	0.34	0.29	0.12		+	0.00		
10	Socials	5	MC	0.41	0.29	0.38		+	0.00		
10	Socials	18	MC	0.32	0.26	0.36		+	0.05		
10	Socials	27	MC	0.34	0.21	0.28		+	0.06		
10	Socials	34	MC	0.55	0.27	0.32		+	0.06		

Table 6-5 Percentage of Students Attempting Last Operational Item in Test

Content	Grade						
	3	4	5	6	7	8	10
English Language Arts	99.76	99.84	99.64	99.79	99.77	99.76	
Mathematics	99.75	99.88	99.88	99.66	99.56	99.40	
Science		99.84				99.85	
Social Studies		99.91				99.80	99.65

Table 6-6 Item Analysis, English Language Arts Grade 3

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TDA	0.34	0.39	0.54
2	MC	0.42	0.27	0.15
3	MC	0.79	0.43	0.25
4	MC	0.61	0.24	0.20
5	MS	0.30	0.12	0.22
6	MC	0.67	0.23	0.18
7	MC	0.52	0.41	0.20
8	MC	0.73	0.36	0.17
9	TE	0.46	0.16	0.37
10	EBSR	0.56	0.54	0.13
11	MC	0.68	0.48	0.23
12	TE	0.50	0.50	0.24
13	MC	0.44	0.36	0.20
14	MC	0.50	0.30	0.23
15	MC	0.58	0.33	0.10
16	MC	0.69	0.46	0.20
17	MC	0.59	0.25	0.18
18	EBSR	0.50	0.37	0.08
19	MC	0.56	0.38	0.20
20	MC	0.45	0.28	0.21
21	MC	0.33	0.30	0.13
22	MC	0.64	0.39	0.23
23	MC	0.49	0.33	0.23
24	MC	0.85	0.45	0.26
25	MC	0.69	0.50	0.21
26	MS	0.58	0.44	0.18
27	TE	0.70	0.48	1.00
28	MC	0.58	0.42	0.21
29	MC	0.36	0.36	0.24
30	EBSR	0.52	0.52	0.13
31	MC	0.74	0.46	0.24
32	MC	0.52	0.29	0.20
33	EBSR	0.48	0.57	0.16
34	MC	0.60	0.52	0.28
35	MC	0.45	0.39	0.31
36	MC	0.48	0.32	0.32
37	MC	0.44	0.33	0.26
38	MC	0.43	0.35	0.24

Table 6-7 Item Analysis, English Language Arts Grade 4

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TDA	0.40	0.53	0.29
2	MC	0.78	0.28	0.07
3	MC	0.47	0.36	0.11
4	EBSR	0.61	0.54	0.06
5	MC	0.82	0.45	0.10
6	MC	0.60	0.33	0.12
7	TE	0.50	0.48	0.13
8	MC	0.53	0.43	0.14
9	MC	0.59	0.37	0.10
10	MC	0.66	0.46	0.10
11	MC	0.43	0.25	0.12
12	TE	0.70	0.38	0.85
13	MC	0.63	0.48	0.13
14	MC	0.31	0.21	0.13
15	MC	0.65	0.32	0.06
16	MC	0.54	0.39	0.15
17	EBSR	0.36	0.33	0.04
18	MC	0.56	0.30	0.09
19	MC	0.69	0.38	0.07
20	EBSR	0.30	0.34	0.05
21	MS	0.46	0.34	0.07
22	MC	0.44	0.21	0.17
23	MC	0.58	0.25	0.21
24	MC	0.62	0.26	0.23
25	MC	0.79	0.50	0.15
26	MC	0.56	0.45	0.26
27	EBSR	0.32	0.44	0.06
28	MC	0.56	0.36	0.12
29	MC	0.71	0.50	0.19
30	MC	0.57	0.41	0.18
31	MC	0.69	0.38	0.23
32	MC	0.61	0.45	0.17
33	TE	0.50	0.48	0.14
34	MC	0.64	0.47	0.18
35	MC	0.68	0.51	0.24
36	TE	0.70	0.52	0.61
37	MC	0.43	0.37	0.21
38	EBSR	0.38	0.40	0.11
39	MC	0.48	0.41	0.16

Table 6-8 Item Analysis, English Language Arts Grade 5

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TDA	0.38	0.50	0.18
2	MC	0.44	0.19	0.04
3	MC	0.28	0.21	0.08
4	MS	0.52	0.42	0.08
5	MC	0.38	0.27	0.06
6	MC	0.60	0.41	0.11
7	TE	0.73	0.35	0.10
8	MC	0.27	0.16	0.07
9	MC	0.47	0.29	0.07
10	MC	0.78	0.33	0.09
11	TE	0.69	0.53	0.08
12	MC	0.70	0.43	0.07
13	TE	0.36	0.40	0.34
14	MC	0.76	0.29	0.06
15	EBSR	0.62	0.41	0.03
16	MC	0.60	0.20	0.11
17	MC	0.63	0.42	0.08
18	MC	0.59	0.37	0.07
19	EBSR	0.57	0.51	0.03
20	MC	0.79	0.38	0.08
21	MC	0.55	0.29	0.17
22	EBSR	0.36	0.40	0.06
23	MC	0.61	0.38	0.19
24	MC	0.51	0.51	0.29
25	MC	0.67	0.47	0.16
26	MC	0.45	0.21	0.50
27	EBSR	0.30	0.17	0.09
28	MC	0.78	0.48	0.19
29	MC	0.79	0.40	0.18
30	MC	0.36	0.22	0.24
31	MC	0.38	0.32	0.24
32	MC	0.31	0.20	0.18
33	MC	0.47	0.43	0.53
34	MC	0.50	0.37	0.28
35	MC	0.63	0.51	0.32
36	MC	0.49	0.40	0.23
37	MC	0.59	0.51	0.28
38	MC	0.60	0.45	0.21
39	EBSR	0.51	0.60	0.13
40	MC	0.58	0.25	0.36

Table 6-9 Item Analysis, English Language Arts Grade 6

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TDA	0.42	0.52	0.26
2	MC	0.67	0.25	0.07
3	MC	0.43	0.35	0.15
4	MC	0.48	0.20	0.13
5	MC	0.81	0.41	0.17
6	EBSR	0.66	0.49	0.05
7	TE	0.28	0.29	0.18
8	MC	0.36	0.25	0.15
9	EBSR	0.23	0.21	0.08
10	TE	0.65	0.46	0.15
11	MC	0.54	0.35	0.19
12	MC	0.42	0.30	0.13
13	TE	0.40	0.27	0.28
14	MC	0.65	0.33	0.08
15	EBSR	0.31	0.33	0.05
16	MC	0.49	0.23	0.10
17	MC	0.66	0.47	0.10
18	MC	0.51	0.33	0.15
19	TE	0.58	0.45	0.12
20	MC	0.65	0.44	0.08
21	MC	0.70	0.36	0.20
22	MC	0.68	0.45	0.25
23	MS	0.65	0.50	0.28
24	EBSR	0.63	0.55	0.10
25	MC	0.71	0.46	0.24
26	MC	0.44	0.29	0.25
27	MC	0.70	0.52	0.28
28	MC	0.53	0.28	0.22
29	MC	0.76	0.52	0.33
30	MC	0.51	0.38	0.29
31	EBSR	0.60	0.50	0.13
32	TE	0.69	0.57	0.25
33	MC	0.68	0.47	0.24
34	MC	0.54	0.46	0.31
35	MC	0.82	0.51	0.26
36	MC	0.45	0.45	0.26
37	TE	0.47	0.31	0.57
38	MC	0.48	0.44	0.21

Table 6-10 Item Analysis, English Language Arts Grade 7

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TDA	0.47	0.53	0.23
2	MC	0.43	0.34	0.08
3	TE	0.49	0.40	0.12
4	EBSR	0.58	0.48	0.05
5	MC	0.80	0.34	0.07
6	TE	0.45	0.25	0.13
7	MC	0.51	0.20	0.20
8	MC	0.56	0.46	0.17
9	MC	0.68	0.36	0.13
10	MC	0.33	0.26	0.12
11	MC	0.41	0.29	0.19
12	TE	0.83	0.32	0.07
13	EBSR	0.47	0.24	0.10
14	EBSR	0.55	0.47	0.02
15	MC	0.76	0.49	0.11
16	MS	0.63	0.52	0.11
17	MC	0.79	0.44	0.11
18	EBSR	0.66	0.54	0.04
19	MC	0.48	0.39	0.08
20	MS	0.36	0.24	0.18
21	MC	0.75	0.41	0.19
22	MC	0.41	0.21	0.42
23	EBSR	0.52	0.53	0.08
24	MC	0.75	0.39	0.27
25	MS	0.61	0.59	0.27
26	MC	0.40	0.18	0.28
27	MC	0.64	0.39	0.36
28	MC	0.78	0.49	0.24
29	MC	0.41	0.31	0.30
30	TE	0.34	0.43	0.74
31	MC	0.71	0.55	0.21
32	MC	0.56	0.38	0.30
33	MC	0.55	0.37	0.46
34	MC	0.73	0.43	0.28
35	TE	0.56	0.44	0.75
36	MC	0.34	0.27	0.24
37	TE	0.45	0.23	0.23

Table 6-11 Item Analysis, English Language Arts Grade 8

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TDA	0.55	0.57	0.29
2	MC	0.30	0.22	0.06
3	MC	0.49	0.22	0.17
4	MC	0.61	0.39	0.17
5	MC	0.79	0.34	0.18
6	MC	0.82	0.42	0.12
7	MC	0.58	0.35	0.14
8	MC	0.66	0.36	0.18
9	EBSR	0.50	0.54	0.08
10	TE	0.29	0.27	0.30
11	TE	0.62	0.50	0.24
12	MC	0.50	0.25	0.18
13	MS	0.53	0.36	0.13
14	MC	0.39	0.16	0.14
15	MC	0.29	0.34	0.06
16	MC	0.35	0.34	0.14
17	MC	0.47	0.34	0.13
18	MS	0.63	0.48	0.11
19	EBSR	0.63	0.53	0.06
20	MC	0.65	0.34	0.16
21	MC	0.73	0.40	0.08
22	EBSR	0.55	0.45	0.07
23	MC	0.87	0.50	0.30
24	MC	0.39	0.42	0.30
25	MS	0.62	0.50	0.29
26	MC	0.74	0.46	0.26
27	MC	0.56	0.28	0.52
28	MC	0.53	0.30	0.26
29	EBSR	0.61	0.62	0.13
30	MC	0.79	0.52	0.27
31	MC	0.43	0.21	0.24
32	MC	0.55	0.35	0.29
33	MC	0.63	0.48	0.20
34	EBSR	0.26	0.18	0.16
35	MC	0.66	0.45	0.27
36	MS	0.57	0.43	0.31
37	MC	0.44	0.40	0.28
38	MC	0.68	0.47	0.27
39	MC	0.47	0.43	0.24

Table 6-12 Item Analysis, Mathematics Grade 3

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	MC	0.80	0.49	0.07
2	TE	0.64	0.58	0.51
3	MC	0.45	0.42	0.13
4	SA	0.54	0.56	0.20
5	MC	0.57	0.30	0.15
6	TE	0.29	0.55	0.17
7	MC	0.52	0.48	0.17
8	MC	0.36	0.42	0.17
9	SA	0.45	0.59	0.28
10	TE	0.66	0.40	0.22
11	MC	0.71	0.44	0.20
12	SA	0.49	0.53	0.39
13	SA	0.38	0.57	0.29
14	MC	0.61	0.38	0.20
15	MC	0.45	0.31	0.28
16	MC	0.69	0.48	0.23
17	TE	0.36	0.41	0.33
18	MC	0.33	0.38	0.25
19	SA	0.55	0.54	0.31
20	SA	0.69	0.59	0.31
21	MC	0.79	0.46	0.22
22	MC	0.71	0.40	0.07
23	MC	0.34	0.27	0.11
24	TE	0.82	0.42	0.14
25	SA	0.34	0.56	0.18
26	MC	0.37	0.56	0.15
27	SA	0.54	0.47	0.22
28	MC	0.34	0.24	0.15
29	MC	0.88	0.37	0.14
30	SA	0.37	0.51	0.17
31	MC	0.68	0.50	0.18
32	SA	0.60	0.60	0.21
33	MC	0.39	0.44	0.18
34	MC	0.59	0.51	0.14
35	SA	0.48	0.55	0.20
36	MC	0.66	0.46	0.19
37	TE	0.45	0.61	0.19
38	MC	0.64	0.50	0.20
39	MC	0.69	0.59	0.16
40	MC	0.61	0.48	0.16
41	TE	0.46	0.54	1.43
42	SA	0.20	0.49	0.25

Table 6-13 Item Analysis, Mathematics Grade 4

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TE	0.34	0.51	0.05
2	MC	0.70	0.44	0.10
3	SA	0.57	0.58	0.12
4	MC	0.79	0.40	0.12
5	TE	0.72	0.50	0.81
6	MC	0.39	0.48	0.10
7	TE	0.73	0.41	0.52
8	MC	0.43	0.38	0.37
9	MC	0.64	0.47	0.14
10	MC	0.62	0.43	0.14
11	SA	0.36	0.63	0.20
12	MC	0.45	0.23	0.12
13	TE	0.54	0.63	0.20
14	MC	0.42	0.49	0.14
15	MC	0.65	0.35	0.12
16	MC	0.41	0.45	0.37
17	MC	0.38	0.31	0.15
18	SA	0.31	0.51	0.19
19	TE	0.34	0.61	0.15
20	MC	0.66	0.40	0.15
21	SA	0.41	0.49	0.13
22	MC	0.49	0.52	0.11
23	MC	0.27	0.31	0.13
24	MC	0.37	0.49	0.06
25	TE	0.70	0.38	0.32
26	MC	0.27	0.38	0.09
27	SA	0.52	0.53	0.14
28	MC	0.38	0.48	0.15
29	TE	0.68	0.38	0.10
30	MC	0.36	0.21	0.11
31	MC	0.53	0.48	0.35
32	MC	0.33	0.34	0.12
33	SA	0.43	0.54	0.18
34	MC	0.33	0.41	0.16
35	TE	0.84	0.45	0.15
36	MC	0.52	0.21	0.13
37	SA	0.78	0.46	0.12
38	MC	0.38	0.57	0.13
39	MC	0.31	0.42	0.54
40	MC	0.57	0.54	0.16

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
41	SA	0.61	0.49	0.15
42	MC	0.50	0.16	0.18
43	MC	0.37	0.35	0.16
44	SA	0.53	0.59	0.17
45	MC	0.86	0.37	0.13
46	TE	0.25	0.53	0.12

Table 6-14 Item Analysis, Mathematics Grade 5

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	MC	0.69	0.49	0.05
2	MC	0.54	0.23	0.09
3	TE	0.51	0.54	0.08
4	TE	0.52	0.53	0.11
5	MC	0.34	0.20	0.09
6	SA	0.29	0.34	0.14
7	MC	0.41	0.45	0.12
8	SA	0.36	0.45	0.30
9	TE	0.22	0.53	0.11
10	MC	0.30	0.27	0.12
11	MC	0.43	0.31	0.13
12	MC	0.35	0.21	0.14
13	TE	0.44	0.44	0.08
14	SA	0.52	0.55	0.26
15	MC	0.55	0.30	0.09
16	SA	0.48	0.56	0.36
17	MC	0.68	0.33	0.15
18	MC	0.42	0.43	0.17
19	TE	0.24	0.48	0.71
20	SA	0.43	0.48	0.17
21	MC	0.63	0.53	0.13
22	SA	0.60	0.58	0.23
23	MC	0.35	0.43	0.12
24	SA	0.59	0.46	0.11
25	MC	0.65	0.27	0.11
26	SA	0.62	0.57	0.14
27	TE	0.46	0.48	0.11
28	MC	0.53	0.33	0.12
29	SA	0.37	0.46	0.24
30	SA	0.45	0.60	0.14
31	MC	0.34	0.17	0.30
32	MC	0.48	0.40	0.12
33	MC	0.63	0.45	0.15
34	MC	0.45	0.31	0.15
35	TE	0.62	0.52	0.33
36	SA	0.43	0.58	0.15
37	MC	0.37	0.56	0.16
38	MC	0.28	0.14	0.13
39	SA	0.49	0.58	0.52
40	MC	0.61	0.58	0.17

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
41	TE	0.41	0.35	0.13
42	MC	0.48	0.28	0.16
43	SA	0.28	0.57	0.22
44	MC	0.31	0.48	0.14
45	SA	0.62	0.53	0.18
46	TE	0.69	0.55	0.12

Table 6-15 Item Analysis, Mathematics Grade 6

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	MC	0.54	0.41	0.08
2	MC	0.58	0.49	0.15
3	MC	0.38	0.49	0.11
4	TE	0.34	0.56	0.11
5	MC	0.33	0.11	0.13
6	SA	0.22	0.56	0.26
7	MC	0.53	0.47	0.14
8	TE	0.79	0.45	0.15
9	MC	0.48	0.52	0.15
10	MC	0.35	0.34	0.17
11	MC	0.31	0.29	0.15
12	MC	0.76	0.40	0.13
13	MC	0.61	0.55	0.18
14	MC	0.52	0.51	0.16
15	TE	0.39	0.65	0.23
16	SA	0.58	0.60	0.28
17	TE	0.80	0.34	0.79
18	SA	0.63	0.52	0.31
19	MC	0.65	0.37	0.14
20	MC	0.39	0.29	0.15
21	MC	0.36	0.49	0.11
22	SA	0.60	0.57	0.20
23	TE	0.54	0.52	0.16
24	MC	0.41	0.19	0.40
25	MC	0.37	0.24	0.39
26	SA	0.22	0.56	0.70
27	TE	0.64	0.51	0.17
28	MC	0.42	0.31	0.29
29	SA	0.38	0.57	0.31
30	SA	0.36	0.22	0.36
31	MC	0.57	0.35	0.18
32	MC	0.34	0.42	0.29
33	MC	0.56	0.50	0.17
34	SA	0.25	0.54	0.57
35	MC	0.69	0.48	0.29
36	TE	0.37	0.52	1.15
37	MC	0.63	0.39	0.25
38	SA	0.67	0.52	0.36
39	MC	0.34	0.11	0.23
40	TE	0.51	0.57	1.42

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
41	SA	0.42	0.54	0.41
42	MC	0.40	0.26	0.25
43	MC	0.34	0.21	0.22
44	TE	0.38	0.58	0.34
45	MC	0.46	0.41	0.31
46	TE	0.66	0.57	0.34

Table 6-16 Item Analysis, Mathematics Grade 7

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	SA	0.77	0.45	0.13
2	MC	0.53	0.52	0.10
3	SA	0.26	0.48	0.29
4	MC	0.49	0.48	0.09
5	TE	0.31	0.49	0.21
6	MC	0.42	0.28	0.14
7	SA	0.26	0.49	0.22
8	MC	0.62	0.50	0.10
9	SA	0.14	0.44	0.50
10	MC	0.53	0.34	0.18
11	MC	0.27	0.35	0.11
12	MC	0.43	0.24	0.15
13	TE	0.36	0.45	0.46
14	MC	0.65	0.53	0.21
15	MC	0.48	0.38	0.22
16	TE	0.47	0.24	0.19
17	SA	0.27	0.44	0.55
18	TE	0.72	0.23	0.22
19	MC	0.48	0.53	0.38
20	SA	0.73	0.53	0.43
21	MC	0.35	0.30	0.43
22	SA	0.53	0.55	0.45
23	MC	0.34	0.24	0.37
24	TE	0.33	0.63	0.70
25	MC	0.47	0.56	0.26
26	MC	0.44	0.26	0.31
27	MC	0.52	0.46	0.33
28	TE	0.46	0.26	0.28
29	MC	0.57	0.23	0.29
30	MC	0.69	0.49	0.29
31	MC	0.31	0.15	0.43
32	MC	0.43	0.34	0.51
33	SA	0.51	0.54	0.85
34	MC	0.36	0.34	0.42
35	MC	0.38	0.13	0.49
36	MC	0.48	0.53	0.37
37	SA	0.37	0.67	0.71
38	TE	0.65	0.54	0.43
39	SA	0.26	0.51	0.74
40	MC	0.50	0.18	0.37

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
41	SA	0.28	0.54	0.60
42	TE	0.45	0.27	0.70
43	MC	0.26	0.20	0.54
44	SA	0.68	0.61	0.66
45	MC	0.59	0.50	0.50
46	MC	0.51	0.41	0.44

Table 6-17 Item Analysis, Mathematics Grade 8

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	SA	0.19	0.47	0.15
2	MC	0.61	0.48	0.09
3	TE	0.58	0.46	0.10
4	MC	0.45	0.24	0.14
5	MC	0.30	0.47	0.11
6	SA	0.38	0.50	0.30
7	MC	0.33	0.32	0.13
8	MC	0.34	0.23	0.14
9	SA	0.44	0.63	0.21
10	MC	0.52	0.48	0.09
11	MC	0.30	0.42	0.14
12	MC	0.46	0.48	0.12
13	MC	0.47	0.58	0.11
14	SA	0.36	0.51	0.55
15	MC	0.61	0.37	0.24
16	MC	0.75	0.37	0.17
17	MC	0.50	0.33	0.27
18	MC	0.42	0.27	0.18
19	SA	0.32	0.55	0.55
20	TE	0.18	0.38	0.18
21	MC	0.52	0.33	0.39
22	MC	0.45	0.27	0.31
23	MC	0.29	0.48	0.34
24	MC	0.40	0.25	0.31
25	MC	0.69	0.44	0.26
26	MC	0.53	0.37	0.26
27	SA	0.70	0.53	0.39
28	TE	0.54	0.30	0.33
29	SA	0.33	0.69	0.56
30	SA	0.29	0.63	0.54
31	MC	0.57	0.32	0.25
32	TE	0.34	0.52	0.47
33	MC	0.44	0.26	0.36
34	MC	0.49	0.53	0.40
35	TE	0.19	0.62	0.43
36	SA	0.24	0.59	1.29
37	MC	0.47	0.54	0.41
38	MC	0.42	0.29	0.35
39	MC	0.42	0.42	0.32
40	TE	0.46	0.38	0.64

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
41	MC	0.52	0.37	0.30
42	TE	0.22	0.50	1.09
43	MC	0.69	0.43	0.34
44	MC	0.37	0.22	0.42
45	MC	0.66	0.44	0.38
46	TE	0.70	0.47	0.60

Table 6-18 Item Analysis, Science Grade 4

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TE	0.78	0.38	0.21
2	MC	0.60	0.45	0.10
3	TE	0.32	0.57	0.32
4	TE	0.85	0.35	0.20
5	MC	0.70	0.53	0.12
6	TE	0.57	0.46	0.17
7	TE	0.54	0.31	0.16
8	TE	0.55	0.41	1.09
9	TE	0.46	0.59	0.14
10	MC	0.67	0.40	0.12
11	MC	0.67	0.51	0.15
12	TE	0.47	0.45	0.27
13	TE	0.43	0.60	0.24
14	TE	0.42	0.40	0.19
15	MC	0.55	0.49	0.15
16	TE	0.56	0.49	0.10
17	EBSR	0.29	0.35	0.09
18	TE	0.19	0.32	0.16
19	EBSR	0.69	0.39	0.13
20	MC	0.53	0.33	0.20
21	TE	0.53	0.51	0.15
22	TE	0.73	0.33	0.11
23	MC	0.50	0.56	0.20
24	MC	0.55	0.44	0.15
25	TE	0.21	0.38	0.12
26	MC	0.65	0.35	0.06
27	MC	0.82	0.32	0.08
28	MC	0.40	0.37	0.17
29	TE	0.39	0.33	0.33
30	TE	0.36	0.15	0.14
31	TE	0.45	0.62	0.21
32	TE	0.35	0.38	0.37
33	MC	0.65	0.43	0.12
34	TE	0.68	0.47	0.11
35	TE	0.89	0.37	0.09
36	TE	0.56	0.36	0.19
37	TE	0.36	0.30	0.18
38	EBSR	0.38	0.43	0.16
39	TE	0.51	0.39	0.13
40	MC	0.53	0.54	0.16

Table 6-19 Item Analysis, Science Grade 8

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TE	0.73	0.47	0.33
2	MS	0.42	0.46	0.14
3	TE	0.73	0.32	0.19
4	TE	0.65	0.34	0.58
5	TE	0.42	0.39	0.42
6	TE	0.22	0.23	0.07
7	TE	0.70	0.33	0.18
8	MC	0.40	0.23	0.22
9	TE	0.49	0.41	0.24
10	MC	0.61	0.39	0.15
11	MC	0.53	0.41	0.15
12	MC	0.46	0.22	0.26
13	TE	0.40	0.50	0.13
14	TE	0.32	0.39	0.15
15	MC	0.47	0.36	0.17
16	MC	0.48	0.15	0.11
17	TE	0.58	0.38	0.15
18	TE	0.48	0.33	0.12
19	MC	0.52	0.42	0.30
20	TE	0.79	0.37	0.17
21	TE	0.51	0.45	0.13
22	EBSR	0.44	0.54	0.11
23	TE	0.68	0.45	0.16
24	TE	0.61	0.46	0.18
25	TE	0.51	0.44	0.14
26	MC	0.70	0.39	0.08
27	TE	0.53	0.37	0.23
28	MC	0.47	0.41	0.18
29	TE	0.74	0.38	0.15
30	TE	0.41	0.37	0.39
31	MC	0.49	0.45	0.16
32	TE	0.60	0.42	0.16
33	EBSR	0.42	0.49	0.09
34	TE	0.30	0.29	0.12
35	MC	0.55	0.50	0.12
36	TE	0.38	0.41	0.13
37	MC	0.60	0.56	0.17
38	TE	0.26	0.33	0.16
39	TE	0.64	0.42	0.18
40	EBSR	0.43	0.34	0.15

Table 6-20 Item Analysis, Social Studies Grade 4

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TE	0.55	0.47	0.05
2	MC	0.72	0.41	0.08
3	TE	0.25	0.21	0.11
4	MC	0.72	0.47	0.13
5	MC	0.72	0.44	0.12
6	MS	0.37	0.51	0.12
7	TE	0.37	0.40	0.11
8	MC	0.54	0.37	0.10
9	MC	0.73	0.49	0.12
10	MC	0.51	0.29	0.13
11	MC	0.65	0.38	0.13
12	MC	0.80	0.49	0.11
13	MS	0.39	0.50	0.08
14	MC	0.48	0.43	0.09
15	MC	0.55	0.44	0.12
16	MC	0.57	0.39	0.14
17	MC	0.34	0.29	0.12
18	TE	0.43	0.34	0.27
19	MC	0.65	0.41	0.10
20	MC	0.39	0.39	0.12
21	TE	0.54	0.47	0.06
22	MC	0.74	0.48	0.11
23	MC	0.71	0.53	0.11
24	TE	0.52	0.53	0.10
25	MC	0.61	0.50	0.12
26	TE	0.68	0.60	0.09
27	TE	0.50	0.37	0.13
28	MC	0.76	0.49	0.07
29	TE	0.55	0.54	0.08
30	TE	0.54	0.36	0.12
31	TE	0.44	0.39	0.11
32	MC	0.62	0.42	0.16
33	TE	0.31	0.35	0.08
34	MC	0.76	0.48	0.07
35	MC	0.45	0.33	0.15
36	MC	0.48	0.35	0.19
37	TE	0.38	0.41	0.09
38	MC	0.38	0.39	0.14
39	MC	0.68	0.41	0.12
40	MC	0.57	0.39	0.09

Table 6-21 Item Analysis, Social Studies Grade 8

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TE	0.69	0.36	0.04
2	MC	0.44	0.35	0.10
3	TE	0.50	0.48	0.14
4	MS	0.25	0.25	0.13
5	MC	0.64	0.47	0.13
6	MC	0.79	0.49	0.19
7	MC	0.42	0.29	0.15
8	MC	0.80	0.53	0.12
9	TE	0.78	0.37	0.12
10	MC	0.79	0.40	0.14
11	MC	0.51	0.40	0.15
12	MC	0.53	0.51	0.18
13	MC	0.38	0.41	0.16
14	MC	0.70	0.54	0.15
15	MC	0.74	0.46	0.18
16	MC	0.71	0.38	0.18
17	MC	0.47	0.45	0.16
18	MC	0.72	0.39	0.17
19	TE	0.52	0.61	0.16
20	MC	0.51	0.45	0.16
21	MS	0.44	0.54	0.05
22	TE	0.52	0.29	0.17
23	MC	0.36	0.43	0.19
24	MC	0.70	0.50	0.17
25	MC	0.55	0.34	0.14
26	MC	0.66	0.57	0.17
27	MC	0.73	0.40	0.11
28	TE	0.36	0.32	0.21
29	MC	0.66	0.46	0.19
30	MC	0.63	0.51	0.19
31	MC	0.47	0.39	0.20
32	MC	0.62	0.45	0.23
33	MC	0.57	0.35	0.17
34	MC	0.67	0.37	0.15
35	MS	0.52	0.52	0.17
36	MC	0.60	0.40	0.22
37	TE	0.60	0.44	0.22
38	MC	0.60	0.41	0.21
39	MC	0.67	0.42	0.15
40	TE	0.40	0.46	0.20

Table 6-22 Item Analysis, Social Studies Grade 10

Item	Item Type	<i>p</i>-Value	Item-Total Test Corr.	Percent Omit
1	TE	0.68	0.44	0.36
2	MC	0.64	0.47	0.15
3	MC	0.72	0.50	0.16
4	MC	0.60	0.39	0.21
5	MC	0.41	0.29	0.38
6	MC	0.47	0.41	0.34
7	TE	0.62	0.42	0.20
8	MS	0.38	0.45	0.16
9	MC	0.62	0.30	0.17
10	MC	0.69	0.47	0.32
11	MC	0.57	0.49	0.26
12	MC	0.57	0.44	0.24
13	MC	0.53	0.53	0.21
14	MC	0.38	0.34	0.23
15	MC	0.57	0.40	0.29
16	MC	0.35	0.38	0.32
17	MC	0.48	0.46	0.34
18	MC	0.32	0.26	0.36
19	MC	0.50	0.28	0.27
20	MS	0.30	0.30	0.25
21	MC	0.59	0.40	0.12
22	MS	0.45	0.53	0.14
23	TE	0.51	0.44	0.24
24	MC	0.65	0.40	0.27
25	MC	0.66	0.57	0.23
26	MS	0.27	0.32	0.23
27	MC	0.34	0.21	0.28
28	MC	0.61	0.55	0.36
29	MC	0.41	0.35	0.31
30	MC	0.56	0.37	0.30
31	TE	0.59	0.48	0.38
32	TE	0.43	0.42	1.01
33	MC	0.46	0.47	0.34
34	MC	0.55	0.27	0.32
35	MC	0.63	0.44	0.37
36	TE	0.37	0.49	0.40
37	MC	0.73	0.45	0.35
38	MC	0.44	0.22	0.38
39	MC	0.47	0.39	0.37
40	MC	0.59	0.49	0.35

Table 6-23 Test-Level Descriptive Statistics

Content	Grade	N Count	Mean Raw Score	Test Difficulty	Raw Score SD	Skewness	Kurtosis	Min Obtained	Max Obtained	Max Possible	Alpha	SEM
English Language Arts	3	58432	26.84	0.54	9.66	0.03	-0.94	1	52	53	0.88	3.34
	4	58914	29.01	0.56	10.22	0.03	-0.86	0	54	56	0.89	3.39
	5	59305	28.78	0.54	9.90	0.03	-0.90	1	54	56	0.88	3.48
	6	59330	29.72	0.56	10.25	-0.14	-0.87	1	55	56	0.89	3.46
	7	60307	30.04	0.56	10.10	-0.14	-0.81	1	56	56	0.88	3.53
	8	62130	30.68	0.56	10.64	-0.14	-0.82	1	56	56	0.89	3.55
Mathematics	3	58691	22.44	0.54	10.15	-0.02	-1.07	0	42	42	0.93	2.67
	4	59120	23.00	0.50	10.35	0.23	-1.00	0	46	46	0.92	2.85
	5	59532	21.39	0.47	10.41	0.16	-0.97	0	46	46	0.92	2.92
	6	59522	22.02	0.48	10.39	0.23	-0.94	0	46	46	0.92	2.88
	7	60485	20.83	0.45	9.83	0.31	-0.84	0	46	46	0.91	2.92
	8	62283	20.36	0.44	10.15	0.48	-0.74	0	46	46	0.92	2.89
Science	4	59090	21.29	0.53	8.76	-0.04	-1.05	0	40	40	0.91	2.68
	8	62181	20.63	0.52	8.44	0.03	-0.94	0	40	40	0.89	2.78
Social Studies	4	59086	21.92	0.55	8.88	-0.09	-1.03	1	40	40	0.91	2.71
	8	62190	23.19	0.58	8.96	-0.15	-1.05	0	40	40	0.91	2.69
	10	61545	20.58	0.52	8.86	0.10	-1.08	0	40	40	0.90	2.79

Table 6-24 Calibration Sample Demographics Compared to Population, English Language Arts

Grade 3	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	58368	95.33	58497	95.54	61228	100.00	0.21	4.67
Gender								
Male	29962	51.33	30027	51.33	31529	51.49	0.00	0.16
Female	28399	48.66	28463	48.66	29692	48.49	0.00	-0.16
Non-binary	7	0.01	7	0.01	7	0.01	0.00	0.00
Race/Ethnicity								
White	37838	64.83	37862	64.72	39463	64.45	-0.10	-0.37
African American	5886	10.08	5948	10.17	6244	10.20	0.08	0.11
Hispanic	8059	13.81	8088	13.83	8592	14.03	0.02	0.23
Asian/Pacific Islander	2733	4.68	2743	4.69	2854	4.66	0.01	-0.02
American Indian	553	0.95	553	0.95	579	0.95	0.00	0.00
Other	3299	5.65	3303	5.65	3496	5.71	-0.01	0.06
LEP								
No	53249	91.23	53346	91.19	55638	90.87	-0.04	-0.36
Yes	5119	8.77	5151	8.81	5590	9.13	0.04	0.36
Disability								
No	50029	85.71	50140	85.71	52053	85.02	0.00	-0.70
Yes	8339	14.29	8357	14.29	9175	14.98	0.00	0.70
SES Disadvantaged								
No	33318	57.08	33351	57.01	34771	56.79	-0.07	-0.29
Yes	25050	42.92	25146	42.99	26457	43.21	0.07	0.29

Grade 4	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	58850	95.40	58996	95.63	61689	100.00	0.24	4.60
Gender								
Male	29904	50.81	29985	50.83	31449	50.98	0.01	0.17
Female	28936	49.17	29001	49.16	30228	49.00	-0.01	-0.17
Non-binary	10	0.02	10	0.02	12	0.02	0.00	0.00
Race/Ethnicity								
White	38298	65.08	38340	64.99	39918	64.71	-0.09	-0.37
African American	5673	9.64	5734	9.72	6042	9.79	0.08	0.15
Hispanic	8263	14.04	8286	14.05	8778	14.23	0.00	0.19
Asian/Pacific Islander	2791	4.74	2803	4.75	2907	4.71	0.01	-0.03
American Indian	575	0.98	576	0.98	608	0.99	0.00	0.01
Other	3250	5.52	3257	5.52	3436	5.57	0.00	0.05
LEP								
No	53667	91.19	53787	91.17	56088	90.92	-0.02	-0.27
Yes	5183	8.81	5209	8.83	5601	9.08	0.02	0.27
Disability								
No	50692	86.14	50813	86.13	52610	85.28	-0.01	-0.86
Yes	8158	13.86	8183	13.87	9079	14.72	0.01	0.86
SES Disadvantaged								
No	34021	57.81	34066	57.74	35494	57.54	-0.07	-0.27
Yes	24829	42.19	24930	42.26	26195	42.46	0.07	0.27

Grade 5	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	59237	95.45	59386	95.69	62058	100.00	0.24	4.55
Gender								
Male	30334	51.21	30408	51.20	31886	51.38	0.00	0.17
Female	28889	48.77	28964	48.77	30158	48.60	0.00	-0.17
Non-binary	14	0.02	14	0.02	14	0.02	0.00	0.00
Race/Ethnicity								
White	38607	65.17	38648	65.08	40201	64.78	-0.09	-0.39
African American	5602	9.46	5652	9.52	5955	9.60	0.06	0.14
Hispanic	8526	14.39	8561	14.42	9056	14.59	0.02	0.20
Asian/Pacific Islander	2747	4.64	2756	4.64	2858	4.61	0.00	-0.03
American Indian	564	0.95	568	0.96	598	0.96	0.00	0.01
Other	3191	5.39	3201	5.39	3390	5.46	0.00	0.08
LEP								
No	54413	91.86	54535	91.83	56765	91.47	-0.03	-0.39
Yes	4824	8.14	4851	8.17	5293	8.53	0.03	0.39
Disability								
No	51423	86.81	51545	86.80	53353	85.97	-0.01	-0.84
Yes	7814	13.19	7841	13.20	8705	14.03	0.01	0.84
SES Disadvantaged								
No	34170	57.68	34217	57.62	35656	57.46	-0.07	-0.23
Yes	25067	42.32	25169	42.38	26402	42.54	0.07	0.23

Grade 6	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	59240	95.07	59412	95.35	62310	100.00	0.28	4.93
Gender								
Male	30213	51.00	30307	51.01	31837	51.09	0.01	0.09
Female	29013	48.98	29091	48.96	30458	48.88	-0.01	-0.09
Non-binary	14	0.02	14	0.02	15	0.02	0.00	0.00
Race/Ethnicity								
White	38974	65.79	39006	65.65	40620	65.19	-0.14	-0.60
African American	5635	9.51	5703	9.60	6125	9.83	0.09	0.32
Hispanic	8486	14.32	8527	14.35	9051	14.53	0.03	0.20
Asian/Pacific Islander	2619	4.42	2626	4.42	2715	4.36	0.00	-0.06
American Indian	565	0.95	579	0.97	622	1.00	0.02	0.04
Other	2961	5.00	2971	5.00	3177	5.10	0.00	0.10
LEP								
No	55213	93.20	55353	93.17	57868	92.87	-0.03	-0.33
Yes	4027	6.80	4059	6.83	4442	7.13	0.03	0.33
Disability								
No	51737	87.33	51870	87.31	53805	86.35	-0.03	-0.98
Yes	7503	12.67	7542	12.69	8505	13.65	0.03	0.98
SES Disadvantaged								
No	34794	58.73	34836	58.63	36265	58.20	-0.10	-0.53
Yes	24446	41.27	24576	41.37	26045	41.80	0.10	0.53

Grade 7	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	60181	94.59	60413	94.95	63623	100.00	0.36	5.41
Gender								
Male	30821	51.21	30933	51.20	32597	51.23	-0.01	0.02
Female	29332	48.74	29452	48.75	30995	48.72	0.01	-0.02
Non-binary	28	0.05	28	0.05	31	0.05	0.00	0.00
Race/Ethnicity								
White	39777	66.10	39840	65.95	41683	65.52	-0.15	-0.58
African American	5741	9.54	5842	9.67	6300	9.90	0.13	0.36
Hispanic	8682	14.43	8729	14.45	9271	14.57	0.02	0.15
Asian/Pacific Islander	2506	4.16	2509	4.15	2603	4.09	-0.01	-0.07
American Indian	588	0.98	597	0.99	643	1.01	0.01	0.03
Other	2887	4.80	2896	4.79	3123	4.91	0.00	0.11
LEP								
No	56081	93.19	56284	93.17	59126	92.93	-0.02	-0.26
Yes	4100	6.81	4129	6.83	4497	7.07	0.02	0.26
Disability								
No	52767	87.68	52950	87.65	55121	86.64	-0.03	-1.04
Yes	7414	12.32	7463	12.35	8502	13.36	0.03	1.04
SES Disadvantaged								
No	35514	59.01	35577	58.89	37200	58.47	-0.12	-0.54
Yes	24667	40.99	24836	41.11	26423	41.53	0.12	0.54

Grade 8	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	61992	93.84	62249	94.23	66060	100.00	0.39	6.16
Gender								
Male	31905	51.47	32033	51.46	33938	51.37	-0.01	-0.09
Female	30036	48.45	30165	48.46	32065	48.54	0.01	0.09
Non-binary	51	0.08	51	0.08	57	0.09	0.00	0.00
Race/Ethnicity								
White	40893	65.96	40964	65.81	43159	65.33	-0.16	-0.63
African American	6069	9.79	6171	9.91	6730	10.19	0.12	0.40
Hispanic	8861	14.29	8917	14.32	9529	14.42	0.03	0.13
Asian/Pacific Islander	2685	4.33	2687	4.32	2797	4.23	-0.01	-0.10
American Indian	592	0.95	606	0.97	667	1.01	0.02	0.05
Other	2892	4.67	2904	4.67	3178	4.81	0.00	0.15
LEP								
No	58132	93.77	58365	93.76	61796	93.55	-0.01	-0.23
Yes	3860	6.23	3884	6.24	4264	6.45	0.01	0.23
Disability								
No	54527	87.96	54730	87.92	57324	86.78	-0.04	-1.18
Yes	7465	12.04	7519	12.08	8736	13.22	0.04	1.18
SES Disadvantaged								
No	36946	59.60	37022	59.47	38944	58.95	-0.12	-0.65
Yes	25046	40.40	25227	40.53	27116	41.05	0.12	0.65

Table 6-25 Calibration Sample Demographics Compared to Population, Mathematics

Grade 3	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	58643	95.78	58722	95.91	61228	100.00	0.13	4.22
Gender								
Male	30112	51.35	30150	51.34	31529	51.49	0.00	0.15
Female	28524	48.64	28565	48.64	29692	48.49	0.00	-0.15
Non-binary	7	0.01	7	0.01	7	0.01	0.00	0.00
Race/Ethnicity								
White	37868	64.57	37895	64.53	39463	64.45	-0.04	-0.12
African American	5899	10.06	5932	10.10	6244	10.20	0.04	0.14
Hispanic	8252	14.07	8259	14.06	8592	14.03	-0.01	-0.04
Asian/Pacific Islander	2778	4.74	2782	4.74	2854	4.66	0.00	-0.08
American Indian	552	0.94	553	0.94	579	0.95	0.00	0.00
Other	3294	5.62	3301	5.62	3496	5.71	0.00	0.09
LEP								
No	53274	90.84	53345	90.84	55638	90.87	0.00	0.03
Yes	5369	9.16	5377	9.16	5590	9.13	0.00	-0.03
Disability								
No	50299	85.77	50365	85.77	52053	85.02	0.00	-0.76
Yes	8344	14.23	8357	14.23	9175	14.98	0.00	0.76
SES Disadvantaged								
No	33442	57.03	33466	56.99	34771	56.79	-0.04	-0.24
Yes	25201	42.97	25256	43.01	26457	43.21	0.04	0.24

Grade 4	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	59088	95.78	59165	95.91	61689	100.00	0.12	4.22
Gender								
Male	30018	50.80	30064	50.81	31449	50.98	0.01	0.18
Female	29060	49.18	29091	49.17	30228	49.00	-0.01	-0.18
Non-binary	10	0.02	10	0.02	12	0.02	0.00	0.00
Race/Ethnicity								
White	38326	64.86	38345	64.81	39918	64.71	-0.05	-0.15
African American	5693	9.63	5724	9.67	6042	9.79	0.04	0.16
Hispanic	8409	14.23	8427	14.24	8778	14.23	0.01	0.00
Asian/Pacific Islander	2837	4.80	2839	4.80	2907	4.71	0.00	-0.09
American Indian	575	0.97	576	0.97	608	0.99	0.00	0.01
Other	3248	5.50	3254	5.50	3436	5.57	0.00	0.07
LEP								
No	53701	90.88	53766	90.87	56088	90.92	-0.01	0.04
Yes	5387	9.12	5399	9.13	5601	9.08	0.01	-0.04
Disability								
No	50933	86.20	50989	86.18	52610	85.28	-0.02	-0.92
Yes	8155	13.80	8176	13.82	9079	14.72	0.02	0.92
SES Disadvantaged								
No	34127	57.76	34147	57.71	35494	57.54	-0.04	-0.22
Yes	24961	42.24	25018	42.29	26195	42.46	0.04	0.22

Grade 5	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	59503	95.88	59577	96.00	62058	100.00	0.12	4.12
Gender								
Male	30456	51.18	30495	51.19	31886	51.38	0.00	0.20
Female	29033	48.79	29068	48.79	30158	48.60	0.00	-0.20
Non-binary	14	0.02	14	0.02	14	0.02	0.00	0.00
Race/Ethnicity								
White	38642	64.94	38666	64.90	40201	64.78	-0.04	-0.16
African American	5621	9.45	5646	9.48	5955	9.60	0.03	0.15
Hispanic	8687	14.60	8703	14.61	9056	14.59	0.01	-0.01
Asian/Pacific Islander	2791	4.69	2793	4.69	2858	4.61	0.00	-0.09
American Indian	564	0.95	565	0.95	598	0.96	0.00	0.02
Other	3198	5.37	3204	5.38	3390	5.46	0.00	0.09
LEP								
No	54453	91.51	54518	91.51	56765	91.47	0.00	-0.04
Yes	5050	8.49	5059	8.49	5293	8.53	0.00	0.04
Disability								
No	51686	86.86	51743	86.85	53353	85.97	-0.01	-0.89
Yes	7817	13.14	7834	13.15	8705	14.03	0.01	0.89
SES Disadvantaged								
No	34284	57.62	34309	57.59	35656	57.46	-0.03	-0.16
Yes	25219	42.38	25268	42.41	26402	42.54	0.03	0.16

Grade 6	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	59481	95.46	59570	95.60	62310	100.00	0.14	4.54
Gender								
Male	30337	51.00	30382	51.00	31837	51.09	0.00	0.09
Female	29130	48.97	29174	48.97	30458	48.88	0.00	-0.09
Non-binary	14	0.02	14	0.02	15	0.02	0.00	0.00
Race/Ethnicity								
White	38995	65.56	39012	65.49	40620	65.19	-0.07	-0.37
African American	5652	9.50	5688	9.55	6125	9.83	0.05	0.33
Hispanic	8657	14.55	8675	14.56	9051	14.53	0.01	-0.03
Asian/Pacific Islander	2641	4.44	2646	4.44	2715	4.36	0.00	-0.08
American Indian	572	0.96	579	0.97	622	1.00	0.01	0.04
Other	2964	4.98	2970	4.99	3177	5.10	0.00	0.12
LEP								
No	55259	92.90	55339	92.90	57868	92.87	0.00	-0.03
Yes	4222	7.10	4231	7.10	4442	7.13	0.00	0.03
Disability								
No	51963	87.36	52031	87.34	53805	86.35	-0.02	-1.01
Yes	7518	12.64	7539	12.66	8505	13.65	0.02	1.01
SES Disadvantaged								
No	34894	58.66	34911	58.61	36265	58.20	-0.06	-0.46
Yes	24587	41.34	24659	41.39	26045	41.80	0.06	0.46

Grade 7	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	60413	94.95	60559	95.18	63623	100.00	0.23	5.05
Gender								
Male	30937	51.21	31017	51.22	32597	51.23	0.01	0.03
Female	29449	48.75	29514	48.74	30995	48.72	-0.01	-0.03
Non-binary	27	0.04	28	0.05	31	0.05	0.00	0.00
Race/Ethnicity								
White	39807	65.89	39852	65.81	41683	65.52	-0.08	-0.38
African American	5766	9.54	5829	9.63	6300	9.90	0.08	0.36
Hispanic	8835	14.62	8857	14.63	9271	14.57	0.00	-0.05
Asian/Pacific Islander	2528	4.18	2533	4.18	2603	4.09	0.00	-0.09
American Indian	594	0.98	597	0.99	643	1.01	0.00	0.03
Other	2883	4.77	2891	4.77	3123	4.91	0.00	0.14
LEP								
No	56131	92.91	56262	92.90	59126	92.93	-0.01	0.02
Yes	4282	7.09	4297	7.10	4497	7.07	0.01	-0.02
Disability								
No	52995	87.72	53103	87.69	55121	86.64	-0.03	-1.08
Yes	7418	12.28	7456	12.31	8502	13.36	0.03	1.08
SES Disadvantaged								
No	35607	58.94	35639	58.85	37200	58.47	-0.09	-0.47
Yes	24806	41.06	24920	41.15	26423	41.53	0.09	0.47

Grade 8	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	62233	94.21	62360	94.40	66060	100.00	0.19	5.79
Gender								
Male	32027	51.46	32094	51.47	33938	51.37	0.00	-0.09
Female	30156	48.46	30215	48.45	32065	48.54	0.00	0.08
Non-binary	50	0.08	51	0.08	57	0.09	0.00	0.01
Race/Ethnicity								
White	40934	65.78	40964	65.69	43159	65.33	-0.09	-0.44
African American	6088	9.78	6153	9.87	6730	10.19	0.08	0.41
Hispanic	9002	14.46	9024	14.47	9529	14.42	0.01	-0.04
Asian/Pacific Islander	2706	4.35	2707	4.34	2797	4.23	-0.01	-0.11
American Indian	604	0.97	606	0.97	667	1.01	0.00	0.04
Other	2899	4.66	2906	4.66	3178	4.81	0.00	0.15
LEP								
No	58222	93.55	58341	93.56	61796	93.55	0.00	-0.01
Yes	4011	6.45	4019	6.44	4264	6.45	0.00	0.01
Disability								
No	54748	87.97	54849	87.96	57324	86.78	-0.02	-1.20
Yes	7485	12.03	7511	12.04	8736	13.22	0.02	1.20
SES Disadvantaged								
No	37037	59.51	37072	59.45	38944	58.95	-0.07	-0.56
Yes	25196	40.49	25288	40.55	27116	41.05	0.07	0.56

Table 6-26 Calibration Sample Demographics Compared to Population, Science

Grade 4	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	59040	95.71	59141	95.87	61689	100.00	0.16	4.29
Gender								
Male	29998	50.81	30053	50.82	31449	50.98	0.01	0.17
Female	29032	49.17	29078	49.17	30228	49.00	-0.01	-0.17
Non-binary	10	0.02	10	0.02	12	0.02	0.00	0.00
Race/Ethnicity								
White	38316	64.90	38347	64.84	39918	64.71	-0.06	-0.19
African American	5670	9.60	5717	9.67	6042	9.79	0.06	0.19
Hispanic	8404	14.23	8413	14.23	8778	14.23	-0.01	0.00
Asian/Pacific Islander	2834	4.80	2838	4.80	2907	4.71	0.00	-0.09
American Indian	573	0.97	575	0.97	608	0.99	0.00	0.02
Other	3243	5.49	3251	5.50	3436	5.57	0.00	0.08
LEP								
No	53652	90.87	53745	90.88	56088	90.92	0.00	0.05
Yes	5388	9.13	5396	9.12	5601	9.08	0.00	-0.05
Disability								
No	50890	86.20	50970	86.18	52610	85.28	-0.01	-0.91
Yes	8150	13.80	8171	13.82	9079	14.72	0.01	0.91
SES Disadvantaged								
No	34107	57.77	34137	57.72	35494	57.54	-0.05	-0.23
Yes	24933	42.23	25004	42.28	26195	42.46	0.05	0.23

Grade 8	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	62051	93.93	62289	94.29	66060	100.00	0.36	6.07
Gender								
Male	31929	51.46	32051	51.46	33938	51.37	0.00	-0.08
Female	30071	48.46	30187	48.46	32065	48.54	0.00	0.08
Non-binary	51	0.08	51	0.08	57	0.09	0.00	0.00
Race/Ethnicity								
White	40880	65.88	40941	65.73	43159	65.33	-0.15	-0.55
African American	6035	9.73	6136	9.85	6730	10.19	0.12	0.46
Hispanic	8963	14.44	9005	14.46	9529	14.42	0.01	-0.02
Asian/Pacific Islander	2703	4.36	2705	4.34	2797	4.23	-0.01	-0.12
American Indian	595	0.96	604	0.97	667	1.01	0.01	0.05
Other	2875	4.63	2898	4.65	3178	4.81	0.02	0.18
LEP								
No	58058	93.56	58277	93.56	61796	93.55	-0.01	-0.02
Yes	3993	6.44	4012	6.44	4264	6.45	0.01	0.02
Disability								
No	54593	87.98	54784	87.95	57324	86.78	-0.03	-1.21
Yes	7458	12.02	7505	12.05	8736	13.22	0.03	1.21
SES Disadvantaged								
No	36962	59.57	37033	59.45	38944	58.95	-0.11	-0.61
Yes	25089	40.43	25256	40.55	27116	41.05	0.11	0.61

Table 6-27 Calibration Sample Demographics Compared to Population, Social Studies

Grade 4	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	59050	95.72	59131	95.85	61689	100.00	0.13	4.28
Gender								
Male	30000	50.80	30047	50.81	31449	50.98	0.01	0.18
Female	29041	49.18	29074	49.17	30228	49.00	-0.01	-0.18
Non-binary	9	0.02	10	0.02	12	0.02	0.00	0.00
Race/Ethnicity								
White	38319	64.89	38340	64.84	39918	64.71	-0.05	-0.18
African American	5668	9.60	5702	9.64	6042	9.79	0.04	0.20
Hispanic	8408	14.24	8422	14.24	8778	14.23	0.00	-0.01
Asian/Pacific Islander	2836	4.80	2840	4.80	2907	4.71	0.00	-0.09
American Indian	575	0.97	576	0.97	608	0.99	0.00	0.01
Other	3244	5.49	3251	5.50	3436	5.57	0.00	0.08
LEP								
No	53660	90.87	53730	90.87	56088	90.92	-0.01	0.05
Yes	5390	9.13	5401	9.13	5601	9.08	0.01	-0.05
Disability								
No	50905	86.21	50965	86.19	52610	85.28	-0.02	-0.92
Yes	8145	13.79	8166	13.81	9079	14.72	0.02	0.92
SES Disadvantaged								
No	34117	57.78	34138	57.73	35494	57.54	-0.04	-0.24
Yes	24933	42.22	24993	42.27	26195	42.46	0.04	0.24

Grade 8	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	62110	94.02	62261	94.25	66060	100.00	0.23	5.98
Gender								
Male	31939	51.42	32021	51.43	33938	51.37	0.01	-0.05
Female	30120	48.49	30189	48.49	32065	48.54	-0.01	0.04
Non-binary	51	0.08	51	0.08	57	0.09	0.00	0.00
Race/Ethnicity								
White	40895	65.84	40943	65.76	43159	65.33	-0.08	-0.51
African American	6040	9.72	6106	9.81	6730	10.19	0.08	0.46
Hispanic	8984	14.46	9008	14.47	9529	14.42	0.00	-0.04
Asian/Pacific Islander	2707	4.36	2708	4.35	2797	4.23	-0.01	-0.12
American Indian	598	0.96	600	0.96	667	1.01	0.00	0.05
Other	2886	4.65	2896	4.65	3178	4.81	0.00	0.16
LEP								
No	58103	93.55	58250	93.56	61796	93.55	0.01	0.00
Yes	4007	6.45	4011	6.44	4264	6.45	-0.01	0.00
Disability								
No	54652	87.99	54777	87.98	57324	86.78	-0.01	-1.22
Yes	7458	12.01	7484	12.02	8736	13.22	0.01	1.22
SES Disadvantaged								
No	37013	59.59	37058	59.52	38944	58.95	-0.07	-0.64
Yes	25097	40.41	25203	40.48	27116	41.05	0.07	0.64

Grade 10	Calibration Sample (1)		Students with Valid Scores (2)		Scheduled to Test (3)		Difference (2) – (1)	Difference (3) – (1)
	N	%	N	%	N	%	%	%
All Students	60960	87.24	61819	88.47	69876	100.00	1.23	12.76
Gender								
Male	31308	51.36	31776	51.40	35825	51.27	0.04	-0.09
Female	29592	48.54	29983	48.50	33982	48.63	-0.04	0.09
Non-binary	60	0.10	60	0.10	69	0.10	0.00	0.00
Race/Ethnicity								
White	42850	70.29	43110	69.74	46815	67.00	-0.56	-3.29
African American	4320	7.09	4602	7.44	6703	9.59	0.36	2.51
Hispanic	8307	13.63	8523	13.79	9854	14.10	0.16	0.48
Asian/Pacific Islander	2470	4.05	2510	4.06	2734	3.91	0.01	-0.14
American Indian	517	0.85	526	0.85	704	1.01	0.00	0.16
Other	2496	4.09	2548	4.12	3066	4.39	0.03	0.29
LEP								
No	57719	94.68	58453	94.56	65879	94.28	-0.13	-0.40
Yes	3241	5.32	3366	5.44	3997	5.72	0.13	0.40
Disability								
No	54605	89.58	55289	89.44	61296	87.72	-0.14	-1.85
Yes	6355	10.42	6530	10.56	8580	12.28	0.14	1.85
SES Disadvantaged								
No	39632	65.01	39909	64.56	43440	62.17	-0.46	-2.85
Yes	21328	34.99	21910	35.44	26436	37.83	0.46	2.85

Table 6-28 Items Flagged Based on Yen's Q1

Content	Grade	Item Number in Calibration	Model	N	Z	Critical Z
ELA	3	30*	2PPC	58368	294.09	155.65
	3	33	2PPC	58368	259.66	155.65
	4	17	2PPC	58849	162.14	156.93
	4	27	2PPC	58849	312.43	156.93
	5	13	2PPC	59237	303.16	157.97
	5	22	2PPC	59237	183.28	157.97
	5	39	2PPC	59237	708.12	157.97
	6	6	2PPC	59240	308.88	157.97
	6	9	2PPC	59240	567.79	157.97
	6	15	2PPC	59240	770.28	157.97
	7	13	2PPC	60180	591.73	160.48
	7	18*	2PPC	60180	378.61	160.48
	7	23	2PPC	60180	817.64	160.48
	7	25	2PPC	60180	319.8	160.48
	8	29	2PPC	61986	181.44	165.30
	8	34*	2PPC	61986	916.02	165.30
Math	3	10	2PPC	58378	163.93	155.67
	4	13*	2PPC	59000	253.23	157.33
	4	26	3PL	59000	185.66	157.33
	5	13	2PPC	59456	201.68	158.55
	5	22	2PPC	59456	219.34	158.55
	5	27	2PPC	59456	200.83	158.55
	6	24	3PL	59389	301.07	158.37
	6	39	3PL	59389	217.46	158.37
	6	44*	2PPC	59389	367.73	158.37
	7	22	2PPC	60376	180.92	161.00
	7	42*	2PPC	60376	377.21	161.00
	8	3	2PPC	62129	198.88	165.68
8	17	3PL	62129	302.71	165.68	
Science	4	12	2PPC	59004	196.23	157.34
	8	1	2PPC	62020	166.13	165.39
	8	4	2PPC	62020	193.53	165.39
	8	8	3PL	62020	178.74	165.39
	8	38*	2PPC	62020	363.31	165.39
Social Studies	4	3*	2PPC	58960	211.92	157.23
	8	1	2PPC	61943	357.36	165.18
	8	28	2PPC	61943	188.84	165.18
	8	35	2PPC	61943	185.84	165.18

Note: An asterisk (*) indicates an anchor item.

Table 6-29 Equating Evaluation Results, Stocking and Lord Method

Content Area	Grade	Number of Anchors	Stocking and Lord TCC Method Results						Equating Constants	
			TCC Results		Parameter Comparison Statistics					
			# of Iterations	F Value	a-Parameter		b-Parameter			
					Corr	# of RMSD Outliers	Corr	# of RMSD Outliers	A	B
ELA	3	17	3	0.1174	0.97	1	0.98	0	0.9838	-1.3181
	4	17	4	0.3405	0.99	2	0.96	1	1.0841	-0.5435
	5	17	4	0.0412	0.96	0	0.99	1	1.0282	-0.3167
	6	16	5	0.1353	0.94	0	0.99	1	1.0752	-0.0630
	7	18	6	0.1558	0.99	1	0.99	1	1.1183	0.4022
	8	16	3	0.2941	0.99	1	0.97	0	1.3086	0.4352
Math	3	17	3	0.0339	0.97	0	0.99	0	1.0387	-1.2366
	4	20	9	0.0830	0.98	0	1.00	1	1.0487	-0.7264
	5	18	10	0.0739	0.98	1	0.99	1	0.9498	-0.1520
	6	19	18	0.1164	0.99	1	0.99	1	1.0710	0.0093
	7	18	11	0.1091	0.95	0	1.00	1	1.1065	0.2473
	8	20	13	0.2398	0.98	0	0.99	1	1.0296	0.7118
Science	4	14	3	0.0562	0.99	1	1.00	0	1.0900	-0.0924
	8	15	3	0.0307	0.98	0	1.00	0	1.0248	-0.1427
Social Studies	4	19	6	0.0260	0.99	1	1	0	1.0288	-0.0635
	8	20	4	0.0432	0.95	1	1	1	1.0499	-0.1179
	10	18	13	0.0438	0.98	0	0.99	1	1.0447	-0.0924

*Equating results obtained in test equating with a reduced anchor set (final)

Table 6-30 Scale Transformation Constants

Content Area	Grade	Scale Transformation Constants	
		M1	M2
ELA	3–8	43.7445	610.4987
Mathematics	3–8	46.4684	612.0818
Science	4	45.0450	500.4505
	8	45.0450	699.5496
Social Studies	4	42.9923	503.8693
	8	40.2577	708.0515
	10	42.0521	803.7847

Table 6-31 Scoring Table for English Language Arts Grade 3

Raw Score	Scale Score	SEM
0	330	97
1	330	97
2	330	97
3	330	97
4	330	97
5	330	97
6	332	96
7	396	50
8	425	37
9	444	31
10	457	27
11	468	23
12	477	21
13	485	19
14	492	18
15	499	17
16	504	16
17	510	15
18	515	15
19	520	14
*20	*524	*14
21	529	14
22	533	13
23	537	13
24	542	13
25	546	13
26	550	13

Raw Score	Scale Score	SEM
27	554	13
28	558	13
29	562	13
30	566	13
*31	*570	*13
32	574	13
33	579	13
34	583	13
35	587	13
36	592	13
37	597	14
38	602	14
39	608	15
40	614	15
41	621	16
*42	*628	*17
43	636	18
44	646	20
45	657	22
46	670	24
47	686	27
48	705	31
49	731	35
50	763	39
51	805	47
52	899	111
53	900	112

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-32 Scoring Table for English Language Arts Grade 4

Raw Score	Scale Score	SEM
0	340	85
1	340	85
2	340	85
3	340	85
4	340	85
5	340	85
6	355	74
7	406	47
8	433	37
9	452	31
10	467	27
11	479	25
12	489	23
13	498	22
14	507	20
15	514	19
16	521	18
17	527	18
18	533	17
19	538	16
20	544	16
*21	*549	*15
22	554	15
23	558	14
24	563	14
25	567	14
26	572	14
27	576	13
28	580	13

Raw Score	Scale Score	SEM
29	585	13
30	589	13
*31	*593	*13
32	597	13
33	602	13
34	606	13
35	610	14
36	615	14
37	620	14
38	625	14
39	630	15
40	635	15
41	640	15
42	646	16
*43	*652	*17
44	658	17
45	666	18
46	673	19
47	682	21
48	692	22
49	703	25
50	717	28
51	733	31
52	755	36
53	783	40
54	818	42
55	864	52
56	930	95

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-33 Scoring Table for English Language Arts Grade 5

Raw Score	Scale Score	SEM
0	350	79
1	350	79
2	350	79
3	350	79
4	350	79
5	350	79
6	350	79
7	393	53
8	428	40
9	450	34
10	468	30
11	482	28
12	494	26
13	505	24
14	514	22
15	522	21
16	530	20
17	537	19
18	543	18
19	549	17
20	555	16
21	560	16
*22	*565	*15
23	570	15
24	574	14
25	579	14
26	583	14
27	587	13
28	591	13

Raw Score	Scale Score	SEM
29	596	13
30	600	13
31	604	13
32	608	13
*33	*612	*13
34	616	13
35	621	13
36	625	13
37	630	14
38	635	14
39	640	15
40	645	15
41	651	16
42	657	16
43	663	17
*44	*670	*18
45	677	18
46	685	19
47	694	21
48	704	22
49	716	24
50	730	27
51	747	31
52	769	35
53	798	40
54	835	44
55	886	58
56	940	97

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-34 Scoring Table for English Language Arts Grade 6

Raw Score	Scale Score	SEM
0	360	78
1	360	78
2	360	78
3	360	78
4	360	78
5	360	78
6	368	73
7	422	44
8	450	36
9	470	31
10	485	28
11	498	26
12	508	24
13	518	22
14	526	21
15	533	20
16	540	18
17	546	17
18	552	16
19	557	16
20	562	15
21	567	15
*22	*572	*14
23	576	14
24	580	14
25	585	14
26	589	14
27	593	14
28	597	13

Raw Score	Scale Score	SEM
29	602	13
30	606	13
31	610	13
32	614	13
33	619	13
*34	*623	*13
35	627	13
36	632	13
37	637	14
38	642	14
39	647	14
40	652	15
41	658	15
42	663	16
43	670	17
*44	*677	*18
45	684	19
46	692	20
47	701	21
48	711	22
49	722	24
50	734	26
51	748	28
52	766	32
53	787	37
54	817	47
55	870	74
56	950	138

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-35 Scoring Table for English Language Arts Grade 7

Raw Score	Scale Score	SEM
0	370	76
1	370	76
2	370	76
3	370	76
4	370	76
5	377	72
6	425	48
7	452	38
8	472	33
9	488	29
10	500	27
11	511	25
12	521	23
13	530	21
14	537	20
15	545	19
16	551	18
17	557	18
18	563	17
19	569	16
20	574	16
21	579	16
22	584	15
*23	*589	*15
24	594	15
25	598	15
26	603	15
27	607	15
28	612	15

Raw Score	Scale Score	SEM
29	617	15
30	621	15
31	626	15
32	631	15
33	636	15
*34	*641	*15
35	646	16
36	651	16
37	656	16
38	661	16
39	667	16
40	673	16
41	678	17
42	685	17
43	691	17
*44	*698	*18
45	705	18
46	713	19
47	722	20
48	731	22
49	742	24
50	755	26
51	770	30
52	790	35
53	817	44
54	858	62
55	938	105
56	960	118

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-36 Scoring Table for English Language Arts Grade 8

Raw Score	Scale Score	SEM
0	380	70
1	380	70
2	380	70
3	380	70
4	380	70
5	380	70
6	380	70
7	405	58
8	441	44
9	464	37
10	482	33
11	497	30
12	510	28
13	521	26
14	530	24
15	539	23
16	547	21
17	554	20
18	561	19
19	567	18
20	573	17
21	579	17
22	584	17
23	590	17
*24	*595	*16
25	600	16
26	606	16
27	611	16
28	616	16

Raw Score	Scale Score	SEM
29	621	16
30	626	16
31	631	16
32	636	16
33	641	16
34	646	16
35	651	16
*36	*656	*16
37	661	16
38	666	16
39	672	16
40	678	17
41	683	17
42	690	17
43	696	18
44	703	18
*45	*710	*19
46	718	20
47	727	21
48	736	23
49	747	24
50	759	27
51	774	30
52	792	34
53	816	41
54	851	55
55	917	91
56	970	133

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-37 Scoring Table for Mathematics Grade 3

Raw Score	Scale Score	SEM
0	360	104
1	360	104
2	360	104
3	360	104
4	360	104
5	429	40
6	454	28
7	469	22
8	481	19
9	490	17
10	497	16
11	504	15
12	510	14
13	516	13
*14	*521	*13
15	525	12
16	530	12
17	534	12
18	538	11
19	542	11
20	546	11
21	550	11

Raw Score	Scale Score	SEM
22	554	11
23	558	11
*24	*562	*11
25	566	11
26	569	11
27	573	11
28	577	11
29	581	11
30	586	11
31	590	11
32	594	11
33	599	12
34	604	12
35	610	13
*36	*616	*13
37	623	14
38	632	16
39	642	18
40	657	23
41	684	36
42	760	104

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-38 Scoring Table for Mathematics Grade 4

Raw Score	Scale Score	SEM
0	405	78
1	405	78
2	405	78
3	405	78
4	405	78
5	405	78
6	405	78
7	459	35
8	482	26
9	497	22
10	509	19
11	519	18
12	527	17
13	534	16
*14	*541	*15
15	547	15
16	553	14
17	558	13
18	563	13
19	568	12
20	572	12
21	577	12
22	581	11
23	585	11

Raw Score	Scale Score	SEM
*24	*589	*11
25	592	11
26	596	10
27	600	10
28	603	10
29	607	10
30	610	10
31	614	10
32	618	10
33	621	10
34	625	10
35	629	10
*36	*633	*10
37	637	11
38	642	11
39	647	12
40	653	13
41	659	14
42	667	15
43	677	18
44	691	22
45	716	33
46	800	107

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-39 Scoring Table for Mathematics Grade 5

Raw Score	Scale Score	SEM
0	430	100
1	430	100
2	430	100
3	430	100
4	430	100
5	461	69
6	509	32
7	528	23
8	541	19
9	551	17
10	558	15
11	565	14
12	570	13
*13	*576	*12
14	580	12
15	585	11
16	589	11
17	593	11
18	596	11
19	600	10
20	604	10
21	607	10
22	610	10
*23	*614	*10

Raw Score	Scale Score	SEM
24	617	10
25	620	10
26	624	10
27	627	10
28	631	10
29	634	10
30	637	10
31	641	10
32	645	10
33	649	11
34	653	11
35	657	11
*36	*661	*11
37	666	12
38	672	13
39	677	13
40	684	14
41	692	16
42	701	18
43	713	21
44	731	27
45	766	46
46	830	102

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-40 Scoring Table for Mathematics Grade 6

Raw Score	Scale Score	SEM
0	440	81
1	440	81
2	440	81
3	440	81
4	440	81
5	440	81
6	490	40
7	516	28
8	532	23
9	544	20
10	553	18
11	562	16
12	569	15
13	575	15
14	581	14
*15	*586	*13
16	591	13
17	596	12
18	601	12
19	605	12
20	609	12
21	613	11
22	617	11
23	621	11

Raw Score	Scale Score	SEM
24	625	11
*25	*629	*11
26	633	11
27	637	11
28	641	11
29	644	11
30	648	11
31	652	11
32	656	11
33	660	11
34	665	11
35	669	11
36	674	12
37	679	12
38	684	12
*39	*690	*13
40	696	13
41	703	14
42	712	16
43	723	19
44	740	25
45	774	45
46	870	137

Note: **Bold** and *represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-41 Scoring Table for Mathematics Grade 7

Raw Score	Scale Score	SEM
0	450	92
1	450	92
2	450	92
3	450	92
4	450	92
5	450	92
6	450	92
7	511	44
8	538	30
9	554	24
10	566	20
11	575	18
12	583	16
13	590	15
14	596	14
15	602	14
*16	*607	*13
17	612	13
18	617	12
19	622	12
20	626	12
21	631	12
22	635	12
23	639	11

Raw Score	Scale Score	SEM
24	643	11
*25	*647	*11
26	652	11
27	656	11
28	660	11
29	664	12
30	668	12
31	673	12
32	677	12
33	682	12
34	687	13
35	692	13
36	698	14
37	704	14
38	710	15
*39	*717	*16
40	726	17
41	735	19
42	747	22
43	762	26
44	784	33
45	824	53
46	880	97

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-42 Scoring Table for Mathematics Grade 8

Raw Score	Scale Score	SEM
0	470	101
1	470	101
2	470	101
3	470	101
4	470	101
5	470	101
6	487	84
7	539	39
8	561	28
9	576	23
10	587	20
11	597	18
12	605	17
13	612	16
14	618	15
*15	*624	*14
16	629	13
17	634	13
18	639	12
19	643	12
20	647	11
21	651	11
22	655	11
23	659	10

Raw Score	Scale Score	SEM
24	663	10
25	666	10
*26	*670	*10
27	673	10
28	677	10
29	680	10
30	684	10
31	687	10
32	691	10
33	695	10
34	699	11
35	703	11
36	708	11
37	712	12
*38	*718	*13
39	723	13
40	730	14
41	737	16
42	747	18
43	758	21
44	775	26
45	805	40
46	890	113

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-43 Scoring Table for Science Grade 4

Raw Score	Scale Score	SEM
0	300	62
1	300	62
2	300	62
3	317	52
4	352	36
5	373	29
6	389	26
7	402	23
8	413	22
9	423	20
10	431	19
11	439	18
12	446	17
*13	*453	*16
14	459	16
15	465	15
16	470	14
17	475	14
18	480	14
19	485	13
20	490	13

Raw Score	Scale Score	SEM
21	495	13
*22	*499	*13
23	504	13
24	509	13
25	514	13
26	519	13
27	524	14
28	529	14
29	535	15
30	541	15
*31	*548	*16
32	556	17
33	564	19
34	574	21
35	586	23
36	600	27
37	619	32
38	645	40
39	693	63
40	725	83

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-44 Scoring Table for Science Grade 8

Raw Score	Scale Score	SEM
0	480	80
1	480	80
2	480	80
3	497	67
4	549	39
5	574	31
6	591	26
7	604	23
8	615	21
9	624	20
10	633	18
11	640	18
12	647	17
*13	*653	*16
14	659	16
15	665	15
16	670	15
17	676	14
18	681	14
19	686	14
20	691	14

Raw Score	Scale Score	SEM
*21	*695	*14
22	700	14
23	705	14
24	710	14
25	715	14
26	720	14
27	725	14
28	731	15
*29	*737	*15
30	743	16
31	749	16
32	757	17
33	765	18
34	774	20
35	785	22
36	798	25
37	815	30
38	839	37
39	881	56
40	945	100

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-45 Scoring Table for Social Studies Grade 4

Raw Score	Scale Score	SEM
0	330	91
1	330	91
2	330	91
3	330	91
4	330	91
5	347	74
6	390	39
7	409	28
8	422	23
9	432	20
10	440	17
11	447	16
12	454	15
13	459	14
*14	*464	*13
15	469	13
16	474	13
17	478	12
18	483	12
19	487	12
*20	*492	*12

Raw Score	Scale Score	SEM
21	496	12
22	500	12
23	505	12
24	509	12
25	513	12
26	518	12
27	523	12
28	527	12
29	532	13
*30	*538	*13
31	543	13
32	549	14
33	556	15
34	563	16
35	572	17
36	582	20
37	596	24
38	616	31
39	653	49
40	700	83

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-46 Scoring Table for Social Studies Grade 8

Raw Score	Scale Score	SEM
0	540	79
1	540	79
2	540	79
3	540	79
4	540	79
5	540	79
6	550	70
7	591	40
8	611	30
9	625	24
10	635	20
11	643	18
12	650	16
13	656	15
14	661	14
*15	*667	*13
16	671	12
17	676	12
18	680	12
19	684	11
20	688	11

Raw Score	Scale Score	SEM
21	692	11
*22	*696	*11
23	700	11
24	704	11
25	708	11
26	712	11
27	716	11
28	720	11
29	725	11
30	729	11
*31	*734	*12
32	739	12
33	745	13
34	751	14
35	759	15
36	768	17
37	779	20
38	796	26
39	829	43
40	860	65

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-47 Scoring Table for Social Studies Grade 10

Raw Score	Scale Score	SEM
0	645	91
1	645	91
2	645	91
3	645	91
4	645	91
5	645	91
6	645	91
7	691	49
8	718	33
9	733	26
10	744	22
11	753	19
12	760	17
13	767	15
*14	*773	*14
15	778	14
16	783	13
17	787	13
18	792	12
19	796	12
20	800	12

Raw Score	Scale Score	SEM
21	804	11
*22	*808	*11
23	812	11
24	816	11
25	820	11
26	825	11
27	829	12
28	833	12
*29	*838	*12
30	843	12
31	848	13
32	854	13
33	860	14
34	867	15
35	875	16
36	884	18
37	896	21
38	913	26
39	941	38
40	980	65

Note: **Bold** and * represents CSEM around cut score (or the next higher scale score if the cut score value is not in the table).

Table 6-48 Numbers and Percentages of Students at LOSS and HOSS

Content	Grade	LOSS	N	Percentage	HOSS	N	Percentage
ELA	3	330	70	0.12	900	1	0.00
	4	340	14	0.02	930	0	0.00
	5	350	29	0.05	940	0	0.00
	6	360	25	0.04	950	0	0.00
	7	370	18	0.03	960	1	0.00
	8	380	46	0.07	970	6	0.01
Mathematics	3	360	741	1.26	760	262	0.45
	4	405	886	1.50	800	86	0.15
	5	430	1054	1.77	830	45	0.08
	6	440	968	1.62	870	90	0.15
	7	450	1011	1.67	880	30	0.05
	8	470	1155	1.85	890	96	0.15
Science	4	300	45	0.08	725	42	0.07
	8	480	51	0.08	945	31	0.05
Social Studies	4	330	511	0.86	700	89	0.15
	8	540	600	0.96	860	167	0.27
	10	645	1930	3.12	980	48	0.08

Figure 6-1 Anchor Set Test Characteristic Curves, English Language Arts Grades 3 through 8

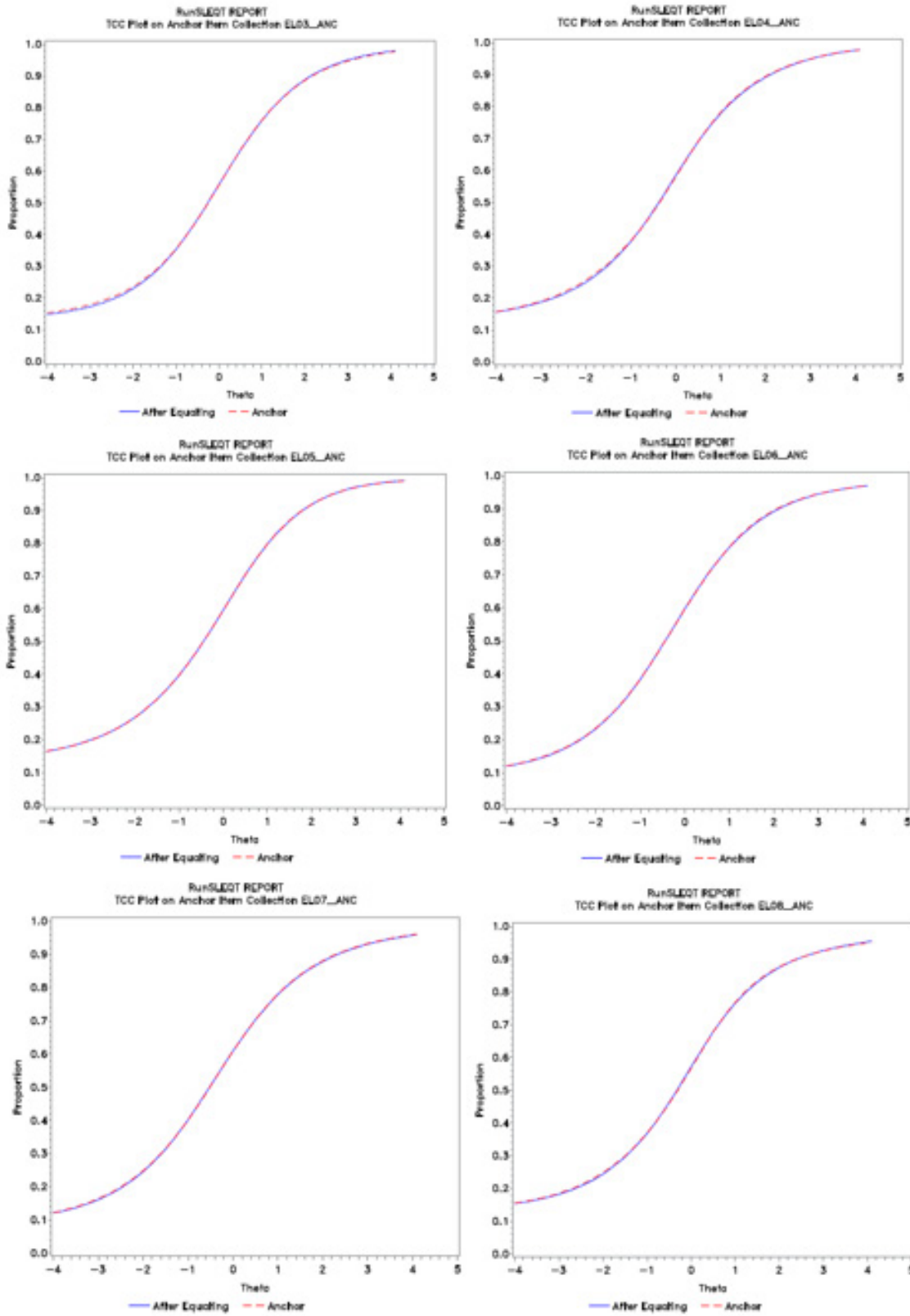


Figure 6-2 Anchor Set Test Characteristic Curves, Mathematics Grades 3 through 8

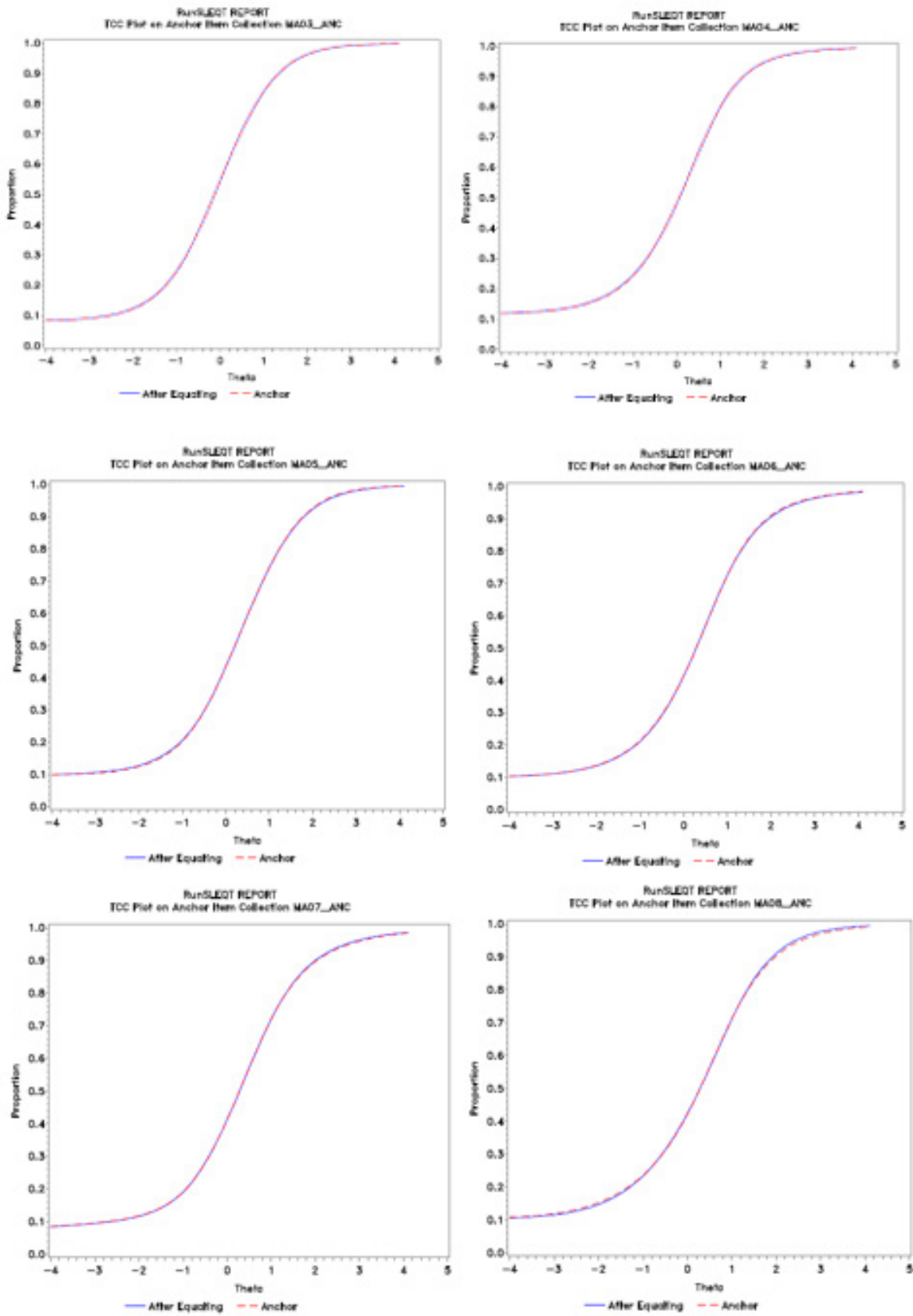


Figure 6-3 Anchor Set Test Characteristic Curves, Science Grades 4 and 8

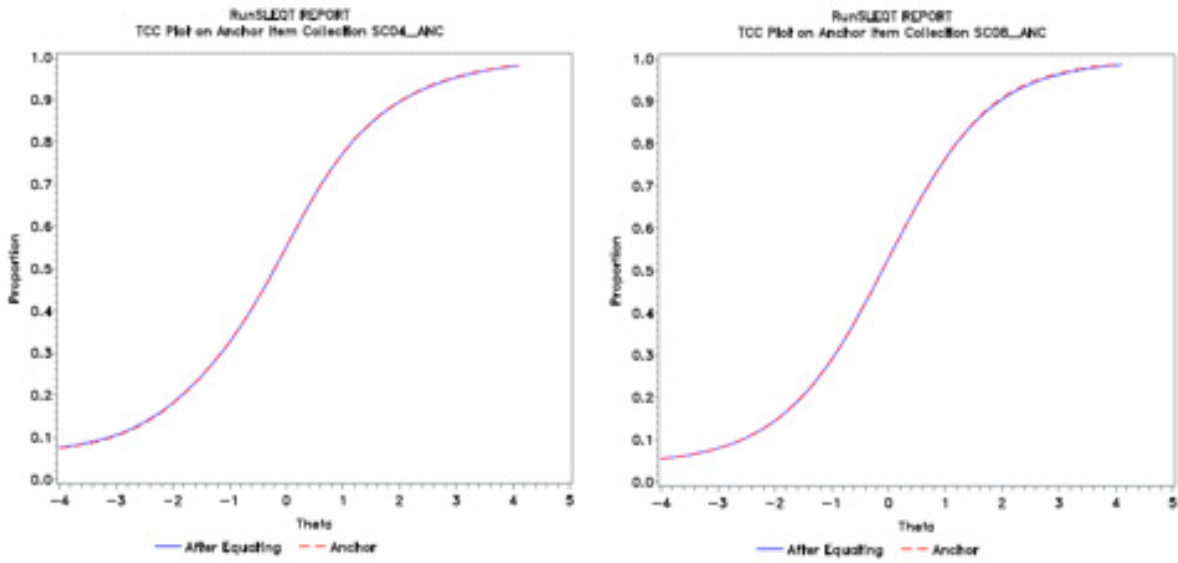


Figure 6-4 Anchor Set Test Characteristic Curves, Social Studies Grades 4, 8, and 10

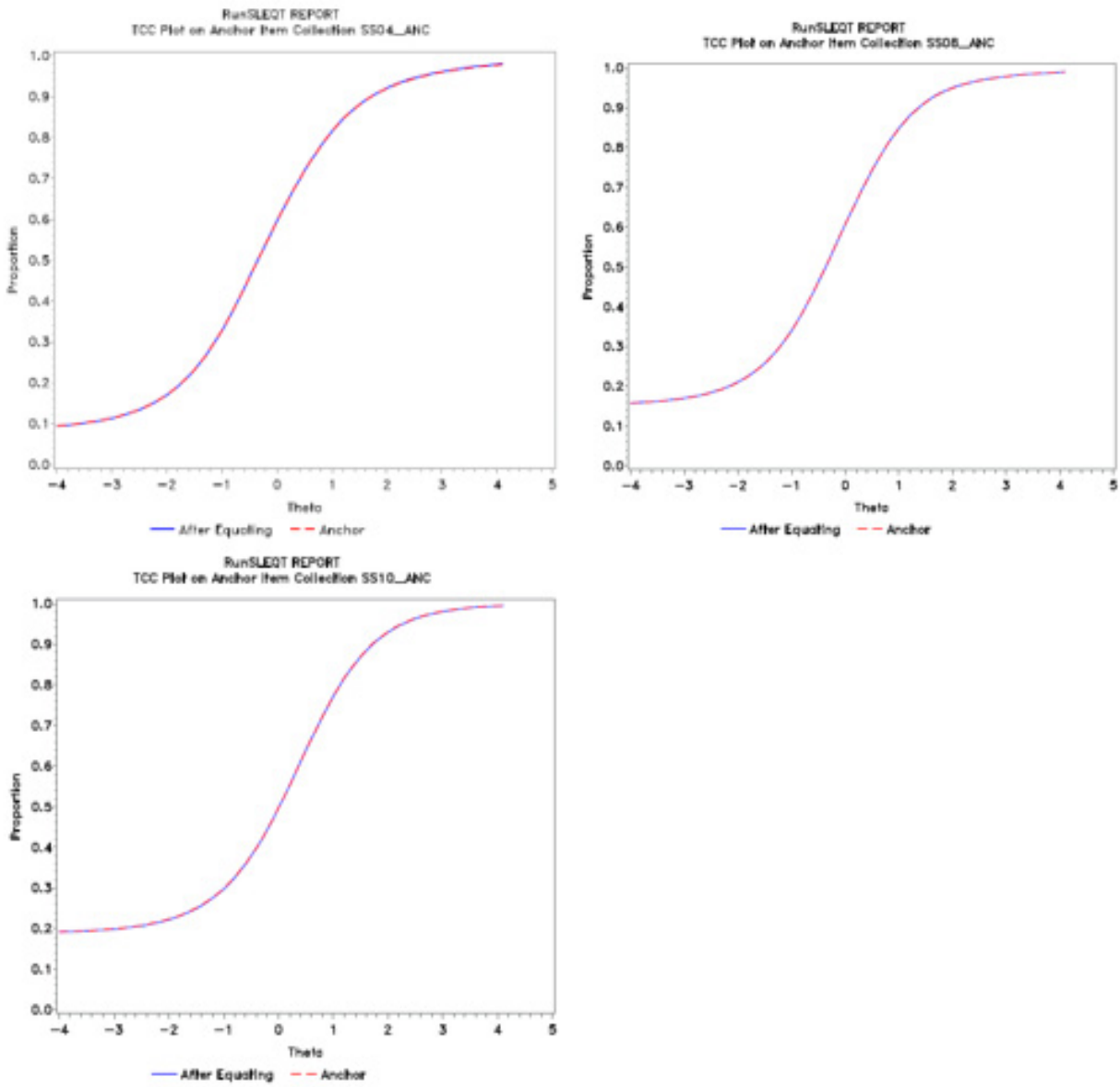


Figure 6-5 Test Characteristic Curves, English Language Arts

Scale Evaluation REPORT
TCC Plots: Ability vs. Proportion

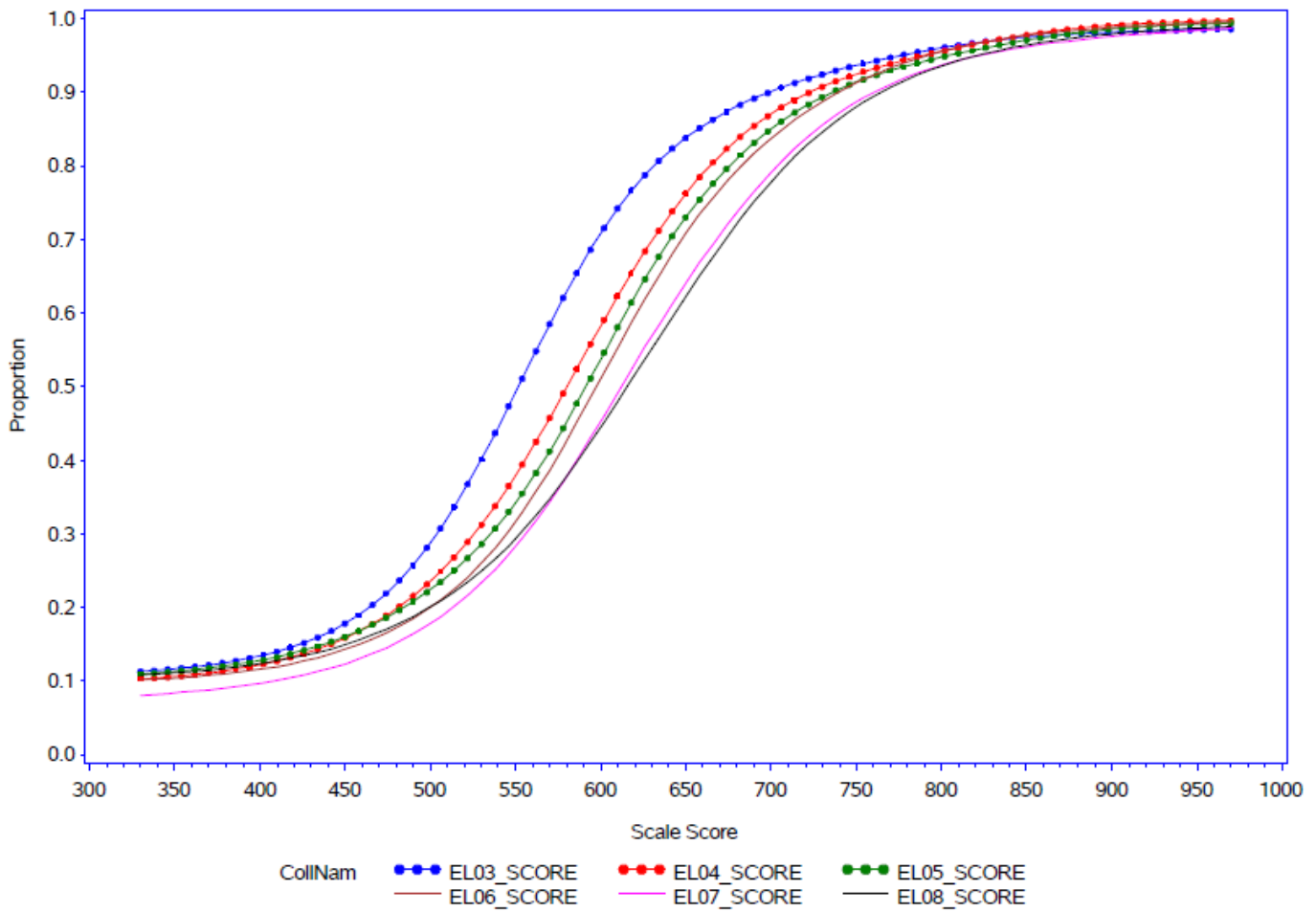


Figure 6-6 Standard Error Curves, English Language Arts

Scale Evaluation REPORT REPORT
Constrained IP SEM Plots: Ability vs. SEM

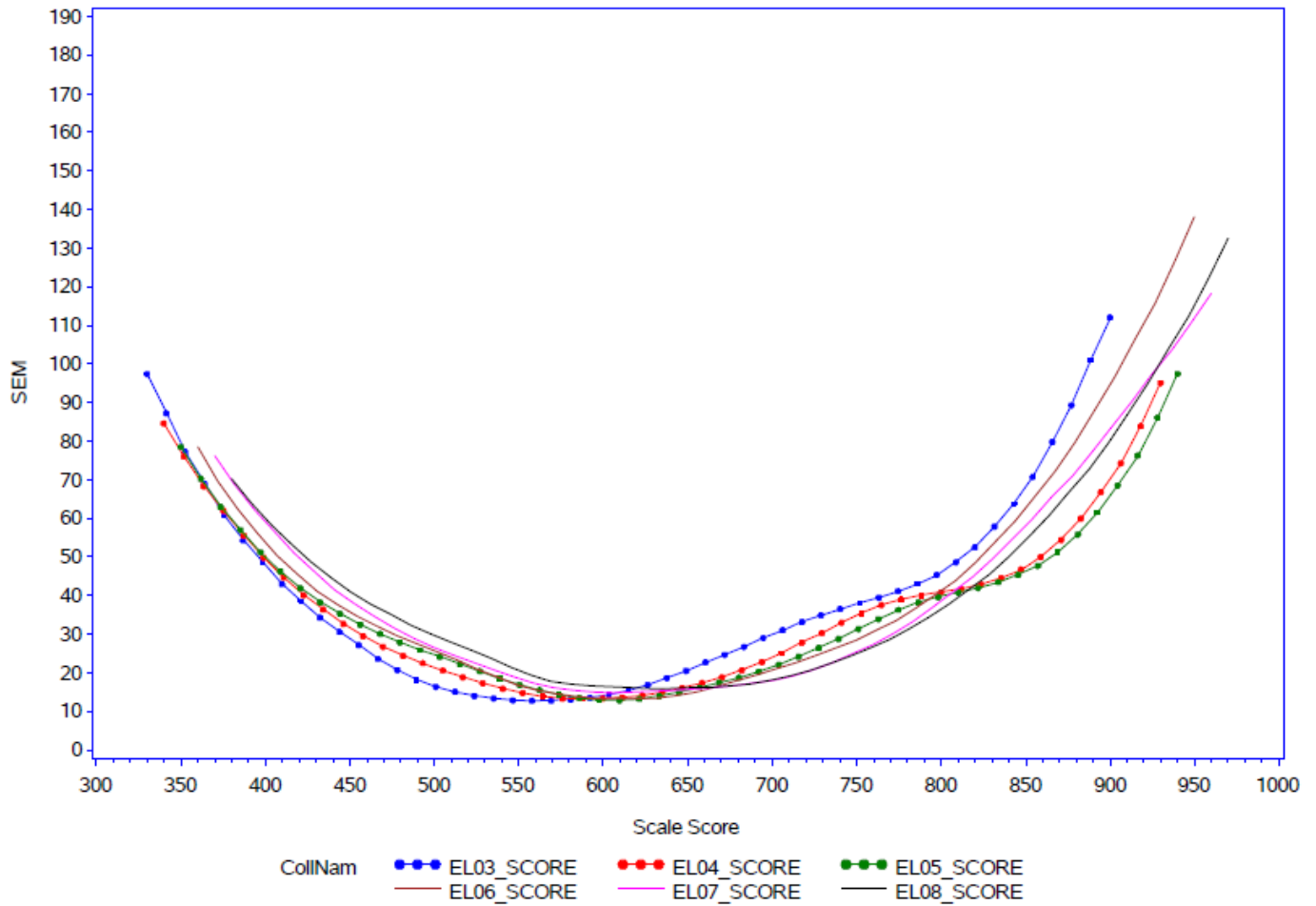


Figure 6-7 Scale Scores and Growth at Quartiles, English Language Arts

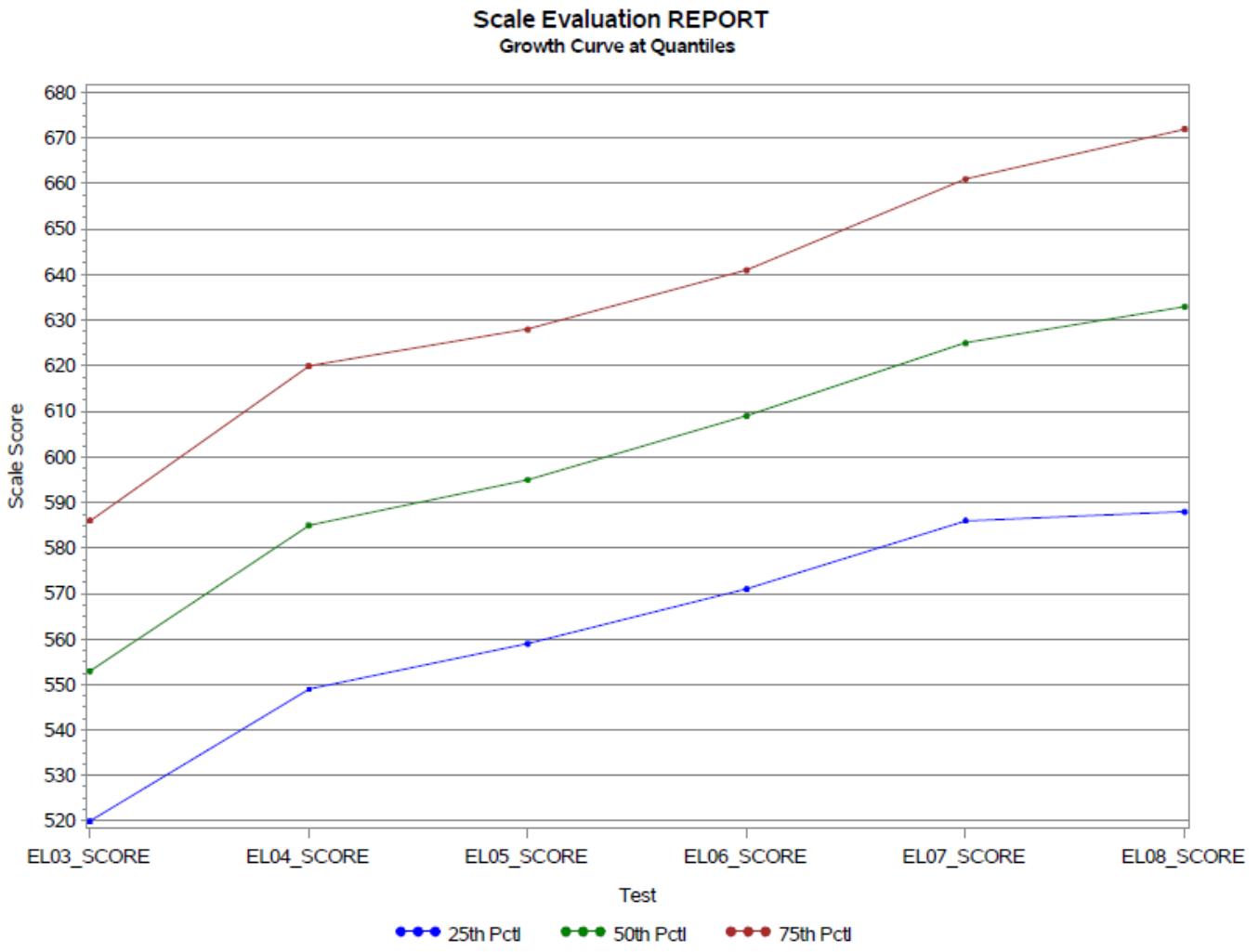


Figure 6-8 Test Characteristic Curves, Mathematics

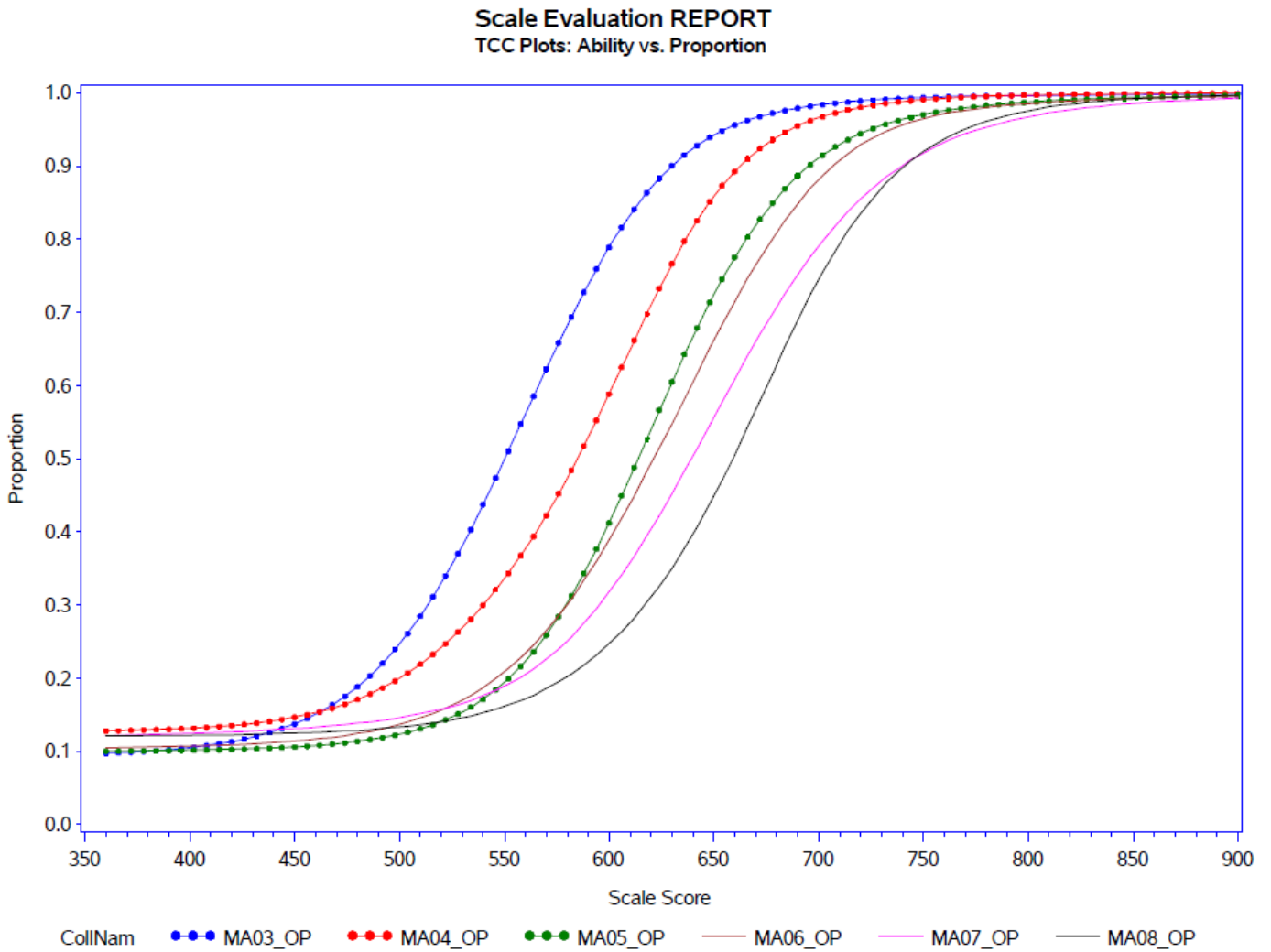


Figure 6-9 Standard Error Curves, Mathematics

Scale Evaluation REPORT REPORT
Constrained IP SEM Plots: Ability vs. SEM

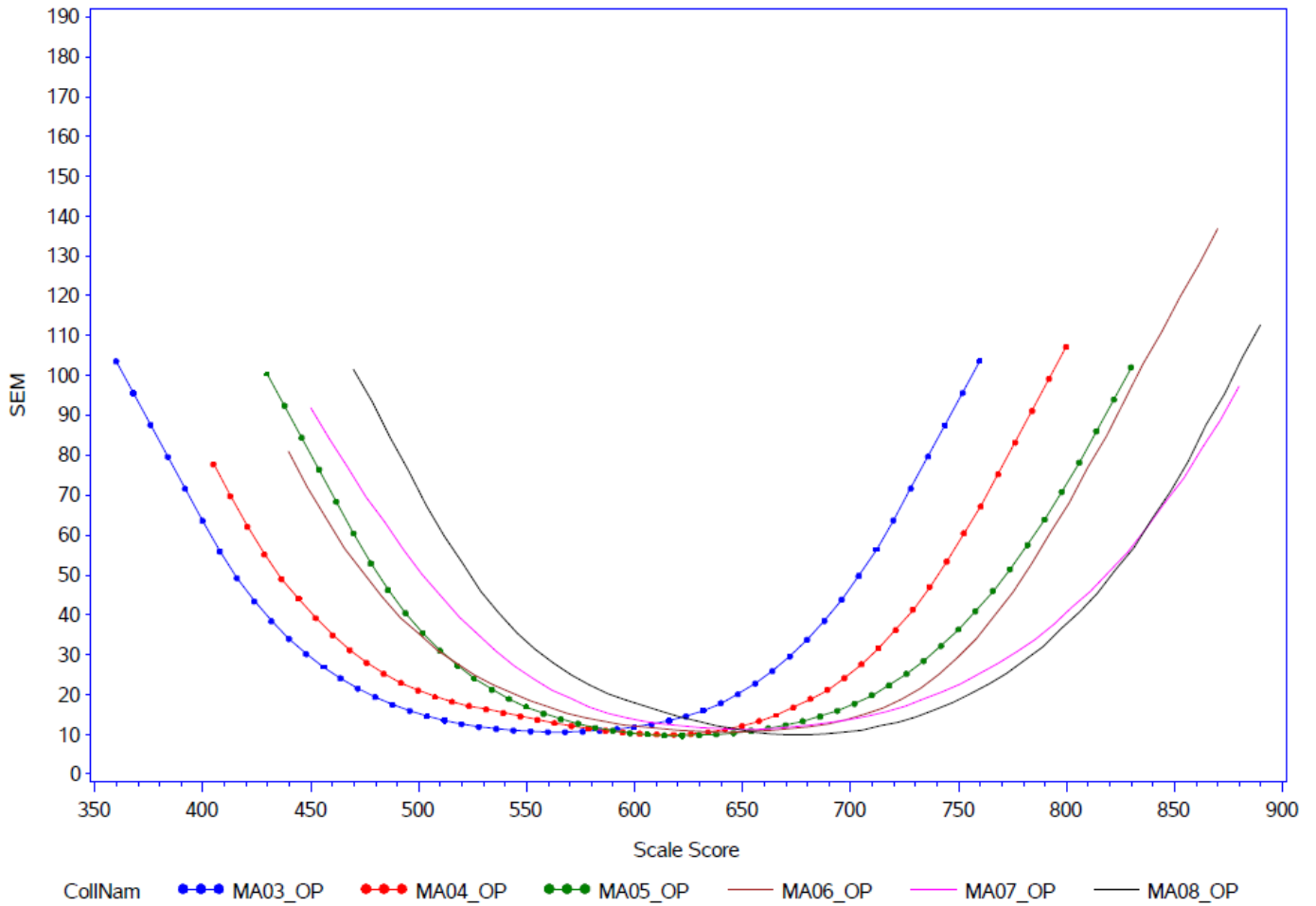


Figure 6-10 Scale Scores and Growth at Quartiles, Mathematics

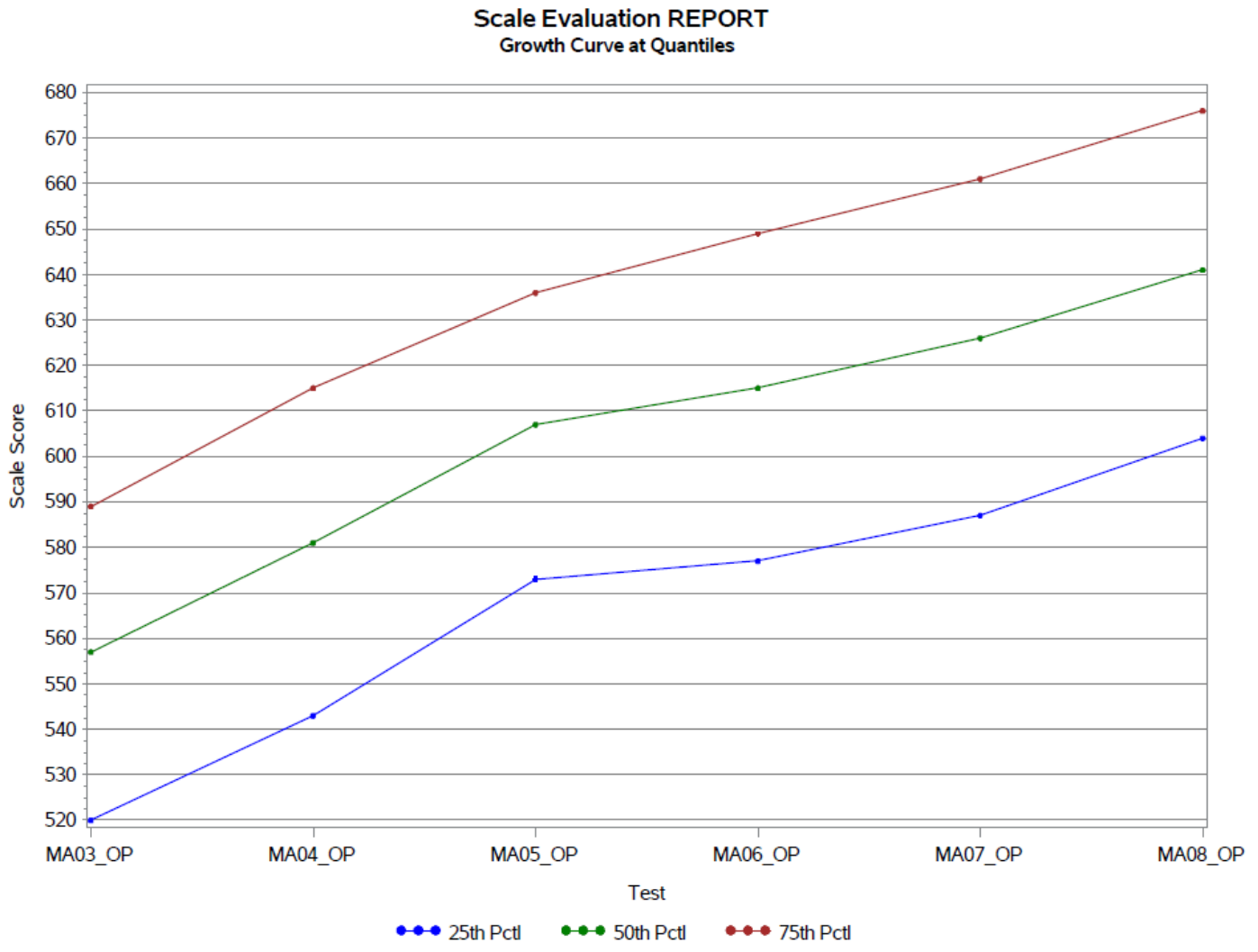


Figure 6-11 Test Characteristic Curves, Science

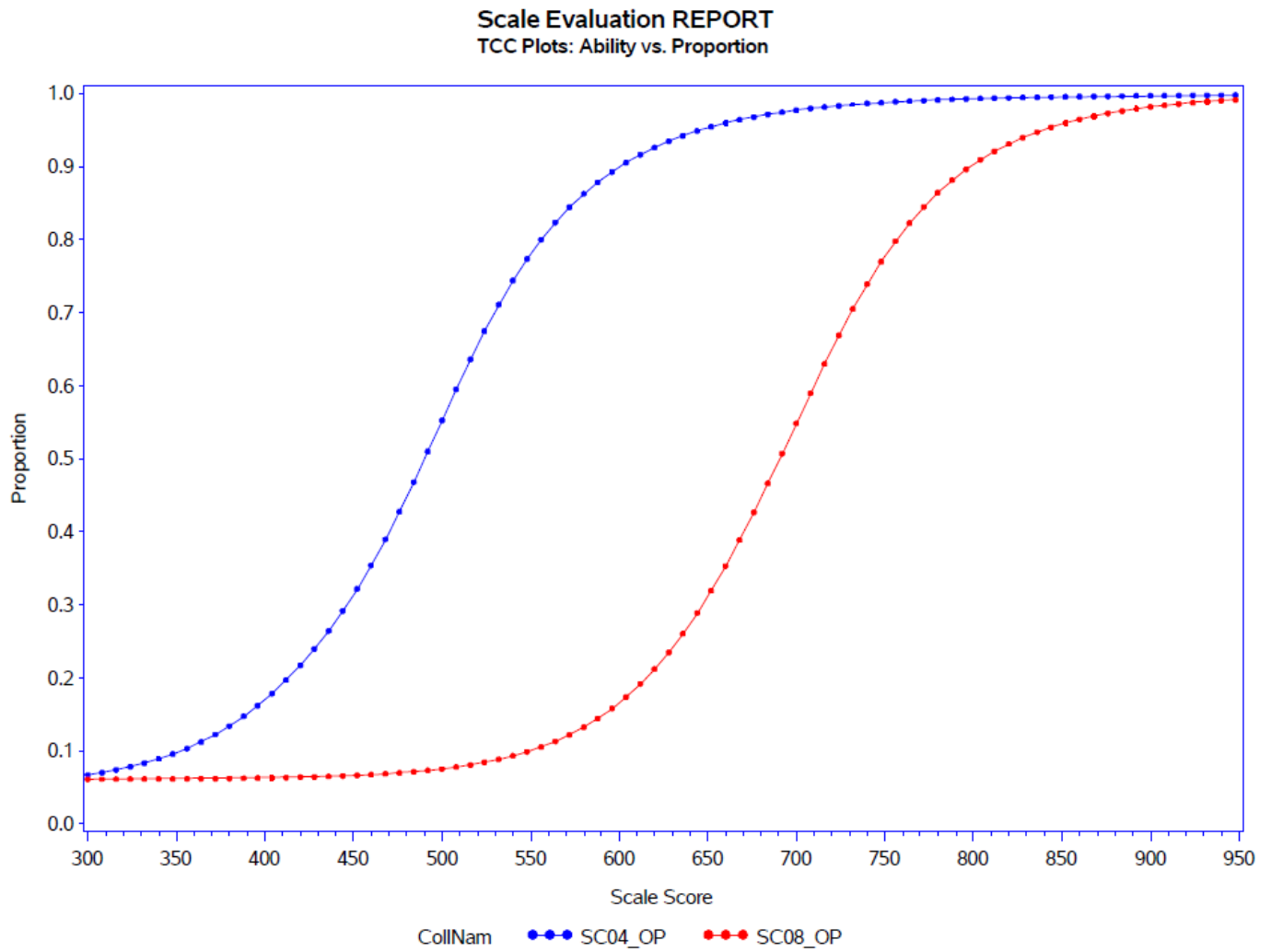


Figure 6-12 Standard Error Curves, Science

Scale Evaluation REPORT REPORT
Constrained IP SEM Plots: Ability vs. SEM

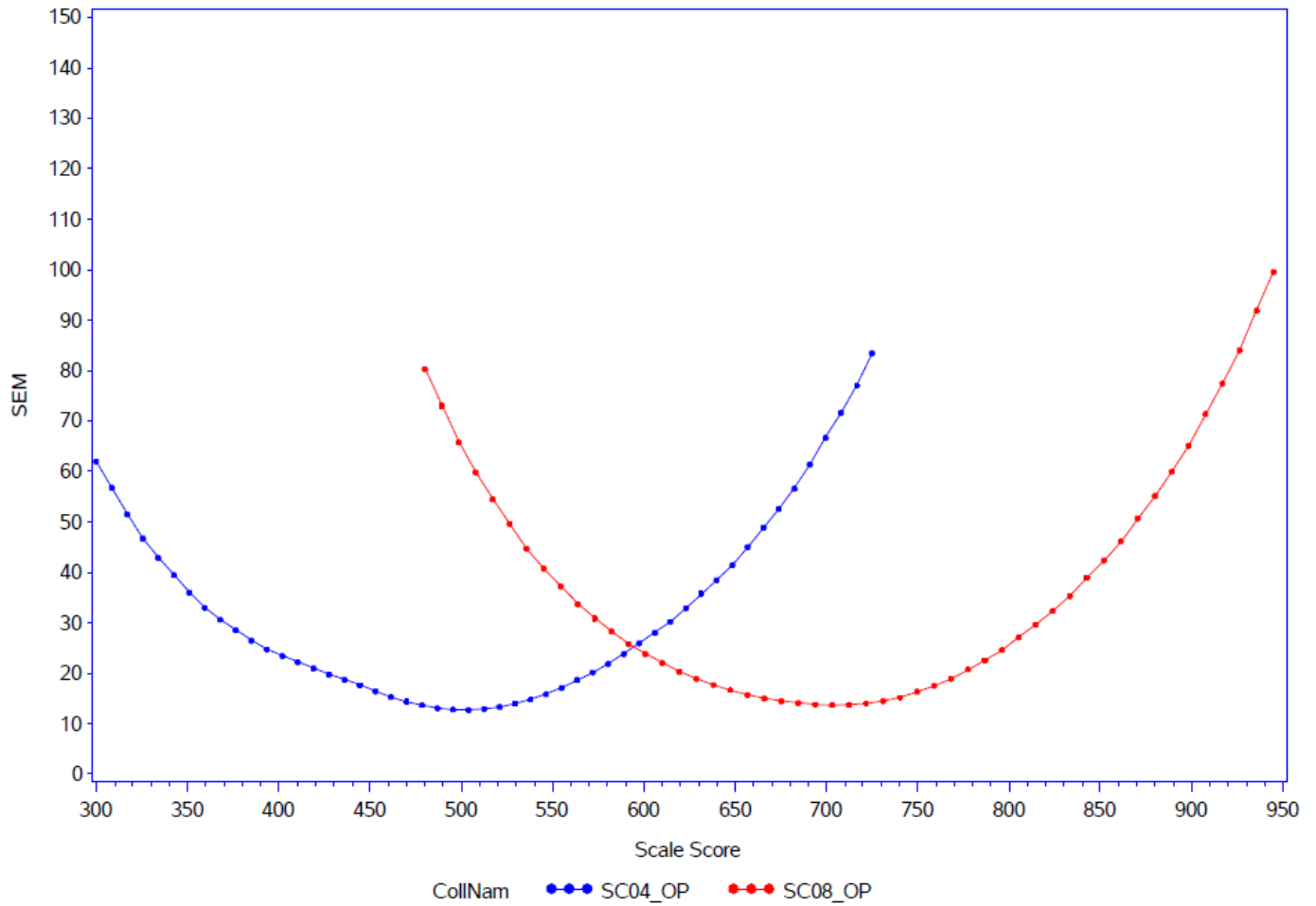


Figure 6-13 Scale Scores at Quartiles, Science

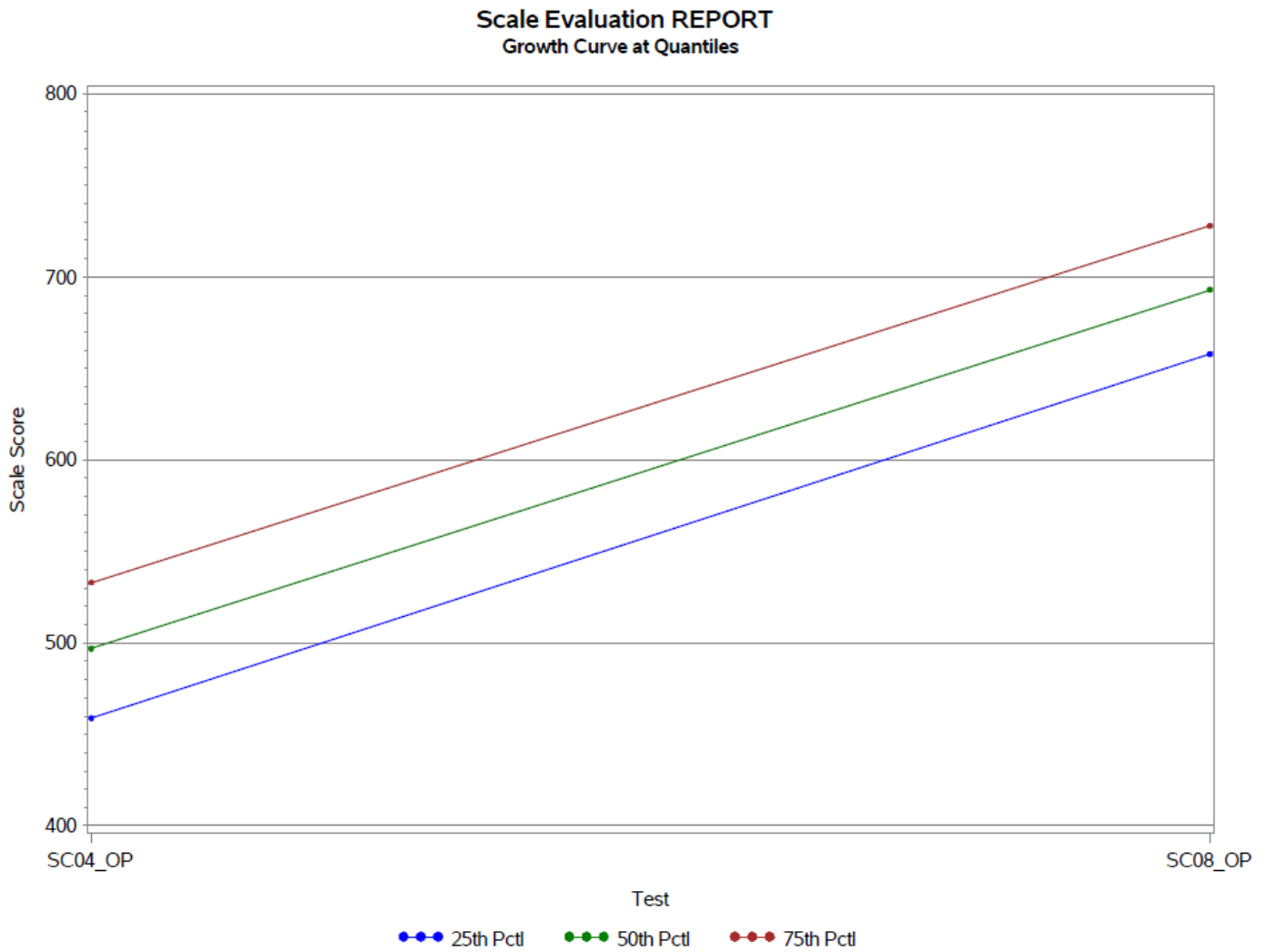


Figure 6-14 Test Characteristic Curves, Social Studies

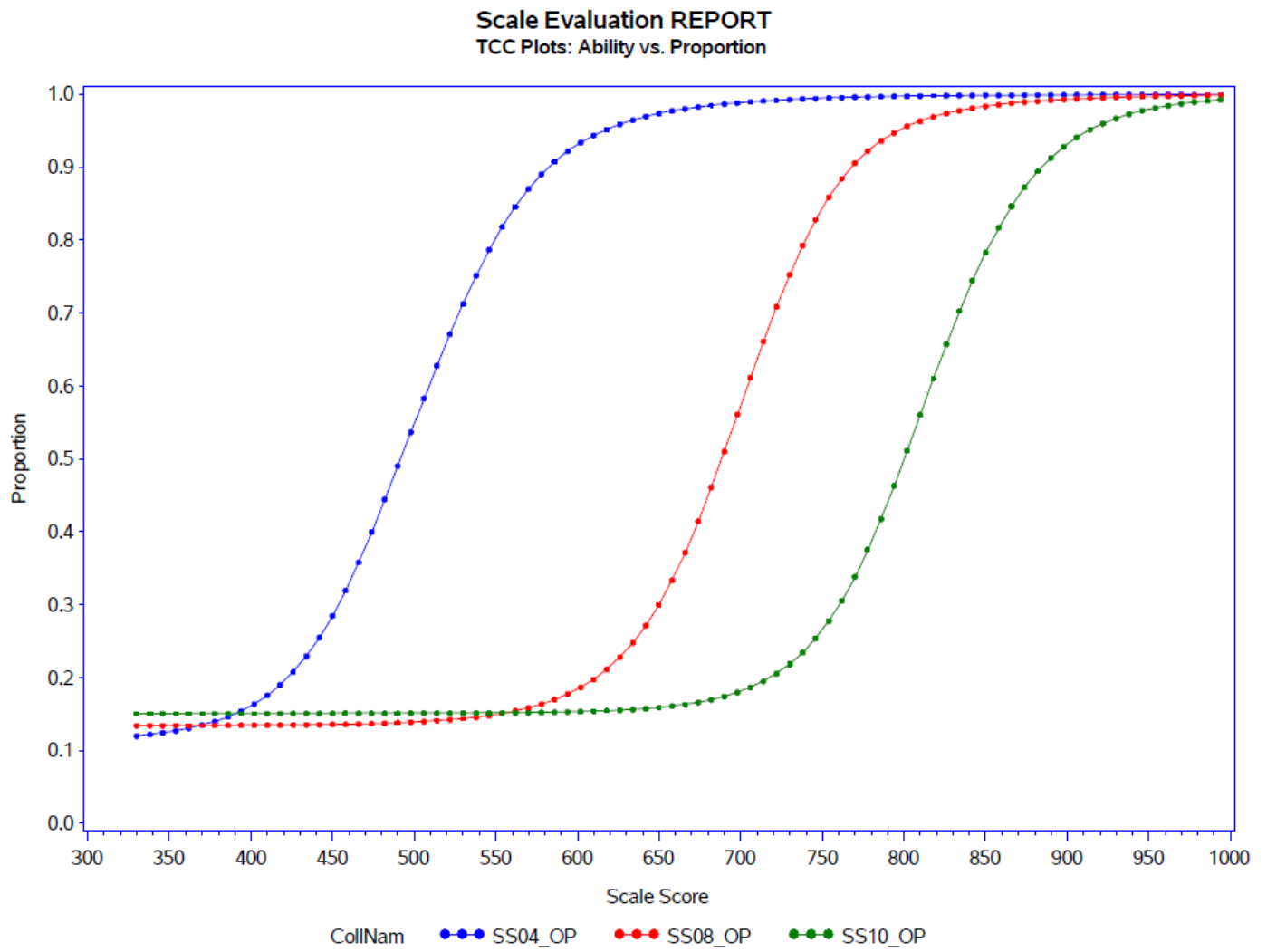


Figure 6-15 Standard Error Curves, Social Studies

Scale Evaluation REPORT REPORT
Constrained IP SEM Plots: Ability vs. SEM

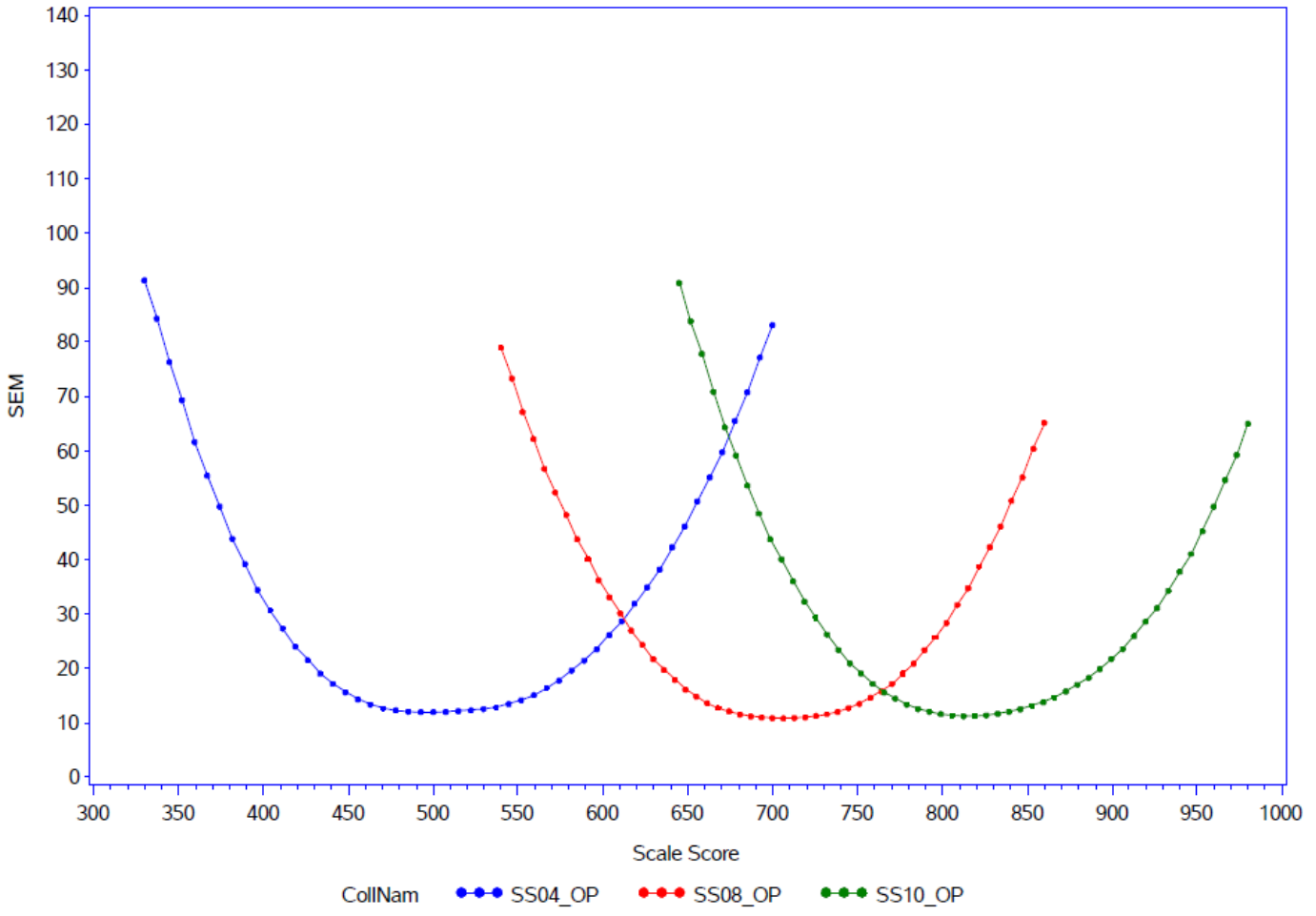
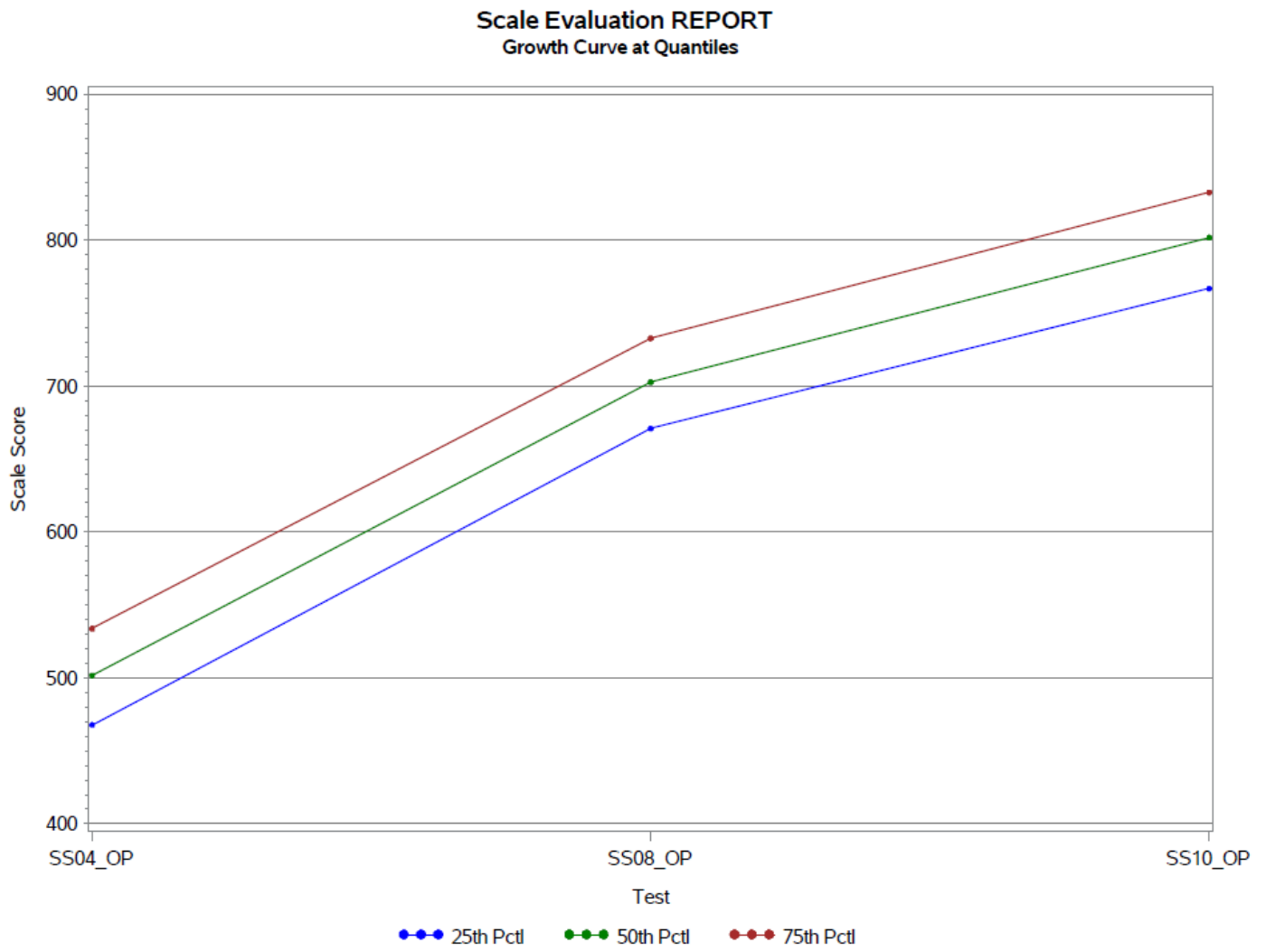


Figure 6-16 Scale Scores at Quartiles, Social Studies



Part 7: Standard Setting

In this part of the report, the standard settings that were conducted for Wisconsin Forward Exams are briefly described. The first standard setting for ELA, Mathematics, Science, and Social Studies occurred in Spring 2016. After the implementation of new Wisconsin Standards for Science, a new standard setting was conducted for Science in Spring 2019. Following the implementation of new Wisconsin Standards for Social Studies, a new standard setting for Social Studies was held in Spring 2022.

The cut scores established during these workshops and the performance level descriptors derived from the standard setting are also presented in this section of the report. The information on the ELA, Mathematics, Science, and Social Studies Spring 2016 standard setting comes from the *Wisconsin Standard Setting Technical Report 2016*. The information on the Science Spring 2019 standard setting comes from the *Wisconsin Forward Exam 2019 Science Standard Setting Technical Report*. The information on the Social Studies Spring 2022 standard setting comes from the *Wisconsin Forward Exam 2022 Social Studies Standard Setting Technical Report*. All three reports are available at <http://dpi.wi.gov/assessment/forward/resources>.

7.1 Background Information

Several changes were made to Wisconsin's statewide tests in recent years. In the 2014–15 school year, the Wisconsin Badger Exam measured students' abilities in ELA and Mathematics using assessments developed by the Smarter Balanced Assessment Consortium (SBAC). Cut scores for the Wisconsin Badger Exam were taken from the national SBAC standard setting, conducted in 2014. For Science and Social Studies, the Wisconsin Knowledge and Concepts Examination (WKCE) was administered. Cut scores for the WKCE were established in 2005.

In the 2015–16 school year, DPI consolidated the Wisconsin Badger Exam and the WKCE into a unified program, the Wisconsin Forward Exam. At the inception of the Wisconsin Forward Exam, DPI indicated that they would no longer use SBAC items or test scales for ELA and Mathematics and that new test scales would be established for the Wisconsin Forward Exam. New test scales and performance levels were established for all four content areas using data from the Spring 2016 administration of the Wisconsin Forward Exam.

Changes to Wisconsin Science standards, test blueprint, and test design were implemented for the Spring 2019 Science operational test administration. New scales were developed, and new performance level cut scores were set for Science tests in Spring 2019.

New test blueprints and test designs were implemented for Spring 2022 Social Studies after adopting new Wisconsin Standards for Social Studies. New scales were developed and new performance level cut scores were set for Social Studies tests in Spring 2022.

7.2 Standard Setting Methodology and Process

Prior to the standard setting workshops in Spring 2016, 2019, and 2022, DPI worked in collaboration with DRC and its other technical advisors to select the methodology to be used at the standard setting. In recognition of its use in Wisconsin and widespread use across the country, the Bookmark Standard Setting Procedure (BSSP) for the Wisconsin Forward Exam was selected for use by DPI. The BSSP was well suited for standard setting for these assessments because (a) the tests are composed of both MC and non-MC items, (b) the items are scaled and can be mapped using item mapping techniques, and (c) the BSSP allows participants to focus on the knowledge, skills, and abilities expected of students in each performance level. The BSSP has been well documented in standard setting literature. Developed in 1996, the BSSP has been implemented in over half of the states in the United States and abroad by DRC and by other major testing firms, making it the most widely used standard setting procedure in K–12 education (Karantonis & Sireci, 2006; Cizek & Bunch, 2007).

7.2.1 Spring 2016 Standard Setting for All Content Areas

On June 14–17, 2016, DPI and DRC conducted the Wisconsin Forward Exam standard setting for grades 3–8 in ELA and Mathematics, grades 4 and 8 in Science, and grades 4, 8, and 10 in Social Studies. The purpose of the standard setting was to develop performance standards for the Wisconsin Forward Exam, including the development of cut scores that divide students into four performance levels: *Below Basic*, *Basic*, *Proficient*, and *Advanced*. During this benchmarked standard setting, DPI developed cut scores on the Wisconsin Forward Exam that reflected these content-based expectations on the tests, as informed by test data from well respected measures of student achievement.

A total of 59 Wisconsin educators and stakeholders worked individually and in committees to recommend performance standards associated with the four performance levels: *Below Basic*, *Basic*, *Proficient*, and *Advanced*. This process yielded performance standards for the 17 tests of the Wisconsin Forward Exam program. The performance standards were approved by the Superintendent of Public Instruction in July 2016. For more information on the ELA, Mathematics, Science, and Social Studies standard setting, refer to *Wisconsin Standard Setting Technical Report 2016*, available at <http://dpi.wi.gov/assessment/forward/resources>.

7.2.2 Spring 2019 Standard Setting for Science

Because the Science test blueprint and design changed for the Spring 2019 administration and new Science reporting scales were developed, a new performance level setting was needed for this content area. On May 29 and 30, 2019, DPI and DRC conducted the Wisconsin Forward Exam standard setting for grades 4 and 8 in Science. The purpose of the standard setting was to develop new performance standards for the Science tests, including the development of cut scores that divided students into the four performance levels: *Below Basic*, *Basic*, *Proficient*, and *Advanced*. During the standard setting, Wisconsin educators made recommendations for cut scores on the Wisconsin Forward Exam that reflected the content-based expectations on the tests, as informed by test data from other measures of student Science achievement.

A total of 27 Wisconsin educators, 13 for grade 4 and 14 for grade 8, working individually and in grade-specific committees, recommended performance standards associated with the four performance levels for the two

Science assessments: *Below Basic*, *Basic*, *Proficient*, and *Advanced*. Participants engaged in three rounds of discussions and judgments to make their cut score recommendations. The cut scores recommended by the committee were approved by the State Superintendent of Public Instruction on June 5, 2019. For more information on the Science standard setting, refer to *Wisconsin Forward Exam 2019 Science Standard Setting Technical Report*, available at <http://dpi.wi.gov/assessment/forward/resources>.

7.2.3 Spring 2022 Standard Setting for Social Studies

The standard setting was required for Social Studies in 2022 because of recent changes to the Social Studies standards. The Spring 2022 Social Studies tests were the first to measure operationally the Wisconsin Standards for Social Studies published in 2018. New reporting scales were developed and new cut scores were needed to align with the new content standards.

On May 24–26, 2022, a committee of 50 Wisconsin educators participated in an online standard setting workshop for the Wisconsin Forward Exam for Social Studies in grades 4, 8, and 10. At the workshop, participating educators recommended cut scores to divide students into four performance levels: *Below Basic*, *Basic*, *Proficient*, and *Advanced*.

During the standard setting workshop, participants were divided into three groups, one per tested grade. Participants engaged in three rounds of discussions and judgments to make their cut score recommendations. The cut scores recommended by the committee were approved by the State Superintendent of Public Instruction in June 2022. For details on the Social Studies standard setting process and results, refer to *Wisconsin Forward Exam 2022 Social Studies Standard Setting Technical Report*, available at <http://dpi.wi.gov/assessment/forward/resources>.

The process of all three standard settings adhered to AERA, APA, & NCME (2014) Standards 5.21 and 5.22, which state the following:

Standard 5.21 When proposed score interpretations involve one or more cut scores, the rationale and procedures used for establishing cut scores should be documented clearly. (p. 107)

Standard 5.22 When cut scores defining pass-fail or proficiency levels are based on direct judgments about the adequacy of item or test performances, the judgmental process should be designed so that the participants providing the judgments can bring their knowledge and experience to bear in a reasonable way. (p. 108)

7.3 Performance Level Descriptors

In terms of the validity of the Wisconsin Forward Exam scores, it is essential to understand that descriptors and cut scores are established in a collaborative and participatory process. The descriptors clearly establish, in plain language, the proper frame of reference for understanding how to interpret test scores, particularly cut scores. Performance level descriptors (PLDs) summarize the knowledge, skills, and abilities expected of students in each performance level. DPI provided policy PLDs for the Wisconsin Forward Exam during the Spring 2016,

Spring 2019, and Spring 2022 standard settings. The brief policy performance level descriptors, shown in Table 7-1, describe DPI’s vision for each performance level. In addition, the standards-based PLDs for the Wisconsin Forward Exam in Science and Social Studies were provided to the standard setting participants in Spring 2019 and Spring 2022. (For detailed standards-based PLDs, refer to *Wisconsin Forward Exam 2019 Science Standard Setting Technical Report* and *Wisconsin Forward Exam 2022 Social Studies Standard Setting Technical Report*.) At the most recent standard settings for Science and Social Studies, Wisconsin educators used the policy PLDs in conjunction with standards-based PLDs to consider the content-based expectations for students in each performance level on each Science test in the Wisconsin Forward Exam program.

7.4 Cut Scores

Table 7-2 shows the cut scores for all grades and content areas. The cut scores reflect the content-based expectations for students and policy-based decisions (i.e., the impact of the cut scores on Wisconsin students as shown through the impact data). The cut scores for ELA and Mathematics, established in Spring 2016, remained unchanged for the 2023 assessments. New cut scores were established for Science after the Spring 2019 test administration, and these cut scores were used for student classification into performance levels in Spring 2019 through 2023. New cut scores reflecting Wisconsin student performance on the new Social Studies assessments were established for Social Studies after the Spring 2022 test administration and were used for student classification into performance levels in Spring 2022 and Spring 2023.

7.5 Summary

Part 7 presented a brief overview of the standard setting process used to establish the Wisconsin Forward Exam cut scores for all content areas after the Spring 2016 test administration, for Science after the Spring 2019 test administration, and for Social Studies after the Spring 2022 administration. The standard settings undertaken by DPI and facilitated by DRC support Standards 5.21 and 5.22 from the *Standards* (AERA, APA, & NCME, 2014).

Table 7-1 Policy Performance Level Descriptors for the Wisconsin Forward Exam

Level	Performance Level Descriptor
<i>Below Basic</i>	Student demonstrates minimal understanding of and ability to apply the knowledge and skills for their grade level that are associated with college content-readiness.
<i>Basic</i>	Student demonstrates partial understanding of and ability to apply the knowledge and skills for their grade level that are associated with college content-readiness.
<i>Proficient</i>	Student demonstrates understanding of and ability to apply the knowledge and skills for their grade level that are associated with college content-readiness.
<i>Advanced</i>	Student demonstrates exemplary understanding of and ability to apply the knowledge and skills for their grade level that are associated with college content-readiness.

Table 7-2 Wisconsin Forward Exam Cut Scores

Content	Grade	Basic	Proficient	Advanced
ELA	3	522	570	624
	4	546	592	650
	5	564	610	670
	6	572	622	671
	7	585	638	697
	8	592	652	708
Mathematics	3	517	560	611
	4	536	588	633
	5	574	611	658
	6	582	626	688
	7	606	647	712
	8	620	667	718
Science	4	447	496	543
	8	653	695	737
Social Studies	4	461	491	537
	8	662	693	734
	10	770	805	837

Part 8: Studies of Reliability

Part 8 of the Technical Report builds upon existing analyses of the summary results by providing additional estimates of the reliability of those results. Reliability can be defined as the consistency of an assessment when the testing procedure is repeated with the same testing target group. A reliable assessment is one that would produce stable scores if the same group of students were to take the same test repeatedly, without any fatigue or memory of the test. As detailed below, the reliability of the Spring 2023 Wisconsin Forward Exam was estimated in three ways:

- Internal consistency was assessed for all items using Cronbach’s alpha (1951).
- Standard error of measurement (SEM) was calculated for raw score and scale score.
- Classification consistency and classification accuracy were estimated for the performance level classifications.

This part of the report addresses AERA, APA, & NCME (2014) Standards 2.0, 2.3, 2.11, 2.13, 2.14, and 2.16, which are cited below:

Standard 2.0 Appropriate evidence of reliability/precision should be provided for the interpretation for each intended score use. (p. 42)

Standard 2.3 For each total score, subscore, or combination of scores that is to be interpreted, estimates of relevant indices of reliability/precision should be reported. (p. 43)

Standard 2.11 Test publishers should provide estimates of reliability/precision as soon as feasible for each relevant subgroup for which the test is recommended. (p. 45)

Standard 2.13 The standard error of measurement, both overall and conditional (if reported), should be provided in units of each reported score. (p. 45)

Standard 2.14 When possible and appropriate, conditional standard errors of measurement should be reported at several score levels unless there is evidence that the standard error is constant across score levels. Where cut scores are specified for selection or classification, the standard errors of measurement should be reported in the vicinity of each cut score. (p. 46)

Standard 2.16 When a test or combination of measures is used to make classification decisions, estimates should be provided of the percentage of test takers who would be classified in the same way on two replications of the procedure. (p. 46)

Standard 2.3 advises providing reliability estimates and the SEM for all total scores and subscores reported, Standard 2.13 advises reporting SEM in both raw score and scale score units, and Standard 2.11 advises assessing reliability and SEM for all population subgroups. This part of the report presents raw score reliability coefficients and SEMs for the four Wisconsin Forward Exam content areas, for each reported content standard

for the total group of examinees, and for the subgroups identified by gender, race/ethnicity, economic status, disability status, and English language proficiency. The scale score CSEMs are provided in Section 8.1.1.

Standard 2.16 advises that when testing measures are used to make categorical decisions, the reliability of those decisions should be estimated. In the present context, Standard 2.16 applies specifically to performance level determinations, such as *Proficient* or *Advanced*. As described below, the Spring 2023 Wisconsin Forward Exam adhered to this standard by applying a detailed analysis of classification consistency and classification accuracy—two related measures used to evaluate the reliability of the performance level classifications used in the test program. This analysis also addresses Standard 2.14 by providing a CSEM for the cut scores that separate the performance levels.

Combined, Cronbach’s alpha, SEM, classification consistency, and classification accuracy provide several forms of evidence related to the reliability of the Wisconsin Forward Exam. Cronbach’s alpha and the SEM operate at the content level. For example, they provide estimates of reliability for student scores in ELA or Mathematics. Classification consistency and classification accuracy operate on the associated performance level classifications. These are of particular interest in the context of the Elementary and Secondary Education Act (ESEA) and the associated accountability requirements. In addition, the Cronbach’s alpha statistics and the SEM were computed for content standards and domains, providing evidence of the reliability and precision of measurement of the Wisconsin Forward Exam subscores. Altogether, the evidence provided in this part of the Technical Report, which is targeted at each intended use of the Wisconsin Forward Exam scores, addresses Standard 2.0.

8.1 Measures of Internal Consistency and Standard Error of Measurement

Cronbach’s alpha is a frequently used measure of internal consistency for tests consisting of MC and CR items. Cronbach’s alpha (α) is computed as

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

where k = number of items, σ_X^2 = the total score variance, and σ_i^2 = the variance of item i (Crocker & Algina, 1986). SEM is defined as

$$SEM = SD\sqrt{1 - R_{xx}},$$

where SD represents the standard deviation of the raw score distribution and R_{xx} represents Cronbach’s alpha.

Cronbach’s alpha and the SEM are shown in Tables 8-1 and 8-2, respectively. These tables include information for all students and for the subgroup categories of gender, race/ethnicity, English language proficiency, disability status, economic status, and accommodation use.

As shown in the “Total” column of Table 8-1, reliability ranges from 0.88 to 0.89 across grades for ELA and from 0.91 to 0.93 across grades for Mathematics. The reliability coefficient was 0.91 for Science grade 4 and

0.89 for Science grade 8. The reliability coefficient was 0.91 for Social Studies in grades 4 and 8 and 0.90 for Social Studies grade 10. All reliability coefficients would ideally be 0.90 or above. However, for relatively short tests that are designed to measure a fairly broad range of content, this is not always a realistic expectation. If 0.90 is considered a conservative criterion for an acceptable level of reliability as measured by Cronbach's alpha, then the ELA assessments in all grades and the Science assessment in grade 8 would not meet this criterion. The reliability coefficients for these tests may be affected by the number of items (and score points) and the diversity of the content being assessed. Applying the Spearman-Brown prophecy formula to these results indicates that to achieve the 0.90 reliability threshold, the current ELA assessments for grades 4, 6, and 8 would need to be increased by approximately 7 to 9 points and the ELA assessments for grades 3, 5, and 7 would need to be increased by approximately 12 to 15 points. For the current Science assessment in grade 8, the increase would need to be at least 4 score points.

Table 8-1 shows that many of the subgroup reliability coefficients were similar to or lower than the total reliability coefficients. Reliability coefficients are particularly sensitive to score distribution and variance, so this result is consistent with the general variability among many of these subgroups.

The test reliability coefficients were comparable within 0.02 for male and female students across all grades and content areas. The reliability coefficients for the non-binary student subgroup were computed only for grades 8 and 10, which had at least 50 non-binary students, and should be interpreted with caution because of the low student counts.

Most differences among the five racial/ethnic groups were small and within 0.05 of one another for all grades in ELA, Science, and Social Studies. In Mathematics, the highest test reliabilities were observed for White and Asian students and the lowest reliability was observed for African American students in all grades.

The differences between test reliability coefficients for economically disadvantaged and not economically disadvantaged students were within 0.03 for all grades and content areas. The test reliability was higher for students who were not economically disadvantaged.

The differences in reliability of test scores between disabled and not disabled students were within 0.03 for most grades and content areas, except for ELA grade 8, Mathematics grades 6 through 8, and Social Studies grade 10, where the differences were between 0.04 to 0.07. The test reliability was found to be higher for students without disabilities.

The greatest differences in test reliability were between fully English proficient and limited English proficient students in all content areas and between students using and not using testing accommodations in Mathematics. The test reliability coefficients for limited English proficient students were consistently lower than the test reliability coefficients for their English proficient peers and ranged from 0.79 to 0.83 for ELA, from 0.80 to 0.90 for Mathematics, from 0.79 to 0.85 for Science, and from 0.76 to 0.86 for Social Studies.

With the exception of grades 3 and 5, the reliability coefficients for students using testing accommodations in Mathematics were below 0.08 and ranged from 0.66 to 0.76 in grades 4, 6, 7, and 8. The reliability coefficients for students using testing accommodations in ELA should be interpreted with caution because of the low

number of students using the accommodations. The reliability coefficients were not computed for students using testing accommodations in Science or Social Studies because the number of students using accommodations in these subject areas was less than 50.

The reliability coefficient is affected by, among other factors, the variability of students' scores. The higher the variability of scores, the higher the reliability coefficient will tend to be. The patterns of the differences in test reliability for different subgroups in Spring 2023 were similar to the patterns observed in Spring 2022.

Table 8-2 presents the raw score SEM for the total population and for the subgroups described above. These values provide important information for raw score interpretation since an individual's obtained score can be expected to fall within two SEMs of the individual's true score approximately 95% of the time. Although there were some observable differences in SEM for the different subgroups, all differences were generally within one-half of a score point. The SEMs for ELA assessments were slightly larger than those for the other content areas. Because these SEMs are in the raw score metric, this result is consistent with the fact that ELA tests have more raw score points and relatively larger raw score SDs than the other content areas. For every grade and content area, the CSEM for individual scale scores is provided in the scoring tables previously discussed in Part 6 (Tables 6-31 through 6-47).

Reliability, as measured by Cronbach's alpha, was also computed for content standards (or reporting categories) within each content area as well as for each language domain in ELA. These data are presented in Table 8-3. The last column presents the reliability for the total test per grade for each content area (with all content standards or domains) for all examinees. It is clear that the reliability per content standard or domain is lower than the reliability for the total test per content area. The number of items or score points has a close relationship with reliability, and a smaller number of items or score points is generally associated with lower reliability. The number of score points for ELA per domain was 7 or 8 in Listening, 22 or 24 in Reading, and 24 in Writing/Language. The number of score points ranged from 4 to 13 per content standard (or reporting category) for ELA, from 7 to 11 per standard for Mathematics, from 8 to 12 per standard for Science, and from 6 to 11 per standard for Social Studies. A lower level of reliability per content standard or domain is therefore expected compared to each full content area assessment. The lower level of reliability per standard or domain is one of the reasons why the information based on the content standards or domains should be used for low-stakes purposes only. (This issue is also discussed in the context of standard performance index scores in Part 10.)

As shown in Table 8-3, the reliability ranges by content standard/domain were as follows:

- For ELA, reliability indices by content standard ranged from 0.34 for the Writing/Language—Language Conventions standard in grade 3 to 0.72 for the Reading—Craft & Structure/Integration of Knowledge & Ideas standard in grade 6. When the ELA domains were considered, the highest reliability, ranging from 0.79 to 0.84 across grades, was found for the Reading domain, followed by the Writing domain, with reliability coefficients ranging from 0.67 to 0.77 across grades, and followed by the Listening domain, which included a smaller number of items, resulting in the reliability coefficients ranging from 0.48 to 0.66 across grades.

- For Mathematics, reliability indices by content standard ranged from 0.55 (for the Geometry standard in grade 7) to 0.78 (for the Operations and Algebraic Thinking standard in grade 3).
- For Science, reliability indices by content standard ranged from 0.65 (for the Earth and Space Science standard in grade 8) to 0.74 (for the Engineering standard also in grade 4).
- For Social Studies, reliability indices by content standard ranged from 0.47 (for the Political Science and Citizenship standard in grade 4) to 0.77 (for the Geography standard also in grade 4).

The SEM associated with each content standard or domain is presented in Table 8-4 by content area and grade level. Some differences in SEM by content standard can be observed. As indicated by the discussion above, these SEMs were smaller than those for the total test and were generally consistent with the number of items within each content standard.

In summary, the reliability indices, as measured by Cronbach’s alpha at the test level, are in a reasonable range given the number of items in each test. As described above, readers should also note that, because reliability is influenced by the number of items, lower reliability for the content standards with fewer items is to be expected.

8.1.1 Conditional Standard Error of Measurement

In contrast to the SEM, the CSEM expresses the degree of measurement error in scale score units and is conditioned on the ability of the student. The CSEM is defined as the reciprocal of the square root of the test information function and can be estimated across all points of the ability continuum (Hambleton & Swaminathan, 1985):

$$\text{CSEM}(\theta_i) = \frac{1}{\sqrt{I(\theta_i)}}$$

where $I(\theta_i)$ is the test information function, computed as a sum of item information functions, obtained as

$$I(\theta_i) = \sum_j \frac{p'_{ij}(\theta_i)^2}{p_{ij}(\theta_i)q_{ij}(\theta_i)},$$

where $p'_{ij}(\theta_i)$ is the derivative of $p_{ij}(\theta_i)$ and $q_{ij}(\theta_i) = 1 - p_{ij}(\theta_i)$.

The CSEM can be used to obtain the range within which a student’s true score is likely to fall (that is, with a certain degree of probability). It is expected that a student’s score obtained from a single testing will fall within one CSEM of that student’s true score 68 percent of the time and that the obtained score will fall within two CSEMs of the true score 95 percent of the time.

Note that the CSEMs vary in magnitude across the entire range of student ability estimates (i.e., scale scores) and are smaller in the middle of the score distribution and larger at the tails. This pattern is seen for all Wisconsin Forward Exam CSEMs and is to be expected when IRT methods are used. In compliance with Standards 2.13 and 2.14, the CSEM of each cut score was presented in the raw score-to-scale score tables (Tables 6-31 through 6-47) for all grades and content areas in Part 6 of this report. In addition, graphical

representations of the CSEM with the cut scores are presented in Figures I-1 through I-17 of Appendix I for all grades and content areas. As shown in Appendix I, the estimates of CSEM tend to be higher at the low and high ends of the scale score range. The CSEM increases when there are few observations at a particular ability level. Generally, there are few students with extreme scores, and these score levels cannot be estimated as accurately as levels toward the middle of the ability range. Figures I-1 through I-17 demonstrate that the CSEM is minimized at the cut scores and in the middle of the scale range, where most students are located.

8.2 Classification Consistency and Accuracy

One of the primary goals of education policy is to improve the performance of all students, with a specific goal of having all students become *Proficient*. Because of this heavy emphasis on moving all students to levels of academic performance at or above each state’s self-defined *Proficient* category, the consistency and accuracy of the classification of students into these performance levels are of particular interest. The following section describes how the consistency and accuracy of these classifications were evaluated and provides evidence that supports the validity of these classifications.

Conceptually, classification consistency is defined as the extent to which two classifications of a single student agree, based either on two independent administrations of the same test or on one administration of two parallel test forms. However, it is difficult to obtain data from repeated administrations of the same form because of the cost, time, and student memory from prior administrations. It is also difficult to construct two psychometrically parallel forms. For these reasons, the common practice is to estimate classification consistency from a single administration.

A contingency table representing the probability of particular classification outcomes under specific scenarios is a convenient way to measure classification consistency. The table below is a contingency table of $(H + 1) \times (H + 1)$, where H is the number of cut scores. Three cut scores yield a 4×4 contingency table, as can be seen below in Table 8-A.

It is common to report two indices of classification consistency: the classification agreement “P” and the coefficient kappa. Hambleton and Novick (1973) proposed P as a measure of classification consistency, where P is defined as the sum of diagonal values of the contingency table:

$$P = P_{11} + P_{22} + P_{33} + P_{44}.$$

Table 8-A Example Contingency Table with Three Cut Scores

Proficiency Level	Level 1	Level 2	Level 3	Level 4	Sum
Level 1	P ₁₁	P ₂₁	P ₃₁	P ₄₁	P _{.1}
Level 2	P ₁₂	P ₂₂	P ₃₂	P ₄₂	P _{.2}
Level 3	P ₁₃	P ₂₃	P ₃₃	P ₄₃	P _{.3}
Level 4	P ₁₄	P ₂₄	P ₃₄	P ₄₄	P _{.4}
Sum	P _{1.}	P _{2.}	P _{3.}	P _{4.}	1.0

To reflect statistical chance agreement, Swaminathan, Hambleton, and Algina (1974) suggest using Cohen's kappa (1960) as

$$\text{kappa} = \frac{P - P_C}{1 - P_C},$$

where P_C is the chance probability of a consistent classification under two completely random assignments. Probability P_C is the sum of the probabilities obtained by multiplying the marginal probability of the first administration and the corresponding marginal probability of the second administration as

$$P_C = (P_{1.} \times P_{.1}) + (P_{2.} \times P_{.2}) + (P_{3.} \times P_{.3}) + (P_{4.} \times P_{.4}).$$

To aid in the interpretation of the kappa statistic, the suggested cutoffs (Landis & Koch, 1977; Altman, 1991) are presented below:

- Kappa of 0 means no agreement.
- Kappa less than 0.20 means poor agreement.
- Kappa from 0.21 to 0.40 means fair agreement.
- Kappa from 0.41 to 0.60 means moderate agreement.
- Kappa from 0.61 to 0.80 means good agreement.
- Kappa from 0.81 to 1.00 means very good agreement.

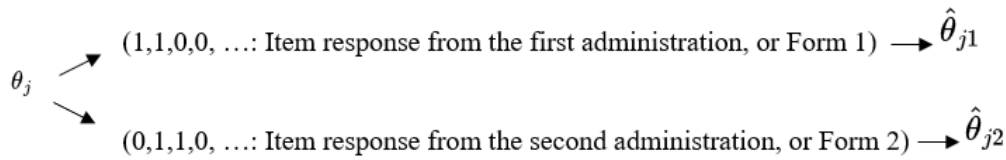
While classification *consistency* refers to the agreement between two observed scores, classification *accuracy* refers to the agreement between the observed score and the true score. Classification accuracy is defined as the extent to which the actual classifications of test takers agree with the classifications that would be made on the basis of their true scores (Livingston & Lewis, 1995). It is common to estimate classification accuracy by assuming the psychometric model to find true scores that correspond to observed scores. For the Wisconsin Forward Exam, the method used to estimate classification accuracy and consistency is the Kolen and Kim method (2004), which is described in the next section of this report (see also Kim, Choi, Um, & Kim, 2006; Kim, Barton, & Kim, 2007).

8.2.1 Kolen and Kim's Method for Pattern Scoring

As stated in Part 6, when IRT is applied to score examinees' responses, two types of scoring are available: number correct scoring and item pattern scoring. The Wisconsin Forward Exam uses item-pattern scoring. Many methods of estimating the consistency and accuracy of classification based on number-correct scoring have been suggested in psychometric literature. However, there have been relatively few studies dealing with item-pattern scoring based on IRT. Kolen and Kim (2004) suggest a simple procedure for pattern scoring (KKM) based on IRT and simulated item responses. The procedure is described below and was implemented with KKCLASS software (Kim, 2005):

Step 1: Obtain item parameters (I) and the ability distribution weight ($\hat{g}(\theta)$) at each quadrature point.

Step 2: Compute two ability estimates at each quadrature point. At a given quadrature point, θ_j , generate two sets of item responses using the item parameters from a test form, assuming that the same test form was administered twice to an examinee with the true ability θ_j .



If two parallel (or alternative) forms (e.g., Form 1 and Form 2) are available, the two response patterns can be generated based on the item parameters from the two forms. Each of the two response patterns will result in an estimated ability $\hat{\theta}_{j1}$ and $\hat{\theta}_{j2}$

Step 3: Construct a classification matrix at each quadrature point. Determine the joint event for the cells in Table 8-B by using the two ability estimates obtained in Step 2. Note that this table is constructed for each quadrature point and replication. One, and only one, cell will have a value of one and zeros elsewhere.

Table 8-B Example Classification Table for One Cut Point (C_1)

Classification		First Administration, or Form 1	
		$\hat{\theta}_{j1} \geq C_1$	$\hat{\theta}_{j1} < C_1$
Second Administration, or Form 2	$\hat{\theta}_{j2} \geq C_1$		
	$\hat{\theta}_{j2} < C_1$		

Step 4: Repeat Steps 2 and 3 R times and get average values over R replications. R should be a large number (e.g., 500) to obtain stable results.

Step 5: Multiply distribution weight ($\hat{g}(\theta)$) by the average values in Step 4 for each quadrature point and sum across all quadrature points. From this, a final contingency table and classification consistency indices, such as kappa, can be computed.

Because the examinees' abilities are estimated at each quadrature point, these quadrature points can be considered the true scores. Therefore, classification accuracy is computed using the examinees' estimated abilities (observed scores) and quadrature points (true scores). Just as 0.90 is generally considered the criterion for acceptable test score reliability, the criterion value of 0.90 is considered to be an acceptably high level of classification accuracy.

In Tables 8-5 through 8-21, there are two tables for each grade and content area. The first table is a contingency table with all three cut scores that was prepared based on the KKM procedure. The rows represent the first administration of an assessment, and the columns represent the second administration of the same assessment to the same students. As mentioned above, in the KKM procedure, the score distributions for the first administration and the second administration are estimated using a simulation. Therefore, the value in each cell

represents the probability of belonging to a particular pair of performance levels in the first administration and the second administration. For example, when considering the first column of data in the ELA grade 3 table, 0.21 represents the probability of belonging in *Below Basic* in both the first and second administrations. The 0.05 value represents the probability of belonging in *Below Basic* in the first administration and *Basic* in the second administration. The probability of belonging in *Below Basic* in the first administration and *Proficient* or *Advanced* in the second administration is 0.00. “Sum” is obtained simply by adding the four row values or the four column values. Because the values displayed have been rounded to two decimal places, this sum is not always identical to the sum of the values shown in the table.

The second table shows indices for classification consistency and classification accuracy. Because there are four performance levels for the Wisconsin Forward Exam, there are three cut scores. The values in “All Cuts” were obtained by applying all three cuts together. In Table 8-5 for ELA grade 3, when all three cuts are used for the computation, classification consistency (P) is 0.74, probability of chance is 0.29, kappa (*k*) is 0.63, and classification accuracy is 0.81. The values for “Cut 1” were obtained by applying only the first cut score. There are two levels whenever only one cut is applied (i.e., performance levels above and below the cut). It is clear that the values for P, *k*, and classification accuracy with all three cuts are smaller than those for any single cut point. The probability of assigning students to the incorrect performance level will increase with the number of cut scores.

Because the *Proficient* cut score is a criterion for accountability reports, the reliability values for this second cut need to be considered carefully. In Table 8-5, for example, the P for the second cut, which establishes the *Proficient* performance level, was 0.90; kappa was 0.78; and classification accuracy was 0.92. The interpretation of the values illustrated for Table 8-5 is the same for Tables 8-6 through 8-21.

As shown in Tables 8-5 through 8-21, when only the *Proficient* cut score was applied, the classification consistency (P) was greater than or equal to 0.88 and the classification accuracy was greater than or equal to 0.91 for all tests. The kappa value was greater than or equal to 0.76 for all tests. According to criteria for kappa (*k*) values (presented earlier in this section of the report in the discussion of classification consistency), all tests showed good or very good agreement based on the cut for the *Proficient* performance level.

In addition, the indices for classification consistency and classification accuracy were computed for the subgroups of students. These data are presented in Appendix J. As seen in Tables J-1 through J-17, when the *Proficient* cut is considered, classification consistency, accuracy coefficients, and kappa values were good or very good for all subgroups, grades, and content areas. Specifically, for ELA, the classification consistency was greater than or equal to 0.85 and the classification accuracy was greater than or equal to 0.89 for all subgroups across all grades. For Mathematics, the classification consistency was greater than or equal to 0.89 and the classification accuracy was greater than or equal to 0.92 for all subgroups across all grades. For Science, the classification consistency was greater than or equal to 0.86 and the classification accuracy was greater than or equal to 0.89 for all subgroups across both grades. For Social Studies, the classification consistency was greater than or equal to 0.86 and the classification accuracy was greater than or equal to 0.91 for all subgroups across all grades.

The kappa values indicated good or very good agreement based on the cut for the *Proficient* performance level for all subgroups across all grades and content areas. The kappa values were greater than or equal to 0.67 for all subgroups in ELA, greater than or equal to 0.65 for all subgroups in Mathematics, greater than or equal to 0.66 for all subgroups in Science, and greater than or equal to 0.68 for all subgroups, except the non-binary group, in Social Studies. The kappa value for the non-binary group in Social Studies grade 8 was 0.51. The indices for classification consistency and classification accuracy were not computed for the non-binary subgroup in grades 3 through 7 across all content areas or for students using testing accommodations in Science and Social Studies. The number of students in these groups was less than 50 per grade. The indices for classification consistency and classification accuracy for the non-binary group in grades 8 and 10 (across all content areas) and for students using testing accommodations in all ELA grades should be interpreted with caution because of the low number of students in these subgroups.

8.3 Inter-rater Reliability for TDA Items

The reliability of scoring of TDA items was measured in two ways: (1) tabulations of exact and adjacent agreement of two scorers and (2) reliability coefficients. Reliability for TDA items was examined by calculating indices of inter-rater agreement, which is the degree of reliability with which the AI engine and a human scorer assign scores to a given student response. Two indices for inter-rater reliability, intraclass correlation and weighted kappa, are presented here.

Notation: To assess reliability, it is necessary to replicate the scoring process for a subset of papers. This is usually done with “blind double-reads.” Suppose that there are N responses, each of which is scored twice. The two scores of response n are denoted by X_{n1} and X_{n2} , where $n = 1, 2, \dots, N$. The resulting data may be presented in two ways: enumeration by response and cross tabulation. Table 8-C shows the enumeration by response data structure, where each row represents a single student response.

Table 8-C Data Structure 1: Enumeration by Response

Response #	Score 1	Score 2	Mean Score
1	X_{11}	X_{12}	$\bar{X}_{.1}$
2	X_{21}	X_{22}	$\bar{X}_{.2}$
...
...
N	X_{N1}	X_{N2}	$\bar{X}_{.N}$
Column Mean	$\bar{X}_{.1}$	$\bar{X}_{.2}$	$\bar{X}_{..}$

Where

$$\bar{X}_{.1} = (X_{11} + X_{12})/2$$

is the mean score for Response 1 (similarly for responses 2, 3, ... N),

$$\bar{X}_{.1} = \frac{1}{N} \sum_{n=1}^N X_{n1} = (X_{11} + X_{21} + \dots + X_{N1})/N$$

is the mean of Score 1 over all responses (similarly for Score 2), and

$$\bar{X}_{..} = \frac{1}{N} \sum_{n=1}^N (X_{n1} + X_{n2})/2$$

is the overall mean score across both scores of all responses.

As an alternative, a square table of counts may be created for each Score 1 by Score 2 (i.e., $X_{n1} \times X_{n2}$) combination. An example of this data structure is presented in Table 8-D.

Table 8-D Data Structure 2: Cross-Tabulation of Score 1 and Score 2

		Score 2				Row Total
		0	1	...	m	
Score 1	0	n_{00}	n_{01}	...	n_{0m}	n_{0+}
	1	n_{10}	n_{11}	...	n_{1m}	n_{1+}

	m	n_{m0}	n_{m1}	...	n_{mm}	n_{m+}
Column Total		n_{+0}	n_{+1}	...	n_{+m}	n_{++}

where m is the maximum score (for a rubric including zero) obtainable for an item, n_{ij} is the number of responses for which Score 1 = i and Score 2 = j , n_{i+} is the number of responses for which Score 1 = i , and n_{+j} is the number of responses for which Score 2 = j .

Formulas for the two reliability coefficients of interest are then given:

1. Intraclass correlation, ρ_{IC} , describes the percentage of overall score variance accounted for by the variance of mean response scores:

$$\rho_{IC} = \frac{\text{Var}_n(\bar{X}_n)}{\text{Var}_n(X_{n1}, X_{n2})} = \frac{\frac{1}{N-1} \sum_{n=1}^N (\bar{X}_n - \bar{X}_{..})^2}{\frac{1}{2(N-1)} \sum_{n=1}^N [(X_{n1} - \bar{X}_{..})^2 + (X_{n2} - \bar{X}_{..})^2]}$$

If agreement is perfect, $\rho_{IC} = 1$. The following is always true: $0 \leq \rho_{IC} \leq 1$.

2. Weighted kappa, k , is used in many contexts as a measure of association in square contingency tables:

$$k = \frac{\sum_{i=0}^m \sum_{j=0}^m w_{ij} \frac{n_{ij}}{n_{++}} - \sum_{i=0}^m \sum_{j=0}^m w_{ij} \frac{n_{i+n_j}}{n_{++}^2}}{1 - \sum_{i=0}^m \sum_{j=0}^m w_{ij} \frac{n_{i+n_j}}{n_{++}^2}}, \text{ where } w_{ij} = 1 - \frac{(i-j)^2}{M^2}.$$

If agreement is perfect, $k = 1$. If agreement is what would be expected by chance, $k = 0$. The following is always true: $0 \leq k \leq 1$.

Ordinal rating scales (e.g., 1, 2, 3, 4) used in scoring TDA items contain a certain level of chance agreement that is expected. Although the intraclass correlation is reported in this report, it does not take into account the possibility of chance agreement between the two raters. Cohen’s kappa does take this into consideration. In general, k will have values equal to or less than the intraclass correlation. If agreement is exact, the value of k is 1.0. If agreement is at chance levels, the value of k is 0. As noted in Section 8.2, values of k greater than 0.81 indicate “very good agreement,” values between 0.61 and 0.80 represent “good agreement,” values between 0.41 and 0.60 represent “moderate agreement,” values between 0.21 and 0.40 represent “fair agreement” beyond chance, and values below 0.20 denote “poor agreement.” Specific criteria for intraclass correlation or weighted k are not established.

Table 8-22 presents the rater agreement statistics for TDA items. The evidence supporting inter-rater reliability is presented in terms of the percentage of agreement between raters (the AI engine and a human rater), two indices of inter-rater reliability, and the distributions of scores across score levels. In the table, “Exact” agreement is defined as scores that are exactly the same. “Adjacent” agreement is defined as scores differing by 1 point. “Discrepant” cases are those cases in which the scores of the two raters differed by more than one raw score point. For example, as shown in Table 8-22, for the grade 3 TDA item, the exact agreement, adjacent agreement, and discrepant agreement rates are 88.69%, 10.99%, and 0.32%, respectively. “Mean” reflects the item mean score from the second reads, which are done by human scorers. “No. of Second Reads” is the number of student responses selected for the purpose of the second read and computing inter-rater reliability. The “Score Frequency” columns represent the scoring outcomes for the student responses based on the raw scores given by the human scorers. The column for “All Codes” reflects the number of students who received the condition codes B, C, N, R, or T (described in detail in Part 5, Table 5-2 of this report).

Overall, the exact rater agreement percentages were high for all TDA items and ranged from 76.18% in grade 8 to 88.69% in grade 3. The combined exact and adjacent agreement percentages were over 99% in grades 3 through 7 and over 98% in grade 8. The intraclass correlation coefficients ranged from 0.86 in grade 3 to 0.93 in grades 6 and 7. The weighted kappa ranged from 0.73 in grade 3 to 0.87 in grade 7, indicating good or very good rater agreement for all TDA items.

8.4 Summary

Overall, the analyses discussed in this section of the report indicated acceptable levels of reliability for the Wisconsin Forward Exam. The internal consistency reliability estimates, as measured by Cronbach’s alpha coefficient, were reasonable given the number of items in each test. The analyses of classification consistency and accuracy indicated acceptable levels of consistency and accuracy of student proficiency level classifications, and the SEM around the *Proficient* cut score was low in every grade and content area. The levels of rater agreement were high, and the discrepancy rates were low, with acceptably high values for the weighted kappa and intraclass correlations. The results of the inter-rater reliability analyses indicated an acceptable degree of reliability for scores on the ELA TDA items in the Wisconsin Forward Exam.

Table 8-1 Cronbach’s Alpha Reliability Coefficients for Total Group and Subgroups

Content	Grade	Total	Gender			Race/Ethnicity						ELP		Disability		Economic Status		Accommodations	
			Female	Male	Non-Binary	White	African American	Hispanic	Asian	American Indian	Two or More	Fully English Proficient	Limited English Proficient	Disabled	Not Disabled	Economically Disadvantaged	Not Economically Disadvantaged	Students with Accommodations	Students without Accommodations
English Language Arts	3	0.88	0.88	0.88		0.87	0.82	0.86	0.88	0.84	0.87	0.88	0.82	0.85	0.88	0.85	0.87	0.90	0.88
	4	0.89	0.89	0.89		0.88	0.84	0.87	0.89	0.85	0.89	0.89	0.83	0.87	0.88	0.87	0.88	0.90	0.89
	5	0.88	0.88	0.87		0.86	0.81	0.85	0.88	0.83	0.87	0.88	0.79	0.84	0.87	0.85	0.86	0.90	0.88
	6	0.89	0.88	0.89		0.87	0.85	0.87	0.88	0.86	0.89	0.88	0.82	0.85	0.88	0.87	0.87	0.88	0.89
	7	0.88	0.87	0.88		0.86	0.85	0.87	0.87	0.85	0.88	0.88	0.82	0.84	0.87	0.86	0.86	0.89	0.88
	8	0.89	0.88	0.89	0.85	0.88	0.86	0.87	0.89	0.86	0.89	0.89	0.82	0.84	0.88	0.87	0.88	0.88	0.88
Mathematics	3	0.93	0.93	0.93		0.92	0.89	0.91	0.93	0.90	0.93	0.93	0.90	0.93	0.93	0.92	0.92	0.84	0.93
	4	0.92	0.92	0.93		0.92	0.85	0.90	0.93	0.89	0.92	0.92	0.87	0.91	0.92	0.90	0.92	0.75	0.92
	5	0.92	0.91	0.93		0.91	0.86	0.90	0.93	0.88	0.92	0.92	0.87	0.91	0.92	0.90	0.91	0.80	0.92
	6	0.92	0.92	0.93		0.91	0.84	0.89	0.93	0.90	0.92	0.92	0.84	0.88	0.92	0.90	0.92	0.76	0.92
	7	0.91	0.91	0.92		0.90	0.84	0.88	0.92	0.87	0.91	0.91	0.82	0.86	0.91	0.88	0.91	0.66	0.91
	8	0.92	0.91	0.92	0.90	0.91	0.82	0.88	0.93	0.87	0.92	0.92	0.80	0.85	0.92	0.89	0.92	0.66	0.92
Science	4	0.91	0.90	0.91		0.89	0.85	0.89	0.90	0.88	0.91	0.91	0.85	0.89	0.90	0.89	0.89		0.91
	8	0.89	0.88	0.90	0.85	0.88	0.83	0.87	0.89	0.86	0.89	0.89	0.79	0.86	0.88	0.87	0.88		0.89
Social Studies	4	0.91	0.90	0.91		0.89	0.86	0.89	0.91	0.87	0.90	0.91	0.86	0.90	0.90	0.89	0.90		0.91
	8	0.91	0.90	0.92	0.85	0.90	0.87	0.89	0.91	0.88	0.91	0.91	0.83	0.88	0.90	0.89	0.90		0.91
	10	0.90	0.89	0.91	0.91	0.90	0.85	0.88	0.90	0.87	0.90	0.90	0.76	0.85	0.90	0.88	0.90		0.90

Note: The reliability coefficients were not computed for non-binary students in grades 3 through 7 or students using testing accommodations in Science or Social Studies because the number of students in these grades and subject areas was less than 50 per grade.

Table 8-2 Standard Error of Measurement for Total Group and Subgroups

Content	Grade	Total	Gender			Race/Ethnicity						ELP		Disability		Economic Status		Accommodations	
			Female	Male	Non-Binary	White	African American	Hispanic	Asian	American Indian	Two or More	Fully English Proficient	Limited English Proficient	Disabled	Not Disabled	Economically Disadvantaged	Not Economically Disadvantaged	Students with Accommodations	Students without Accommodations
English Language Arts	3	3.34	3.33	3.34		3.29	3.37	3.38	3.37	3.39	3.36	3.33	3.40	3.37	3.32	3.39	3.27	3.35	3.34
	4	3.39	3.39	3.39		3.37	3.37	3.40	3.40	3.41	3.38	3.39	3.39	3.36	3.39	3.40	3.36	3.42	3.39
	5	3.48	3.47	3.47		3.45	3.45	3.49	3.49	3.48	3.49	3.47	3.46	3.43	3.46	3.48	3.44	3.50	3.48
	6	3.46	3.44	3.46		3.42	3.48	3.49	3.48	3.50	3.47	3.45	3.51	3.44	3.44	3.49	3.40	3.49	3.46
	7	3.53	3.52	3.51		3.48	3.55	3.57	3.53	3.57	3.52	3.52	3.53	3.44	3.50	3.55	3.47	3.60	3.53
	8	3.55	3.53	3.56	3.38	3.49	3.57	3.61	3.50	3.57	3.57	3.54	3.59	3.49	3.52	3.60	3.48	3.66	3.55
Mathematics	3	2.67	2.68	2.64		2.65	2.61	2.69	2.66	2.72	2.69	2.66	2.70	2.65	2.67	2.70	2.63	2.52	2.67
	4	2.85	2.86	2.84		2.85	2.78	2.86	2.83	2.86	2.85	2.85	2.86	2.83	2.85	2.86	2.83	2.76	2.85
	5	2.92	2.93	2.91		2.95	2.71	2.87	2.90	2.91	2.91	2.93	2.83	2.77	2.94	2.88	2.94	2.62	2.93
	6	2.88	2.89	2.87		2.89	2.76	2.88	2.86	2.87	2.87	2.89	2.81	2.79	2.89	2.88	2.87	2.68	2.89
	7	2.92	2.92	2.92		2.94	2.76	2.89	2.91	2.88	2.92	2.93	2.82	2.77	2.94	2.89	2.93	2.65	2.93
	8	2.89	2.90	2.87	2.94	2.91	2.74	2.85	2.85	2.83	2.86	2.89	2.77	2.72	2.90	2.85	2.90	2.63	2.90
Science	4	2.68	2.69	2.66		2.67	2.62	2.70	2.69	2.71	2.68	2.67	2.71	2.66	2.68	2.70	2.65		2.68
	8	2.78	2.79	2.76	2.74	2.77	2.75	2.80	2.75	2.80	2.80	2.78	2.76	2.74	2.78	2.80	2.76		2.78
Social Studies	4	2.71	2.71	2.70		2.68	2.75	2.76	2.72	2.78	2.73	2.70	2.79	2.73	2.70	2.77	2.66		2.71
	8	2.69	2.69	2.67	2.50	2.65	2.79	2.76	2.64	2.78	2.71	2.68	2.81	2.75	2.67	2.77	2.62		2.69
	10	2.79	2.80	2.77	2.65	2.76	2.79	2.83	2.79	2.83	2.79	2.79	2.79	2.75	2.79	2.83	2.75		2.79

Note: The SEMs were not computed for non-binary students in grades 3 through 7 or students using testing accommodations in Science or Social Studies because the number of students in these grades and subject areas was less than 50 per grade.

Table 8-3 Cronbach’s Alpha Reliability Coefficients for Content Standards and Domains

English Language Arts

Grade	Alpha per Content Standard and Domain									
	A	B	C	D	E	F	Listening	Reading	Writing	Total Test
3	0.69	0.55	0.50	0.53	0.46	0.34	0.48	0.81	0.71	0.88
4	0.59	0.59	0.58	0.55	0.57	0.46	0.48	0.81	0.77	0.89
5	0.64	0.56	0.50	0.40	0.47	0.45	0.56	0.80	0.70	0.88
6	0.65	0.72	0.49	0.38	0.39	0.40	0.51	0.84	0.67	0.89
7	0.55	0.54	0.64	0.40	0.37	0.43	0.66	0.79	0.67	0.88
8	0.59	0.61	0.60	0.46	0.50	0.35	0.59	0.81	0.72	0.89

ELA standards: A = Reading—Key Ideas and Details; B = Reading—Craft & Structure/Integration of Knowledge & Ideas; C = Reading—Vocabulary Use; D = Writing/Language—Text Types and Purposes; E = Writing/Language—Research; F = Writing/Language—Language Conventions

Mathematics

Grade	Alpha per Content Standard										
	A	B	C	D	E	F	G	H	I	J	Total Test
3	0.78	0.77	0.72	0.76	0.64						0.93
4	0.77	0.74	0.76	0.69	0.62						0.92
5	0.74	0.77	0.63	0.75	0.64						0.92
6					0.71	0.73	0.75	0.77	0.61		0.92
7					0.55	0.77	0.67	0.72	0.69		0.91
8					0.69		0.72	0.69	0.71	0.74	0.92

Mathematics standards: A = Operations and Algebraic Thinking; B = Number and Operations in Base Ten; C = Number and Operations—Fractions; D = Measurement and Data; E = Geometry; F = Ratios and Proportional Relationships; G = The Number System; H = Expressions and Equations; I = Statistics and Probability; J = Functions

Science

Grade	Alpha per Content Standard				Total Test
	A	B	C	D	
4	0.73	0.71	0.66	0.74	0.91
8	0.69	0.70	0.65	0.66	0.89

Science standards: A = Life Science; B = Physical Science; C = Earth and Space Science; D = Engineering.

Social Studies

Grade	Alpha per Content Standard					Total Test
	A	B	C	D	E	
4	0.77	0.68	0.47	0.61	0.70	0.91
8	0.65	0.69	0.72	0.66	0.64	0.91
10	0.72	0.68	0.64	0.54	0.62	0.90

Social Studies standards: A = Geography; B = History; C = Political Science; D = Economics; E = Behavioral Sciences

Table 8-4 Standard Error of Measurement per Content Standards and Domains

English Language Arts

Grade	SEM per Content Standard and Domain									
	A	B	C	D	E	F	Listening	Reading	Writing	Total Test
3	1.65	1.18	0.82	1.14	1.10	1.06	1.27	2.21	1.92	3.34
4	1.30	1.39	1.05	1.15	1.04	1.01	1.39	2.18	1.86	3.39
5	1.58	1.49	0.80	1.05	1.00	1.18	1.41	2.32	1.88	3.48
6	1.40	1.41	0.82	1.21	1.21	1.09	1.38	2.15	2.03	3.46
7	1.45	1.39	0.86	1.43	1.20	0.94	1.26	2.18	2.09	3.53
8	1.59	1.23	0.94	1.29	1.22	1.12	1.23	2.21	2.10	3.55

ELA standards: A = Reading—Key Ideas and Details; B = Reading—Craft & Structure/Integration of Knowledge & Ideas; C = Reading—Vocabulary Use; D = Writing/Language—Text Types and Purposes; E = Writing/Language—Research; F = Writing/Language—Language Conventions

Mathematics

Grade	SEM per Content Standard										
	A	B	C	D	E	F	G	H	I	J	Total Test
3	1.18	1.12	1.15	1.34	1.12						2.67
4	1.29	1.21	1.33	1.37	1.10						2.85
5	1.27	1.24	1.34	1.32	1.32						2.92
6					1.08	1.03	1.42	1.40	1.41		2.88
7					1.43	1.14	1.06	1.35	1.45		2.92
8					1.34		1.17	1.38	1.18	1.32	2.89

Mathematics standards: A = Operations and Algebraic Thinking; B = Number and Operations in Base Ten; C = Number and Operations—Fractions; D = Measurement and Data; E = Geometry; F = Ratios and Proportional Relationships; G = The Number System; H = Expressions and Equations; I = Statistics and Probability; J = Functions

Science

Grade	SEM per Content Standard				Total Test
	A	B	C	D	
4	1.45	1.38	1.30	1.19	2.68
8	1.48	1.55	1.28	1.22	2.78

Science standards: A = Life Science; B = Physical Science; C = Earth and Space Science; D = Engineering

Social Studies

Grade	SEM per Content Standard					Total Test
	A	B	C	D	E	
4	1.36	1.28	1.12	1.17	1.09	2.71
8	1.29	1.35	1.08	1.14	1.11	2.69
10	1.36	1.22	1.26	1.19	1.16	2.79

Social Studies standards: A = Geography; B = History; C = Political Science; D = Economics; E = Behavioral Sciences

Table 8-5 Classification Consistency and Classification Accuracy for English Language Arts Grade 3

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.21	0.05	0.00	0.00	0.27
Basic	0.05	0.26	0.06	0.00	0.36
Proficient	0.00	0.05	0.22	0.03	0.29
Advanced	0.00	0.00	0.03	0.05	0.08
Sum	0.26	0.36	0.30	0.08	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.90	0.90	0.94	0.74
Probability of Chance	0.61	0.53	0.86	0.29
Kappa (<i>k</i>)	0.73	0.78	0.61	0.63
Classification Accuracy	0.93	0.92	0.96	0.81

Table 8-6 Classification Consistency and Classification Accuracy for English Language Arts Grade 4

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.19	0.05	0.00	0.00	0.24
Basic	0.05	0.21	0.05	0.00	0.31
Proficient	0.00	0.05	0.25	0.03	0.34
Advanced	0.00	0.00	0.03	0.08	0.11
Sum	0.23	0.31	0.34	0.12	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.90	0.90	0.94	0.74
Probability of Chance	0.64	0.50	0.80	0.28
Kappa (<i>k</i>)	0.73	0.79	0.68	0.63
Classification Accuracy	0.93	0.93	0.95	0.81

Table 8-7 Classification Consistency and Classification Accuracy for English Language Arts Grade 5

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.23	0.06	0.00	0.00	0.28
Basic	0.05	0.23	0.05	0.00	0.33
Proficient	0.00	0.05	0.23	0.03	0.31
Advanced	0.00	0.00	0.03	0.05	0.08
Sum	0.28	0.34	0.31	0.07	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.89	0.90	0.95	0.74
Probability of Chance	0.60	0.53	0.86	0.29
Kappa (<i>k</i>)	0.74	0.79	0.60	0.63
Classification Accuracy	0.92	0.92	0.96	0.81

Table 8-8 Classification Consistency and Classification Accuracy for English Language Arts Grade 6

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.22	0.04	0.00	0.00	0.26
Basic	0.04	0.23	0.06	0.00	0.34
Proficient	0.00	0.05	0.20	0.04	0.30
Advanced	0.00	0.00	0.04	0.07	0.11
Sum	0.26	0.33	0.30	0.11	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.91	0.89	0.92	0.72
Probability of Chance	0.62	0.52	0.81	0.28
Kappa (<i>k</i>)	0.77	0.76	0.60	0.62
Classification Accuracy	0.94	0.92	0.95	0.80

Table 8-9 Classification Consistency and Classification Accuracy for English Language Arts Grade 7

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.20	0.04	0.00	0.00	0.24
Basic	0.04	0.24	0.05	0.00	0.33
Proficient	0.00	0.07	0.23	0.03	0.33
Advanced	0.00	0.00	0.03	0.06	0.10
Sum	0.24	0.35	0.31	0.09	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.91	0.88	0.94	0.73
Probability of Chance	0.63	0.51	0.83	0.29
Kappa (<i>k</i>)	0.76	0.76	0.64	0.62
Classification Accuracy	0.94	0.91	0.95	0.80

Table 8-10 Classification Consistency and Classification Accuracy for English Language Arts Grade 8

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.23	0.04	0.00	0.00	0.27
Basic	0.04	0.26	0.05	0.00	0.34
Proficient	0.00	0.06	0.18	0.03	0.27
Advanced	0.00	0.00	0.04	0.07	0.11
Sum	0.27	0.36	0.27	0.10	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.92	0.89	0.93	0.74
Probability of Chance	0.61	0.53	0.81	0.28
Kappa (<i>k</i>)	0.79	0.77	0.65	0.64
Classification Accuracy	0.94	0.93	0.95	0.82

Table 8-11 Classification Consistency and Classification Accuracy for Mathematics Grade 3

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.19	0.04	0.00	0.00	0.23
Basic	0.04	0.21	0.05	0.00	0.30
Proficient	0.00	0.05	0.25	0.03	0.33
Advanced	0.00	0.00	0.03	0.11	0.14
Sum	0.23	0.29	0.34	0.14	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.93	0.90	0.94	0.76
Probability of Chance	0.64	0.50	0.76	0.27
Kappa (<i>k</i>)	0.79	0.80	0.73	0.67
Classification Accuracy	0.95	0.93	0.96	0.83

Table 8-12 Classification Consistency and Classification Accuracy for Mathematics Grade 4

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.19	0.04	0.00	0.00	0.23
Basic	0.05	0.26	0.03	0.00	0.34
Proficient	0.00	0.04	0.19	0.02	0.25
Advanced	0.00	0.00	0.02	0.16	0.18
Sum	0.23	0.34	0.24	0.18	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.91	0.93	0.96	0.80
Probability of Chance	0.65	0.51	0.70	0.26
Kappa (<i>k</i>)	0.75	0.85	0.87	0.73
Classification Accuracy	0.94	0.94	0.96	0.83

Table 8-13 Classification Consistency and Classification Accuracy for Mathematics Grade 5

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.22	0.04	0.00	0.00	0.26
Basic	0.04	0.18	0.05	0.00	0.27
Proficient	0.00	0.05	0.27	0.03	0.35
Advanced	0.00	0.00	0.03	0.09	0.12
Sum	0.26	0.27	0.35	0.12	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.92	0.90	0.94	0.76
Probability of Chance	0.61	0.50	0.79	0.28
Kappa (<i>k</i>)	0.80	0.80	0.72	0.67
Classification Accuracy	0.94	0.93	0.96	0.83

Table 8-14 Classification Consistency and Classification Accuracy for Mathematics Grade 6

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.24	0.05	0.00	0.00	0.28
Basic	0.05	0.21	0.05	0.00	0.30
Proficient	0.00	0.04	0.28	0.02	0.34
Advanced	0.00	0.00	0.02	0.06	0.08
Sum	0.28	0.30	0.34	0.08	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.91	0.91	0.96	0.78
Probability of Chance	0.59	0.51	0.86	0.29
Kappa (<i>k</i>)	0.77	0.81	0.71	0.68
Classification Accuracy	0.93	0.94	0.97	0.84

Table 8-15 Classification Consistency and Classification Accuracy for Mathematics Grade 7

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.32	0.04	0.00	0.00	0.36
Basic	0.05	0.20	0.05	0.00	0.29
Proficient	0.00	0.04	0.24	0.02	0.30
Advanced	0.00	0.00	0.02	0.04	0.05
Sum	0.37	0.28	0.30	0.05	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.91	0.91	0.96	0.79
Probability of Chance	0.54	0.54	0.90	0.30
Kappa (<i>k</i>)	0.80	0.81	0.65	0.69
Classification Accuracy	0.93	0.93	0.98	0.84

Table 8-16 Classification Consistency and Classification Accuracy for Mathematics Grade 8

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.30	0.05	0.00	0.00	0.35
Basic	0.06	0.24	0.04	0.00	0.33
Proficient	0.00	0.03	0.19	0.02	0.24
Advanced	0.00	0.00	0.02	0.05	0.07
Sum	0.35	0.33	0.25	0.07	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.89	0.92	0.96	0.78
Probability of Chance	0.54	0.57	0.87	0.30
Kappa (<i>k</i>)	0.76	0.82	0.71	0.68
Classification Accuracy	0.92	0.94	0.97	0.84

Table 8-17 Classification Consistency and Classification Accuracy for Science Grade 4

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.15	0.05	0.00	0.00	0.20
Basic	0.04	0.19	0.05	0.00	0.29
Proficient	0.00	0.06	0.21	0.04	0.31
Advanced	0.00	0.00	0.04	0.16	0.21
Sum	0.19	0.30	0.30	0.21	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.91	0.89	0.91	0.71
Probability of Chance	0.68	0.50	0.67	0.26
Kappa (<i>k</i>)	0.71	0.79	0.73	0.61
Classification Accuracy	0.94	0.92	0.93	0.79

Table 8-18 Classification Consistency and Classification Accuracy for Science Grade 8

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.18	0.05	0.00	0.00	0.23
Basic	0.05	0.18	0.05	0.00	0.28
Proficient	0.00	0.06	0.18	0.04	0.28
Advanced	0.00	0.00	0.04	0.17	0.21
Sum	0.23	0.29	0.27	0.21	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.90	0.89	0.91	0.70
Probability of Chance	0.65	0.50	0.67	0.25
Kappa (<i>k</i>)	0.71	0.78	0.74	0.60
Classification Accuracy	0.93	0.92	0.93	0.78

Table 8-19 Classification Consistency and Classification Accuracy for Social Studies Grade 4

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.17	0.04	0.00	0.00	0.22
Basic	0.04	0.11	0.05	0.00	0.19
Proficient	0.00	0.05	0.25	0.06	0.35
Advanced	0.00	0.00	0.05	0.19	0.24
Sum	0.22	0.19	0.34	0.25	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.91	0.90	0.90	0.71
Probability of Chance	0.66	0.52	0.63	0.26
Kappa (<i>k</i>)	0.74	0.79	0.72	0.61
Classification Accuracy	0.94	0.93	0.93	0.80

Table 8-20 Classification Consistency and Classification Accuracy for Social Studies Grade 8

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.16	0.04	0.00	0.00	0.21
Basic	0.04	0.12	0.05	0.00	0.20
Proficient	0.00	0.05	0.24	0.05	0.34
Advanced	0.00	0.00	0.05	0.21	0.25
Sum	0.20	0.21	0.33	0.26	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.92	0.90	0.91	0.73
Probability of Chance	0.68	0.52	0.62	0.26
Kappa (<i>k</i>)	0.74	0.80	0.75	0.63
Classification Accuracy	0.94	0.93	0.93	0.80

Table 8-21 Classification Consistency and Classification Accuracy for Social Studies Grade 10

Contingency Table with All Cut Scores

Performance Level	Below Basic	Basic	Proficient	Advanced	Sum
Below Basic	0.22	0.05	0.00	0.00	0.28
Basic	0.05	0.15	0.05	0.00	0.25
Proficient	0.00	0.05	0.15	0.05	0.25
Advanced	0.00	0.00	0.04	0.19	0.23
Sum	0.27	0.25	0.24	0.24	1.00

Indices for Classification Consistency and Classification Accuracy

Indices	Cut 1	Cut 2	Cut 3	All Cuts
Classification Consistency (P)	0.89	0.90	0.91	0.71
Probability of Chance	0.60	0.50	0.64	0.25
Kappa (<i>k</i>)	0.74	0.79	0.74	0.61
Classification Accuracy	0.93	0.93	0.93	0.79

Table 8-22 Inter-Rater Reliability, English Language Arts

Grade	Item No.	Max	Percentage of Agreement			Intra. Corr.	Weighted Kappa	Mean	Score Frequency					
			Exact	Adjacent	Discrepant				No. of Second	1	2	3	4	All Codes
3	1	4	88.69%	10.99%	0.32%	0.86	0.73	1.25	15709	7436	1260	154	22	6837
4	1	4	81.96%	17.19%	0.86%	0.89	0.79	1.49	15001	5375	2058	615	104	6849
5	1	4	82.27%	17.13%	0.61%	0.88	0.77	1.46	9759	4649	1997	401	53	2659
6	1	4	85.16%	14.61%	0.23%	0.93	0.85	1.50	11368	5231	2524	817	150	2646
7	1	4	84.88%	14.76%	0.36%	0.93	0.87	1.79	15813	4739	5234	1714	212	3914
8	1	4	76.18%	22.34%	1.48%	0.91	0.81	1.83	13383	4582	3024	1333	260	4184

Note: The sum of the modes of agreement and codes may not equal exactly 100% due to rounding.

Note: TDA item scores presented in this table reflect a 1–4-point scoring rubric (before application of a weight of 2).

Part 9: Studies of Construct-Related Validity

As stated in Part 2 of this Technical Report, validity is the overarching component of the Wisconsin Forward Exam program. The following excerpt is from the *Standards* (AERA, APA, & NCME, 2014):

Ultimately, the validity of an intended interpretation of test scores relies on all the available evidence relevant to the technical quality of a testing system. Different components of validity evidence . . . include evidence of careful test construction; adequate score reliability; appropriate test administration and scoring; accurate score scaling, equating, and standard setting; and careful attention to fairness for all test takers, as appropriate to the test interpretation in question. (p. 22)

Part 9 addresses four additional issues related to the evidence of the validity of an intended interpretation of test scores: test fairness, evidence of validity based on the internal structure of the test, evidence of validity based on the relationship between test scores and other variables, and test integrity. In the subsequent pages, Part 9 will, as stated, present additional metrics with which to evaluate the validity of an intended interpretation of test scores of the Wisconsin Forward Exam program.

As described below, the Wisconsin Forward Exam program formally assessed the issue of test fairness through an analysis of differential item functioning (DIF). It is possible for items to function differently across different population groups, and it is also possible that results for an item do not reflect student ability but instead reflect irrelevant information influenced by demographic factors. The DIF analysis provided below serves to determine whether that possibility occurred and, if so, to what degree, item by item, for each of the categories of gender, race/ethnicity, economic status, disability status, accommodation use, and English language proficiency.

This part is particularly relevant to AERA, APA, & NCME (2014) Standards 3.1, 3.2, 3.3, and 3.6. Each of these standards and the way in which the standard is addressed will be presented in this part.

Standard 3.6 Where credible evidence indicates that test scores may differ in meaning for relevant subgroups in the intended examinee population, test developers and/or users are responsible for examining the evidence for validity of score interpretations for intended uses for individuals from those subgroups. What constitutes a significant difference in subgroup scores and what actions are taken in response to such differences may be defined by applicable laws. (p. 65)

There is no particular research on the Wisconsin Forward Exam showing that the test scores of examinee subgroups differ in meaning; however, this is an ongoing concern in any large-scale testing program. To lessen the possibility of differences in test score meaning, DRC follows multiple best practices of the testing industry in item development and selection, as is explained in Part 3. These practices adhere to AERA, APA, & NCME (2014) Standards 3.1, 3.2, and 3.3:

Standard 3.1 Those responsible for test development, revision, and administration should design all steps of the testing process to promote valid score interpretations for intended score uses for the widest possible range of individuals and relevant subgroups in the intended population. (p. 63)

Standard 3.2 Test developers are responsible for developing tests that measure the intended construct and for minimizing the potential for tests' being affected by construct irrelevant characteristics, such as linguistic, communicative, cognitive, cultural, physical, or other characteristics. (p. 64)

Standard 3.3 Those responsible for test development should include relevant subgroups in validity, reliability/precision, and other preliminary studies used when constructing the test. (p. 64)

DRC conducted DIF studies following the operational administration of the Wisconsin Forward Exam. Items are first evaluated for possible DIF in the field test phase of test development and again after their operational administration. Items flagged for DIF are further examined for possible bias. All items on the Spring 2023 test forms were field-tested or previously administered operationally in Wisconsin, and DIF analyses were already performed on these items after their first administration. The DIF analyses were repeated in Spring 2023. Items flagged for DIF were again evaluated by DRC content experts for potential bias. Section 9.1 of this part of the Technical Report explains the steps taken to evaluate the Wisconsin Forward Exam items through the use of DIF.

Section 9.2 of the report provides the evidence of the validity of an intended interpretation of test scores related to test construct. Two measures of the test internal structure are provided: correlations between content area reporting category (standard) scores and principal component analysis. Both of these measures are provided to demonstrate the existence of a single, underlying trait or ability for each content area, such as ELA ability or Mathematics ability. The presence of a single, underlying trait is a fundamental issue when scaling and analyzing results through IRT models. Therefore, these analyses are essential elements in assessing the validity of the Wisconsin Forward Exam.

In Section 9.3, the relationship between the Wisconsin Forward Exam scores and other variables is explored in order to support the evidence of the validity of an intended interpretation of test scores. These measures include evaluation of the correlations of the content area scores with other content area scores for the total population and by subgroups. They also include comparisons of student performance on the Wisconsin Forward Exam with performance on the NAEP.

In addition, Section 9.4 provides an overview of the forensic analysis procedures that were employed to ensure the integrity of test scores by identifying schools and individual students that might have engaged in inappropriate behaviors during testing.

9.1 Differential Item Functioning

An empirical DIF approach was used to examine potential item bias and to determine whether item performance differences between identifiable subgroups were due to extraneous or construct-irrelevant information, making the items unfairly difficult for a particular subgroup in the student population. An item was flagged for DIF when there was a significant difference in the scores between a focal group of students and a reference group of students, with both groups at the same overall ability level. Thus, an item flagged for DIF is more difficult for a particular group of students than would be expected based on their total test scores (Camilli & Shepard, 1994; Green, 1975).

DIF analyses were conducted based on gender, race/ethnicity, economic status, disability status, English language proficiency, and accommodation use groups. The reference and focal groups are as follows:

- **Gender**—reference group: male students; focal group: female students
- **Race/Ethnicity**—reference group: White students; focal groups: African American, Asian, Hispanic, American Indian students
- **English language proficiency**—reference group: fully English proficient students; focal group: students of limited English proficiency
- **Disability status**—reference group: students without disabilities; focal group: students with disabilities
- **Economic status**—reference group: not economically disadvantaged students; focal group: economically disadvantaged students
- **Accommodation use**—reference group: students not using testing accommodations; focal group: students using testing accommodations

Two DIF statistics that are commonly used for this purpose are the Mantel-Haenszel (MH) statistic (1959) and the Standardized Mean Difference (SMD) between the reference and focal groups, proposed by Dorans and Schmitt (1991).

The MH statistic is computed as follows (Zwick, Donoghue, & Grima, 1993):

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

where F_k is the sum of scores for the focal group at the k level of the matching variable. Note that the MH statistic is sensitive to N such that larger sample sizes increase the value of the chi square.

In addition to the MH chi-square statistic, the delta statistic (MH-D DIF) was computed for all items. The delta statistic was developed by Educational Testing Service (Holland & Thayer, 1985, 1986). To compute delta, alpha (the odds ratio) is first computed:

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

where N_{r1k} is the number of correct responses in the reference group at ability level k , N_{f0k} is the number of incorrect responses in the focal group at ability level k , N_k is the total number of responses, N_{f1k} is the number of correct responses in the focal group at ability level k , and N_{r0k} is the number of incorrect responses in the reference group at ability level k . **MH-D DIF** is then computed:

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

For selected response items, the $MH(\chi_{MH}^2)$ statistic was used to evaluate potential DIF items. In the MH procedure, subgroups are matched by their raw total test score using a contingency table with k ability levels. When applying the MH procedure, the log-odds ratio α is assumed to be constant across the k matched levels. Then the χ_{MH}^2 estimates a pooled common odds ratio. Taking the natural logarithm of the common-odds ratio and its confidence limits and multiplying these with the constant -2.35 , the resulting values may then be placed on the MH delta metric (Δ_{MH}) for interpretive purposes. Items were flagged for DIF using the following criteria:

- Moderate DIF: Significant MH chi-square statistic ($p < 0.05$) and $1.0 < |MH-D DIF| < 1.5$
- Large DIF: Significant MH chi-square statistic ($p < 0.05$) and $|MH-D DIF| \geq 1.5$

For or non-MC items, an effect size (ES) statistic based on the MH chi-square was used. The ES is obtained by dividing the SMD statistics by the standard deviation (SD) of the item. The SMD is an ES index of DIF, which is relatively easy to interpret (Zwick et al., 1993). The SMD compares the means of the reference and focal groups, adjusting for the distribution of the reference and focal group members on the conditioning variable (Zwick et al., 1993), which for these analyses is the Wisconsin Forward Exam raw score. SMD is computed as follows (Zwick et al., 1993):

$$SMD = p_{Fk} \left(\sum_k m_{Fk} - \sum_k m_{Rk} \right),$$

where p_{Fk} = the proportion of the focal group members at the k th level of the matching variable, $m_{Fk} = 1/N_{F1k}$, and $m_{Rk} = 1/N_{R1k}$. Items are flagged using the same rules that are used in the NAEP:

- Moderate DIF: If the MH statistic is significant ($p < 0.05$) and $|ES|$ is between 0.17 and 0.25
- Large DIF: If the MH statistic is significant ($p < 0.05$) and $|ES| \geq 0.25$

A positive DIF value indicates that the item favors the focal group, while a negative value indicates that the item disadvantages the focal group. Tables 9-1 through 9-4 show the DIF results for ELA, Mathematics, Science, and Social Studies, respectively, and include items flagged for different subgroups.

A negative SMD value implies that the focal group has a lower mean item score than the reference group, whereas a positive value implies that the focal group has a higher mean item score than the reference group, conditioned on the matching test score.

The minimum case count for the focal group was set at 200, and the minimum case count for the reference group was set at 400. The DIF analyses were not performed for subgroups of fewer than 200 students. In these cases, the statistical procedures do not have sufficient power to detect differences should they exist.

Tables 9-1 through 9-4 show items that were flagged for moderate or large DIF based on the criteria described above. The B flag represents a lower threshold for DIF. Each table includes the DIF category, the reference and focal groups, the grade, the item number on the test, and the item type. The tables present the MH SMD

statistics and the Mantel-Haenszel statistics (Δ_{MH}). After specifying these statistics for each item, the final column provides a flag status. The flag is based on SMD statistics and on MH (Δ_{MH}) statistics.

In Table 9-1, looking at the gender category, six items were flagged for moderate (B flag) and one item was flagged for large (C flag) DIF in ELA grades 5, 6, 7, and 8. Of these items, four were flagged in favor of the focal group (females) and three were flagged against the focal group. When ethnicity DIF was considered across all grades, four items were flagged in favor of and five items were flagged against Asian students, one item was flagged in favor of and two items were flagged against African American students, and three items were flagged against American Indian students. One item in grade 3 and one item in grade 7 were flagged against students with limited English proficiency (moderate DIF). A total of twelve items were flagged against students with disabilities across all grades. One item in grade 5 was flagged in favor of students with disabilities. In addition, one item in ELA grade 8 was flagged against economically disadvantaged students. Of all items flagged in ELA, twenty-seven items displayed moderate DIF (Flag B) and ten items displayed large DIF (Flag C).

The other DIF results in Tables 9-2 through 9-4 can be understood in the same fashion. Note that a single item can be flagged for multiple subgroup categories, such as for ethnicity and language proficiency.

When looking at DIF results by item type, it was observed that most of the flagged items were MC and TE items across all content areas and subgroups. In addition, ELA TDA items were flagged for DIF by disability status (against students with disabilities in all grades), by ethnicity status (in favor of Asian students in grade 3, against American Indian students in grades 6 through 8, and against African American students in grades 4 and 8), by gender (in favor of female students in grades 5, 6, and 7), and by economic status (against economically disadvantaged students in grade 8).

Combined, the DIF statistical analyses discussed above and the expert reviews provide an appropriate set of tools with which to minimize the extraneous or construct-irrelevant information associated with item bias or DIF in the Wisconsin Forward Exam. It should be noted that in large-scale assessments, such as the Wisconsin Forward Exam, it is expected that some items will show DIF. All items flagged for DIF are annotated as such in the item pool so that content experts would be able to reevaluate these items in future item selection activities. Items with DIF (particularly items flagged for large DIF) are to be avoided in future selections.

9.2 Validity Evidence Based on Internal Test Structure

Construct-related evidence of the validity of an intended interpretation of test scores can be defined as the extent to which tests measure the skills or constructs they intend to measure and is the central concept underlying the Wisconsin Forward Exam validation process. Evidence for construct-related validity is comprehensive and integrates evidence from both content- and criterion-related validity. The Wisconsin Forward Exam development process included specifications, item writing, review, and test construction.

Threats to construct-related validity include the unintended measurement of variables unrelated to the desired constructs and multidimensionality of the tests. To ensure that the test items are focused on the desired constructs, standardized procedures are employed to select items with sound statistical properties, to align the

items to content standards, and to ensure that each test form meets the Wisconsin Forward Exam blueprint. A test can be said to be unidimensional when all of the items in the test measure the same underlying ability or trait. For example, Mathematics items should measure Mathematics ability and not Reading skills. Standard 1.13 of the *Standards* (AERA, APA, & NCME, 2014) states the following:

If the rationale for a test score interpretation for a given use depends on premises about the relationships among test items or among parts of the test, evidence concerning the internal structure of the test should be provided. (pp. 26–27)

9.2.1 Correlations between Content Standards

Analyses of the internal structure of a test can indicate the extent to which the relationships between test items and components conform to the construct the test purports to measure. For educational assessments that are designed to measure a single construct or content domain, the correlations between content standards within a test can be expected to be relatively high. Table 9-5 shows the correlations between the main test domains for ELA, and Tables 9-6 through 9-9 show the correlations between content standards for each Wisconsin Forward Exam content area. The correlation coefficients here reflect the degree of linear relationship and direction between any two given content standards. The correlation can range from +1 to -1. A correlation of +1 indicates a perfect positive linear relationship between two content standards, and a correlation of -1 indicates a perfect negative linear relationship between two content standards. A correlation of zero means there is no linear relationship. In general, the size of the correlation coefficient is influenced by the number of items or score points and by the score variance. Readers are cautioned not to confuse correlation with causation. The presence of a high correlation between two content standards should not be taken as an indication that there is a causal relationship between them.

As may be observed in Table 9-5, the correlations between the ELA main test domains of Reading, Writing, and Listening are moderate to high and range from 0.56 to 0.76 across all grades. Lower correlations (at or lower than 0.61) were observed between the Listening and Writing domains and the Listening and Reading domains in grades 3 through 6, while higher correlations (at least 0.72 or higher) were observed between the Reading and Writing domains across all grades. The correlations between ELA content standards (see Table 9-6) are typically moderate for all grades and all standard pairs and range from 0.36 to 0.71. It should be noted that the number of items associated with each content standard was smaller than the number of items associated with each ELA domain, resulting in lower correlations at the standard level compared to the correlations at the ELA domain level.

As indicated in Table 9-7, the correlations between Mathematics content standards are moderate to high and range from 0.56 to 0.77. The correlations between Science content standards range from 0.65 to 0.71 (see Table 9-8), and the correlations between Social Studies content standards range from 0.54 to 0.73 (as shown in Table 9-9). Overall, the correlations for all content areas are within the moderate to high range.

Although it may be tempting to try to interpret the differences in magnitude within and across content areas, it is important to note that these correlations are highly dependent upon the numbers of items and the score variance for the different standards. The important finding is that within each content area, the correlations between

content standards are low enough to indicate that the standards are, as intended, somewhat distinct from one another but high enough to indicate that the individual standards are measuring related components of a single content area.

9.2.2 Principal Component Analysis

Wisconsin Forward Exam items are calibrated using unidimensional IRT models, which suggests that the test items are measuring an essentially unidimensional construct. To assess the dimensionality of the Wisconsin Forward Exam, a principal components analysis was conducted for each content area and grade. A principal components analysis is a statistical technique commonly used to evaluate dimensionality by detecting patterns of relationships among items. This method is useful in determining whether the observed scores on a test can be explained largely or entirely in terms of a much smaller number of components. For example, if answering the Mathematics items in a Mathematics test required a high level of reading ability, the Mathematics test would be measuring not only mathematics ability but also reading ability. Such a test would be said to be multidimensional rather than essentially unidimensional. One way of evaluating the dimensions detected in the analysis is by examining the eigenvectors and eigenvalues. In a principal components analysis, the eigenvectors correspond to factors, and the eigenvalues correspond to the variance explained by these factors. The sum of the eigenvalues is equal to the number of items in the test. The eigenvalues can be ordered from first to last in terms of the amount of common variance that each explains. Data are generally considered to be unidimensional if the second eigenvalue is less than or equal to 1.0. Previous research shows that an examination of the ratio of the first two (i.e., the two largest) eigenvalues can be useful in determining the existence of dominant factors. Specifically, where large ratios exist between the first and second eigenvalues, a single dominant factor can be said to exist. Although the definition of “large” in the present context is subjective, the results in Table 9-10 show that the eigenvalue of the first factor is more than five times as large as the eigenvalue of the second factor.

As can be seen in Table 9-10, the ratios of the first two eigenvalues range from 5.89 to 7.68. The eigenvalues are proportional to the amount of common variance explained by each component, indicating that the variance explained by the first component alone is between approximately six and eight times greater than the variance explained by the second component. The eigenvalue ratios range from 6.06 to 6.96 in ELA, from 5.89 to 7.68 in Mathematics, from 6.03 to 7.53 in Science, and from 6.80 to 7.49 in Social Studies. These ratios suggest that the unidimensionality of each of the Wisconsin Forward Exam content assessments is sufficient to meet the requirements of a unidimensional IRT calibration model.

Overall, these results provide support for the construct validity of the Wisconsin Forward Exam assessments. The correlations between content standards and the presence of a single dominant factor for each test confirm that the content standards are sufficiently unidimensional to be combined into a single score.

9.3 Validity Evidence Based on Relationship with Other Variables

The relationship between the Wisconsin Forward Exam scores and other variables was examined to further support the validity of the intended score interpretation. This was done using two measures: evaluation of

correlations between the Wisconsin Forward Exam content area scores and comparisons of the percentages of students classified in different performance levels (impact data) on the State assessment and on the NAEP assessment.

9.3.1 Correlations between Content Area Test Scores

The test score relationship with other variables can be assessed by the extent to which measures of constructs that theoretically should not be related to each other are, in fact, observed as not related to each other. Typically, correlation coefficients between measures of unrelated or distantly related constructs are examined in support of the validity evidence based on the relationship of the test scores with other variables.

To assess the relationships between the Wisconsin Forward Exam content area scores, the correlations between the ELA, Mathematics, Science, and Social Studies scale scores for students who took more than one subject area test in 2022 were computed and examined for the total student population and for each subgroup. Table 9-11 shows the correlations between the content area scores for the total population of Wisconsin students. The correlations between ELA and Mathematics scores ranged from 0.75 to 0.77 across all grades. The correlations between ELA and Science scores and between ELA and Social Studies scores ranged from 0.80 to 0.83 for grades 4 and 8. The correlations between Mathematics and Science scores and between Mathematics and Social Studies scores ranged from 0.73 to 0.77 for grades 4 and 8. Lastly, the correlations between Science and Social Studies scale scores were 0.83 for grade 4 and 0.82 for grade 8.

Tables 9-12 through 9-16 show correlation coefficients between the content area scores by gender, ethnicity, English language proficiency status, economic status, and disability status, respectively. As seen in Table 9-12, the correlations between the content area scores for male or female groups ranged from 0.73 to 0.84 and were comparable for the two gender groups for each pair of correlated scores in each grade. The correlations between the content area scores for the non-binary group in grade 8 ranged from 0.55 to 0.75 across all pairs of scores. Recall that the non-binary results should be interpreted with caution due to a low number of students in this group. The correlations between the content area scores for different ethnic groups ranged from 0.60 to 0.84 (see Table 9-13). The highest correlations by ethnic group were observed for Asian students. Correlations between the content area scores for the African American student subgroup were generally lower than the correlations for other subgroups. As shown in Table 9-14, the correlations between the content area scores by English proficiency status ranged from 0.54 to 0.75 for limited English proficiency students and from 0.73 to 0.83 for fully English proficient students across all grade levels and all pairs of correlated scores. The correlations between the content area scores by student economic status are presented in Table 9-15. These correlations ranged from 0.72 to 0.82 for students who were not economically disadvantaged and from 0.68 to 0.80 for economically disadvantaged students across all grades and pairs of correlated scores. The correlations between the content area scores by student disability status are shown in Table 9-16. These correlations ranged from 0.72 to 0.83 for students without disabilities and from 0.59 to 0.79 for students with disabilities across all grades and pairs of correlated scores. In all grade levels, the correlations between each pair of scores were, in most cases, lower for the groups of students classified as English language learners, economically disadvantaged, or disabled compared to the groups of students classified as fully English proficient, not economically disadvantaged, or not disabled. In all grade levels and for each subgroup, the correlations between

Mathematics test scores and other content area scores tended to be lower than the correlations between ELA scores and other content area scores. The correlation coefficients between the content area scores were not computed by accommodation use, because the accommodation use status was not consistent across content areas for the same students (e.g., students who used accommodations in one content area did not necessarily use accommodations in another content area).

Overall, the correlations between the content area scores for the total population of students were found to be highly related. As mentioned, correlations between Mathematics test scores and other content area scores were found to be consistently lower than the correlations between ELA scores and other content area scores, suggesting larger differences between the Mathematics constructs and constructs measured by other content areas. The correlations between the content area scores for the subgroups of students were found to be moderately to highly related. Despite relatively high correlations, the tests are not perfectly related to one another, suggesting that different constructs are being tapped; however, if the test scores are highly related to one another, they may be tapping into a similar knowledge base or general underlying ability. This outcome is not unexpected for the new generations of large-scale assessments based on new CCR standards, such as the Wisconsin State Standards that emphasize teaching and learning various content skills across content areas. All assessments are intended to be aligned to performance expectations that are more reflective of the current knowledge and skill demands of postsecondary education and careers. Given the cross-content complexity and cognitive demands of the assessments, the relationship between student scores from different tests has necessarily become strengthened, and larger correlations between the content areas are expected.

Partial Correlations

In addition to the simple correlations between the content area scores, partial correlations, which are measures of the strength of the relationship between the content area scores while controlling for the student demographic characteristics (gender, ethnicity, English proficiency status, disability status, and economic status), were also computed. Partial correlations allow for the evaluation of the relationship between two content area scores with the effect of the student demographic characteristics removed (or held constant). The partial correlations between the ELA, Mathematics, Science, and Social Studies test scores for the total population of students and at each grade level are presented in Table 9-17. The partial correlations between ELA and Mathematics scores ranged from 0.66 to 0.70 across all grades. The partial correlations between ELA and Science scores and ELA and Social Studies scores ranged from 0.74 to 0.78. The partial correlations between Mathematics and Science scores and Mathematics and Social Studies scores ranged from 0.64 to 0.69. The partial correlations between Science and Social Studies scores were 0.78 for grade 4 and 0.76 for grade 8.

Although the magnitude of these correlations is considered to be strong, the partial correlations between the content area scores were lower than the corresponding simple correlations, indicating that the student demographic characteristics did contribute to the strength of the relationship between the content area test scores. The differences between the simple correlation coefficients and corresponding partial correlation coefficients suggested that there may be some effect of the student demographic characteristics on the relationships between the ELA, Mathematics, Science, and Social Studies test scores.

9.3.2 Comparison of the Wisconsin Forward Exam and Wisconsin NAEP Impact Data

The NAEP is the largest nationally representative and continuing assessment of what America’s students know and can do in various content areas. Assessments in several content areas, including Reading, Mathematics, and Science, are administered to students in grades 4, 8, and 12 and conducted periodically. Representative samples of students from different states, including Wisconsin, participated in the latest NAEP assessment, which occurred in Spring 2022.

The main NAEP assessments are constructed using detailed frameworks that result from a comprehensive national process in which teachers, curriculum experts, policymakers, and members of the general public work to create a unified vision of how a particular subject ought to be assessed. This vision is based on current educational research on achievement and its measurement as well as good educational practices. The frameworks are devised through a development process that ensures they meet current educational requirements. (For details, refer to <https://nces.ed.gov/nationsreportcard/assessments/frameworks.aspx>.)

The NAEP results are reported for all assessed content areas and for all participating grades at the national level. At the state level, the results for Reading, Mathematics, Science, and Writing are reported for grades 4 and 8. The results may also be reported at the district level (within a state) for these four content areas. No results are reported at the student level.

Wisconsin students participated in the latest Science NAEP assessment in 2019 and in Reading and Mathematics assessments in Spring 2022. The percentages of Wisconsin students classified in different performance levels on the Wisconsin Forward Exam in Spring 2023 and the corresponding NAEP assessments for Reading and Mathematics from Spring 2022 are presented in Table 9-18. Because the Spring 2019 Science results by state are not available at the time of preparation of this report, the Spring 2015 Science NAEP results for Wisconsin students are presented (also in Table 9-18). With three exceptions, the percentages of students classified in the different performance levels on the NAEP assessments and on the Wisconsin Forward Exam were comparable and within 10% of each other for every performance level across both grades and all three content areas. The exceptions were lower percentages of students classified in the *Below Basic* performance level on the Wisconsin Forward ELA test in grade 4 and higher percentages of students classified in the *Advanced* performance level on the Wisconsin Forward Science test for grades 4 and 8 compared to the percentages of students classified in the respective performance levels on the corresponding NAEP assessments.

Looking at the percentages of students classified at or above *Proficient*, higher proportions of students were classified in this combined category of performance levels on the Wisconsin Forward Exams in ELA grades 4 and 8, Mathematics grade 4, and Science grades 4 and 8 compared to the corresponding NAEP Reading and Mathematics assessments in Spring 2022 or the Science assessments in Spring 2015. These differences ranged from approximately 2% for Mathematics grade 4 to 12% for ELA grade 4. A lower proportion of students was classified at or above *Proficient* on the Wisconsin Forward Exams in Mathematics grade 8 compared to the NAEP assessment in Mathematics grade 8 in Spring 2022 (difference of less than 2%). These comparisons should be made with caution and in the context of recovery from disrupted learning in the 2020–21 academic year.

It should be noted that the Spring 2015 Reading and Mathematics Wisconsin NAEP impact data were used as benchmarks during the Wisconsin Forward Exam standard setting for ELA and Mathematics after the Spring 2016 test administration. The Spring 2015 Science Wisconsin NAEP impact data were also shown to the participants for reference and guidance in performance level setting during the Spring 2019 standard setting. While the standard setting participants were free to deviate from the NAEP impact data while placing their bookmarks in the ordered item booklets in consideration of the Wisconsin performance level descriptors (PLDs), the final Wisconsin impact data achieved after the standard setting were generally aligned with the Wisconsin state-level NAEP data. When considering the Wisconsin content standards and impact data articulation across grades, the Wisconsin Forward Exam cut scores for ELA, Mathematics, and Science remained in most cases aligned with the *Proficient* benchmarks, further supporting the evidence of the relationship between the state and the national assessments in these content areas.

9.4 Test Integrity: Data Forensic Analyses

With the high-stakes nature of large-scale statewide assessment programs, there can be situations in which student responses, and hence their scores, may not be a true representation of student ability. Various activities may take place, such as a student copying from another student's paper, a student receiving inappropriate assistance before or during testing, or a student's responses being altered during or after testing. To maintain the integrity of the Wisconsin Forward Exam and the validity of the results, it is important that any such instances be discovered.

Two studies were conducted to evaluate the Wisconsin Forward Exam student data for any indicators of possible inappropriate testing behavior. The first study examines incorrect student responses to MC items on the Spring 2023 Wisconsin Forward Exam in ELA, Mathematics, Science, and Social Studies that were changed to correct responses. These answer changes are referred to as wrong-to-right answer changes. Inordinate numbers of wrong-to-right answer changes in a specifically identifiable testing administration group may indicate inappropriate student behavior or intervention by an educator during the testing session.

The second study evaluates the time spent on the test and individual test items by students. These analyses serve to inform of any events in which students (within one school) spent a very short or very long time on the test or specific items. Inordinate numbers of unusual test or item response times may indicate inappropriate pre-knowledge of the items or other interventions during the testing session.

The results of the two studies are provided to DPI for evaluation. We emphasize that the results from these studies may be used in conjunction with other information to investigate whether inappropriate interventions may have taken place. The statistical results by themselves may simply be coincidental and do not necessarily indicate inappropriate behavior.

9.5 Summary

In summary, the overall purpose of Part 9 was to provide additional evidence of the validity of an intended interpretation of test scores related to test construct. Through the measures of correlations between content area reporting category scores and principal components analysis, the existence of a single underlying trait or ability

for each content area was demonstrated. Next, the relationship between the Wisconsin Forward Exam scores and other variables was explored and validated through the evaluation of correlations between content area scores for the total population and by subgroups. In addition, student performance on the Wisconsin Forward Exam was compared with student performance on the NAEP assessment. The forensic analysis procedures that were employed to ensure the integrity of test scores by identifying schools and individual students that might have engaged in inappropriate behaviors during testing were also described in this part of the report.

Table 9-1 Items Flagged for DIF in English Language Arts

DIF Category	Reference Group	Focal Group	Grade	Item Number	Item Type	MH SMD Statistic	MH Delta Statistic	DIF Flag	
Gender	Male	Female	5	1	TDA	0.14		B	
				35	MC	-0.07	-1.10	B-	
			6	1	TDA	0.20		C	
			7	1	TDA	0.23		B	
			8	5	MC	0.08	1.39	B	
				24	MC	-0.08	-1.00	B-	
Race/Ethnicity	White	Asian	3	1	TDA	0.12		B	
				25	MC	-0.14	-2.00	C-	
			4	35	MC	-0.08	-1.18	B-	
				5	8	MC	0.11	1.34	B
					23	MC	-0.10	-1.20	B-
			7	25	MC	-0.08	-1.17	B-	
		12		TE	0.06	1.43	B		
		8	34	MC	0.07	1.06	B		
			38	MC	-0.07	-1.01	B-		
			African American	4	1	TDA	-0.16		B-
		8		1	TDA	-0.32		C-	
			20	MC	0.09	1.05	B		
		American Indian	6	1	TDA	-0.18		B-	
			7	1	TDA	-0.19		B-	
			8	1	TDA	-0.32		C-	
Limited English Proficiency	No	Yes	3	25	MC	-0.10	-1.20	B-	
			7	12	TE	-0.06	-0.80	B-	
Disability Status	No	Yes	3	1	TDA	-0.17		C-	
			4	1	TDA	-0.18		B-	
				5	MC	-0.08	-1.14	B-	
			5	1	TDA	-0.22		C-	
				10	MC	-0.10	-1.13	B-	
				12	MC	-0.09	-1.03	B-	
				18	MC	0.09	1.01	B	
			6	1	TDA	-0.25		C-	
			7	1	TDA	-0.34		C-	
				12	TE	-0.10	-1.28	C-	
8	1	TDA	-0.37		C-				
	5	MC	-0.11	-1.28	B-				
Economically Disadvantaged	No	Yes	8	1	TDA	-0.18		B-	

Table 9-2 Items Flagged for DIF in Mathematics

DB39:I76IF Category	Reference Group	Focal Group	Grade	Item Number	Item Type	MH SMD Statistic	MH Delta Statistic	DIF Flag
Gender	Male	Female	3	39	MC	-0.06	-1.05	B-
			4	10	MC	-0.08	-1.05	B-
				40	MC	-0.08	-1.12	B-
				45	MC	0.05	1.15	B
				9	TE	-0.09	-1.96	B-
			5	27	TE	-0.09	-1.13	B-
				6	12	MC	-0.07	-1.15
			14		MC	-0.09	-1.16	B-
			7	1	SA	-0.10	-2.02	B-
				14	MC	0.07	1.10	B
				19	MC	-0.14	-1.94	C-
				25	MC	-0.13	-1.92	C-
				33	SA	-0.09	-1.29	B-
				45	MC	0.08	1.06	B
Race/Ethnicity	White	Asian	3	31	MC	-0.07	-1.01	B-
			4	40	MC	-0.10	-1.37	B-
			6	44	TE	-0.09	-1.54	B-
			7	2	MC	0.09	1.33	B
		African American	4	5	TE	-0.12	-1.58	C-
				37	SA	-0.09	-1.24	C-
				41	SA	-0.08	-1.01	B-
			6	17	TE	-0.09	-0.94	B-
				38	SA	-0.09	-1.24	B-
				7	1	SA	-0.11	-1.31
		8	27	SA	-0.14	-2.03	C-	
		American Indian	4	5	TE	-0.12	-1.58	C-
		Disability Status	No	Yes	4	4	MC	-0.09
35	TE					-0.06	-1.07	B-
Accommodation Use	No	Yes	3	29	MC	-0.11	-1.32	B-
				4	4	MC	-0.11	-1.19
			24		MC	0.07	1.08	B
			35		TE	-0.08	-1.06	B-
			6	8	TE	-0.09	-1.04	B-
			7	1	SA	-0.08	-1.00	B-
				20	SA	-0.09	-1.42	B-
				44	SA	-0.08	-1.70	B-
8	27	SA	-0.09	-1.33	B-			

Table 9-3 Items Flagged for DIF in Science

DIF Category	Reference Group	Focal Group	Grade	Item Number	Item Type	MH SMD Statistic	MH Delta Statistic	DIF Flag
Gender	Male	Female	8	5	TE	-0.12	-1.42	B-
				21	TE	-0.10	-1.19	B-
				38	TE	-0.08	-1.28	B-
Race/Ethnicity	White	Asian	8	5	TE	-0.10	-1.18	B-
Disability Status	No	Yes	4	4	TE	-0.07	-1.02	B-
				15	MC	0.08	1.10	B
			8	1	TE	-0.10	-1.34	B-

Table 9-4 Items Flagged for DIF in Social Studies

DIF Category	Reference Group	Focal Group	Grade	Item Number	Item Type	MH SMD Statistic	MH Delta Statistic	DIF Flag
Gender	Male	Female	4	31	TE	-0.10	-1.12	B-
				38	MC	-0.08	-1.01	B-
			8	5	MC	-0.09	-1.25	B-
				9	TE	-0.09	-1.53	B-
				30	MC	0.08	1.17	B
			10	3	MC	-0.07	-1.19	B-
				8	MS	-0.08	-1.10	B-
39	MC	0.09		1.04	B			
Race/Ethnicity	White	Asian	4	1	TE	0.09	1.16	B
			8	3	TE	-0.12	-1.59	B-
				14	MC	-0.09	-1.51	C-
				10	1	TE	-0.13	-1.68
			2		MC	-0.08	-1.16	B-
			25		MC	0.06	1.12	B
		African American	8	5	MC	-0.11	-1.40	B-
				9	TE	-0.09	-1.07	B-
			10	11	MC	0.09	1.15	B
		32		TE	-0.09	-1.44	B-	
		Hispanic	8	5	MC	-0.09	-1.05	B-
		American Indian	10	28	MC	0.09	1.17	B
40	MC			0.17	2.11	C		
Limited English Proficiency	No	Yes	10	1	TE	-0.11	-1.22	B-
Disability Status	No	Yes	8	9	TE	-0.09	-1.05	B-
			10	1	TE	-0.09	-1.01	B-

Table 9-5 Correlations between English Language Arts Test Domains

Grade	ELA Domain	Listening	Reading
3	Reading	0.61	
	Writing	0.56	0.74
4	Reading	0.58	
	Writing	0.57	0.76
5	Reading	0.60	
	Writing	0.56	0.72
6	Reading	0.60	
	Writing	0.56	0.73
7	Reading	0.66	
	Writing	0.64	0.72
8	Reading	0.65	
	Writing	0.64	0.74

Table 9-6 Correlations between Content Standards, English Language Arts

Grade	Standard Code	A	B	C	D	E	F
3	B	0.62					
	C	0.60	0.50				
	D	0.60	0.52	0.53			
	E	0.55	0.48	0.48	0.50		
	F	0.44	0.38	0.40	0.42	0.38	
	G	0.57	0.48	0.48	0.49	0.47	0.36
4	B	0.61					
	C	0.58	0.59				
	D	0.56	0.58	0.56			
	E	0.53	0.55	0.53	0.55		
	F	0.50	0.51	0.50	0.53	0.48	
	G	0.48	0.51	0.49	0.50	0.48	0.42
5	B	0.64					
	C	0.53	0.49				
	D	0.53	0.49	0.44			
	E	0.53	0.50	0.45	0.45		
	F	0.49	0.46	0.41	0.46	0.45	
	G	0.55	0.50	0.48	0.45	0.48	0.42
6	B	0.71					
	C	0.57	0.58				
	D	0.53	0.55	0.42			
	E	0.50	0.53	0.43	0.43		
	F	0.51	0.53	0.43	0.45	0.42	
	G	0.54	0.56	0.44	0.46	0.44	0.43
7	B	0.57					
	C	0.55	0.57				
	D	0.51	0.54	0.51			
	E	0.48	0.47	0.46	0.46		
	F	0.44	0.46	0.48	0.48	0.41	
	G	0.54	0.56	0.59	0.55	0.49	0.49
8	B	0.62					
	C	0.59	0.63				
	D	0.56	0.58	0.58			
	E	0.53	0.56	0.54	0.55		
	F	0.41	0.43	0.43	0.46	0.42	
	G	0.55	0.58	0.56	0.56	0.54	0.43

ELA standards: A = Reading—Key Ideas and Details; B = Reading—Craft & Structure/Integration of Knowledge & Ideas; C = Reading—Vocabulary Use; D = Writing/Language—Text Types and Purposes; E = Writing/Language—Research; F = Writing/Language—Language Conventions; G = Listening

Table 9-7 Correlations between Content Standards, Mathematics

Grade	Standard Code	A	B	C	D	E	F	G	H	I
3	B	0.76								
	C	0.68	0.68							
	D	0.77	0.76	0.69						
	E	0.68	0.66	0.65	0.68					
4	B	0.76								
	C	0.74	0.71							
	D	0.71	0.70	0.71						
	E	0.58	0.56	0.59	0.59					
5	B	0.74								
	C	0.67	0.69							
	D	0.69	0.70	0.66						
	E	0.66	0.66	0.60	0.65					
6	F					0.66				
	G					0.69	0.74			
	H					0.71	0.71	0.75		
	I					0.60	0.61	0.64	0.64	
7	F					0.61				
	G					0.59	0.68			
	H					0.62	0.73	0.69		
	I					0.62	0.70	0.65	0.70	
8	G					0.60				
	H					0.66		0.68		
	I					0.65		0.61	0.67	
	J					0.66		0.64	0.70	0.70

Note: Standard Codes are as follows: A = Operations and Algebraic Thinking; B = Number and Operations in Base Ten; C = Number and Operations—Fractions; D = Measurement and Data; E = Geometry; F = Ratios and Proportional Relationships; G = The Number System; H = Expressions and Equations; I = Statistics and Probability; J = Functions

Table 9-8 Correlations between Content Standards, Science

Grade	Standard Code	A	B	C
4	B	0.71		
	C	0.68	0.68	
	D	0.71	0.71	0.70
8	B	0.69		
	C	0.67	0.66	
	D	0.67	0.68	0.65

Note: Standard Codes are as follows: A = Life Science; B = Physical Science; C = Earth and Space Science;

Table 9-9 Correlations between Content Standards, Social Studies

Grade	Standard Code	A	B	C	D
4	B	0.70			
	C	0.59	0.56		
	D	0.67	0.64	0.54	
	E	0.73	0.68	0.57	0.64
8	B	0.67			
	C	0.69	0.70		
	D	0.65	0.66	0.70	
	E	0.65	0.68	0.67	0.64
10	B	0.70			
	C	0.68	0.64		
	D	0.63	0.60	0.59	
	E	0.68	0.63	0.63	0.58

Note: Standard Codes are as follows: A = Geography; B = History; C = Political Science; D = Economics; E = Behavioral Sciences

Table 9-10 Principal Components Analysis

Content Area	Grade	First Eigenvalue	Second Eigenvalue	Ratio of First Two Eigenvalues
ELA	3	7.49	1.22	6.12
	4	8.02	1.15	6.96
	5	7.54	1.24	6.10
	6	7.95	1.30	6.13
	7	7.52	1.24	6.06
	8	8.06	1.32	6.08
Mathematics	3	11.33	1.48	7.68
	4	11.00	1.87	5.89
	5	10.88	1.42	7.66
	6	11.17	1.79	6.24
	7	10.23	1.49	6.84
	8	10.70	1.44	7.41
Science	4	8.96	1.19	7.53
	8	7.95	1.32	6.03
Social Studies	4	8.95	1.31	6.83
	8	9.16	1.35	6.80
	10	8.56	1.14	7.49

Table 9-11 Correlations between Content Area Scale Scores

Grade	ELA & Mathematics	ELA & Science	ELA & Social Studies	Mathematics & Science	Mathematics & Social Studies	Science & Social Studies
3	0.77					
4	0.77	0.83	0.83	0.77	0.76	0.83
5	0.75					
6	0.77					
7	0.77					
8	0.76	0.80	0.82	0.77	0.73	0.82

Table 9-12 Correlations between Content Area Scale Scores by Gender

Grade	Demographic Group	ELA & Mathematics	ELA & Science	ELA & Social Studies	Mathematics & Science	Mathematics & Social Studies	Science & Social Studies
3	Female	0.78					
	Male	0.77					
4	Female	0.78	0.83	0.84	0.77	0.76	0.83
	Male	0.78	0.83	0.84	0.77	0.76	0.84
5	Female	0.76					
	Male	0.76					
6	Female	0.78					
	Male	0.78					
7	Female	0.78					
	Male	0.79					
8	Female	0.77	0.81	0.82	0.77	0.73	0.81
	Male	0.76	0.80	0.82	0.77	0.73	0.82
	Non-Binary	0.67	0.75	0.72	0.72	0.55	0.71

Note: Correlations for the non-binary group were not computed for grades 3 through 7 due to student counts < 50.

Table 9-13 Correlations between Content Area Scale Scores by Ethnicity/Race

Grade	Demographic Group	ELA & Mathematics	ELA & Science	ELA & Social Studies	Mathematics & Science	Mathematics & Social Studies	Science & Social Studies
3	White	0.73					
	African American	0.65					
	Hispanic	0.73					
	Asian	0.77					
	American Indian	0.70					
	Two or More	0.76					
4	White	0.73	0.80	0.82	0.73	0.72	0.81
	African American	0.69	0.73	0.72	0.65	0.64	0.71
	Hispanic	0.74	0.80	0.80	0.74	0.71	0.80
	Asian	0.78	0.83	0.84	0.78	0.77	0.84
	American Indian	0.69	0.78	0.78	0.70	0.68	0.78
	Two or More	0.77	0.81	0.83	0.75	0.75	0.83
5	White	0.73					
	African American	0.61					
	Hispanic	0.68					
	Asian	0.76					
	American Indian	0.69					
	Two or More	0.74					
6	White	0.74					
	African American	0.68					
	Hispanic	0.72					
	Asian	0.78					
	American Indian	0.73					
	Two or More	0.76					
7	White	0.75					
	African American	0.69					
	Hispanic	0.73					
	Asian	0.77					
	American Indian	0.71					
	Two or More	0.77					
8	White	0.74	0.78	0.80	0.74	0.71	0.81
	African American	0.63	0.71	0.73	0.63	0.60	0.71
	Hispanic	0.69	0.77	0.79	0.71	0.67	0.78
	Asian	0.79	0.81	0.84	0.78	0.76	0.83
	American Indian	0.70	0.76	0.75	0.68	0.65	0.75
	Two or More	0.76	0.80	0.82	0.76	0.72	0.80

Table 9-14 Correlations between Content Area Scale Scores by English Proficiency Status

Grade	Demographic Group	ELA & Mathematics	ELA & Science	ELA & Social Studies	Mathematics & Science	Mathematics & Social Studies	Science & Social Studies
3	Fully English Proficient	0.77					
	Limited English Proficient	0.67					
4	Fully English Proficient	0.77	0.83	0.83	0.77	0.75	0.83
	Limited English Proficient	0.67	0.74	0.74	0.68	0.65	0.75
5	Fully English Proficient	0.75					
	Limited English Proficient	0.58					
6	Fully English Proficient	0.77					
	Limited English Proficient	0.63					
7	Fully English Proficient	0.77					
	Limited English Proficient	0.65					
8	Fully English Proficient	0.76	0.80	0.81	0.76	0.73	0.81
	Limited English Proficient	0.57	0.65	0.69	0.59	0.54	0.67

Table 9-15 Correlations between Content Area Scale Scores by Economic Status

Grade	Demographic Group	ELA & Mathematics	ELA & Science	ELA & Social Studies	Mathematics & Science	Mathematics & Social Studies	Science & Social Studies
3	Not Economically Disadvantaged	0.74					
	Economically Disadvantaged	0.73					
4	Not Economically Disadvantaged	0.74	0.81	0.82	0.74	0.72	0.82
	Economically Disadvantaged	0.74	0.80	0.80	0.74	0.71	0.80
5	Not Economically Disadvantaged	0.73					
	Economically Disadvantaged	0.70					
6	Not Economically Disadvantaged	0.74					
	Economically Disadvantaged	0.73					
7	Not Economically Disadvantaged	0.75					
	Economically Disadvantaged	0.73					
8	Not Economically Disadvantaged	0.75	0.78	0.80	0.75	0.72	0.81
	Economically Disadvantaged	0.70	0.77	0.78	0.72	0.68	0.78

Table 9-16 Correlations between Content Area Scale Scores by Disability Status

Grade	Demographic Group	ELA & Mathematics	ELA & Science	ELA & Social Studies	Mathematics & Science	Mathematics & Social Studies	Science & Social Studies
3	Not Disabled	0.76					
	Disabled	0.73					
4	Not Disabled	0.76	0.82	0.83	0.76	0.74	0.83
	Disabled	0.74	0.79	0.78	0.74	0.71	0.79
5	Not Disabled	0.74					
	Disabled	0.65					
6	Not Disabled	0.75					
	Disabled	0.68					
7	Not Disabled	0.75					
	Disabled	0.68					
8	Not Disabled	0.74	0.79	0.80	0.75	0.72	0.81
	Disabled	0.62	0.72	0.73	0.66	0.59	0.73

Table 9-17 Partial Correlations between Content Area Scale Scores

Grade	ELA & Mathematics	ELA & Science	ELA & Social Studies	Mathematics & Science	Mathematics & Social Studies	Science & Social Studies
3	0.70					
4	0.70	0.78	0.78	0.69	0.67	0.78
5	0.66					
6	0.69					
7	0.69					
8	0.68	0.74	0.76	0.68	0.64	0.76

Table 9-18 Comparison of Most Recent Wisconsin NAEP and Spring 2023 Wisconsin Forward Exam Impact Data

Content	Grade	Wisconsin NAEP Percentages of Students							Wisconsin Forward Exam Spring 2023 Percentages of Students					
		NAEP Year	Below Basic	Basic	Proficient	Advanced	At or Above Proficient	At or Above Basic	Below Basic	Basic	Proficient	Advanced	At or Above Proficient	At or Above Basic
Reading/ ELA	4	2022	37	31	25	8	33	63	23.21	31.78	34.74	10.27	45.01	76.79
Reading/ ELA	8	2022	28	39	29	3	32	72	26.84	35.54	28.24	9.38	37.62	73.16
Math	4	2022	21	36	33	10	43	79	21.00	33.70	30.96	14.34	45.30	79.00
Math	8	2022	30	37	25	8	33	70	34.69	34.12	24.83	6.36	31.19	65.31
Science	4	2015	21	38	40	1	41	79	18.82	29.93	31.85	19.40	51.25	81.18
Science	8	2015	25	35	38	2	40	75	21.89	29.25	29.13	19.74	48.86	78.11

Note: The NAEP assessed student knowledge and skills in Reading, while the Wisconsin Forward Exam assessed student knowledge and skills in ELA, which included Reading, Listening, and Writing.

Note: NAEP data are from <https://nces.ed.gov/nationsreportcard>.

Part 10: Test Results

Part 10 of the Technical Report provides short descriptions of the Wisconsin Forward Exam score reports and interpretive guide. It also presents a summary of student test results for the Spring 2023 Wisconsin Forward Exam administration. The summary results are presented for all Wisconsin students and cover four types of reported scores: total test scale scores; total test performance levels; scores based on each of the content standards within each content area, which are called standard performance index (SPI) scores; and performance levels based on SPI scores. The four types of scores offer the reader several points from which to understand and evaluate the performance of Wisconsin students on the Wisconsin Forward Exam. In addition, the longitudinal test participation rates and test results are presented in this part of the report. The AERA, APA, & NCME (2014) Standards addressed in Part 8 include 5.1, 6.10, 7.0, 7.1, and 12.18.

10.1 Types of Reports

Score reports are the primary means of communicating test scores to relevant district personnel (e.g., district assessment coordinators, superintendents), teachers, and parents. AERA, APA, & NCME (2014) Standard 6.10 states the following:

When test score information is released, those responsible for testing programs should provide interpretations appropriate to the audience. The interpretations should describe in simple language what the test covers, what scores represent, the precision/reliability of the scores, and how scores are intended to be used. (p. 119)

Standard 5.1 is related in that it states the following:

Test users should be provided with clear explanations of the characteristics, meaning, and intended interpretation of scale scores, as well as their limitations. (p. 102)

Interpretations related to the test scores are disseminated in two ways: (1) the individual score report and (2) the *User's Guide to Interpreting Reports* (DRC, 2022).

In addition to providing an explanation of the intended interpretation of test scores, a testing program must also ensure that the information related to the test scores is understandable by the target audience. Standards 7.0 and 7.1 of the *Standards* (AERA, APA, & NCME, 2014) state the following:

Standard 7.0 Information relating to tests should be clearly documented so that those who use tests can make informed decisions regarding which test to use for a specific purpose, how to administer the chosen test, and how to interpret test scores. (p. 125)

Standard 7.1 The rationale for a test, recommended uses of the test, support for such uses, and information that assists in score interpretation should be documented. When particular misuses of a test can be reasonably anticipated, cautions against such misuses should be specified. (p. 125)

In support of Standards 7.0 and 7.1, the *User’s Guide to Interpreting Reports* is accessible to parents, teachers, and the general public at <https://dpi.wi.gov/assessment/forward/data#resources>.

In the 2022–23 administration year, DPI reported the Wisconsin Forward Exam results in WISEdash Public and WISEdash for Districts. These dashboards provide comprehensive data analysis for statewide assessments, attendance, graduation, coursework and other data of interest to district, school and other statewide data users. DRC reported the Wisconsin Forward Exam results through the Wisconsin Forward Exam Reporting System, which is a browser-based system designed to deliver online interactive reporting to authorized users at the state and district levels for Wisconsin schools.

10.1.1 Description of Each Type of Report

In this section, descriptions of the following reports are provided: Individual Student Report (ISR), Student Roster, Summary by Subject, and Summary by Reporting Category. In compliance with AERA, APA, & NCME (2014) Standard 12.18, the Wisconsin Forward Exam score reports provide clear information about the achievements of individual students and groups of students. Standard 12.18 states the following:

In educational settings, score reports should be accompanied by a clear presentation of information on how to interpret the scores, including the degree of measurement error associated with each score or classification level, and by supplementary information related to group summary scores. In addition, dates of test administration and relevant norming studies should be included in score reports. (p. 200)

Individual Student Report

The Individual Student Report (ISR) is one of the types of reports available through the Wisconsin Forward Exam Reporting System. The ISR is the primary means for sharing student test results with parents and guardians. It is a stand-alone document, giving parents or guardians relevant information that enables them to understand their child’s test score. The ISRs are provided to schools to be sent home to parents or guardians. The ISR consists of three to four pages (depending on how many tests a student took). On the upper-left side of the first page, the student’s identifying information is provided. Underneath the student information, there is a short description of the Wisconsin Forward Exam and the purpose of the report. On the upper-right side of the first page there is a table with information on the student’s performance level in each content area as well as the student’s percentile rank in that content. Under the table summarizing the student’s proficiency classification, a short description of the Wisconsin Forward Exam performance levels is presented.

The second page of the report includes the presentation of the total test scale scores and the performance levels for ELA and Mathematics. This information is presented in the form of a bar graph, and the student’s scale score for a given content area is shown, along with the performance level associated with that scale score. The total test results are followed by the reporting category results presented for that content area. These results include number of points obtained, number of points possible, SPI score, and the SPI proficiency classification for each reporting category. The third page of the report includes the total test and reporting category results for Science and Social Studies for students in grades 4, 8, and 10 (the grades at which these assessments are administered).

The last page of the ISR includes information on accommodations and designated support use, Wisconsin Academic Standards, and a short explanation of the reported scores. An example of an ISR can be found in the *User's Guide to Interpreting Reports*.

Roster

Another report available from the Wisconsin Forward Exam Reporting System is an online interactive Roster that displays a list of students based on the specific report filter options selected, such as test administration, grade, school, district, gender, race/ethnicity, disability status, and English proficiency status. Total test scale scores and performance level indicators, as well as the reporting category performance levels, are displayed in a table format for the content area chosen.

Subject Summary Report

The Subject Summary online interactive report contains performance level information for the school, district, and state. It includes a mean scale score along with numeric and graphic representations of the performance level summary for the subject and grade. This report also includes the number and percentage of students in each of the performance levels.

Reporting Category Summary Report

The online interactive Reporting Category Summary Report includes performance level information based on the individual reporting category selected for the school, district, and state. By selecting any one of the available reporting categories from the Reporting Category filter, the table and chart data will be based on that category for the subject and grade chosen. Along with numeric and graphic representations of the performance level by reporting category, subject, and grade, the report includes the number and percentage of students in each of the performance levels for the reporting category selected.

Demographic Summary Report

The online interactive Demographic Summary Report provides at-a-glance comparisons of performance between various demographic subgroups. The percentages of students in each performance level, means scale scores, as well as the test participation rates are presented graphically and numerically for each subgroup within each demographic category. The summary table shows the number of students tested, the mean scale scores, and the percentage of students in each performance level for each demographic subgroup combination.

Examples of Roster, Subject Summary, Reporting Category Summary, and Demographic Summary reports are presented in the *User's Guide to Interpreting Reports* available on DPI's website.

10.1.2 Interpreting Test Results

A student's correct responses to the assessment questions are used to derive that student's Wisconsin Forward Exam scale score. The scale score describes performance on a continuum that spans the complete range of grades 3–8 for ELA and Mathematics. These scores range in value from 330 to 970 for ELA and from 360 to 890 for Mathematics. Because ELA and Mathematics assessments are on vertical scales, scores from adjacent

grades may be compared within a content area. For example, it is appropriate to compare a student's grade 5 Mathematics scale score with the student's grade 6 Mathematics scale score in a subsequent administration year. ELA and Mathematics scale scores can also be compared within a content area across the administrations from Spring 2016 to Spring 2023.

Science scale scores range from 300 to 725 for grade 4 and from 480 to 945 for grade 8. Science scores can be compared within a grade level, but since Science assessments are not on a vertical scale, the scale scores cannot be compared across grades. Because new reporting scales were developed for the Science assessments in Spring 2019, the Science scale scores from the current administration can only be compared with the scores from the Spring 2019, 2021, and 2022 administrations.

Social Studies scale scores range from 330 to 700 for grade 4, from 540 to 860 for grade 8, and from 645 to 980 for grade 10. Social Studies scores can be compared within a grade level, but since the assessments are not on a vertical scale, the scale scores cannot be compared across grades. New Social Studies scales were established in Spring 2022 and became a new baseline for longitudinal comparison. As such, the Spring 2023 scale scores can only be compared with the Social Studies scores from the Spring 2022 administration.

Scale scores cannot be compared across content areas. For example, it is not appropriate to compare a student's Mathematics and ELA scores as they do not represent comparable achievement.

The Wisconsin Forward Exam scale scores determine a student's performance level. Student performance is reported in terms of four performance levels that describe a pathway to proficiency and college and career readiness. Each performance level represents standards of performance for each assessed content area. Performance level scores provide a description of what students can do in terms of the content and skills assessed, as described in the Wisconsin Academic Standards.

In addition to the total test score, students receive scores in each reporting category of the test taken. The reporting category scores are SPI scores and performance levels. The SPI is an estimate of the number of questions that a student could be expected to answer correctly if there had been 100 such questions measuring that content standard on the test in a given administration year. More information on the SPI scores is provided in Section 10.4 of this report.

Last but not least, state percentile ranks are computed for each student based on the student's total test scale score. The state percentile ranks, ranging from 1 to 99, provide information that compares the student's achievement with that of a larger reference group, the state. The percentile rank tables for the most recent test administration can be found on DPI's website at <https://dpi.wi.gov/assessment/forward/data>.

Information on score interpretation is included in the *User's Guide to Interpreting Reports*, which was written for Wisconsin teachers and administrators who received score reports from the 2022–23 administration of the Wisconsin Forward Exam. The *Guide* was developed collaboratively by DRC and DPI staff.

10.2 Scale Scores Summary Statistics

The primary scores reported in Wisconsin Forward Exam program reports are scale scores. The scale score of a student in a given content area represents the student's level of performance in that content area. Higher scale scores indicate higher levels of performance, and lower scale scores indicate lower levels of performance. Scale scores are based on the entire set of scored operational items per grade and content area.

Summary descriptive statistics based on the scale score results are described below. Table 10-1 is the summary scale score table based on the Spring 2023 census data. The table shows the following: mean scale score, standard deviation of the scale scores, skewness and kurtosis, minimum and maximum obtained scale scores, and lowest and highest obtainable scale scores (LOSS and HOSS, respectively) for all content areas and grades based on the census data (all students with valid test scores). The LOSS and HOSS, as discussed in Part 6, identify the lower and upper limits of the scale score range. These values were established when the current scales were developed and do not change from one administration to another.

English Language Arts

- Mean scale score increased as grade level increased, ranging from 552.33 for grade 3 to 628.49 for grade 8. This mean scale score pattern supports the ELA vertical scale properties.
- Standard deviations ranged from 48.85 to 64.03 scale score points across grades.
- Student scores spanned the full-scale score range from the LOSS to the HOSS in grades 3, 7, and 8. No student reached the HOSS in grades 4, 5, or 6.

Mathematics

- Mean scale score increased as grade level increased, ranging from 553.02 for grade 3 to 638.57 for grade 8. This mean scale score pattern supports the Mathematics vertical scale properties.
- Standard deviations ranged from 52.62 to 60.98 scale score points across grades.
- In each grade level, student scores spanned the full-scale score range from the LOSS to the HOSS.

Science

- Mean scale scores were 496.06 and 692.80 for grades 4 and 8, respectively.
- Standard deviations were 54.65 and 51.93 scale score points for grades 4 and 8, respectively.
- In each grade level, student scores spanned the full-scale score range from the LOSS to the HOSS.

Social Studies

- Mean scale scores were 499.74, 700.49, and 796.91 for grades 4, 8, and 10, respectively.
- Standard deviations ranged from 48.94 to 53.02 scale score points.
- In each grade level, student scores spanned the full-scale score range from the LOSS to the HOSS.

10.2.1 Subgroup Performance Patterns in Scale Score Results

In addition to the evaluation of performance of all students, scale score means were computed and compared for subgroups of students. The mean differences were evaluated using Cohen's d effect size measure (Cohen, 1988). For the purpose of the mean scale score difference interpretation, the d statistic: $|d|$ equal to or larger than 0.10 (one-tenth of a pooled standard deviation) was considered to be an important and substantial difference in performance between subgroups. The scale score means, standard deviations, lowest and highest obtained scores, and the effect size of mean differences for subgroups of students are reported in Tables 10-2 through 10-5 for ELA, Mathematics, Science, and Social Studies, respectively.

The scale score means for subgroups traditionally considered disadvantaged or minority were compared against the performance of the reference group in each demographic category. The reference groups were male students, White students, fully English proficient students, students without disabilities, students not considered economically disadvantaged, and students not using testing accommodations. The positive value of Cohen's d statistics indicates that the mean score for the reference group was higher than the mean score for the focal group and vice versa. The scale score statistics were computed based on the census data.

English Language Arts

- Female students outperformed male students in each grade with differences in mean scale scores between approximately 7 points in grades 3 and 4 to over 16 points in grade 8. These differences were approximately a quarter of the pooled standard deviation or less.
- White students as a group had the highest mean scale scores, followed by Asian students, Hispanic students, American Indian students, and African American students. The differences between mean scale scores of White students and mean scale scores of Asian students were approximately 17 points (about one-third of the pooled standard deviation) in grades 3 and 4 and were found to be decreasing as the grade level increased. The difference between mean scale scores of White students and mean scale scores of Asian students was approximately 2 scale score points in grade 8. The differences between mean scale scores of White students and mean scale scores of Hispanic students ranged from approximately 31 points in grade 3 to 37 scale score points in grade 8. These differences were more than half of the pooled standard deviation in each grade. The differences between mean scale scores of White students and mean scale scores of American Indian students ranged from approximately 36 points in grade 3 to over 43 scale score points in grade 8. These differences were approximately three-quarters of the pooled standard deviation or larger across all grades. The differences between mean scale scores of White students and mean scale scores of African American students ranged from approximately 54 points in grade 3 to 62 scale score points in grade 8. These differences were more than one pooled standard deviation in each grade.
- Fully English proficient students outperformed their limited English proficiency peers in each grade with mean differences ranging from approximately 31 scale score points (about two-thirds of the pooled standard deviation) to about 55 scale score points (close to one pooled standard deviation).

- Students without disabilities outperformed students with disabilities by approximately 35 scale score points (in grade 3) to 68 points (in grade 8). These differences were from three-quarters to over one pooled standard deviation.
- Students considered not economically disadvantaged outperformed economically disadvantaged students in every grade with mean differences ranging from approximately 35 scale score points in grade 3 to over 43 scale score points in grade 8. These differences were approximately three-quarters of the pooled standard deviation in each grade.
- Students not using testing accommodations, as a group, had higher mean scale scores than students using testing accommodations in each grade. However, due to low numbers of students using testing accommodations in ELA, the performance of these two groups cannot be reliably compared.

Mathematics

- Male students showed higher mean scale scores in grades 3 through 7 (differences between 3 and 8 scale score points) and a slightly lower mean scale score in grade 8 (difference of approximately 1.5 scale score points) than female students. All differences were approximately one-tenth of the pooled standard deviation or less.
- Similar to the pattern of mean scale scores observed for ELA, White students outperformed students in all other ethnic groups in all grades. The differences between mean scale scores of White students and mean scale scores of Asian students ranged from just over 2 points (less than one-tenth of the pooled standard deviation) in grade 8 to close to 15 points in grade 3 (about one-quarter of the pooled standard deviation). The differences between mean scale scores of White students and mean scale scores of Hispanic students ranged from approximately 37 points to 42 scale score points, and the differences between mean scale scores of White students and mean scale scores of American Indian students ranged from approximately 38 points to over 47 scale score points across all grade levels. These differences were approximately three-quarters of the pooled standard deviation or larger across all grades. The differences between mean scale scores of White students and mean scale scores of African American students ranged from approximately 62 points to 74 scale score points across all grades. These differences were more than one pooled standard deviation in each grade.
- Fully English proficient students outperformed their limited English proficiency peers in each grade with mean differences ranging from approximately 33 scale score points (about two-thirds of the pooled standard deviation) to about 51 scale score points (close to one pooled standard deviation) across all grades.
- Students without disabilities outperformed students with disabilities by approximately 41 scale score points (in grade 3) to 63 points (in grade 7). These differences were from three-quarters to over one pooled standard deviation.
- Students considered not economically disadvantaged outperformed economically disadvantaged students in every grade with mean differences ranging from approximately 40 scale score points in

grade 5 to over 45 scale score points in grade 7. These differences were larger than three-quarters of the pooled standard deviation in each grade.

- Students not using testing accommodations outperformed students using testing accommodations in each grade. The mean scale score differences ranged from approximately 67 points in grade 5 to 82 points in grade 7. These differences were larger than one pooled standard deviation in each grade.

Science

- The mean scale scores were only slightly lower for female students, with a difference of about 3 scale score points (one-twentieth of the pooled standard deviation) in grade 4 and no practical difference in grade 8.
- White students outperformed their Asian, Hispanic, American Indian, and African American peers in both Science grades. The difference between mean scale scores of White students and mean scale scores of Asian students was about 22 points (less than half of the pooled standard deviation) in grade 4 and about 7 points in grade 8. The difference between mean scale scores of White students and mean scale scores of Hispanic students was approximately 39 scale score points (three-quarters of the pooled standard deviation) in grade 4 and about 33 points (two-thirds of the pooled standard deviation) in grade 8. The difference between mean scale scores of White students and mean scale scores of American Indian students was approximately 35 points (about three-quarters of the pooled standard deviation) in both Science grades. The difference between mean scale scores of White students and mean scale scores of African American students was approximately 65 points in grade 4 and about 55 points in grade 8. These differences were more than one pooled standard deviation in each grade.
- Fully English proficient students outperformed their limited English proficiency peers in both grades with mean differences of about 39 points in grade 4 (about three-quarters of the pooled standard deviation) and approximately 45 points (close to one pooled standard deviation) in grade 8.
- Students without disabilities outperformed students with disabilities by approximately 35 scale score points (two-thirds of the pooled standard deviation) in grade 4 and 46 points (close to one pooled standard deviation) in grade 8.
- Students considered not economically disadvantaged outperformed economically disadvantaged students in both grades with a mean difference of approximately 39 scale score points (three-quarters of the pooled standard deviation) in grade 4 and a mean difference of approximately 34 scale score points (two-thirds of the pooled standard deviation) in grade 8.
- Students not using testing accommodations, as a group, had higher mean scale scores than students using testing accommodations in each grade. However, due to low numbers of students using testing accommodations in Science, the performance of these two groups cannot be reliably compared.

Social Studies

- There was no practical difference between female and male student mean scale scores in grade 4. Female students performed better than male students in grades 8 and 10, with a difference between

mean scale scores of approximately 3 points (less than one-tenth of the pooled standard deviation) in these grades.

- Similar to the patterns observed for ELA, Mathematics, and Science, White students outperformed their Asian, Hispanic, American Indian, and African American peers in all Social Studies grades. The differences between mean scale scores of White students and mean scale scores of Asian students ranged from about 19 points (less than half of the pooled standard deviation) in grade 4 to about 4 points (less than one-tenth of the pooled standard deviation) in grade 8. The differences between mean scale score of White students and mean scale score of Hispanic students ranged from 27 scale score points in grade 8 to approximately 33 points in grade 4. These differences were more than one-half of the pooled standard deviation in each grade. The differences between mean scale scores of White students and mean scale scores of American Indian students were approximately 32 points in grade 4 and about 29 points in grades 8 and 10. These differences were more than one-half of the pooled standard deviation in each grade. The differences between mean scale scores of White students and mean scale scores of African American ranged from approximately 46 points in grade 8 to about 57 points in grade 4. These differences were approximately one pooled standard deviation or larger in each grade.
- Fully English proficient students outperformed their limited English proficiency peers in all grades with mean differences ranging from 34 points in grade 4 (about three-quarters of the pooled standard deviation) and approximately 53 points (one pooled standard deviation) in grade 10.
- Students without disabilities outperformed students with disabilities by approximately 36 scale score points (two-thirds of the pooled standard deviation) in grade 4 to 51 points (close to one pooled standard deviation) in grade 10.
- Students considered not economically disadvantaged outperformed economically disadvantaged students in all grades with the mean differences ranging from approximately 32 scale score points in grade 8 to over 36 scale score points in grade 4. These differences were approximately two-thirds of the pooled standard deviation or larger in all grades.
- Students not using testing accommodations, as a group, had higher mean scale scores than students using testing accommodations in each grade. However, due to low numbers of students using testing accommodations in Social Studies, the performance of these two groups cannot be reliably compared.

10.3 Performance Level Classifications

Student performance on the Wisconsin Forward Exam is reported in terms of four performance categories: *Below Basic*, *Basic*, *Proficient*, and *Advanced*. These performance categories are established through cut scores.

Standard 5.21 of the *Standards* (AERA, APA, & NCME, 2014) indicates that “when proposed score interpretations involve one or more cut scores, the rationale and procedures used for establishing cut scores should be documented clearly” (p. 107).

In terms of the validity of the Wisconsin Forward Exam, it is essential to understand that cut scores and PLDs are established in a collaborative and participatory process. The descriptors clearly establish, in plain language, the proper frame of reference for understanding how to interpret test scores, particularly cut scores. PLDs summarize the knowledge, skills, and abilities expected of students in each achievement level. As stated in Part 7, DPI provided policy PLDs for the Wisconsin Forward Exam assessments. At the standard setting, DPI used the policy PLDs in conjunction with the content standards to consider the content-based expectations for students in each achievement level on each test in the Wisconsin Forward Exam program.

Tables 10-6 through 10-9 provide the scale score ranges that define performance levels together with the percentage of students in each performance level. The results for each content area and grade are summarized below.

English Language Arts

- Between approximately 37% (grade 3) and 45% (grade 4) of students were either *Proficient* or *Advanced*.
- Between about 6% and 10% of students were classified as *Advanced*, depending on the grade level.
- Across all grade levels, more than 50% of students were below *Proficient*. These percentages ranged from approximately 55% below *Proficient* in grade 4 to approximately 63% below *Proficient* in grade 8.

Mathematics

- Between approximately 31% (grade 8) and 48% (grade 3) of students were either *Proficient* or *Advanced*.
- The percentage of students who were *Advanced* ranged from approximately 4.5% (grade 7) to over 14% (grade 4).
- Across all grade levels, the percentages of students below *Proficient* ranged from approximately 52% in grade 3 to 69% in grade 8.

Science

- Approximately 51% of students in grade 4 and about 49% of students in grade 8 were either *Proficient* or *Advanced*.
- The percentage of students classified as *Advanced* was approximately 19% in grade 4 and about 20% in grade 8.
- The percentage of students classified as below *Proficient* was approximately 49% in grade 4 and about 51% in grade 8.

Social Studies

- Approximately 59% of students in grades 4 and 8 and about 47% of students in grade 10 were classified as *Proficient* or *Advanced*.
- Between 22.5% and 25% of students were classified as *Advanced* across the three grades.
- The percentage of students classified as below *Proficient* ranged from approximately 41% in grade 4 and 8 to 53% in grade 10.

10.3.1 Subgroup Patterns in Performance Level Results

The performance level results varied by subgroup: gender, race/ethnicity, English language proficiency, disability status, economic status, and accommodation use. The main subgroup performance patterns are described below. These comparisons are based on Tables 10-10 through 10-13.

In terms of gender, higher percentages of female students were classified as *Proficient* or above in all ELA grades. The differences in the percentages of male and female students in the *Proficient* or above categories ranged from approximately 5% (grade 4) to about 9% (grades 6 and 8) for ELA. A reversed trend was observed for Mathematics. Between approximately 2% (grade 8) and 8% (grade 4) more male students than female students were classified as *Proficient* or above in Mathematics. In Science, more male students than female students were classified as *Proficient* or above in grade 4 (with a difference of approximately 2%), and there was no practical difference between gender groups in grade 8. In Social Studies, the differences in percentages of male and female students classified as *Proficient* or above were about 4% in grade 8 and 3% in grade 10, with more female students being classified in the two highest performance categories in these grades. There was no practical difference in the percentages of male and female students classified as *Proficient* or above in grade 4. The performance of non-binary students cannot be reliably compared with the performance of male or female students due to the small number of the non-binary group members in each grade.

There were some consistent patterns in performance by ethnicity across grades and content areas. This pattern followed the pattern of the mean scale score differences. In terms of the *Proficient* or above categories, the prevailing tendency was that there were higher percentages of White students as a group, followed by Asian students, Hispanic students, American Indian students, and African American students. The inverse sequence was found at the *Below Basic* performance level.

Performance level results showed that there were higher percentages of fully English proficient students who were classified as *Proficient* or above compared to students who were of limited English proficiency in every grade and content area. These differences ranged from approximately 26% in ELA grade 3 to 40% in Social Studies grade 10.

Performance level results showed a similar pattern in comparisons for students without disabilities who were classified as *Proficient* or above compared to students with disabilities, with differences ranging from approximately 25% (ELA grade 3) to over 41% (Social Studies grade 8) depending on grade level and content area.

There were consistent differences in performance between economically disadvantaged students and not economically disadvantaged students. In every grade and content area, between approximately 27% (ELA and Mathematics grade 8 and Social Studies grade 10) and about 33% (Mathematics grades 4 through 6) more students who were not economically disadvantaged were classified as *Proficient* or above compared to their economically disadvantaged peers.

Performance level results showed that there were higher percentages of students not using testing accommodations who were classified as *Proficient* or above compared to students using testing accommodations. These differences ranged from approximately 9% to 20% across ELA grades and from 32% to 43% across Mathematics grades. The differences in the percentages of students in different performance levels between groups of students using and not using testing accommodations should be interpreted with caution for ELA due to a small number of students per grade who used ELA testing accommodations. The comparisons of percentages of students using and not using testing accommodations for Science and Social Studies are less reliable due to the fact that fewer than 40 students per grade used testing accommodations in these content areas.

10.4 Standard Performance Index for Content Standards

In addition to raw scores and scale scores, teachers and educational decision makers frequently need diagnostic information to inform instructional strategies. Diagnostic information also helps to identify individual student strengths and needs. This kind of information can be derived from scores on subsets of test items that estimate how much a student knows in a clearly defined skill domain. These skill domains are called content standards or standards and they reflect Wisconsin Forward Exam reporting categories. Scores on subsets of test items at the content standard level are called standard performance index (SPI) scores. The purpose of reporting SPI scores on the Wisconsin Forward Exam is to show the relationship between the overall achievement being measured (represented by the test score) and the skills within each of the content standards associated with the overall content area. Teachers may use the SPI scores for individual students as indicators of strengths and weaknesses, but the SPI scores are best corroborated by other evidence, such as homework, class participation, diagnostic test scores, or observation. District and school administrators may compare their results by content standard and grade level with the state mean percentage to better understand their strengths and weaknesses within a particular content area and grade level.

An SPI score can be interpreted as an estimate of the number of items a student would be expected to answer correctly if there had been 100 similar items for a given reporting category. For example, an SPI score of 77 for a given reporting category means that, if the student were given 100 similar items, the student would be expected to answer 77 of them correctly. This is a criterion-referenced score because it estimates how much a student knows in a clearly defined skill domain (i.e., the criterion). Technical readers can refer to Appendix K of this report for more details.

This approach—identifying student proficiency on each content standard—relates to the ELA and Mathematics Wisconsin Academic Standards and Wisconsin Standards for Science and Social Studies. SPI scores provide a more reliable estimate of student achievement on each content standard than is possible by simply reporting the

percentage correct. However, *SPI scores should be used for low-stakes purposes because these scores cannot be considered stable for any content standard with a small number of items.*

Readers should note that the average difficulty of items will vary across content standards and grades. Content standards vary in their complexity, level of abstraction, and cognitive demand. Some standards may be intrinsically more difficult than others, and the difficulty of individual items is determined, in part, by the difficulty of the content domain being measured. The current test blueprints do not specify the average difficulty level of items for each content standard within grades or across grades. If the difficulty of the items varies across years, grades, or content standards, the mean SPI scores will be affected by differences in item difficulty as well as differences in student ability. *Thus, differences in SPI scores across years, grades, or content standards should not be seen as reliable indicators of differences in student ability since these differences may be explained in whole or in part by differences in the difficulty of the items themselves.*

In general, comparisons across years, grades, or content standards are appropriate for assessing the relative difficulty of the items. Comparisons of individual student scores or group mean scores can provide useful information about the *relative* strengths and needs of individual students or groups on these standards in a given administration year.

Tables 10-14 through 10-18 identify the content standards or domains, the number of MC and non-MC items within each standard or domain, the total number of possible points per standard or domain, the mean raw score, the mean *p*-value, the standard deviation of the raw scores, the mean SPI score, and the standard deviation of SPI scores for all content areas across grades. The results from Tables 10-14 through 10-28 are summarized below. Tables 10-19 through 10-22 identify the SPI cut scores for each content area, reporting category, and grade level.

English Language Arts

Tables 10-14 and 10-15 present mean *p*-values and SPI scores for ELA across content standards and across domains, respectively. Results show that the mean ELA SPI scores across grades ranged from approximately 40 to 71 for content standards (including Listening) and from about 45 to 61 for domains (Reading or Writing), indicating that the items were, on average, moderately difficult for examinees. Content standard C (Reading—Vocabulary Use) was found to be easier than other standards with mean SPI values ranging from 61 to 71 across grades. Content standard D (Writing/Language—Text Types and Purposes) was found to be more difficult than other standards in grades 3 through 6 with the SPI means ranging from 40 to 44 in these grades. No particular pattern of difficulty was detected across other standards.

The Writing domain was more difficult than the Reading and Listening domains for students in all grades except grade 4 (with mean SPI scores ranging from 45 to 53 for the Writing domain), and the Listening domain was most difficult in grade 4 (with the mean SPI score of approximately 47). The Reading domain was moderately difficult across grades (with mean SPI scores ranging from 52 to 61).

Mathematics

Table 10-16 presents mean p -values and SPI scores for Mathematics across grades and content standards. Results show that the mean p -values and SPI scores varied across standards in all grades. Mean SPI scores across all content standards ranged from approximately 40 for content standard E (Geometry) in grades 6 through 8 and content standard G (The Number System) in grade 8 to about 57 for content standard A (Operations and Algebraic Thinking) and content standard G (Geometry) in grade 3. The Mathematics items were more challenging in higher grades than lower grades. No particular pattern of difficulty was detected across Mathematics standards within a grade.

Science

Table 10-17 presents mean p -values and SPI scores for Science across grades and content standards. The mean Science SPI scores across both grades and all content standards ranged from approximately 44 to 60, indicating that the test items were of medium difficulty, on average. Mean SPI scores indicated that content standard C (Earth and Space Science) was more difficult than other standards in both grades (with the mean SPI scores of approximately 44 in grade 4 and about 46 in grade 8).

Social Studies

Table 10-18 presents mean p -values and SPI scores for Social Studies across grades and content standards. The mean Social Studies SPI scores across all grades and content standards ranged from approximately 46 to 62, indicating that the test items were, on average, of medium difficulty. The mean SPI scores indicated that the most difficult content standard varied across the three Social Studies grades. The most difficult standard was content standard C (Political Science and Citizenship) in grade 4, content standard D (Economics) in grade 8, and content standard E (Behavioral Sciences) in grade 10.

Summary of Student Performance Indicator Results

Overall, the mean SPI scores across grades and content standards ranged in difficulty from approximately 40 (standard E, Geometry in Mathematics grades 6 through 8; standard G, The Number System in Mathematics grade 8; standard D, Writing/Language—Text Types and Purposes in ELA grades 3 and 6) to 71 (standard C, Reading—Vocabulary Use in ELA grade 7).

It is important to note that some variation in the difficulty of the items across content standards within and across grades and test forms is inevitable and that some of that variation is independent of any intrinsic differences in the difficulty of the standards themselves (e.g., variations in the difficulty of the particular items that were selected for the test forms). For this reason, SPI scores should be interpreted with caution and should not be used to make comparisons of student performance across testing years or grade levels.

10.5 Longitudinal Comparisons of Test Scores

It is often desirable to examine the scores of students across time and to monitor group performance. This is possible if the test content and the construct measured by the test are comparable from year to year and if the scores are reported on the same scale in multiple years.

Seven years of test scores on the same reporting scales are available for the Wisconsin Forward Exam assessments in ELA and Mathematics. The state-level mean scale scores and standard deviations for the 2016 through 2023 administrations are presented for ELA and Mathematics in Tables 10-23 and 10-24, respectively. New scales were established for the Science assessments after the Spring 2019 test administration. Because the new Science assessments were not linked to the previous scales, the Spring 2023 scale scores are comparable only with the Spring 2019, 2021, and 2022 scale scores and not with the scores prior to the Spring 2019 administration. Therefore, only four years of scale score data are presented for Science in Table 10-25. New scales were established for the Social Studies assessments after the Spring 2022 test administration. Because the new Social Studies assessments were not linked to the previous Social Studies scales, the Spring 2023 scale scores are comparable only with the Spring 2022 administration scores. Therefore, only two years of scale score data are presented for Social Studies in Table 10-26. The historical Social Studies data can be found in the *Wisconsin Forward Exam Spring 2021 Technical Report* available at DPI’s website: <https://dpi.wi.gov/assessment/forward/resources#documentation>.

The statistics presented in Tables 10-23 through 10-26 are based on the total population of Wisconsin students, including students attending public, choice, and private schools. The test participation rates for each year are also included in Tables 10-23 through 10-26. The participation rates are computed as the percentage of students who received a valid scale score given the number of students expected to take the test in each grade and content area. The “Enrolled” column shows the total number of students expected to take the test each year. The “Number Tested” and “Percent Tested” columns show the number and percentage of students who participated in the assessment and received a valid scale score.

It should be noted that students who participated in the Spring 2021 test administration were not fully representative of the Wisconsin student population. As such, the interpretation of the state-level performance trend changes between Spring 2019 and Spring 2021 and between Spring 2021 and subsequent administrations should be done with caution and in the context of disrupted learning due to the COVID-19 pandemic in the 2020–21 academic year. In this report we focus on the differences in student performance between Spring 2022 and Spring 2023.

It was observed that the mean scale scores were higher for ELA grades 3, 4, 6, 7, and 8 (see Table 10-23), all Mathematics grades (see Table 10-24), and Social Studies grade 8 (see Table 10-26) in Spring 2023 compared to Spring 2022. These score increases ranged from about 2 to 6 scale score points across ELA grades and from 1 to 6 scale score points across Mathematics grades. The mean scale score for Social Studies grade 8 in 2023 was about 1 scale score higher than the mean scale score for the same grade in 2022. There was no practical difference in the year-to-year mean scale score for Science grade 4 and Social Studies grade 4 (refer to Tables 10-25 and 10-26, respectively). Small mean scale score decreases from Spring 2022 to Spring 2023 were noted

for ELA grade 5, Science grade 8, and Social Studies grade 10. These differences were between 1 and 2 scale score points.

When student performance in Spring 2023 was compared to the pre-pandemic performance of students in Spring 2019, it was observed that the mean scale scores for ELA were still lower in Spring 2023 compared to Spring 2019 in grades 3, 5, 6, and 7 (difference of 2 to 5 scale score points). The Mathematics mean scale scores were lower in Spring 2023 compared to Spring 2019 in grades 7 and 8 (difference of approximately 5 to 6 scale score points). The Science mean scale scores were lower in Spring 2023 in both grades (difference of approximately 4 scale score points in grade 4 and about 7 scale score points in grade 8.) The Spring 2023 mean scale scores for ELA grades 4 and 8 and Mathematics grades 3 through 5 were either slightly higher than or comparable to the Spring 2019 mean scale scores in the corresponding grades.

Tables 10-27 through 10-30 show the percentages of students in each achievement level in the Spring 2016 through 2023 test administrations for ELA and Mathematics, in the Spring 2019 through 2023 test administrations for Science, and in the Spring 2022 and 2023 test administrations for Social Studies. The results presented in these tables are based on all Wisconsin students who participated in the assessments in a given year, including students attending public, choice, and private schools. The pattern of student performance classification change between Spring 2022 and Spring 2023 was consistent with the pattern of the mean scale score change between the two test administrations.

For ELA, an increase in the percentage of students at or above the *Proficient* cut was observed for all ELA grades except for grade 5. The increase ranged from approximately 2% in grade 6 to 4% in grade 8. The decrease in the percentage of students at or above the *Proficient* cut score was less than 2% for grade 5 (see Table 10-27). For Mathematics, the increase in the percentage of students classified as *Proficient* or above ranged from less than 1% in grades 3 and 8 to about 3% in grades 6 (see Table 10-28).

For Science (see Table 10-29), approximately the same percentage of students were classified as *Proficient* or above in Spring 2023 compared to Spring 2022 in both grades (with year-to year difference of less than half a percent). As stated earlier in the report, new performance level cut scores were established for Science after the Spring 2019 test administration, which is considered to be a new baseline for longitudinal comparisons. Therefore, no impact data from administrations prior to 2019 are presented in this report. The historical Science data can be found in the *Wisconsin Forward Exam 2018 Technical Report* available at DPI's website: <https://dpi.wi.gov/assessment/forward/resources#documentation>.

Similarly, for Social Studies (see Table 10-30), approximately the same percentage of students were classified as *Proficient* or above in Spring 2023 compared to Spring 2022 in all three grades (with year-to year difference of less than half a percent). As stated earlier in the report, new performance level cut scores were established for Social Studies after the Spring 2022 test administration. Therefore, the Spring 2023 Social Studies impact data can only be directly compared with the Spring 2022 impact data. The historical Social Studies data can be found in the *Wisconsin Forward Exam Spring 2021 Technical Report* available at DPI's website: <https://dpi.wi.gov/assessment/forward/resources#documentation>.

Overall, the percentages of students classified in the *Proficient* or above categories in Spring 2023 were found to be higher than or comparable to Spring 2022 for each grade and content area, except for ELA grade 5 (where there was a small year-to-year decline). These results likely reflect a continued recovery from the effects of the COVID-19 pandemic on student learning in the last two years.

10.6 Summary

In the Wisconsin Forward Exam, the purpose of the ELA, Mathematics, Science, and Social Studies assessments is to demonstrate student achievement through test scores in the respective content areas. The results presented in Part 10, together with the reliability and validity evidence presented in Parts 8 and 9, indicate that the scale scores and performance levels reported in the Wisconsin Forward Exam program are valid and reliable evidence of student achievement in the tested content areas and grades in Spring 2023. However, due to the circumstances related to the COVID-19 pandemic and lower participation in the assessment than in a typical year, we recommend that the results of the Spring 2021 test administration be treated and interpreted with caution. More reliable comparisons of student performance can be made between other administration years.

Classroom teachers may use the Spring 2023 scores as evidence of student achievement for students who participated in the assessment. District and school administrators may use this information for activities such as planning teaching activities in the next school year.

Table 10-1 Scale Score Descriptive Statistics for Total Population

Content	Grade	N Count	Mean	SD	Skewness	Kurtosis	Min	Max	LOSS	HOSS
English Language Arts	3	58497	552.33	48.85	-0.21	0.74	330	900	330	900
	4	58996	583.88	53.27	-0.06	0.27	340	801	340	930
	5	59386	592.89	51.66	-0.22	0.45	350	824	350	940
	6	59412	604.39	53.00	-0.34	0.42	360	864	360	950
	7	60413	622.59	55.27	-0.14	0.23	370	960	370	960
	8	62249	628.49	64.03	-0.23	0.41	380	970	380	970
Mathematics	3	58722	553.02	56.43	-0.35	1.47	360	760	360	760
	4	59165	576.62	55.40	-0.45	0.74	405	800	405	800
	5	59577	601.51	52.62	-0.63	1.39	430	830	430	830
	6	59570	610.58	57.27	-0.30	0.98	440	870	440	870
	7	60559	620.62	60.98	-0.49	0.68	450	880	450	880
	8	62360	638.57	56.40	-0.25	1.04	470	890	470	890
Science	4	59141	496.06	54.65	0.04	0.24	300	725	300	725
	8	62289	692.80	51.93	0.03	0.47	480	945	480	945
Social Studies	4	59131	499.74	50.82	-0.22	0.81	330	700	330	700
	8	62261	700.49	48.94	-0.28	0.73	540	860	540	860
	10	61819	796.91	53.02	-0.53	0.76	645	980	645	980

Table 10-2 Scale Score Descriptive Statistics by Subgroup, English Language Arts

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
3	Gender	Male	30027	548.79	48.86	330	900	
		Female	28463	556.06	48.56	330	871	-0.15
		Non-Binary	7	615.00	28.85	578	661	
	Race/Ethnicity	White	37862	564.14	44.26	330	871	
		African American	5948	510.45	46.44	330	690	1.20
		Hispanic	8088	532.77	46.93	330	735	0.70
		Asian	2743	547.00	49.60	330	900	0.38
		American Indian	553	527.74	46.10	330	662	0.82
		Two or More	3303	548.84	47.95	330	732	0.34
		Limited English Proficiency	No	53346	555.11	48.57	330	900
	Yes	5151	523.62	42.05	330	686	0.66	
	Disability Status	No	50140	557.39	47.22	330	900	
		Yes	8357	521.98	47.45	330	734	0.75
	Economically Disadvantaged	No	33351	567.22	45.30	330	900	
Yes		25146	532.59	46.32	330	753	0.76	
Accommodation Use	No	58416	552.37	48.83	330	900		
	Yes	81	524.77	56.31	371	640	0.57	
4	Gender	Male	29985	580.40	53.34	340	798	
		Female	29001	587.47	52.95	340	801	-0.13
		Non-Binary	10	604.20	90.80	379	698	
	Race/Ethnicity	White	38340	596.66	48.78	340	801	
		African American	5734	536.73	48.62	340	743	1.23
		Hispanic	8286	562.29	50.19	340	798	0.70
		Asian	2803	579.38	54.11	340	781	0.35
		American Indian	576	559.86	46.83	363	694	0.75
		Two or More	3257	579.57	52.57	394	798	0.35
		Limited English Proficiency	No	53787	587.07	52.89	340	801
	Yes	5209	550.96	45.49	340	741	0.69	
	Disability Status	No	50813	589.70	51.26	340	801	
		Yes	8183	547.74	51.28	340	776	0.82
	Economically Disadvantaged	No	34066	600.28	49.51	340	801	
Yes		24930	561.48	49.99	340	798	0.78	
Accommodation Use	No	58900	583.94	53.24	340	801		
	Yes	96	548.18	60.69	405	691	0.67	

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
5	Gender	Male	30408	587.62	51.52	350	782	
		Female	28964	598.41	51.22	350	824	-0.21
		Non-Binary	14	622.07	43.18	526	695	
	Race/Ethnicity	White	38648	604.86	47.31	350	824	
		African American	5652	549.00	47.36	350	728	1.18
		Hispanic	8561	571.56	49.44	350	759	0.70
		Asian	2756	592.94	53.09	350	807	0.25
		American Indian	568	567.44	46.98	350	730	0.79
		Two or More	3201	587.40	51.68	350	779	0.37
		Limited English Proficiency	No	54535	596.22	50.99	350	824
		Yes	4851	555.44	43.67	350	722	0.81
	Disability Status	No	51545	599.23	48.71	350	824	
		Yes	7841	551.20	51.10	350	735	0.98
	Economically Disadvantaged	No	34217	608.92	47.43	350	824	
		Yes	25169	571.11	49.13	350	773	0.79
Accommodation Use	No	59292	592.95	51.61	350	824		
	Yes	94	557.28	67.26	395	764	0.69	
6	Gender	Male	30307	598.19	53.43	360	863	
		Female	29091	610.82	51.76	360	864	-0.24
		Non-Binary	14	633.71	53.64	506	703	
	Race/Ethnicity	White	39006	616.18	48.09	360	864	
		African American	5703	559.26	51.88	360	786	1.17
		Hispanic	8527	584.41	51.38	360	764	0.65
		Asian	2626	606.38	52.39	360	819	0.20
		American Indian	579	575.21	54.43	360	724	0.85
		Two or More	2971	597.39	54.38	360	802	0.39
		Limited English Proficiency	No	55353	607.49	52.11	360	864
	Yes		4059	562.01	46.44	360	756	0.88
	Disability Status	No	51870	611.03	49.95	360	864	
		Yes	7542	558.72	50.77	360	741	1.05
	Economically Disadvantaged	No	34836	620.34	48.16	360	864	
		Yes	24576	581.77	51.29	360	803	0.78
Accommodation Use	No	59319	604.43	52.99	360	864		
	Yes	93	574.72	50.33	473	697	0.56	

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
7	Gender	Male	30933	616.41	55.05	370	945	
		Female	29452	629.05	54.74	370	960	-0.23
		Non-Binary	28	647.93	55.41	494	759	
	Race/Ethnicity	White	39840	633.84	51.08	370	945	
		African American	5842	578.79	52.96	370	765	1.07
		Hispanic	8729	602.02	53.73	370	809	0.62
		Asian	2509	628.53	55.01	443	960	0.10
		American Indian	597	596.01	50.60	370	791	0.74
		Two or More	2896	618.55	56.89	370	804	0.30
		Limited English Proficiency	No	56284	625.82	54.43	370	960
		Yes	4129	578.52	47.18	370	765	0.88
	Disability Status	No	52950	629.78	52.00	370	960	
		Yes	7463	571.58	50.73	370	762	1.12
	Economically Disadvantaged	No	35577	638.45	51.26	370	960	
		Yes	24836	599.87	52.81	370	835	0.74
Accommodation Use	No	60286	622.66	55.23	370	960		
	Yes	127	588.12	62.81	460	740	0.63	
8	Gender	Male	32033	620.57	64.69	380	970	
		Female	30165	636.84	62.24	380	970	-0.26
		Non-Binary	51	661.49	51.26	566	806	
	Race/Ethnicity	White	40964	641.44	59.38	380	970	
		African American	6171	579.52	61.67	380	783	1.04
		Hispanic	8917	604.72	61.92	380	907	0.61
		Asian	2687	639.06	62.09	428	970	0.04
		American Indian	606	597.90	58.98	380	762	0.73
		Two or More	2904	619.36	66.53	380	970	0.37
		Limited English Proficiency	No	58365	631.91	63.15	380	970
		Yes	3884	577.07	54.45	380	759	0.88
	Disability Status	No	54730	636.72	60.27	380	970	
		Yes	7519	568.53	58.44	380	846	1.14
	Economically Disadvantaged	No	37022	646.03	59.52	380	970	
		Yes	25227	602.74	61.67	380	970	0.72
Accommodation Use	No	62127	628.58	63.97	380	970		
	Yes	122	580.32	74.44	380	730	0.75	

Table 10-3 Scale Score Descriptive Statistics by Subgroup, Mathematics

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
3	Gender	Male	30150	556.76	58.31	360	760	
		Female	28565	549.07	54.11	360	760	0.14
		Non-Binary	7	594.71	27.99	559	636	
	Race/Ethnicity	White	37895	567.54	49.80	360	760	
		African American	5932	499.98	54.98	360	760	1.34
		Hispanic	8259	529.12	52.76	360	760	0.76
		Asian	2782	552.58	60.77	360	760	0.30
		American Indian	553	524.29	51.55	360	760	0.87
		Two or More	3301	546.68	56.10	360	760	0.41
		Limited English Proficiency	No	53345	556.09	56.00	360	760
	Yes	5377	522.60	51.48	360	760	0.60	
	Disability Status	No	50365	558.92	53.10	360	760	
		Yes	8357	517.49	62.55	360	760	0.76
	Economically Disadvantaged	No	33466	570.69	50.86	360	760	
Yes		25256	529.62	54.92	360	760	0.78	
Accommodation Use	No	58478	553.35	56.21	360	760		
	Yes	244	476.25	57.75	360	623	1.37	
4	Gender	Male	30064	580.41	57.58	405	800	
		Female	29091	572.69	52.76	405	800	0.14
		Non-Binary	10	594.30	73.80	444	674	
	Race/Ethnicity	White	38345	591.52	47.99	405	800	
		African American	5724	520.54	54.09	405	684	1.45
		Hispanic	8427	551.44	52.91	405	800	0.82
		Asian	2839	578.05	56.10	405	800	0.28
		American Indian	576	550.86	51.46	405	709	0.85
		Two or More	3254	568.19	55.50	405	800	0.48
		Limited English Proficiency	No	53766	579.77	54.98	405	800
	Yes	5399	545.28	49.53	405	688	0.63	
	Disability Status	No	50989	582.82	51.97	405	800	
		Yes	8176	537.96	60.31	405	800	0.84
	Economically Disadvantaged	No	34147	594.57	48.80	405	800	
Yes		25018	552.12	54.49	405	800	0.83	
Accommodation Use	No	57134	578.95	54.20	405	800		
	Yes	2031	510.96	48.29	405	642	1.26	

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
5	Gender	Male	30495	603.43	54.50	430	830	
		Female	29068	599.49	50.49	430	830	0.07
		Non-Binary	14	622.79	41.97	560	708	
	Race/Ethnicity	White	38666	615.14	45.48	430	830	
		African American	5646	550.30	51.76	430	699	1.40
		Hispanic	8703	576.90	52.15	430	830	0.82
		Asian	2793	607.44	53.70	430	830	0.17
		American Indian	565	577.11	48.99	430	707	0.84
		Two or More	3204	593.32	52.53	430	753	0.47
	Limited English Proficiency	No	54518	604.81	51.67	430	830	
		Yes	5059	565.97	49.45	430	717	0.75
	Disability Status	No	51743	607.59	49.10	430	830	
		Yes	7834	561.37	57.21	430	830	0.92
	Economically Disadvantaged	No	34309	618.50	45.39	430	830	
Yes		25268	578.45	52.98	430	830	0.82	
Accommodation Use	No	57160	604.22	50.97	430	830		
	Yes	2417	537.58	50.50	430	682	1.31	
6	Gender	Male	30382	612.05	58.90	440	870	
		Female	29174	609.04	55.49	440	870	0.05
		Non-Binary	14	625.36	54.63	535	732	
	Race/Ethnicity	White	39012	625.57	50.51	440	870	
		African American	5688	554.38	53.47	440	747	1.40
		Hispanic	8675	584.59	53.24	440	870	0.80
		Asian	2646	613.99	59.70	440	870	0.23
		American Indian	579	580.21	57.55	440	870	0.90
		Two or More	2970	600.00	58.65	440	870	0.50
	Limited English Proficiency	No	55339	614.09	56.33	440	870	
		Yes	4231	564.62	49.09	440	736	0.89
	Disability Status	No	52031	617.68	53.64	440	870	
		Yes	7539	561.53	57.49	440	870	1.04
	Economically Disadvantaged	No	34911	629.02	51.67	440	870	
Yes		24659	584.47	54.61	440	870	0.84	
Accommodation Use	No	56632	614.27	55.17	440	870		
	Yes	2938	539.42	50.15	440	734	1.36	

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
7	Gender	Male	31017	622.13	62.09	450	880	
		Female	29514	619.03	59.74	450	880	0.05
		Non-Binary	28	627.86	64.58	450	739	
	Race/Ethnicity	White	39852	635.86	53.40	450	880	
		African American	5829	561.47	59.62	450	773	1.37
		Hispanic	8857	593.68	58.73	450	841	0.78
		Asian	2533	628.54	63.25	450	880	0.14
		American Indian	597	588.26	57.80	450	759	0.89
		Two or More	2891	612.09	61.94	450	880	0.44
	Limited English Proficiency	No	56262	624.26	59.81	450	880	
		Yes	4297	573.03	56.02	450	770	0.86
	Disability Status	No	53103	628.36	56.80	450	880	
		Yes	7456	565.52	61.34	450	770	1.10
	Economically Disadvantaged	No	35639	639.42	54.35	450	880	
Yes		24920	593.74	59.87	450	880	0.81	
Accommodation Use	No	57581	624.64	58.68	450	880		
	Yes	2978	542.92	51.83	450	727	1.40	
8	Gender	Male	32094	637.81	58.97	470	890	
		Female	30215	639.35	53.52	470	890	-0.03
		Non-Binary	51	654.39	53.43	528	890	
	Race/Ethnicity	White	40964	651.62	51.51	470	890	
		African American	6153	589.84	51.08	470	770	1.20
		Hispanic	9024	614.60	52.14	470	890	0.72
		Asian	2707	649.37	59.50	470	890	0.04
		American Indian	606	610.35	50.32	470	827	0.80
		Two or More	2906	628.03	58.88	470	890	0.45
	Limited English Proficiency	No	58341	641.51	55.77	470	890	
		Yes	4019	595.94	47.49	470	890	0.82
	Disability Status	No	54849	645.34	53.37	470	890	
		Yes	7511	589.18	53.22	470	890	1.05
	Economically Disadvantaged	No	37072	655.14	52.40	470	890	
Yes		25288	614.29	53.17	470	890	0.77	
Accommodation Use	No	59473	641.79	54.82	470	890		
	Yes	2887	572.36	46.83	470	736	1.27	

Table 10-4 Scale Score Descriptive Statistics by Subgroup, Science

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
4	Gender	Male	30053	497.45	55.71	300	725	
		Female	29078	494.61	53.47	300	725	0.05
		Non-Binary	10	551.30	86.67	367	665	
	Race/Ethnicity	White	38347	510.38	50.06	300	725	
		African American	5717	444.99	46.23	300	645	1.32
		Hispanic	8413	471.65	50.29	300	725	0.77
		Asian	2838	487.92	53.18	300	725	0.45
		American Indian	575	475.50	46.55	301	612	0.70
		Two or More	3251	490.94	54.01	300	725	0.39
	Limited English Proficiency	No	53745	499.58	54.30	300	725	
		Yes	5396	460.97	44.85	300	635	0.72
	Disability Status	No	50970	500.96	53.39	300	725	
		Yes	8171	465.52	52.50	300	695	0.67
	Economically Disadvantaged	No	34137	512.54	51.01	300	725	
Yes		25004	473.57	51.34	300	725	0.76	
Accommodation Use	No	59103	496.09	54.64	300	725		
	Yes	38	446.71	52.81	368	569	0.90	
8	Gender	Male	32051	693.28	54.13	480	945	
		Female	30187	692.26	49.47	480	945	0.02
		Non-Binary	51	718.47	43.48	599	836	
	Race/Ethnicity	White	40941	704.60	48.72	480	945	
		African American	6136	649.32	44.83	480	811	1.15
		Hispanic	9005	671.44	47.57	480	945	0.68
		Asian	2705	697.39	49.77	480	945	0.15
		American Indian	604	669.83	47.13	480	851	0.71
		Two or More	2898	685.16	52.75	480	945	0.40
	Limited English Proficiency	No	58277	695.67	51.42	480	945	
		Yes	4012	651.18	40.01	480	786	0.88
	Disability Status	No	54784	698.32	49.94	480	945	
		Yes	7505	652.57	48.34	480	945	0.92
	Economically Disadvantaged	No	37033	706.59	49.44	480	945	
Yes		25256	672.60	48.78	480	899	0.69	
Accommodation Use	No	62260	692.82	51.92	480	945		
	Yes	29	657.17	60.42	515	755	0.69	

Table 10-5 Scale Score Descriptive Statistics by Subgroup, Social Studies

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
4	Gender	Male	30047	500.24	52.73	330	700	
		Female	29074	499.21	48.75	330	700	0.02
		Non-Binary	10	525.30	60.30	404	611	
	Race/Ethnicity	White	38340	512.00	46.41	330	700	
		African American	5702	455.10	47.54	330	676	1.22
		Hispanic	8422	479.17	47.80	330	700	0.70
		Asian	2840	493.15	51.34	330	700	0.40
		American Indian	576	479.59	44.33	330	652	0.70
		Two or More	3251	495.97	50.47	330	700	0.34
		Limited English Proficiency	No	53730	502.85	50.39	330	700
	Yes	5401	468.73	44.24	330	603	0.68	
	Disability Status	No	50965	504.76	48.68	330	700	
		Yes	8166	468.35	52.59	330	700	0.74
	Economically Disadvantaged	No	34138	515.09	46.73	330	700	
Yes		24993	478.76	48.63	330	700	0.76	
Accommodation Use	No	59093	499.77	50.79	330	700		
	Yes	38	450.37	66.65	330	636	0.97	
8	Gender	Male	32021	699.03	51.18	540	860	
		Female	30189	701.98	46.39	540	860	-0.06
		Non-Binary	51	729.53	36.35	661	844	
	Race/Ethnicity	White	40943	710.15	45.81	540	860	
		African American	6106	663.86	47.52	540	860	1.01
		Hispanic	9008	682.79	46.33	540	860	0.60
		Asian	2708	705.99	48.55	540	860	0.09
		American Indian	600	680.78	44.13	540	844	0.64
		Two or More	2896	695.01	50.07	540	860	0.33
		Limited English Proficiency	No	58250	703.07	48.38	540	860
	Yes	4011	662.99	40.94	540	812	0.84	
	Disability Status	No	54777	706.21	46.07	540	860	
		Yes	7484	658.61	48.91	540	860	1.03
	Economically Disadvantaged	No	37058	713.33	45.71	540	860	
		Yes	25203	681.60	47.37	540	860	0.68
	Accommodation Use	No	62232	700.50	48.93	540	860	
		Yes	29	677.48	51.78	580	761	0.47

Grade	Category	Subgroup	N Count	Mean	SD	Min.	Max.	Cohen's <i>d</i>
10	Gender	Male	31776	795.34	55.10	645	980	
		Female	29983	798.53	50.65	645	980	-0.06
		Non-Binary	60	820.52	56.04	645	956	
	Race/Ethnicity	White	43110	806.21	49.19	645	980	
		African American	4602	753.34	53.81	645	951	1.06
		Hispanic	8523	775.65	52.69	645	980	0.61
		Asian	2510	799.09	53.06	645	951	0.14
		American Indian	526	777.68	49.62	645	915	0.58
		Two or More	2548	791.25	54.43	645	980	0.30
	Limited English Proficiency	No	58453	799.79	51.83	645	980	
		Yes	3366	746.89	48.36	645	876	1.02
	Disability Status	No	55289	802.27	50.32	645	980	
		Yes	6530	751.58	53.64	645	980	1.00
	Economically Disadvantaged	No	39909	808.67	49.02	645	980	
		Yes	21910	775.50	53.32	645	980	0.66
	Accommodation Use	No	61798	796.92	53.02	645	980	
		Yes	21	762.62	44.45	659	861	0.65

Table 10-6 Score Ranges and Associated Impact Data, English Language Arts

Grade	Score Range				Impact Data				
	Below Basic	Basic	Proficient	Advanced	Below Basic	Basic	Proficient	Advanced	Proficient + Advanced
3	330–521	522–569	570–623	624–900	25.93	36.89	30.89	6.30	37.18
4	340–545	546–591	592–649	650–930	23.21	31.78	34.74	10.27	45.01
5	350–563	564–609	610–669	670–940	27.81	33.75	32.35	6.09	38.44
6	360–571	572–621	622–670	671–950	25.57	34.38	31.11	8.94	40.05
7	370–584	585–637	638–696	697–960	24.40	34.70	32.72	8.18	40.90
8	380–591	592–651	652–707	708–970	26.84	35.54	28.24	9.38	37.62

Table 10-7 Score Ranges and Associated Impact Data, Mathematics

Grade	Score Range				Impact Data				
	Below Basic	Basic	Proficient	Advanced	Below Basic	Basic	Proficient	Advanced	Proficient + Advanced
3	360–516	517–559	560–610	611–760	23.17	29.02	34.84	12.97	47.81
4	405–535	536–587	588–632	633–800	21.00	33.70	30.96	14.34	45.30
5	430–573	574–610	611–657	658–830	25.53	27.30	35.66	11.52	47.18
6	440–581	582–625	626–687	688–870	27.81	30.43	34.81	6.94	41.75
7	450–605	606–646	647–711	712–880	35.62	29.20	30.66	4.51	35.18
8	470–619	620–666	667–717	718–890	34.69	34.12	24.83	6.36	31.19

Table 10-8 Score Ranges and Associated Impact Data, Science

Grade	Score Range				Impact Data				
	Below Basic	Basic	Proficient	Advanced	Below Basic	Basic	Proficient	Advanced	Proficient + Advanced
4	300–446	447–495	496–542	543–725	18.82	29.93	31.85	19.40	51.25
8	480–652	653–694	695–736	737–945	21.89	29.25	29.13	19.74	48.86

Table 10-9 Score Ranges and Associated Impact Data, Social Studies

Grade	Score Range				Impact Data				
	Below Basic	Basic	Proficient	Advanced	Below Basic	Basic	Proficient	Advanced	Proficient + Advanced
4	330–460	461–490	491–536	537–700	21.01	19.87	36.14	22.98	59.12
8	540–661	662–692	693–733	734–860	20.03	20.91	34.36	24.70	59.06
10	645–769	770–804	805–836	837–980	27.39	25.36	24.72	22.53	47.24

Table 10-10 Percentage of Students in Each Performance Level by Subgroup, English Language Arts

Grade	Performance Level and Student Count	Gender			Race/Ethnicity						ELP		Disability		Economic Status		Accommodations	
		Female	Male	Non-Binary	White	African American	Hispanic	Asian	American Indian	Two or More	Fully English Proficient	Limited English Proficient	Not Disabled	Disabled	Not Economically Disadvantaged	Economically Disadvantaged	No Accommodations	Using Accommodations
3	BB	23.36	28.37	0.00	16.60	60.57	40.70	30.95	43.22	27.22	23.81	47.85	21.66	51.50	15.26	40.07	25.89	49.38
	B	36.55	37.22	0.00	37.27	29.93	38.23	38.61	40.33	39.72	36.69	38.94	37.54	32.98	35.25	39.07	36.91	22.22
	P	32.79	29.07	71.43	38.01	8.78	18.60	23.55	15.19	27.79	32.67	12.37	33.77	13.59	40.09	18.67	30.89	27.16
	A	7.30	5.35	28.57	8.11	0.72	2.47	6.89	1.27	5.27	6.83	0.83	7.03	1.93	9.40	2.18	6.31	1.23
	N Count	28463	30027	7	37862	5948	8088	2743	553	3303	53346	5151	50140	8357	33351	25146	58416	81
4	BB	21.21	25.16	10.00	14.31	58.89	37.04	25.19	37.50	25.82	21.24	43.56	18.89	50.08	12.81	37.42	23.17	48.96
	B	31.44	32.11	20.00	30.98	28.65	35.23	35.32	38.72	33.53	31.15	38.26	31.96	30.62	29.18	35.32	31.79	22.92
	P	35.93	33.58	40.00	41.55	10.95	23.58	29.47	21.01	31.75	36.47	16.89	37.70	16.33	42.98	23.47	34.75	25.00
	A	11.42	9.15	30.00	13.16	1.50	4.15	10.02	2.78	8.90	11.14	1.29	11.45	2.97	15.02	3.78	10.29	3.13
	N Count	29001	29985	10	38340	5734	8286	2803	576	3257	53787	5209	50813	8183	34066	24930	58900	96
5	BB	24.44	31.03	14.29	18.58	62.97	43.43	28.96	45.07	31.43	25.26	56.52	22.73	61.20	16.36	43.37	27.77	54.26
	B	33.22	34.26	7.14	34.28	27.64	34.91	33.67	36.97	34.43	33.71	34.12	34.90	26.13	32.66	35.22	33.75	27.66
	P	34.81	29.99	71.43	39.36	8.67	19.41	29.75	17.08	29.15	34.44	8.86	35.49	11.72	41.85	19.44	32.38	12.77
	A	7.53	4.72	7.14	7.79	0.73	2.24	7.62	0.88	5.00	6.59	0.49	6.87	0.94	9.12	1.97	6.09	5.32
	N Count	28964	30408	14	38648	5652	8561	2756	568	3201	54535	4851	51545	7841	34217	25169	59292	94
6	BB	21.79	29.20	14.29	17.17	59.74	38.45	24.60	44.21	30.46	23.24	57.40	20.42	60.97	14.92	40.66	25.54	46.24
	B	33.76	34.99	14.29	34.32	28.72	37.59	35.76	37.13	35.07	34.46	33.28	35.34	27.78	32.78	36.65	34.39	32.26
	P	33.31	28.99	57.14	37.22	10.19	20.15	29.47	16.93	26.79	32.76	8.70	34.20	9.85	39.14	19.73	31.13	19.35
	A	11.14	6.82	14.29	11.29	1.35	3.81	10.17	1.73	7.67	9.55	0.62	10.03	1.41	13.15	2.96	8.95	2.15
	N Count	29091	30307	14	39006	5703	8527	2626	579	2971	55353	4059	51870	7542	34836	24576	59319	93

Grade	Performance Level and Student Count	Gender			Race/Ethnicity						ELP		Disability		Economic Status		Accommodations	
		Female	Male	Non-Binary	White	African American	Hispanic	Asian	American Indian	Two or More	Fully English Proficient	Limited English Proficient	Not Disabled	Disabled	Not Economically Disadvantaged	Economically Disadvantaged	No Accommodations	Using Accommodations
7	BB	20.66	27.98	14.29	16.72	55.70	37.15	20.84	42.04	27.94	22.21	54.32	19.00	62.70	14.57	38.49	24.35	48.03
	B	34.19	35.21	14.29	34.78	30.37	36.60	37.47	36.52	33.91	34.66	35.31	35.81	26.83	32.91	37.27	34.71	31.50
	P	35.01	30.50	60.71	38.41	12.50	22.77	30.49	18.93	29.94	34.40	9.76	36.00	9.39	40.80	21.14	32.75	14.96
	A	10.14	6.31	10.71	10.09	1.44	3.47	11.20	2.51	8.22	8.74	0.61	9.18	1.09	11.73	3.10	8.19	5.51
	N Count	29452	30933	28	39840	5842	8729	2509	597	2896	56284	4129	52950	7463	35577	24836	60286	127
8	BB	22.19	31.25	7.84	19.13	56.55	39.88	21.59	44.39	33.61	24.73	58.63	21.43	66.25	17.03	41.23	26.79	51.64
	B	35.48	35.59	39.22	35.79	31.81	36.97	37.22	38.28	33.47	35.65	33.88	36.90	25.64	34.52	37.03	35.55	29.51
	P	30.96	25.67	39.22	33.37	10.24	19.39	27.80	15.18	24.55	29.65	7.03	31.14	7.18	35.10	18.18	28.26	17.21
	A	11.36	7.50	13.73	11.72	1.39	3.76	13.40	2.15	8.37	9.97	0.46	10.54	0.93	13.34	3.56	9.39	1.64
	N Count	30165	32033	51	40964	6171	8917	2687	606	2904	58365	3884	54730	7519	37022	25227	62127	122

Note: The abbreviation “BB” is for the Below Basic performance level, “B” is for the Basic performance level, “P” is for the Proficient performance level, and “A” is for the Advanced performance level.

Table 10-11 Percentage of Students in Each Performance Level by Subgroup, Mathematics

Grade	Performance Level and Student Count	Gender			Race/Ethnicity						ELP		Disability		Economic Status		Accommodations	
		Female	Male	Non-Binary	White	African American	Hispanic	Asian	American Indian	Two or More	Fully English Proficient	Limited English Proficient	Not Disabled	Disabled	Not Economically Disadvantaged	Economically Disadvantaged	No Accommodations	Using Accommodations
3	BB	24.88	21.56	0	13.21	60.82	38.67	25.02	40.69	26.60	21.14	43.31	19.12	47.59	12.42	37.42	22.96	73.77
	B	31.10	27.05	14.29	27.99	26.77	33.32	30.63	36.71	31.45	28.54	33.77	29.19	27.98	25.81	33.27	29.05	21.31
	P	33.30	36.29	57.14	41.93	11.13	23.25	30.01	19.71	31.60	36.31	20.18	37.48	18.89	42.73	24.37	34.96	4.51
	A	10.73	15.10	28.57	16.87	1.28	4.76	14.34	2.89	10.36	14.01	2.73	14.21	5.54	19.03	4.94	13.03	0.41
	N Count	28565	30150	7	37895	5932	8259	2782	553	3301	53345	5377	50365	8357	33466	25256	58478	244
4	BB	22.22	19.81	20.00	11.20	59.94	36.51	20.75	34.55	25.63	19.10	39.93	16.82	47.03	10.54	35.27	19.31	68.44
	B	36.65	30.86	10.00	32.68	29.68	38.46	36.00	42.36	37.03	32.98	40.90	33.95	32.18	30.25	38.42	33.93	27.47
	P	29.47	32.39	40.00	37.60	8.89	19.70	27.23	19.10	26.03	32.44	16.17	33.41	15.68	38.36	20.85	31.92	3.79
	A	11.66	16.93	30.00	18.53	1.48	5.33	16.03	3.99	11.31	15.48	3.00	15.82	5.11	20.85	5.45	14.84	0.30
	N Count	29091	30064	10	38345	5724	8427	2839	576	3254	53766	5399	50989	8176	34147	25018	57134	2031
5	BB	25.78	25.29	14.29	15.06	66.47	43.21	23.38	41.24	30.77	23.11	51.61	20.78	56.85	13.52	41.82	23.37	76.42
	B	29.43	25.27	28.57	26.75	23.08	30.82	28.14	33.45	29.93	26.86	32.00	27.92	23.21	25.25	30.08	27.66	18.70
	P	35.30	36.00	35.71	43.34	9.65	21.97	33.51	22.30	30.18	37.57	15.10	38.54	16.65	44.39	23.80	36.98	4.55
	A	9.49	13.44	21.43	14.85	0.80	4.00	14.97	3.01	9.11	12.47	1.28	12.76	3.29	16.84	4.29	11.99	0.33
	N Count	29068	30495	14	38666	5646	8703	2793	565	3204	54518	5059	51743	7834	34309	25268	57160	2417
6	BB	28.12	27.53	21.43	16.95	69.08	45.68	27.32	47.67	35.93	25.12	63.03	22.58	63.93	15.66	45.02	25.05	81.18
	B	32.06	28.87	28.57	30.96	22.45	32.93	30.57	32.12	31.08	30.66	27.46	31.37	23.94	29.06	32.38	31.17	16.13
	P	33.73	35.85	35.71	43.13	8.07	19.42	31.78	18.13	27.71	36.79	8.91	38.30	10.74	44.74	20.76	36.49	2.55
	A	6.09	7.75	14.29	8.97	0.40	1.96	10.32	2.07	5.29	7.43	0.59	7.75	1.38	10.54	1.85	7.29	0.14
	N Count	29174	30382	14	39012	5688	8675	2646	579	2970	55339	4231	52031	7539	34911	24659	56632	2938

Grade	Performance Level and Student Count	Gender			Race/Ethnicity						ELP		Disability		Economic Status		Accommodations	
		Female	Male	Non-Binary	White	African American	Hispanic	Asian	American Indian	Two or More	Fully English Proficient	Limited English Proficient	Not Disabled	Disabled	Not Economically Disadvantaged	Economically Disadvantaged	No Accommodations	Using Accommodations
7	BB	36.26	35.03	25.00	24.62	76.41	54.69	34.35	60.97	42.51	32.89	71.45	30.15	74.57	22.30	54.67	32.79	90.50
	B	30.53	27.93	32.14	31.41	17.10	28.18	27.91	23.95	28.43	29.81	21.18	30.95	16.72	30.51	27.33	30.28	8.36
	P	29.27	31.98	39.29	38.22	6.18	16.00	29.53	14.57	25.04	32.46	7.14	33.84	8.02	40.45	16.66	32.19	1.07
	A	3.94	5.06	3.57	5.74	0.31	1.13	8.21	0.50	4.01	4.84	0.23	5.05	0.68	6.74	1.34	4.74	0.07
	N Count	29514	31017	28	39852	5829	8857	2533	597	2891	56262	4297	53103	7456	35639	24920	57581	2978
8	BB	33.35	35.98	17.65	24.03	73.30	53.52	30.14	58.42	43.98	32.26	69.94	29.30	74.04	22.09	53.16	32.14	87.25
	B	36.28	32.07	49.02	36.92	21.62	31.62	32.58	28.88	31.42	34.75	24.91	36.10	19.69	35.72	31.77	35.21	11.67
	P	24.74	24.91	25.49	31.01	4.47	13.08	25.71	11.39	19.30	26.21	4.75	27.50	5.31	32.81	13.12	25.99	0.97
	A	5.63	7.05	7.84	8.04	0.62	1.78	11.56	1.32	5.30	6.77	0.40	7.10	0.96	9.38	1.95	6.67	0.10
	N Count	30215	32094	51	40964	6153	9024	2707	606	2906	58341	4019	54849	7511	37072	25288	59473	2887

Note: The abbreviation “BB” is for the Below Basic performance level, “B” is for the Basic performance level, “P” is for the Proficient performance level, and “A” is for the Advanced performance level.

Table 10-12 Percentage of Students in Each Performance Level by Subgroup, Science

Grade	Performance Level and Student Count	Gender			Race/Ethnicity						ELP		Disability		Economic Status		Accommodations	
		Female	Male	Non-Binary	White	African American	Hispanic	Asian	American Indian	Two or More	Fully English Proficient	Limited English Proficient	Not Disabled	Disabled	Not Economically Disadvantaged	Economically Disadvantaged	No Accommodations	Using Accommodations
4	BB	18.56	19.07	10.00	10.04	54.28	31.97	22.20	27.83	21.44	16.88	38.14	15.65	38.55	9.95	30.92	18.79	55.26
	B	31.36	28.55	10.00	27.59	31.71	36.68	33.79	38.26	32.11	28.92	40.05	29.23	34.32	25.51	35.96	29.94	23.68
	P	32.30	31.41	20.00	37.30	11.58	23.10	28.82	26.26	29.50	33.20	18.35	33.93	18.87	37.52	24.10	31.86	13.16
	A	17.78	20.96	60.00	25.07	2.43	8.25	15.19	7.65	16.95	21.00	3.47	21.19	8.26	27.01	9.01	19.41	7.89
	N Count	29078	30053	10	38347	5717	8413	2838	575	3251	53745	5396	50970	8171	34137	25004	59103	38
8	BB	20.73	23.01	1.96	13.86	54.01	34.95	18.60	36.26	26.81	19.90	50.70	17.53	53.72	13.46	34.25	21.87	51.72
	B	30.54	28.04	29.41	27.53	31.34	34.28	30.61	35.26	30.92	28.74	36.69	29.38	28.29	25.96	34.07	29.26	13.79
	P	30.65	27.69	35.29	33.58	11.65	22.10	29.21	20.86	26.78	30.37	11.04	31.38	12.70	33.92	22.11	29.13	24.14
	A	18.09	21.27	33.33	25.04	3.00	8.67	21.59	7.62	15.49	20.99	1.57	21.72	5.29	26.66	9.58	19.74	10.34
	N Count	30187	32051	51	40941	6136	9005	2705	604	2898	58277	4012	54784	7505	37033	25256	62260	29

Note: The abbreviation “BB” is for the Below Basic performance level, “B” is for the Basic performance level, “P” is for the Proficient performance level, and “A” is for the Advanced performance level.

Table 10-13 Percentage of Students in Each Performance Level by Subgroup, Social Studies

Grade	Performance Level and Student Count	Gender			Race/Ethnicity						ELP		Disability		Economic Status		Accommodations	
		Female	Male	Non-Binary	White	African American	Hispanic	Asian	American Indian	Two or More	Fully English Proficient	Limited English Proficient	Not Disabled	Disabled	Not Economically Disadvantaged	Economically Disadvantaged	No Accommodations	Using Accommodations
4	BB	20.10	21.90	10.00	12.22	55.93	34.47	25.00	31.77	23.16	19.00	40.99	17.24	44.59	11.21	34.41	20.98	65.79
	B	20.81	18.95	10.00	17.98	22.03	24.86	23.24	28.30	20.92	19.07	27.83	19.36	23.01	16.47	24.50	19.87	7.89
	P	37.68	34.65	30.00	40.70	17.94	29.60	32.01	30.73	35.74	37.16	25.98	38.26	22.90	40.45	30.25	36.15	13.16
	A	21.41	24.49	50.00	29.09	4.10	11.07	19.75	9.20	20.18	24.77	5.20	25.14	9.50	31.87	10.85	22.99	13.16
	N Count	29074	30047	10	38340	5702	8422	2840	576	3251	53730	5401	50965	8166	34138	24993	59093	38
8	BB	17.44	22.50	3.92	13.60	46.72	30.10	16.69	31.50	24.07	18.31	45.05	15.54	52.89	12.23	31.51	20.02	44.83
	B	21.25	20.61	7.84	18.67	25.63	26.89	21.53	27.67	21.89	20.10	32.61	20.40	24.59	17.26	26.27	20.91	20.69
	P	37.30	31.57	47.06	37.27	21.59	30.45	34.53	30.83	32.87	35.39	19.42	36.83	16.29	37.53	29.70	34.37	13.79
	A	24.02	25.32	41.18	30.45	6.06	12.57	27.25	10.00	21.17	26.20	2.92	27.22	6.24	32.99	12.52	24.70	20.69
	N Count	30189	32021	51	40943	6106	9008	2708	600	2896	58250	4011	54777	7484	37058	25203	62232	29
10	BB	25.12	29.58	8.33	20.43	61.71	42.39	25.10	42.02	32.26	25.18	65.89	23.00	64.56	18.88	42.91	27.38	57.14
	B	26.08	24.68	26.67	25.10	21.90	27.56	27.09	29.28	26.14	25.39	24.81	25.87	21.06	24.42	27.07	25.36	28.57
	P	26.64	22.91	23.33	27.55	11.13	19.18	24.86	17.30	21.31	25.69	7.84	26.56	9.16	27.98	18.77	24.72	9.52
	A	22.16	22.83	41.67	26.92	5.26	10.86	22.95	11.41	20.29	23.74	1.46	24.57	5.22	28.72	11.25	22.53	4.76
	N Count	29983	31776	60	43110	4602	8523	2510	526	2548	58453	3366	55289	6530	39909	21910	61798	21

Note: The abbreviation “BB” is for the Below Basic performance level, “B” is for the Basic performance level, “P” is for the Proficient performance level, and “A” is for the Advanced performance level.

Table 10-14 Summary Statistics for Content Standards Raw and SPI Scores, English Language Arts

Grade	Content Standard	Standard	No. of Items		Total Score Points	Mean	Mean <i>p</i> -Value	SD	SPI	
			MC	CR					Mean	SD
3	A	Reading—Key Ideas and Details	3	4	11	6.26	0.57	2.96	56.84	24.11
	B	Reading—Craft & Structure/Integration of Knowledge & Ideas	7	0	7	3.32	0.48	1.77	47.94	20.15
	C	Reading—Vocabulary Use	4	0	4	2.54	0.64	1.16	63.42	22.57
	D	Writing/Language—Text Types and Purposes	2	2	12	4.81	0.56	2.18	40.20	15.12
	E	Writing/Language—Research	3	2	6	3.16	0.54	1.49	52.86	18.78
	F	Writing/Language—Language Conventions	4	1	6	2.87	0.52	1.31	47.99	13.90
	G	Listening	5	1	7	3.86	0.56	1.76	54.97	19.29
4	A	Reading—Key Ideas and Details	3	3	9	4.64	0.55	2.03	51.64	18.67
	B	Reading—Craft & Structure/Integration of Knowledge & Ideas	5	2	9	4.72	0.53	2.18	52.50	20.29
	C	Reading—Vocabulary Use	6	0	6	3.66	0.61	1.63	60.91	22.15
	D	Writing/Language—Text Types and Purposes	2	2	12	5.26	0.56	2.45	44.06	17.32
	E	Writing/Language—Research	4	1	6	3.37	0.58	1.59	56.30	20.75
	F	Writing/Language—Language Conventions	4	1	6	3.56	0.58	1.38	59.33	17.31
	G	Listening	4	2	8	3.76	0.52	1.93	47.46	18.18
5	A	Reading—Key Ideas and Details	6	2	10	4.97	0.51	2.63	49.96	23.02
	B	Reading—Craft & Structure/Integration of Knowledge & Ideas	8	1	10	4.59	0.48	2.25	46.32	18.66
	C	Reading—Vocabulary Use	4	0	4	2.81	0.70	1.14	69.38	22.01
	D	Writing/Language—Text Types and Purposes	2	2	12	5.02	0.45	1.94	42.20	13.02
	E	Writing/Language—Research	1	2	5	2.70	0.55	1.37	53.89	20.83
	F	Writing/Language—Language Conventions	5	1	7	3.70	0.53	1.59	53.01	16.85
	G	Listening	4	2	8	4.97	0.63	2.13	61.44	21.23
6	A	Reading—Key Ideas and Details	4	3	10	5.89	0.59	2.38	58.75	21.06
	B	Reading—Craft & Structure/Integration of Knowledge & Ideas	6	2	10	6.09	0.60	2.69	60.84	24.49
	C	Reading—Vocabulary Use	4	0	4	2.72	0.68	1.15	67.64	21.38
	D	Writing/Language—Text Types and Purposes	2	2	12	4.73	0.43	2.20	39.94	14.65
	E	Writing/Language—Research	1	3	6	2.81	0.47	1.55	46.89	18.54
	F	Writing/Language—Language Conventions	4	1	6	3.36	0.54	1.41	56.09	16.85
	G	Listening	4	2	8	4.10	0.54	1.98	51.63	18.83

Grade	Content Standard	Standard	No. of Items		Total Score Points	Mean	Mean <i>p</i> -Value	SD	SPI	
			MC	CR					Mean	SD
7	A	Reading—Key Ideas and Details	3	4	10	4.40	0.43	2.17	44.44	18.04
	B	Reading—Craft & Structure/Integration of Knowledge & Ideas	5	2	9	4.77	0.53	2.05	53.12	18.80
	C	Reading—Vocabulary Use	5	0	5	3.58	0.72	1.43	70.94	24.08
	D	Writing/Language—Text Types and Purposes	1	3	13	5.79	0.44	2.62	44.96	16.71
	E	Writing/Language—Research	2	2	6	3.11	0.52	1.52	52.14	18.21
	F	Writing/Language—Language Conventions	4	1	5	3.14	0.63	1.24	62.69	17.73
	G	Listening	2	3	8	5.22	0.68	2.17	64.70	23.46
8	A	Reading—Key Ideas and Details	5	3	11	5.66	0.54	2.48	51.67	18.74
	B	Reading—Craft & Structure/Integration of Knowledge & Ideas	6	1	8	4.48	0.55	1.96	56.19	20.94
	C	Reading—Vocabulary Use	3	1	5	3.49	0.72	1.49	69.19	25.54
	D	Writing/Language—Text Types and Purposes	3	2	12	6.25	0.55	2.69	52.21	18.98
	E	Writing/Language—Research	2	2	6	3.33	0.56	1.73	55.67	22.73
	F	Writing/Language—Language Conventions	4	1	6	3.18	0.53	1.39	53.27	15.89
	G	Listening	4	2	8	4.26	0.50	1.91	53.48	19.68

Table 10-15 Summary Statistics for Domain Raw and SPI Scores, English Language Arts

Grade	Domain	No. of Items		Total Score Points	Mean	Mean <i>p</i> -Value	SD	SPI	
		MC	CR					Mean	SD
3	Listening	5	1	7	3.86	0.56	1.76	54.97	19.29
	Reading	14	4	22	12.12	0.55	5.08	55.16	22.18
	Writing	9	5	24	10.84	0.54	3.97	45.29	15.28
4	Listening	4	2	8	3.76	0.52	1.93	47.46	18.18
	Reading	14	5	24	13.02	0.56	4.99	54.26	19.85
	Writing	10	4	24	12.20	0.57	4.52	50.92	17.79
5	Listening	4	2	8	4.97	0.63	2.13	61.44	21.23
	Reading	18	3	24	12.37	0.54	5.16	51.66	20.48
	Writing	8	5	24	11.42	0.51	3.92	47.76	15.14
6	Listening	4	2	8	4.10	0.54	1.98	51.63	18.83
	Reading	14	5	24	14.70	0.62	5.48	61.17	22.23
	Writing	7	6	24	10.90	0.49	4.10	45.65	15.51
7	Listening	2	3	8	5.22	0.68	2.17	64.70	23.46
	Reading	13	6	24	12.75	0.54	4.78	53.22	18.96
	Writing	7	6	24	12.04	0.54	4.36	50.39	16.67
8	Listening	4	2	8	4.26	0.50	1.91	53.48	19.68
	Reading	14	5	24	13.62	0.58	5.13	56.83	20.39
	Writing	9	5	24	12.76	0.54	4.78	53.29	18.51

Table 10-16 Summary Statistics for Content Standards Raw and SPI Scores, Mathematics

Grade	Content Standard	Standard	No. of Items		Total Score Points	Mean	Mean p -Value	SD	SPI	
			MC	CR					Mean	SD
3	A	Operations and Algebraic Thinking	4	5	9	5.11	0.57	2.50	56.72	25.58
	B	Number and Operations in Base Ten	3	5	8	3.92	0.49	2.32	49.23	26.19
	C	Number and Operations—Fractions	6	2	8	4.42	0.55	2.16	55.20	23.26
	D	Measurement and Data	5	5	10	5.02	0.50	2.76	50.32	25.25
	E	Geometry	5	2	7	3.97	0.57	1.87	56.58	22.38
4	A	Operations and Algebraic Thinking	6	4	10	5.18	0.52	2.66	51.94	24.36
	B	Number and Operations in Base Ten	5	4	9	4.50	0.50	2.37	50.08	23.64
	C	Number and Operations—Fractions	7	3	10	4.98	0.50	2.71	49.98	24.34
	D	Measurement and Data	6	4	10	4.58	0.46	2.48	46.02	21.90
	E	Geometry	4	3	7	3.76	0.54	1.80	53.42	19.98
5	A	Operations and Algebraic Thinking	4	5	9	4.41	0.49	2.51	48.98	25.26
	B	Number and Operations in Base Ten	3	6	9	4.38	0.49	2.58	48.57	26.20
	C	Number and Operations—Fractions	5	4	9	3.85	0.43	2.21	42.99	20.89
	D	Measurement and Data	6	4	10	4.58	0.46	2.62	45.82	23.30
	E	Geometry	5	4	9	4.17	0.47	2.19	46.50	20.51
6	E	Geometry	3	4	7	2.76	0.40	2.00	39.80	24.92
	F	Ratios and Proportional Relationships	3	4	7	3.74	0.54	2.00	53.20	25.43
	G	The Number System	6	5	11	5.59	0.51	2.87	50.76	23.76
	H	Expressions and Equations	6	5	11	4.90	0.45	2.90	44.77	24.18
	I	Statistics and Probability	8	2	10	5.03	0.50	2.26	50.36	18.58
7	E	Geometry	7	3	10	3.96	0.40	2.15	39.96	17.29
	F	Ratios and Proportional Relationships	4	4	8	4.54	0.57	2.37	56.48	27.31
	G	The Number System	3	4	7	3.03	0.43	1.84	43.38	22.50
	H	Expressions and Equations	5	5	10	4.65	0.47	2.55	46.60	22.77
	I	Statistics and Probability	7	4	11	4.63	0.42	2.60	42.29	20.59
8	E	Geometry	5	5	10	3.94	0.40	2.41	39.71	20.60
	G	The Number System	5	3	8	3.17	0.40	2.21	39.92	23.52
	H	Expressions and Equations	7	3	10	4.55	0.46	2.49	45.73	21.98
	I	Statistics and Probability	6	2	8	4.32	0.54	2.20	53.60	24.43
	J	Functions	6	4	10	4.36	0.44	2.59	43.73	23.14

Table 10-17 Summary Statistics for Content Standards Raw and SPI Scores, Science

Grade	Content Standard	Standard	No. of Items		Total Score Points	Mean	Mean <i>p</i> -Value	SD	SPI	
			MC	CR					Mean	SD
4	A	Life Science	3	9	12	7.18	0.60	2.81	59.63	21.40
	B	Physical Science	4	7	11	5.97	0.54	2.56	54.28	20.94
	C	Earth and Space Science	3	6	9	4.13	0.46	2.23	46.27	21.54
	D	Engineering	3	5	8	3.99	0.50	2.35	50.09	26.37
8	A	Life Science	5	6	11	5.86	0.53	2.64	53.35	21.56
	B	Physical Science	3	9	12	6.26	0.52	2.82	52.17	21.17
	C	Earth and Space Science	1	8	9	3.92	0.44	2.17	43.81	21.06
	D	Engineering	3	5	8	4.57	0.57	2.08	57.11	22.66

Table 10-18 Summary Statistics for Content Standards Raw and SPI Scores, Social Studies

Grade	Content Standard	Standard	No. of Items		Total Score Points	Mean	Mean <i>p</i> -Value	SD	SPI	
			MC	CR					Mean	SD
4	A	Geography	8	3	11	6.89	0.63	2.83	62.46	23.80
	B	History	7	2	9	4.83	0.54	2.27	53.90	22.03
	C	Political Science	2	4	6	2.79	0.47	1.54	46.81	19.48
	D	Economics	5	2	7	3.49	0.50	1.86	50.27	22.34
	E	Behavioral Sciences	3	4	7	3.91	0.56	1.99	55.81	25.03
8	A	Geography	5	4	9	5.42	0.60	2.18	60.18	21.15
	B	History	8	2	10	5.59	0.56	2.42	56.10	21.18
	C	Political Science	6	1	7	3.92	0.56	2.05	56.15	26.21
	D	Economics	5	2	7	3.88	0.56	1.95	55.71	23.93
	E	Behavioral Sciences	5	2	7	4.36	0.63	1.85	62.11	22.71
10	A	Geography	7	3	10	5.57	0.56	2.61	55.74	23.54
	B	History	7	1	8	4.39	0.55	2.14	54.68	23.58
	C	Political Science	6	2	8	4.05	0.51	2.10	50.80	22.54
	D	Economics	6	1	7	3.31	0.48	1.75	47.52	20.18
	E	Behavioral Sciences	4	3	7	3.22	0.47	1.89	46.39	22.80

Table 10-19 SPI Cut Scores, English Language Arts

Content Standard/Domain	Performance Level	Grade 3		Grade 4		Grade 5	
		Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound
Reading—Key Ideas and Details	1	0	37	0	35	0	31
	2	38	68	36	54	32	57
	3	69	90	55	74	58	84
	4	91	100	75	100	85	100
Reading—Craft & Structure	1	0	31	0	34	0	31
	2	32	52	35	54	32	51
	3	53	81	55	79	52	74
	4	82	100	80	100	75	100
Reading—Vocabulary Use	1	0	47	0	41	0	57
	2	48	73	42	66	58	81
	3	74	94	67	88	82	93
	4	95	100	89	100	94	100
Writing/Language—Text Types and Purposes	1	0	29	0	30	0	33
	2	30	45	31	45	34	43
	3	46	60	46	66	44	62
	4	61	100	67	100	63	100
Writing/Language—Research	1	0	39	0	38	0	40
	2	40	59	39	59	41	61
	3	60	79	60	83	62	82
	4	80	100	84	100	83	100
Writing/Language—Language Conventions	1	0	37	0	46	0	42
	2	38	51	47	62	43	57
	3	52	67	63	79	58	76
	4	68	100	80	100	77	100
Listening	1	0	40	0	32	0	48
	2	41	62	33	48	49	70
	3	63	81	49	69	71	89
	4	82	100	70	100	90	100
Reading	1	0	37	0	36	0	35
	2	38	64	37	57	36	59
	3	65	88	58	79	60	81
	4	89	100	80	100	82	100
Writing	1	0	33	0	36	0	37
	2	34	50	37	53	38	51
	3	51	67	54	73	52	70
	4	68	100	74	100	71	100

Content Standard/Domain	Performance Level	Grade 6		Grade 7		Grade 8	
		Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound
Reading—Key Ideas and Details	1	0	43	0	29	0	37
	2	44	67	30	47	38	58
	3	68	83	48	70	59	75
	4	84	100	71	100	76	100
Reading—Craft & Structure	1	0	41	0	37	0	40
	2	42	70	38	58	41	64
	3	71	92	59	78	65	83
	4	93	100	79	100	84	100
Reading—Vocabulary Use	1	0	53	0	53	0	54
	2	54	78	54	83	55	84
	3	79	90	84	95	85	95
	4	91	100	96	100	96	100
Writing/Language—Text Types and Purposes	1	0	29	0	32	0	40
	2	30	42	33	47	41	59
	3	43	58	48	67	60	74
	4	59	100	68	100	75	100
Writing/Language—Research	1	0	33	0	37	0	39
	2	34	53	38	57	40	64
	3	54	68	58	75	65	84
	4	69	100	76	100	85	100
Writing/Language—Language Conventions	1	0	45	0	50	0	43
	2	46	60	51	68	44	57
	3	61	76	69	84	58	72
	4	77	100	85	100	73	100
Listening	1	0	37	0	48	0	40
	2	38	57	49	75	41	59
	3	58	74	76	91	60	79
	4	75	100	92	100	80	100
Reading	1	0	44	0	37	0	42
	2	45	70	38	58	43	65
	3	71	88	59	78	66	82
	4	89	100	79	100	83	100
Writing	1	0	34	0	37	0	40
	2	35	49	38	54	41	60
	3	50	65	55	72	61	76
	4	66	100	73	100	77	100

Table 10-20 SPI Cut Scores, Mathematics

Content Standard/Domain	Performance Level	Grade 3		Grade 4		Grade 5	
		Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound
Operations and Algebraic Thinking	1	0	33	0	27	0	27
	2	34	60	28	54	28	52
	3	61	86	55	81	53	80
	4	87	100	82	100	81	100
Number and Operations in Base Ten	1	0	24	0	28	0	26
	2	25	50	29	51	27	51
	3	51	82	52	78	52	81
	4	83	100	79	100	82	100
Number and Operations—Fractions	1	0	34	0	26	0	25
	2	35	54	27	50	26	42
	3	55	84	51	80	43	69
	4	85	100	81	100	70	100
Measurement and Data	1	0	26	0	25	0	26
	2	27	52	26	44	27	45
	3	53	80	45	73	46	76
	4	81	100	74	100	77	100
Geometry	1	0	37	0	36	0	29
	2	38	59	37	56	30	46
	3	60	82	57	74	47	71
	4	83	100	75	100	72	100

Content Standard/Domain	Performance Level	Grade 6		Grade 7		Grade 8	
		Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound
Geometry	1	0	19	0	30	0	27
	2	20	39	31	43	28	46
	3	40	83	44	73	47	75
	4	84	100	74	100	76	100
Ratios and Proportional Relationships*	1	0	35	0	45		
	2	36	60	46	73		
	3	61	90	74	93		
	4	91	100	94	100		
The Number System	1	0	32	0	31	0	23
	2	33	57	32	50	24	51
	3	58	84	51	84	52	78
	4	85	100	85	100	79	100
Expressions and Equations	1	0	25	0	32	0	32
	2	26	48	33	55	33	55
	3	49	84	56	86	56	82
	4	85	100	87	100	83	100
Statistics and Probability	1	0	37	0	30	0	41
	2	38	53	31	49	42	68
	3	54	78	50	79	69	90
	4	79	100	80	100	91	100
Functions**	1					0	28
	2					29	52
	3					53	85
	4					86	100

* Content standard in grades 6 and 7 only.

** Content standard in grade 8 only.

Table 10-21 SPI Cut Scores, Science

Content Standard/Domain	Performance Level	Grade 4		Grade 8	
		Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound
Life Science	1	0	37	0	33
	2	38	61	34	52
	3	62	79	53	74
	4	80	100	75	100
Physical Science	1	0	32	0	33
	2	33	54	34	52
	3	55	74	53	72
	4	75	100	73	100
Earth and Space Science	1	0	23	0	24
	2	24	43	25	43
	3	44	67	44	63
	4	68	100	64	100
Engineering	1	0	20	0	36
	2	21	48	37	58
	3	49	77	59	79
	4	78	100	80	100

Table 10-22 SPI Cut Scores, Social Studies

Content Standard/Domain	Performance Level	Grade 4		Grade 8		Grade 10	
		Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound	Score Lower Bound	Score Upper Bound
Geography	1	0	39	0	39	0	36
	2	40	59	40	54	37	58
	3	60	83	55	77	59	77
	4	84	100	78	100	78	100
History	1	0	31	0	34	0	36
	2	32	47	35	50	37	57
	3	48	72	51	73	58	75
	4	73	100	74	100	76	100
Political Science	1	0	29	0	28	0	32
	2	30	40	29	49	33	49
	3	41	61	50	78	50	70
	4	62	100	79	100	71	100
Economics	1	0	29	0	31	0	31
	2	30	41	32	47	32	46
	3	42	68	48	75	47	63
	4	69	100	76	100	64	100
Behavioral Sciences	1	0	31	0	39	0	27
	2	32	50	40	59	28	46
	3	51	77	60	81	47	65
	4	78	100	82	100	66	100

Table 10-23 Longitudinal Comparison of State-Level Participation Rates and Scale Score Means, English Language Arts

Grade	Year	Enrolled	Number Tested	Percent Tested	Scale Score Mean	Scale Score SD
3	2016	65793	64107	97.44	560.57	47.31
	2017	65340	63946	97.87	559.12	46.93
	2018	64693	63194	97.68	556.70	46.66
	2019	62646	61091	97.52	554.59	45.54
	2021	60785	52930	87.08	550.13	46.61
	2022	60956	58275	95.60	549.75	47.41
	2023	61228	58497	95.54	552.33	48.85
4	2016	64361	62609	97.28	582.71	49.41
	2017	66001	64423	97.61	585.26	52.44
	2018	65885	64354	97.68	580.90	51.81
	2019	65222	63528	97.40	582.01	51.05
	2021	61127	52706	86.22	577.65	51.31
	2022	61465	58773	95.62	578.32	52.20
	2023	61689	58996	95.63	583.88	53.27
5	2016	64045	62300	97.28	599.62	51.11
	2017	64624	62995	97.48	603.24	51.00
	2018	66542	64903	97.54	600.78	48.35
	2019	66250	64654	97.59	595.58	48.77
	2021	62405	54010	86.55	593.12	49.01
	2022	61735	59044	95.64	593.86	52.07
	2023	62058	59386	95.69	592.89	51.66
6	2016	64594	62728	97.11	610.36	52.16
	2017	64446	62754	97.37	614.59	49.82
	2018	65363	63600	97.30	609.61	50.18
	2019	67243	65386	97.24	607.00	50.15
	2021	64925	55511	85.50	604.02	50.12
	2022	63048	60112	95.34	602.26	51.02
	2023	62310	59412	95.35	604.39	53.00
7	2016	64044	62084	96.94	623.84	54.85
	2017	65049	63091	96.99	626.80	59.14
	2018	64975	63140	97.18	627.43	56.56
	2019	65904	63878	96.93	627.70	54.88
	2021	66361	56295	84.83	625.36	55.21
	2022	65386	61871	94.62	618.80	56.30
	2023	63623	60413	94.95	622.59	55.27
8	2016	63861	61486	96.28	637.23	57.27
	2017	64265	62109	96.65	637.69	61.61
	2018	65638	63248	96.36	630.98	59.94
	2019	65355	63056	96.48	629.06	59.84
	2021	67572	56756	83.99	628.22	58.76
	2022	66789	62684	93.85	622.78	61.00
	2023	66060	62249	94.23	628.49	64.03

Note: Caution should be exercised when interpreting the Spring 2021 statewide data due to participation rates lower than in a typical administration year, resulting in overrepresentation of some subgroups and underrepresentation of other subgroups.

Table 10-24 Longitudinal Comparison of State-Level Participation Rates and Scale Score Means, Mathematics

Grade	Year	Enrolled	Number Tested	Percent Tested	Scale Score Mean	Scale Score SD
3	2016	65793	64194	97.57	554.28	46.47
	2017	65340	64066	98.05	555.03	48.63
	2018	64693	63314	97.87	555.94	50.87
	2019	62646	61210	97.71	555.78	53.50
	2021	60785	52892	87.01	548.97	56.39
	2022	60956	58449	95.89	552.25	56.04
	2023	61228	58722	95.91	553.02	56.43
4	2016	64361	62674	97.38	573.45	56.15
	2017	66001	64533	97.78	574.33	54.92
	2018	65885	64462	97.84	576.76	52.99
	2019	65222	63630	97.56	577.09	51.78
	2021	61127	52658	86.15	571.50	53.33
	2022	61465	58931	95.88	573.54	57.00
	2023	61689	59165	95.91	576.62	55.40
5	2016	64045	62368	97.38	599.57	50.19
	2017	64624	63152	97.72	599.73	51.00
	2018	66542	65021	97.71	598.82	56.65
	2019	66250	64728	97.70	601.48	53.14
	2021	62405	53932	86.42	594.26	56.01
	2022	61735	59187	95.87	599.38	52.67
	2023	62058	59577	96.00	601.51	52.62
6	2016	64594	62772	97.18	612.67	53.00
	2017	64446	62847	97.52	612.93	54.81
	2018	65363	63669	97.41	611.97	57.64
	2019	67243	65470	97.36	610.77	58.31
	2021	64925	55462	85.42	602.06	57.74
	2022	63048	60234	95.54	604.64	59.36
	2023	62310	59570	95.60	610.58	57.27
7	2016	64044	62144	97.03	627.49	57.40
	2017	65049	63200	97.16	627.48	58.65
	2018	64975	63218	97.30	622.82	65.55
	2019	65904	63973	97.07	625.25	60.69
	2021	66361	56247	84.76	620.00	60.06
	2022	65386	61969	94.77	617.55	62.13
	2023	63623	60559	95.18	620.62	60.98
8	2016	63861	61551	96.38	640.79	57.54
	2017	64265	62175	96.75	641.11	59.36
	2018	65638	63318	96.47	644.24	60.78
	2019	65355	63108	96.56	644.53	57.85
	2021	67572	56726	83.95	638.33	56.94
	2022	66789	62762	93.97	634.95	60.04
	2023	66060	62360	94.40	638.57	56.40

Note: Caution should be exercised when interpreting the Spring 2021 statewide data due to participation lower than in a typical administration year, resulting in overrepresentation of some subgroups and underrepresentation of other subgroups.

Table 10-25 Longitudinal Comparison of State-Level Participation Rates and Scale Score Means, Science

Grade	Year	Enrolled	Number Tested	Percent Tested	Scale Score Mean	Scale Score SD
4	2019	65222	63611	97.53	499.88	50.24
	2021	61127	52417	85.75	497.39	50.30
	2022	61465	58880	95.79	495.81	52.11
	2023	61689	59141	95.87	496.06	54.65
8	2019	65355	63062	96.49	699.70	50.55
	2021	67572	56485	83.59	697.31	49.86
	2022	66789	62644	93.79	693.86	51.11
	2023	66060	62289	94.29	692.80	51.93

Note 1: New reporting scales were established for Science in Spring 2019.

Note 2: Caution should be exercised when interpreting the Spring 2021 statewide data due to participation rates lower than in a typical administration year, resulting in overrepresentation of some subgroups and underrepresentation of other subgroups.

Table 10-26 Longitudinal Comparison of State-Level Participation Rates and Scale Score Means, Social Studies

Grade	Year	Enrolled	Number Tested	Percent Tested	Scale Score Mean	Scale Score SD
4	2022	61465	58833	95.72	499.85	50.27
	2023	61689	59131	95.85	499.74	50.82
8	2022	66789	62606	93.74	699.24	50.80
	2023	66060	62261	94.25	700.49	48.94
10	2022	67427	59391	88.08	798.49	50.78
	2023	69876	61819	88.47	796.91	53.02

Note 1: New reporting scales were established for Social Studies in Spring 2022.

Table 10-27 Longitudinal Comparison of State-Level Impact Data, English Language Arts

Grade	Year	N	Below Basic	Basic	Proficient	Advanced	Prof. & Adv.
3	2016	64107	21.99	34.88	34.29	8.84	43.13
	2017	63946	21.45	36.72	33.81	8.02	41.83
	2018	63194	22.78	37.47	32.58	7.17	39.75
	2019	61091	23.28	38.04	33.21	5.48	38.69
	2021	52930	27.71	37.74	29.48	5.08	34.56
	2022	58275	27.55	37.52	29.61	5.32	34.93
	2023	58497	25.93	36.89	30.89	6.30	37.18
4	2016	62609	22.81	33.88	34.77	8.54	43.30
	2017	64423	21.14	32.14	37.00	9.71	46.72
	2018	64354	24.04	32.06	35.72	8.19	43.91
	2019	63528	23.88	33.14	34.10	8.89	42.98
	2021	52706	27.17	32.71	32.67	7.45	40.12
	2022	58773	26.66	31.74	33.63	7.96	41.60
	2023	58996	23.21	31.78	34.74	10.27	45.01
5	2016	62300	23.17	34.37	34.55	7.91	42.47
	2017	62995	20.36	33.22	37.88	8.54	46.42
	2018	64903	21.53	34.30	37.40	6.77	44.17
	2019	64654	26.11	33.83	34.34	5.72	40.06
	2021	54010	27.67	34.80	32.41	5.12	37.52
	2022	59044	27.99	31.73	34.21	6.06	40.28
	2023	59386	27.81	33.75	32.35	6.09	38.44
6	2016	62728	21.12	36.30	31.67	10.91	42.58
	2017	62754	18.23	36.52	33.51	11.75	45.26
	2018	63600	22.06	35.08	32.73	10.12	42.86
	2019	65386	23.56	35.48	31.87	9.09	40.96
	2021	55511	25.18	36.36	30.65	7.80	38.45
	2022	60112	26.19	35.83	30.71	7.27	37.98
	2023	59412	25.57	34.38	31.11	8.94	40.05
7	2016	62084	23.11	34.91	34.09	7.89	41.98
	2017	63091	22.27	34.10	33.52	10.11	43.63
	2018	63140	21.29	33.57	35.72	9.43	45.15
	2019	63878	21.88	33.25	35.36	9.51	44.87
	2021	56295	23.08	33.99	34.12	8.81	42.92
	2022	61871	26.71	35.03	30.70	7.55	38.26
	2023	60413	24.40	34.70	32.72	8.18	40.90
8	2016	61486	21.24	37.21	31.26	10.30	41.56
	2017	62109	21.66	37.22	29.19	11.93	41.12
	2018	63248	24.66	38.01	27.93	9.40	37.33
	2019	63056	25.94	37.04	28.80	8.23	37.03
	2021	56756	26.05	38.29	28.02	7.64	35.66
	2022	62684	29.81	36.72	26.45	7.02	33.46
	2023	62249	26.84	35.54	28.24	9.38	37.62

Table 10-28 Longitudinal Comparison of State-Level Impact Data, Mathematics

Grade	Year	N	Below Basic	Basic	Proficient	Advanced	Prof. & Adv.
3	2016	64194	18.59	33.41	38.90	9.10	48.00
	2017	64066	18.90	33.06	37.84	10.20	48.03
	2018	63314	18.68	31.48	38.47	11.37	49.83
	2019	61210	19.28	31.28	37.17	12.27	49.44
	2021	52892	23.78	31.23	34.46	10.54	44.99
	2022	58449	22.76	30.04	35.02	12.18	47.20
	2023	58722	23.17	29.02	34.84	12.97	47.81
4	2016	62674	19.59	36.22	33.33	10.86	44.20
	2017	64533	19.13	37.37	32.67	10.83	43.50
	2018	64462	18.37	37.17	32.71	11.74	44.46
	2019	63630	18.87	36.09	32.82	12.23	45.05
	2021	52658	22.32	36.61	30.84	10.23	41.07
	2022	58931	22.24	34.03	30.95	12.79	43.73
	2023	59165	21.00	33.70	30.96	14.34	45.30
5	2016	62368	25.94	29.98	34.14	9.94	44.08
	2017	63152	24.97	30.57	34.58	9.88	44.46
	2018	65021	24.73	29.32	35.05	10.90	45.95
	2019	64728	24.22	29.20	35.09	11.49	46.58
	2021	53932	29.46	28.94	32.12	9.47	41.59
	2022	59187	26.82	28.38	34.28	10.52	44.80
	2023	59577	25.53	27.30	35.66	11.52	47.18
6	2016	62772	25.51	31.66	36.78	6.05	42.84
	2017	62847	24.70	31.68	37.50	6.11	43.61
	2018	63669	24.78	31.27	37.78	6.18	43.96
	2019	65470	26.72	30.79	35.80	6.69	42.49
	2021	55462	32.61	31.82	30.90	4.66	35.57
	2022	60234	30.96	30.28	33.27	5.49	38.76
	2023	59570	27.81	30.43	34.81	6.94	41.75
7	2016	62144	30.45	30.28	34.81	4.45	39.26
	2017	63200	30.80	29.92	34.53	4.75	39.29
	2018	63218	31.36	29.67	34.33	4.64	38.97
	2019	63973	32.18	28.99	34.05	4.78	38.83
	2021	56247	34.99	30.17	31.38	3.47	34.84
	2022	61969	38.11	28.16	29.58	4.15	33.73
	2023	60559	35.62	29.20	30.66	4.51	35.18
8	2016	61551	28.66	37.48	28.12	5.74	33.86
	2017	62175	28.43	36.95	28.33	6.29	34.62
	2018	63318	27.95	35.44	28.71	7.90	36.61
	2019	63108	28.55	35.60	27.83	8.01	35.85
	2021	56726	32.47	37.53	23.62	6.38	30.00
	2022	62762	35.96	33.73	24.05	6.26	30.31
	2023	62360	34.69	34.12	24.83	6.36	31.19

Table 10-29 Longitudinal Comparison of State-Level Impact Data, Science

Grade	Year	N	Below Basic	Basic	Proficient	Advanced	Prof. & Adv.
4	2019	63611	14.98	32.25	33.29	19.49	52.78
	2021	52417	16.18	32.67	33.03	18.13	51.16
	2022	58880	18.20	30.95	32.32	18.53	50.85
	2023	59141	18.82	29.93	31.85	19.40	51.25
8	2019	63062	17.76	28.29	31.50	22.45	53.95
	2021	56485	18.56	29.97	30.65	20.82	51.47
	2022	62644	20.77	30.33	28.99	19.91	48.90
	2023	62289	21.89	29.25	29.13	19.74	48.86

Note: New cut scores were used to classify students into performance levels after the Spring 2019 Science test administration.

Table 10-30 Longitudinal Comparison of State-Level Impact Data, Social Studies

Grade	Year	N	Below Basic	Basic	Proficient	Advanced	Prof. & Adv.
4	2022	58833	21.02	20.09	36.79	22.10	58.89
	2023	59131	21.01	19.87	36.14	22.98	59.12
8	2022	62606	19.53	21.74	35.08	23.64	58.73
	2023	62261	20.03	20.91	34.36	24.70	59.06
10	2022	59391	25.49	26.84	25.73	21.94	47.67
	2023	61819	27.39	25.36	24.72	22.53	47.24

Note: New cut scores were used to classify students into performance levels after the Spring 2022 Social Studies test administration.

Part 11: Summary and Recommendations

Results and key findings of the Spring 2023 Wisconsin Forward Exam administration are presented throughout the body of this report. This last section of the report presents some recommendations for DPI consideration.

The 2023 test administration was the seventh administration of the Forward assessments. Since Spring 2016, the assessment results have been reported on the same scales and students have been classified into performance levels using the same cut scores, allowing for longitudinal tracking of student performance in ELA and Mathematics. New test scales and new performance level cut scores were set for the Science assessments after the Spring 2019 test administration. The 2023 Wisconsin Forward Exam administration was the fourth administration of the Science assessments that measured the new Science standards and were reported on the new scales. For Social Studies, new test scales and new performance level cut scores were set after the Spring 2022 test administration. The 2023 Wisconsin Forward Exam administration was the second administration of the Social Studies assessments that measured the new Social Studies standards.

In keeping with the field-testing of new test items in the past administrations, DRC recommends that, in the future, all new items continue to be field-tested in Wisconsin prior to their operational test administration to provide accurate information on how students may perform on these items once they are administered operationally. DRC also recommends continuing to embed field test items in each operational test administration for all content areas in order to build a high-quality Wisconsin item bank for future form development.

DRC recommends continuing to use an artificial intelligence (AI) engine in the scoring of the short- or long-response writing items for efficiency and accuracy. As indicated in Part 5 and Part 8 of this report, the AI scores were in good or very good agreement with scores by trained human scorers.

In Mathematics grades 4 through 8, substantially fewer students received the lowest obtainable scale score (LOSS) in Spring 2023 compared to Spring 2022. While the Mathematics assessments continue to be difficult for some students, including more and easier non-multiple-choice items in the Spring 2023 assessments contributed to lowering the percentage of students at the LOSS to less than 2% in each grade. As explained in Part 6 of this report, for students to receive a scale score above the LOSS, they need to correctly answer more items, including some non-MC items. DRC continues to recommend that some easier non-MC items be included in the future forms of Mathematics tests. In addition, approximately 3% of students scored at LOSS in Social Studies grade 10. Students who scored at LOSS in Social Studies grade 10 did not answer any of the non-MC items correctly. It is recommended that more non-MC items, including easier items, be included in the Social Studies grade 10 assessment in the future.

From the psychometric perspective and per the peer review recommendation for Science, the magnitude of the conditional standard error, particularly at the lower and upper end of the scale, will continue to be monitored. Efforts will be made to include more items of higher discrimination that span the entire range of the difficulty scale with a goal to improve the precision of measurement, particularly for students of lower and higher ability.

References

- Allen, M. J., & Yen, W. M. (1979). *Introduction to measurement theory*. Monterey, CA: Brooks/Cole.
- Altman, D. G. (1991). *Practical statistics for medical research*. Chapman and Hall.
- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for educational and psychological testing*. Washington, DC: American Educational Research Association.
- Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: An application of an EM algorithm. *Psychometrika*, 46, 443–459.
- Burket, G. R. (2002). PARDUX, Version 1.54 [Computer program]. CTB/McGraw-Hill.
- Camilli, G., & Shepard, A. L. (1994). *Methods for identifying biased test items*. Thousand Oaks, CA: Sage Publications.
- Candell, G. L., & Drasgow, F. (1988). An iterative procedure for linking metrics and assessing bias in item response theory. *Applied Psychological Measurement*, 12(3), 253–260.
- Cizek, G. J., & Bunch, M. B. (2007). *Standard setting: A guide to establishing and evaluating performance standards on tests*. Thousand Oaks, CA: Sage Publications.
- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Lawrence Erlbaum Associates.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20, 37–46.
- Crocker, L., & Algina, J. (1986). *Introduction to classical and modern test theory*. Belmont, CA: Wadsworth Group/Thomson Learning.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16(3), 297–334.
- CTB/McGraw-Hill. (1997). *TerraNova* (1st ed.). Monterey, CA: Author.
- CTB/McGraw-Hill. (2000). *TerraNova* (2nd ed.). Monterey, CA: Author.
- CTB/McGraw-Hill. (2009). *TerraNova 3rd Edition Technical Addendum: Forms E and F*. Monterey, CA: Author.
- Dorans, N. J., & Schmitt, A. P. (1991). *Constructed response and differential item functioning: A pragmatic approach*. Princeton, NJ: Educational Testing Service.

- Fitzpatrick, A. R. (1991). *Status report on the results of preliminary analyses of dichotomous and multi-level items using the PARMATE program*. Monterey, CA: CTB/McGraw-Hill.
- Fitzpatrick, A. R., & Julian, M. W. (1996). *Two studies comparing the parameter estimates produced by PARDUX and PARSCALE*. Unpublished manuscript.
- Green, D. R. (December 1975). *Procedures for assessing bias in achievement tests*. Paper presented at the National Institute of Education Conference on Test Bias, Annapolis, MD.
- Hambleton, R. K., & Novick, M. R. (1973). Toward an integration of theory and method for criterion-referenced tests. *Journal of Educational Measurement*, 10(3), 159–170.
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Hingham, MA: Kluwer-Nijhoff Publishing.
- Holland, P. W., & Thayer, D. T. (1985). *An alternate definition of the ETS delta scale of item difficulty* (Research Report No. RR-85-43). Princeton, NJ: Educational Testing Service.
- Holland, P. W., & Thayer, D. T. (April 1986). *Differential item performance and the Mantel-Haenszel procedure*. Paper presented at the annual meeting of the American Educational Research Association Annual Meeting, San Francisco, CA.
- Karantonis, A., & Sireci, S. G. (2006). The Bookmark standard-setting method: A literature review. *Educational Measurement: Issues and Practice*, 25(1), 4–12.
- Kim, D. (2005). KKCLASS [Computer program]. Unpublished.
- Kim, D., Barton, K., & Kim, J. (April 2007). *Estimating classification consistency and classification accuracy with pattern scoring*. Paper presented at the annual meeting of the American Educational Research Association, Chicago, IL.
- Kim, D., Choi, S., Um, K., & Kim, J. (April 2006). *A comparison of methods for estimating classification consistency*. Paper presented at the annual meeting of the National Council on Measurement in Education, San Francisco, CA.
- Kolen, M. J., & Brennan, R. L. (1995). *Test equating: Methods and practices*. New York, NY: Springer-Verlag.
- Kolen, M., & Kim, D. (2004). [Personal correspondence].
- Landis, J. R., & Koch, G. G. (1977). The measurement of observer agreement for categorical data. *Biometrics*, 33(1), 159–174.
- Linn, R. L. (Ed.) (1989). *Educational measurement* (3rd ed.). New York, NY: Macmillan.

- 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*. Reading, MA: Addison-Wesley.
- 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.
- McCallin, R. C. (2006). Test administration. In S. M. Downing & T. M. Haladyna (Eds.), *Handbook of test development* (pp. 625–652). Mahwah, NJ: Lawrence Erlbaum Associates.
- Michaelides, M. P., & Haertel, E. H. (2004). *Sampling of common items: An unrecognized source of error in test equating*. Los Angeles, CA: Center for the Study of Evaluation.
- Muraki, E. (1990). Fitting a polytomous item response model to Likert-type data. *Applied Psychological Measurement*, 14(1), 59–71.
- Muraki, E. (1992). A generalized partial credit model: Application of an EM algorithm. *Applied Psychological Measurement*, 16(2), 159–176.
- Muraki, E., & Bock, R. D. (1991). PARSCALE: Parameter scaling of rating data [Computer program]. Chicago, IL: Scientific Software, Inc.
- National Research Council (2012). *A Framework for K-12 Science Education: Practices, Crosscutting Concepts, and Core Ideas*. Washington, DC: The National Academies Press.
<https://doi.org/10.17226/13165>.
- Shu, L. (2020). PARDUX, Version 1.69 [Computer program]. Data Recognition Corporation.
- Stocking, M. L., & Lord, F. M. (1983). Developing a common metric in item response theory. *Applied Psychological Measurement*, 7(2), 201–210.
- Swaminathan, H., Hambleton, R. K., & Algina, J. (1974). Reliability of criterion referenced tests: A decision theoretic formulation. *Journal of Educational Measurement*, 11(4), 263–267.
- Swineford, F. (1956). *Technical manual for users of test analysis* (Statistical Report No. 56 42). Princeton, NJ: Educational Testing Service.
- Thissen, D. (1982). Marginal maximum-likelihood estimation for the one-parameter logistic model. *Psychometrika*, 47(2), 175–186.
- Thissen, D. (1990). MULTILOG: Multiple categorical item analysis and test scoring (Version 6) [Computer program]. Chicago, IL: Scientific Software, Inc.

- Thompson, S. J., Johnstone, C. J., & Thurlow, M. L. (2002). *Universal design applied to large scale assessment* (Synthesis Report No. 44). Minneapolis, MN: University of Minnesota, National Center on Educational Outcomes.
- Tong, Y., Wu, S., & Xu, M. (2008). *A comparison of pre-equating and post-equating using large-scale assessment data*. Paper presented at the annual meeting of the American Educational Research Association, New York, NY.
- van der Linden, W. J., & Hambleton, R. K. (Eds.) (1997). *Handbook of modern item response theory*. New York, NY: Springer.
- Wright, B. D., & Linacre, J. M. (1992). BIGSTEPS Rasch analysis [Computer program]. Chicago, IL: MESA Press.
- Yen, W. M. (1981). Using simulation results to choose a latent trait model. *Applied Psychological Measurement*, 5(2), 245–262.
- Yen, W. M. (1984). Obtaining maximum likelihood trait estimates from number-correct scores for the three-parameter logistic model. *Journal of Educational Measurement*, 21(2), 93–111.
- Yen, W. M. (1993). Scaling performance assessments: Strategies for managing local item dependence. *Journal of Educational Measurement*, 30(3), 187–213.
- Yen, W. M., & Burket, G. R. (1997). Comparison of item response theory and Thurstone methods of vertical scaling. *Journal of Educational Measurement*, 34(4), 293–313.
- Yen, W. M., & Candell, G. L. (1991). Increasing score reliability with item-pattern scoring: An empirical study in five score metrics. *Applied Measurement in Education*, 4(3), 209–228.
- Zwick, R., Donoghue, J. R., & Grima, A. (1993). Assessment of differential item functioning for performance tasks. *Journal of Educational Measurement*, 30(3), 233–251.

APPENDIX A
ITEM REVIEW TRAINING SLIDES



Wisconsin Forward Exam Science Virtual Item Review

August 16–18, 2021

Meeting Overview

- Virtual Item Review and using a virtual platform
- Brief overview of the Forward Exam
- Item review process and training
- Break into breakout rooms
- Review items for placement on exam

Roles & Responsibilities

- **Participants**
 - Review Items
- **DRC Facilitators**
 - Lead the group through the agenda
 - Encourage interaction
 - Lead discussions
 - Monitor electronic submission of documents
- **DPI and DRC**
 - Answer questions

Critical Importance of Security and Confidentiality

- All item review participants must complete a security/nondisclosure agreement
- All passage and item content are secure
- Note-taking policy
- Cellphone and personal computer use– phones not allowed to be used during meeting time
- Alexa turned off during meeting
- Communication following the meeting–take the process back with you, not the content

Tips for a Remote Meeting

- Keep your audio on Mute unless you need to speak to the group.
- Stay logged in to Zoom for the entire meeting time, even during breaks (just leave Mute on).
- Do not exit out of Zoom when opening another window in your browser.
- If you encounter technical difficulties, call DRC Customer Service at: **1-833-867-5680**.

ZOOM Meeting Features



Meeting Website

WISCONSIN DEPARTMENT OF PUBLIC INSTRUCTION
DRC

Wisconsin Item Review

ELA
August 16th-19th, 2021

Science
August 16th-18th, 2021

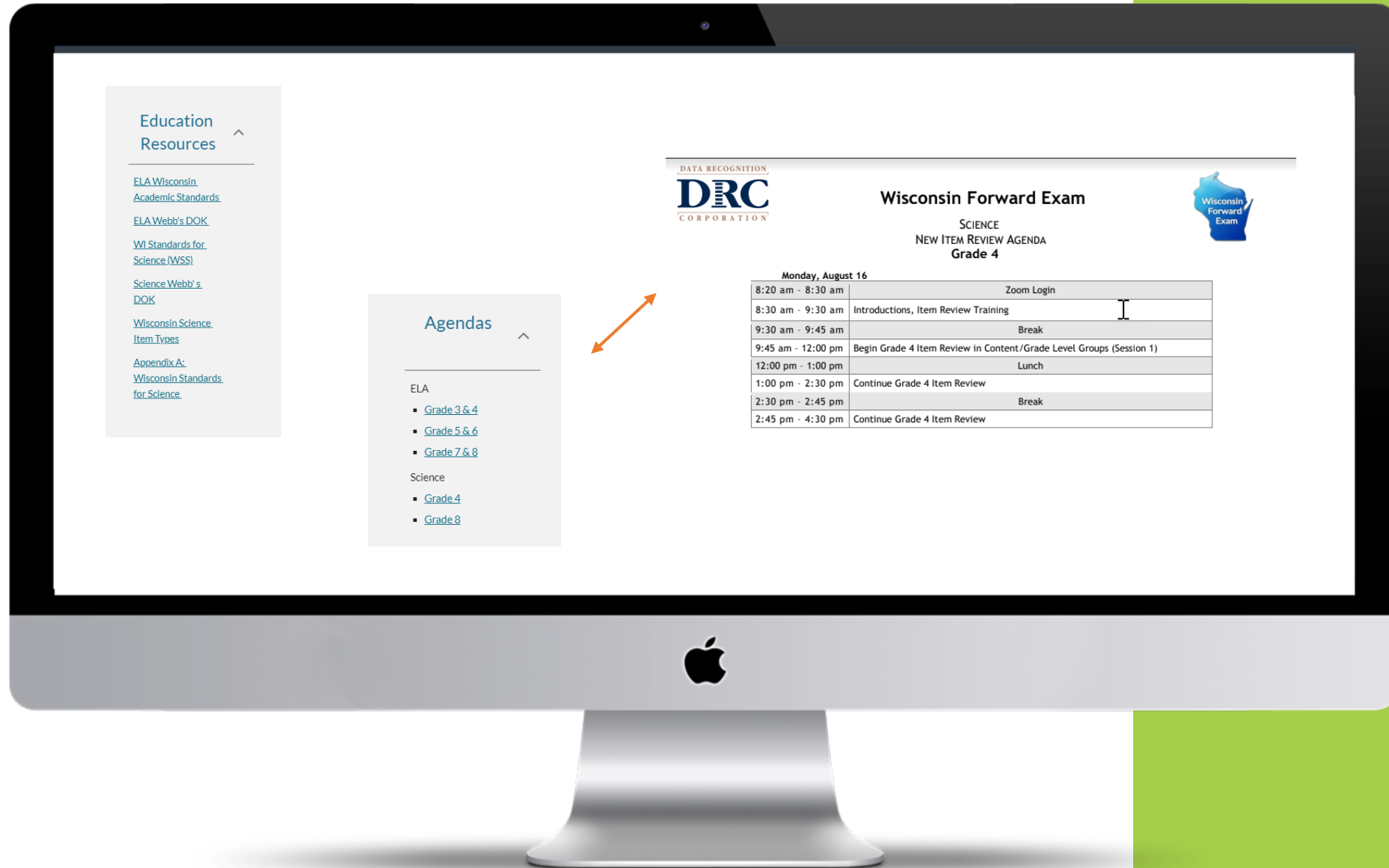
[General Meeting Information](#)

[Agendas](#)

[Test Environment](#)

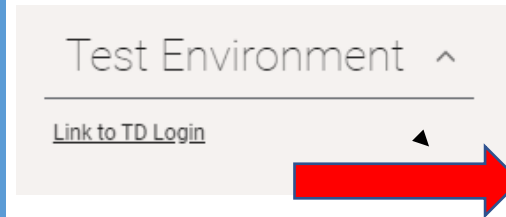
Reminder: This confidential website is for participants only.

Meeting Website

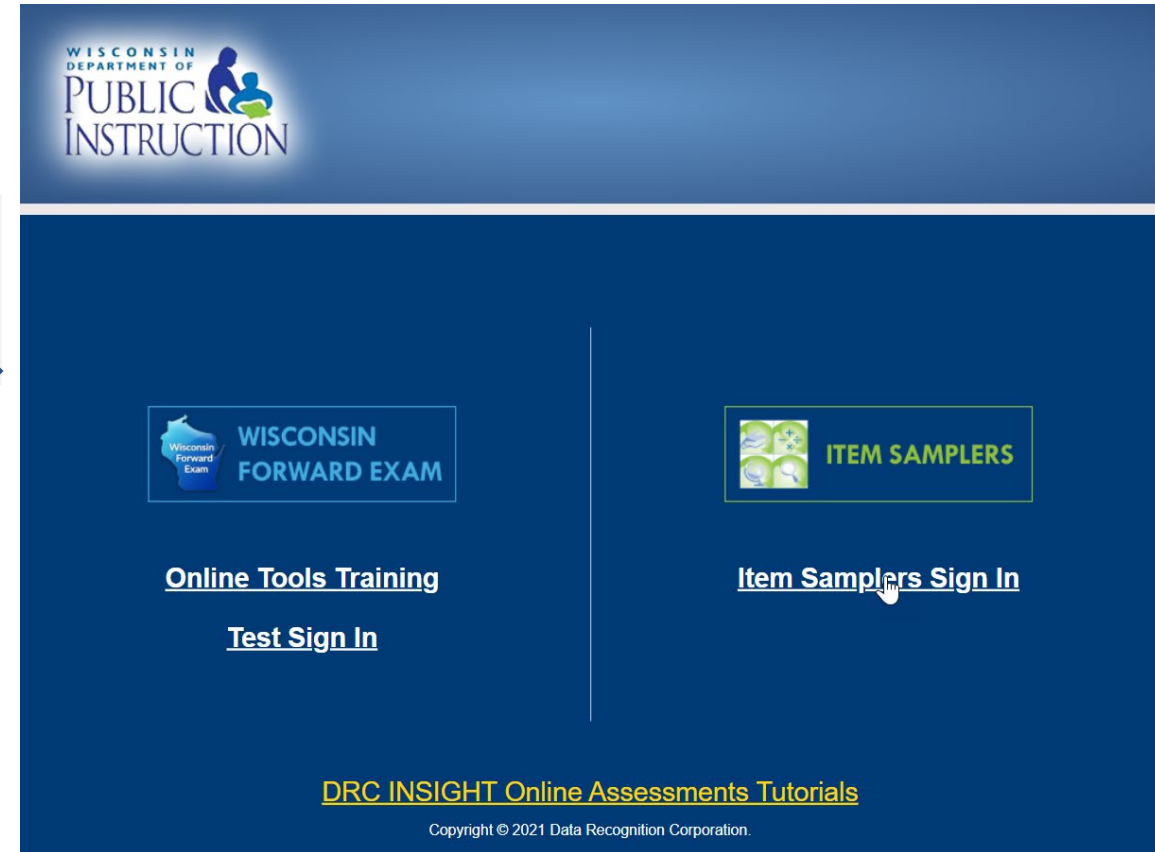


Each link in the Home drop-down menu reveals the subject area resource and agenda links.

Meeting Website – Test Environment link



Click the “Test Environment” link to launch the TD environment. Then select “Test Sign In.”



Remote Meeting Website— Reviewer Feedback Documents

Reviewer
Feedback



Session #	Seq #	Item ID	Scenario Title	WSS Performance Indicator	SEP	CCC	Item Type	Points	Key(s)
1	1	1064623		SCI.ESS2.D.m	Developing Models	Systems and System Models	MC	1	C
1	2	1064620		SCI.LS2.C.m		Cause and Effect	TE	1	decrease; decrease
1	3	1064621		SCI.LS2.D.m	Argue from Evidence		EBSR	1	Part A: A; Part B: C,E
1	4	1064648		SCI.ESS3.C.m	Argue from Evidence	Energy and Matter	TE	1	X to Y; marbles & table to air; kinetic to sound
1	5	1064646		SCI.PS3.D.m	Construct an Explanation	Energy and Matter	TE	1	oxygen; sugars
1	6	1064638		SCI.ESS3.B.m	Analyze and Interpret Data	Patterns	TE	1	oceanic; continental; younger
1	7	1064654	Paper Airplanes	SCI.PS1.A.m	Analyze and Interpret Data		MS	1	R, C; (Danver, Las Vegas)



Reviewer Feedback ^

- [ELA Grade 3](#)
- [ELA Grade 4](#)
- [ELA Grade 5](#)
- [ELA Grade 6](#)
- [ELA Grade 7](#)
- [ELA Grade 8](#)
- [Science Grade 4](#)
- [Science Grade 8](#)

From the drop-down menu, launch Google-based reviewer feedback sheets by clicking the links under the **Reviewer Feedback** heading.

Do not share the materials from this document with anyone.

Do not cut or copy text from one cell and paste it in another. It will throw the hidden formulas off.

Wisconsin's Vision for Science Learning

“[By] the end of 12th grade, all students have some appreciation of the beauty and wonder of science; possess sufficient knowledge of science and engineering to engage in public discussions on related issues; are careful consumers of scientific and technological information related to their everyday lives; are able to continue to learn about science outside school; and have the skills to enter careers of their choice, including (but not limited to) careers in science, engineering, and technology.”

Wisconsin's 3D Standards for Science

1-PS4 Waves and their Applications in Technologies for Information Transfer

1-PS4-1. Plan and conduct investigations to provide evidence that vibrating materials can make sound and that sound can make materials vibrate. (Clarification Statement: Examples of vibrating materials that make sound could include tuning forks and plucking a stretched string. Examples of how sound can make matter vibrate could include holding a piece of paper near a speaker making sound and holding an object near a vibrating tuning fork.)

1-PS4-2. Make observations to construct an evidence-based account that objects can be seen only when illuminated. (Clarification Statement: Examples of observations could include those made in a completely dark room, a white box, and a video of a cave explorer with a flashlight. Illumination could be from an external light source or by an object giving off its own light.)

1-PS4-3. Plan and conduct an investigation to determine the effect of placing objects made with different materials in the path of a beam of light. (Clarification Statement: Examples of materials could include those that are transparent (such as clear plastic), translucent (such as wax paper), opaque (such as cardboard), and reflective (such as a mirror.) [Assessment Boundary: Assessment does not include the speed of light.]

1-PS4-4. Use tools and materials to design and build a device that uses light or sound to solve the problem of communicating over a distance.* (Clarification Statement: Examples of devices could include a light source to send signals, paper cup and string telephone,* and a system of alarm bells.) [Assessment Boundary: Assessment does not include technological design for non-communicative devices such as the microwave oven.]

The performance expectations above were developed using the following elements from the DRC document *A Framework for K-12 Science Education*.

Science and Engineering Practices	Disciplinary Core Ideas	Crosscutting Concepts
<p>Plan and Carry Out Investigations</p> <p>Planning and carrying out investigations to answer questions or test solutions to problems in K-2 builds on prior experiences and progresses to simple investigations, based on fair tests which provide data to support explanations or design solutions.</p> <ul style="list-style-type: none"> Plan and conduct investigations collaboratively to produce data to serve as the basis for evidence to answer a question. (1-PS4-1)(1-PS4-2) <p>Constructing Explanations and Designing Solutions</p> <p>Constructing explanations and designing solutions in K-2 builds on prior experiences and progresses to the use of evidence and ideas in constructing evidence-based accounts of natural phenomena and designing solutions.</p> <ul style="list-style-type: none"> Make observations (firsthand or from media) to construct an evidence-based account for natural phenomena. (1-PS4-2) Use tools and materials provided to design a device that solves a specific problem. (1-PS4-4) 	<p>PS4.A: Wave Properties</p> <ul style="list-style-type: none"> Sound can make matter vibrate, and vibrating matter can make sound. (1-PS4-1) <p>PS4.B: Electromagnetic Radiation</p> <ul style="list-style-type: none"> Objects can be seen if light is available to illuminate them or if they give off their own light. (1-PS4-2) Some materials allow light to pass through them, others allow only some light through and others block all the light and create a dark shadow on any surface beyond them, where the light cannot reach. Mirrors can be used to redirect a light beam. (Boundary: The idea that light travels from place to place is developed through experiences with light sources, mirrors, and shadows, but no attempt is made to discuss the speed of light.) (1-PS4-3) <p>PS4.C: Information Technologies and Instrumentation</p> <ul style="list-style-type: none"> People also use a variety of devices to communicate (and receive information) over long distances. (1-PS4-4) 	<p>Cause and Effect</p> <ul style="list-style-type: none"> Simple tests can be designed to gather evidence to support or refute student ideas about causes. (1-PS4-1)(1-PS4-2)(1-PS4-3) <p>Connections to Engineering, Technology, and Applications of Science</p> <p>Influence of Engineering, Technology, and Science, on Society and the Natural World</p> <ul style="list-style-type: none"> People depend on various technologies in their lives. Future life would be very different without technology. (1-PS4-4)

Connections to Nature of Science

Scientific Investigations Use a Variety of Methods

- Science investigations begin with a question. (1-PS4-1)
- Scientists use different ways to study the world. (1-PS4-1)

Connections to other K-12 in their grade band

Application of 3D's across grade bands: 1-PS4-1 (1-PS4-1), 1-PS4-2 (1-PS4-2), 1-PS4-3 (1-PS4-3), 1-PS4-4 (1-PS4-4)

Common Core State Standards Connections:

ELA/Literacy

- W.1.1 Write informative/explanatory texts in which they name a topic, supply some facts about the topic, and provide some series of claims. (1-PS4-2)
- W.1.2 Participate in shared research and writing projects (e.g., explore a number of "how-to" topics on a given topic and use them to write a sequence of instructions). (1-PS4-2)
- W.1.3 (1-PS4-2)(1-PS4-3)(1-PS4-4)
- W.1.8 With guidance and support from adults, recall information from experiences or gather information from provided sources to answer a question. (1-PS4-1)(1-PS4-2)(1-PS4-3)

SL.1.1 Participate in collaborative conversations with diverse partners about grade 1 topics and texts with peers and adults in small and larger groups. (1-PS4-1)(1-PS4-2)(1-PS4-3)

Mathematics

- Use appropriate tools strategically. (1-PS4-4)
- 1.MD.1.1 Order three objects by length; compare the lengths of two objects indirectly by using a third object. (1-PS4-4)
- 1.MD.1.2 Express the length of an object as a whole number of length units, by laying multiple copies of a shorter object (the length unit) end to end; understand that the length measurement of an object is the number of same-size length units that span it with no gaps or overlaps. (1-PS4-4)

Standard SCI.5.EPS: Students use mathematics and computational thinking, in conjunction with using and disciplinary core ideas, to make sense of phenomena and solve problems.

Learning Element	Performance Indicators (By Grade Band)
SCI.5.EPS.A: Qualitative and Quantitative Data	<p>K-2</p> <p>SCI.5.EPS.A.2 Students recognize that mathematics can be used to describe the natural and designed world. This includes the following:</p> <ul style="list-style-type: none"> Use counting and numbers to identify and describe patterns in the natural and designed worlds. Describe, measure, or compare quantitative attributes of different objects and display the data using simple graphs. Use qualitative and/or quantitative data to compare two alternative solutions to a problem.
	<p>3-5</p> <p>SCI.5.EPS.A.3 Students extend quantitative measurements to a variety of physical properties, using computation and mathematics to analyze data and compare alternative design solutions. This includes the following:</p> <ul style="list-style-type: none"> Organize simple data sets to reveal patterns that suggest relationships. Describe, measure, estimate, and/or graph quantities such as area, volume, weight, and time to address scientific and engineering questions and problems. Create and use graphs or charts generated from simple algorithms to compare alternative solutions to an engineering problem.
	<p>6-8</p> <p>SCI.5.EPS.A.m Students identify patterns in large data sets and use mathematical concepts to support explanations and arguments. This includes the following:</p> <ul style="list-style-type: none"> Decide when to use qualitative vs. quantitative data. Use digital tools (e.g., computers) to analyze very large data sets for patterns and trends. Use mathematical representations to describe and support scientific conclusions and design solutions. Create algorithms (a series of ordered steps) to solve a problem. Apply mathematical concepts and processes (such as ratio, rate, percent, basic operations, and simple algebra) to scientific and engineering questions and problems. Use digital tools and mathematical concepts and arguments to test and compare proposed solutions to an engineering design problem.
	<p>9-12</p> <p>SCI.5.EPS.A.m Students use mathematical concepts to support explanations and arguments. This includes the following:</p> <ul style="list-style-type: none"> Use mathematical representations to describe and support scientific conclusions and design solutions. Use digital tools and mathematical concepts and arguments to test and compare proposed solutions to an engineering design problem.

Wisconsin Standards for Science

LIFE SCIENCE

Standard SCI.LE: Students use science and engineering practices, crosscutting concepts, and an understanding of structures and processes (on a scale from molecules to organisms) to make sense of phenomena and solve problems.

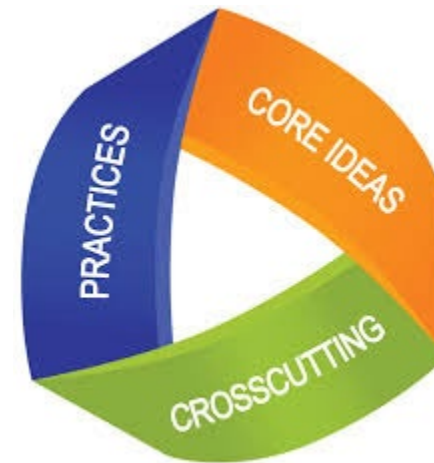
Learning Element	Performance Indicators (By Grade Band)
SCI.LE.A: Structures and Processes	<p>K-2</p> <p>SCI.LE.A.2 All organisms have systems that function and external mechanisms that allow for growth, survival, behavior and reproduction.</p>
SCI.LE.B: Growth and Development of Organisms	<p>K-2</p> <p>SCI.LE.B.2 Growth and offspring often engage in behaviors that help the offspring survive.</p>
SCI.LE.C: Organization for Matter and Energy Flow in Organisms	<p>K-2</p> <p>SCI.LE.C.2 Animals obtain food they need from plants or other animals. Plants need water and light.</p>
SCI.LE.D: Information Systems and Communication	<p>K-2</p> <p>SCI.LE.D.2 Different forms of communication are used to coordinate activities.</p>
	<p>3-5</p> <p>SCI.LE.A.3 All life forms are made up of cells. In organisms, cells work together to form tissues and organs that are essential for vertebrate body functions.</p> <p>SCI.LE.B.3 Reproduction is essential to every level of organism. Organisms have unique and diverse life cycles.</p> <p>SCI.LE.C.3 Plants provide animals with the materials and energy they need for basic needs, growth, survival, and motion. Plants acquire material for growth chiefly from air, water, and inorganic matter, and obtain energy from sunlight which is used to synthesize molecules necessary for survival.</p> <p>SCI.LE.D.3 Different forms of communication are essential for vertebrate levels of organization, including the brain.</p>
	<p>6-8</p> <p>SCI.LE.A.4 Cells and tissues are made up of cells. In organisms, cells work together to form tissues and organs that are essential for vertebrate body functions.</p> <p>SCI.LE.B.4 Growth and division of cells in organisms occur by mitosis and differentiation for specific cell types.</p> <p>SCI.LE.C.4 Plants use the energy from light to make sugars through photosynthesis. While additional oxygen, food is broken down through a series of chemical reactions that release energy and release energy.</p> <p>SCI.LE.D.4 The molecules produced through photosynthesis are used to make amino acids and other molecules that can be assembled into proteins on the ribosome.</p>
	<p>9-12</p> <p>SCI.LE.A.4 The molecules produced through photosynthesis are used to make amino acids and other molecules that can be assembled into proteins on the ribosome.</p>

Standard SCI.CC.3: Students use science and engineering practices, disciplinary core ideas, and an understanding of scale, proportion and quantity to make sense of phenomena and solve problems.

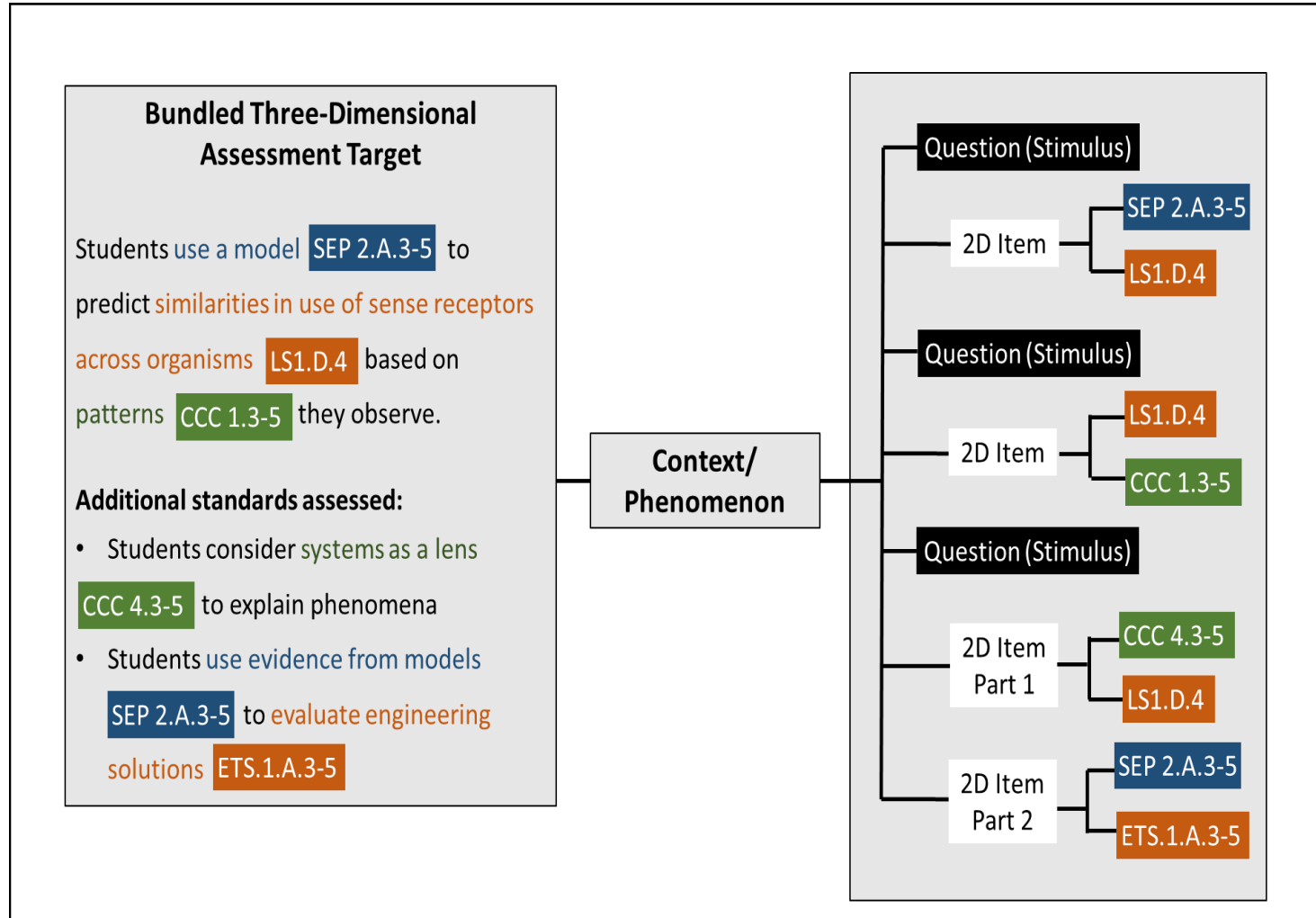
Learning Element	Performance Indicators (By Grade Band)
SCI.CC.3: Scale, Proportion, and Quantity	<p>K-2</p> <p>SCI.CC.3.K-2 Students use relative scales (e.g., bigger and smaller; hotter and colder; faster and slower) to describe objects. They use standard units to measure length.</p>
	<p>3-5</p> <p>SCI.CC.3.3-5 Students recognize natural objects and observable phenomena exist from the very small to the immensely large. They use standard units to measure and describe physical quantities such as mass, time, temperature, and volume.</p>
	<p>6-8</p> <p>SCI.CC.3.6-8 Students observe time, space, and energy phenomena at various scales using models to study systems that are too large or too small. They understand phenomena observable at one scale may not be observable at another scale, and the function of natural and designed systems may change with scale. They use proportional relationships (e.g., speed as the ratio of distance traveled to time taken) to gather information about the magnitude of properties and processes. They use algebraic thinking to examine scientific data and predict the effect of a change in one variable on another (e.g., linear growth vs. exponential growth).</p>
	<p>9-12</p> <p>SCI.CC.3.9-12 Students understand the significance of a phenomenon is dependent on the scale, proportion, and quantity at which it occurs. They recognize patterns observable at one scale may not be observable at another scale, and the function of natural and designed systems may change with scale. They use proportional relationships (e.g., speed as the ratio of distance traveled to time taken) to gather information about the magnitude of properties and processes. They use algebraic thinking to examine scientific data and predict the effect of a change in one variable on another (e.g., linear growth vs. exponential growth).</p>

Core Science Standards Statement

“Students use science and engineering practices, disciplinary core ideas, and crosscutting concepts to make sense of phenomena and solve problems.”



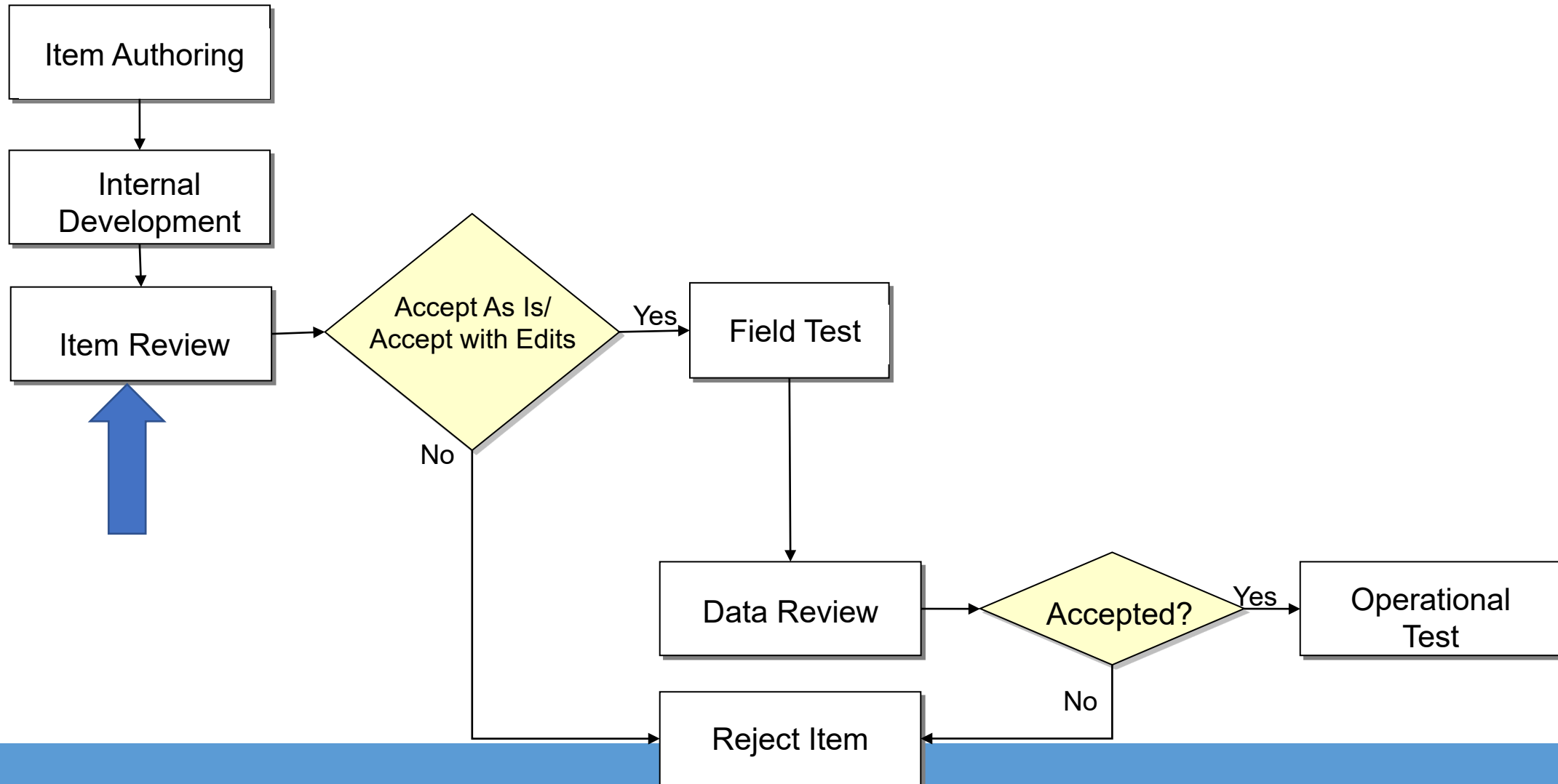
Targets and items



Wisconsin Forward Exam

- Provides a measure of whether students are proficient in the skills and abilities identified in the Wisconsin Academic Standards
- All exam items are aligned to the standards:
 - English Language Arts and Mathematics tested in grades 3-8
 - Science tested in grades 4 and 8
 - Social Studies tested in grades 4, 8, and 10

Life Cycle of an Item



Forward Science Exam Item Types

- Selected Response
 - Multiple Choice (MC)
 - Evidence-Based Selected Response (EBSR)
- Technology Enhanced (TE)

Multiple Choice (MC)

- All MC items have 4 answer choices
 - 3 distractors and 1 correct answer
- Used in all content areas
- Can be linked to a passage or stimulus or used as a “stand-alone MC”
- May have graphs, tables, or other information to support the stem

MC Sample

WBTE Preview Albert Einstein 912355 //

Question 1 Item ID ?

(Click the Enlarge button to learn more about this item.)

[Enlarge](#)

A student examines the data table shown.

Planet	Distance from the Sun (x 1,000,000 km)
Jupiter	778.6
Saturn	1,433.5
Uranus	2,872.5

These data allow the student to produce which type of diagram?

- (a) a diagram comparing the compositions of planets
- (b) a diagram comparing the surface features and temperatures of planets
- (c) a diagram comparing the locations of planets in the solar system
- (d) a diagram comparing the sizes of planets

Evidence-Based Selected Response (EBSR)

2-Part Item (Worth 1 point)

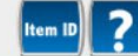
- Part A—one or more correct answers;
 - can be MC or TE
- Part B—one or more correct answers;
 - can be MC or TE

EBSR Sample

WBTE Preview

Albert Einstein 795696 //

Question 1
Page 1 of 2



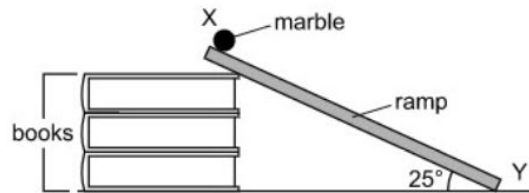
Model Slide

(Practice Hint 1: Use the Line Guide tool to focus on a single line of text.)

(Practice Hint 2: Use the Sticky Note tool to leave a reminder of important information.)

A group of four students worked on a science project. They wanted to learn what angle of a slide on a playground would cause an object to slide down the fastest. They measured the angle between the slide on their playground and the ground and determined that it was 25 degrees. The students set up their experiment with a ramp and books to make a model of their playground slide. They performed the first set of trials with an angle of 25 degrees. During each trial, they measured the time it took a marble to travel from point X to point Y. The students then changed the setup so that the angle of the slide was 30 degrees and repeated their tests.

Experiment Setup

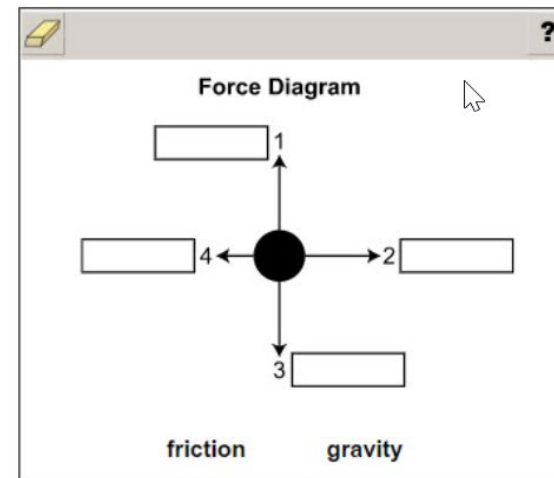


Use the diagram in the scenario to answer the questions.

Part A

One student drew the forces on the marble at point Y, just after it left the ramp. Drag and drop the names of the two forces below into the correct boxes in the diagram.

(Practice Hint 1: Select and drag the words below to the correct boxes to label the forces in the diagram. To change an answer drag the word out of the box.)



(Practice Hint 2: Select Next to move on to Part B.)

EBSR Sample – (cont.)

WBTE Preview

Albert Einstein 795696 //

Question 1
Page 2 of 2



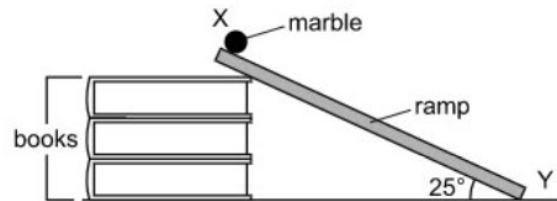
Model Slide

(Practice Hint 1: Use the Line Guide tool to focus on a single line of text.)

(Practice Hint 2: Use the Sticky Note tool to leave a reminder of important information.)

A group of four students worked on a science project. They wanted to learn what angle of a slide on a playground would cause an object to slide down the fastest. They measured the angle between the slide on their playground and the ground and determined that it was 25 degrees. The students set up their experiment with a ramp and books to make a model of their playground slide. They performed the first set of trials with an angle of 25 degrees. During each trial, they measured the time it took a marble to travel from point X to point Y. The students then changed the setup so that the angle of the slide was 30 degrees and repeated their tests.

Experiment Setup



Use the diagram in the scenario to answer the questions.

Part B

The student wants to create another drawing to show the forces on the marble at point X, before it was released. Use the drop-down menus to complete the sentence.

(Practice Hint 3: Select the drop-down menus to view the possible options.)

In this drawing, arrow 4 should be to show forces. arrow 2

balanced
unbalanced

shorter than
longer than
the same length as

Technology Enhanced (TE)

- TE items present in all content areas
- Interactive
- Wide variety, including: drop-down menu, drag and drop, matching, graphing, highlighting text

TE Sample Item

WBTE Preview

Albert Einstein 795695 //

Question 1



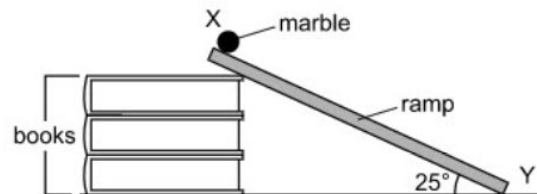
Model Slide

(Practice Hint 1: Use the Line Guide tool to focus on a single line of text.)

(Practice Hint 2: Use the Sticky Note tool to leave a reminder of important information.)

A group of four students worked on a science project. They wanted to learn what angle of a slide on a playground would cause an object to slide down the fastest. They measured the angle between the slide on their playground and the ground and determined that it was 25 degrees. The students set up their experiment with a ramp and books to make a model of their playground slide. They performed the first set of trials with an angle of 25 degrees. During each trial, they measured the time it took a marble to travel from point X to point Y. The students then changed the setup so that the angle of the slide was 30 degrees and repeated their tests.

Experiment Setup



During all of the experiments, the students released the marble at point X.

Select one box in each row of the table to show the changes in the potential energy and kinetic energy of the marble as it traveled from point X to point Y.

(Practice Hint: Select one of the two boxes next to each type of energy in the table to add a check mark.)

	Increased	Decreased
kinetic energy	<input type="checkbox"/>	<input type="checkbox"/>
potential energy	<input type="checkbox"/>	<input type="checkbox"/>

Item Review Process

Participants will view items online using INSIGHT– the same testing engine students use

- Allows interaction with item functionality, particularly useful for technology-enhanced items
- Facilitator will provide specific directions for logging in to begin reviews

Item Review Process – (cont.)

Reviews will be completed in groups and individually. Items will be reviewed for

- Standard alignment
- Grade-level appropriateness
- Correct answer key(s)
- Correct Content
- Depth of Knowledge (DOK) level
- Bias and sensitivity concerns
- Is the wording and technical requirements of the item clear and easy to understand?

Step 1: Standard Alignment

After reading item ask yourself:

Does the standard listed match the state standard?

- Each member will have access to the standards
- Match item to appropriate standard as noted on item rating sheet
- Indicate agreement of alignment on item rating sheet or note the recommended standard and discuss with committee when group reviews item

Step 2: Check the Answer(s)

- Is the answer (or answers) listed correct?
 - If yes, move on to step 3
 - If no, note new answer(s) and discuss with committee when group reviews item

Step 3: Confirm the Depth of Knowledge Level

- Is the DOK level listed correct?
 - If yes, move on to step 4.
 - If no, mark your thinking and discuss with when group reviews item.

We will go into detail about DOK a little later in this presentation.

Step 4: Check for Bias and Sensitivity

- Stereotyping
- Gender
- Regional or geographical
- Ethnic or cultural
- Socioeconomic class
- Persons with a disability
- Ageism
- Religious

Steps 5 and 6: Mark Comments

In spreadsheet, mark column noting the following:

- Accept “A”
 - Item is OK as is
- Accept with Revisions “AR”
 - Accept but apply recommended edits

Dissenting View “DV”

- Additional comments as needed

Step 7: Indicate Item Preference

- Rank item on a scale of 1–5 (with 5 being highest), on your preference for having this item appear on the Wisconsin Forward Exam.

NOTE: This ranking will be used internally and not necessarily discussed as a committee for consensus.

What's DOK?



**Webb's
Depth-of-Knowledge
Levels**

Cognitive Level vs Difficulty

DOK is used by item writers to gauge the *cognitive level* of item, it **does not** correlate to the *difficulty* of the item.

DIFFICULTY ≠ COMPLEXITY

DIFFICULTY	COMPLEXITY
<p>How much effort is needed to answer a question, address a problem, or accomplish a task?</p> <p>How many people can answer a question, address a problem, or accomplish a task correctly or successfully?</p> <p>Easy or Hard</p>	<p>What kind of thinking, action, or knowledge must be demonstrated and communicated to answer a question, address a problem, or accomplish a task?</p> <p>How many different ways can a question be answered, a problem be addressed, or a task be accomplished?</p> <p>Simple or Complex</p>

Definition of DOK

The degree or complexity of knowledge that the content curriculum standards and expectations require.

- Includes four levels, from lowest (basic recall) to highest (extended thinking)
- Focuses on how well the students need to know the content before they can respond to a given item
- Used by item writers to gauge the ***cognitive level*** of item, **does not** correlate to the ***difficulty*** of the item

DOK 1 Recall and Reproduction (not on Forward Exam)

DOK 2 Skills and Concepts

DOK 3 Strategic Thinking and Reasoning

DOK 4 Extended Thinking

(rarely on standardized assessments — more “project-like” or on performance assessments)

DOK 1: Recall and Reproduction

- Students demonstrate a rote response, use a well-known formula, or follow a simple procedure.
- A “simple” procedure is well defined and typically involves only **one** step.

Key words: identify, recall, recognize facts, use, measure, solve a one-step problem

DOK 2: Skills and Concepts

- Students make some decisions regarding how to approach the question or problem.
- Requires deeper knowledge than just giving a definition, such as explaining how or why
- It may involve two or more steps, however two steps does not automatically make a DOK 2.

Key words: explain, categorize, use context clues, select a procedure, compare/contrast

DOK 2 - (cont.)

Activities may include:

- Making observations/collecting information
- Classifying/comparing information
- Organizing/displaying data or information in tables and graphs

Note: Some action verbs, such as “explain,” “describe,” or “interpret,” could be classified at different DOK levels, depending on the complexity of the action.

Sample of a DOK 2

WBTE Preview Albert Einstein 912352 //

Question 1 Item ID ?

(Click the Enlarge button to learn more about this item.)

[Enlarge](#)

In a population of plants, a mutation allows one plant to grow taller than usual. As a result, the plant receives more sunlight than other plants in the area.

Drag the terms into the boxes in the model to show the energy transfer process in which this plant's height gives it an advantage.

?

carbon dioxide

electromagnetic waves

oxygen

sugars

DOK 3: Strategic Thinking and Reasoning

- Students demonstrate deep understanding through planning, using evidence, and exhibiting higher levels of cognitive reasoning.

Key Words: connect ideas, explain thinking, cite evidence, analyze, apply a concept

Activities may include the following:

- Use concepts to solve non-routine problems
- Describe how word choice, point of view or bias, may help the readers' interpretation of text
- Apply a concept in a new context
- Cite evidence and develop a logical argument for concepts
- Compare information within or across data sets

Sample of an Easy DOK 3 Item

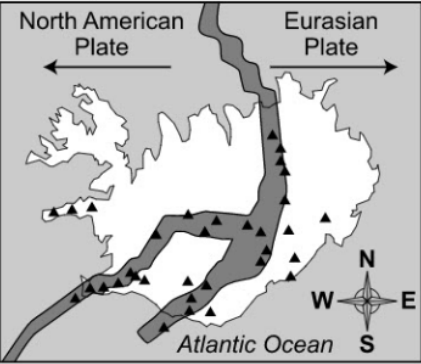
WBTE Preview Albert Einstein 905026 //

Question 1 Item ID ?

Iceland

Iceland provides an amazing geologic laboratory for scientists. It is one of the few places on Earth where scientists can study seafloor spreading—above sea level. This is because Iceland is a product of volcanic activity along the Mid-Atlantic Ridge. The map below shows the path of the Mid-Atlantic Ridge through Iceland and the location of the country's major volcanoes.

Major Volcanoes in Iceland



Key
■ Mid-Atlantic Ridge


Iceland is positioned along a divergent boundary where the North American and Eurasian Plates are moving away from one another. This divergence occurs as a result of convection in Earth's mantle. New, hot magma rises through Earth's mantle and escapes through cracks in Earth's crust. When the magma cools, it forms new crust.

[More Text Below](#)

(Click the Enlarge button to learn more about this item.)

[Enlarge](#)

Compare the two locations on the map of Iceland below with the Major Volcanoes in Iceland map in the scenario.



[Enlarge](#)

Which statement best compares the locations on the map?

- (a) Location 1 is more likely than location 2 to experience a volcanic eruption because it is surrounded by water.
- (b) Location 1 is more likely than location 2 to experience an earthquake because it is on a plate moving to the west.
- (c) Location 2 is more likely than location 1 to experience a volcanic eruption because it is closer to the Mid-Atlantic Ridge.
- (d) Location 2 is more likely than location 1 to experience an earthquake because it has a larger surrounding landmass.

Things to Keep in Mind...

Standards and Difficulty

Items need to measure what students should know and be able to do at their grade level based on the academic standards. This may be different from what your personal experience is with students.

Questions to ask during review:

- Does the item provide for an optimal standard assessment of all students?
- Are there items written to ALL ability levels? It is OK to have easy items.

Things to Keep in Mind...

Technical Design

- Is the item free from being confusing or tricky?
- Does the item meet requirements for technical quality?
- Do graphics/visuals complement and support item?
- Does the stem provide a complete, clear and concise question/problem and directions?
- Is the stem free from clueing the correct answer(s)?
- Are correct answer(s) clear and accurate?
 - Distractors (or incorrect options) may contain common misperceptions or processes

Things to Keep in Mind...

Principles of Universal Design

Items should respect the diversity of the assessment population.

- Every student must be able to access the information.
- Items must measure what is intended.

Items should have

- A clear format for text
- Clear pictures and graphics
- Concise and readable text

The Google Sheet

A Google sheet has been created for each grade. Each reviewer will have their own colored tab located at the bottom. There are also tabs for a DPI reviewer, a BVI reviewer, and the DRC facilitator.

- A link to the Google sheet is located on the website.
- All information (metadata) for each item can be found on the Google sheet.
- All acceptance and reviewer comments will be recorded on the Google sheet.
- Please do not cut and paste in Google sheet.

Item Review Tally Sheet

Step 2

Step 4

Step 6

Session #	Seq #	Item ID	Scenario Title	WSS Performance Indicator	SEP	CCC	Item Type	Points	Key(s)	Proposed DOK	Bias/Sensitivity Comments	Accept (A) Accept with Revisions (AR) Dissenting View (DV)	Comments	Item Preference
1	1	1064623		SCI.ESS2.D.m	Developing Models	Systems and System Models	MC	1	C	2				
1	2	1064620		SCI.LS2.C.m		Cause and Effect	TE	1	decrease; decrease	3				
1	3	1064621		SCI.LS2.D.m	Argue from Evidence		EBSR	1	Part A: A; Part B: C,E	3				
1	4	1064648		SCI.ESS3.C.m	Argue from Evidence	Energy and Matter	TE	1	X to Y; marbles & table to air; kinetic to sound	2				
1	5	1064646		SCI.PS3.D.m	Construct an Explanation	Energy and Matter	TE	1	oxygen; sugars	2				

Step 1

Step 3

Step 5

Step 7

Evaluating an Item: Grade 4 Science

WBTE Preview

795692 // Albert Einstein

Question 1



More Text Above

(Practice Hint 3: Use the scroll bar on the right to access the information below.)

Experiment Setup



Experiment Results



During week 1 of the experiment, the students placed the plant in a

(Practice Hint: Eliminate answer choices by using the Cross-Off tool.)

Which statement best explains why parts of the leaves that were yellow became green after the paper was removed?

- a The green color of the leaves makes seeds so the plant can reproduce.
- b The green color of the leaves attracts insects to help pollinate the plant.
- c The green color of the leaves helps to turn sunlight into food for the plant.
- d The green color of the leaves protects the plant from absorbing too much sunlight.

When to Edit an Item

Reasons to edit an item include, but are not limited to the following:

- If the subject matter is above grade level or out of scope for the standard.
- If there is an opportunity to make the item/passage/stimulus easier for students to understand.
- If the topic or language is inappropriate, controversial, or inflammatory.
- If assigned DOK is not appropriate.

What if I Disagree with the Committee?

- Speak up! It's possible that another committee member has the same concern or you may have noticed something that other committee members have not.
- Record your dissenting view on the item review tracking sheet. Discussion by all is encouraged; however, if you choose not to share your opinion, your facilitator can voice your concern for you.
- DRC and DPI will reconcile any major disagreements/concerns noted on tracking sheet following the meeting.
- **A consensus is not always needed.**

Committee Members

- ✓ Invest yourself in the process.
- ✓ Share your opinions.
- ✓ Listen to your colleagues.
- ✓ Think about *all* Wisconsin students.
- ✓ Remember all passages and items are secure. Integrity during the review process is essential.



Item Review Process: Summary

- ✓ Standard alignment
- ✓ Key(s)
- ✓ DOK levels
- ✓ Grade-level appropriateness
- ✓ Bias and sensitivity

Meeting Evaluation

General Meeting Information ▾

Item Review Evaluation

https://forms.gle/d5eB8...



Item Review Meeting Evaluation

Your email address (YOUR EMAIL WILL APPEAR HERE) will be recorded when you submit this form. Not you? [Switch account](#)

* Required

Please tell us which session you attended. *

Questions?



Thank You!



APPENDIX B
DATA REVIEW TRAINING SLIDES

Wisconsin Forward Exam Item Data Review

August 9, 2022



WISCONSIN DEPARTMENT OF
Public Instruction

Purpose

- Establish a robust pool of items for use in new test development to ensure proper representation:
 - New content standards
 - Test design
- General statistical guidelines are presented
 - Item flags are not created equal
 - Items with statistical flags are not necessarily poor items
 - Item content needs to be considered as well
 - Long-lasting effect of the pandemic on student performance should be taken into consideration
 - Approving an item does not guarantee its appearance on a future test, but rather maximizes the size of the pool for item selection during test development.

Key Objectives

- Review and understand item card layout
- Understand and interpret item statistics
- Review item cards for a few field test items with different statistics
- Apply knowledge of item statistics to evaluate the remaining field test items

Sample Item Card

Assessed standard

Item ID

Content Area

Standard

Stem

Answer choices

Grade

Item Type

Key(s)

Standard: Analyze how a person's local actions can have global consequences, and how global patterns and processes can affect seemingly unrelated local actions.		WI - Data Card
1. Read the information in the box.		Item ID
<div style="border: 1px solid black; padding: 5px; margin-bottom: 10px;"> <p>In 1988, a biologist secretly made a video that showed large numbers of dolphins dying in nets that fishermen used to catch tuna. The release of the video resulted in a consumer boycott of tuna, led in part by school children who urged their parents to not buy tuna. Two years later, one of the largest tuna companies announced it would no longer buy tuna caught by methods that threatened dolphins. Other major producers of canned tuna followed suit. Tuna fishermen changed their fishing methods so that dolphins would no longer accidentally get caught in their nets.</p> </div> <p>Which conclusion is <u>best</u> supported by the information in the box?</p> <p>A. Fishing nets are dangerous and should be outlawed. B. Individuals acting locally can make an impact far away. C. Large companies do not care about fish or mammals. D. Evidence made secretly cannot be used in a court of law.</p>		1035047
		Content Area
		Social Studies
		Course
		Passage ID
		Passage Title
		Grade
		8
		Standards
		WMAS: BH.3.a.r
		Item Type
		Multiple Choice
		Points
		1
		Depth of Knowledge
		3
		Est Difficulty
		Key
		B

Sample Item Card (cont.)

Administration(s)

Form Name	Use Function	Seq	Period	Year	Session	Calc	Model/Ext	Grade
FT1	FT	13	Spring	2021	1	No	3PL/3PL	8

Administration information

Traditional Statistics

N	P-Val	Mean	Item Total Corr
10723	0.52		0.42

Classical statistics

Fit Statistics

Outfit t	Infit t	Outfit MnSq	Infit MnSq	Chi-sq	Deg Free	Item Fit	Fit
						1.32	

Item Fit

IRT Statistics

Label	Final	Final S.E.	Preliminary	Preliminary S.E.	Displ
Slope	0.98				
Location	-0.19				
Asymptote	0.25				

IRT statistics

Distractor/Step Specific

Part	Label	Freq	Proportion	Corr	Avg Meas	Step Meas
	A	3069	0.29	-0.14		
	B	5618	0.52	0.42		
	C	1371	0.13	-0.27		
	D	665	0.06	-0.23		
	OMITS	28	0.00			

Distractor analysis

DIF Analysis

Category	Bias Code	Num Value	N - Ref	N - Focal
ACC			3802	18
DISAB	A-	-0.34	6300	4380
ECODISAD	A	-0.01	5830	4919
LEP	A	-0.16	9056	1417
MALEFEMALE	A	0.35	6181	4568
WHITEAMIN			6238	179
WHITEASIAN	A	0.45	6785	444
WHITEBLACK	A	0.18	6746	998
WHITEHISPANIC	A	-0.13	6799	1899
WHITEMULTI	A	0.38	6768	431

DIF Index

Classical Statistics: Item Difficulty

Difficulty

- “P-Value” : proportion of students who answered an item correctly (or a percent of maximum points possible for polytomously scored items)
 - 0.0 means all students answered incorrectly
 - 1.0 means all students answered correctly
 - The higher the p-value, the easier the item
 - P-value 0.30 – 30% of students answered item correctly
 - P-value 0.50 – half of the students answered item correctly
 - P-value 0.90 – 90% of students answered item correctly

Traditional Statistics

N	P-Val	Mean	Item Total Corr
4349	0.73		0.49

Traditional Statistics

N	P-Val	Mean	Item Total Corr
2803	0.61	1.21	0.36

Item Difficulty: Considerations

Targeted Range

- P-Value: 0.20 to 0.90
- Items outside of target range may be approved if content is appropriate

Content Consideration

- We need to build tests with a wide range of p-values in order to effectively place students into the four performance categories
 - Hard items to distinguish between Proficient/Advanced
 - Easy items to distinguish between Below Basic/Basic
- Why did most students answer this item correctly or incorrectly?
- Are there any reasons other than item difficulty to support a decision to ACCEPT or REJECT this item?

Classical Statistics: Item Discrimination

Discrimination

- Measures item's ability to differentiate between high and low performers
- Item-Total Test Correlation (or point biserial for dichotomously scored items) is the correlation of the examinees' raw scores on a single item with their raw scores on all remaining test items (-1.0 to +1.0)
 - Positive—high achievers outperformed low achievers (targeted).
 - Negative—low achievers outperformed high achievers (unexpected).
 - Around zero—high and low achievers performed about the same on an item (not desired).

Traditional Statistics

N	P-Val	Mean	Item Total Corr
4349	0.73		0.49

Item Discrimination: Considerations

Targeted Range

- above 0.15
 - Higher discrimination is always better
 - Items with negative should be rejected
 - Items with low discrimination often should be rejected
 - Higher items discrimination leads to smaller standard error around student scores

Content Consideration

- Why is this item less able to differentiate between high and low achievers?
- Is the low discrimination associated with extreme low or high P-Values (item difficulty)?
- Are there any other reasons other than item discrimination to support your decision on ACCEPTING or REJECTING this item?

Distractor Specific Analysis (MC Items)

Distractor/Step Specific

Label	Proportion	Corr	Avg Meas	Threshold
A	0.05	-0.22		
B	0.10	-0.26		
C	0.12	-0.28		
D*	0.73	0.49		
MULTS	0.00			
OMITS	0.00			

Guideline

•MC items:

- Correlations for the distractors should be negative.
- Correlations for the distractors should never be higher than correlation for the correct answer
- Proportion of distractor < proportion of key

Content Consideration

- Is the correlation of selecting any incorrect option greater than 0? If yes, why does this option distract more high achievers than low achievers?
- Is the proportion of selecting any incorrect option greater than the proportion of selecting the key? If yes, why?

Score Point Analysis (non-MC Items)

Distractor/Step Specific

Part	Label	Freq	Proportion	Corr
	0	408	0.16	-0.37
	1	2107	0.84	0.37
	BL	3	0.00	

Guideline

- Non-MC items:
 - Worth 1 or more points
 - Correlations for the score 0 expected to be negative
 - Correlation for the highest score should be positive

Content Consideration

- Is the proportion of students at each a score point reasonable?
- Is the pattern of item score correlation as expected?

Option Analysis for EBSR and MS items

Distractor/Step Specific

Part	Label	Freq	Proportion	Corr	Avg Meas	Step Meas
	0	1727	0.63	-0.19		
	1	465	0.17	-0.08		
	2	545	0.20	0.31		
	BL	5	0.00			
1	A	987	0.36	-0.08		
1	B	1010	0.37	0.25		
1	C	405	0.15	-0.12		
1	D	332	0.12	-0.12		
2	A	830	0.30	0.18		
2	B	305	0.11	-0.20		
2	C	1143	0.42	0.02		
2		456	0.17	-0.06		

Distractor/Step Specific

Part	Label	Freq	Proportion	Corr	Avg Meas	Step Meas
	0	5634	0.81	-0.34		
	1	1258	0.18	0.37		
	BL	41	0.01			
1	A	1326	0.19	0.03		
1	B	912	0.13	-0.15		
1	C	2635	0.38	0.29		
1	D	4222	0.61	0.30		
1	E	881	0.13	-0.20		

Guidelines

- Correlations for correct options should be positive; for incorrect options – negative
- Proportions of students at correct options expected to be higher than for incorrect options
- Is the pattern of option proportions and correlations as expected?

IRT Statistics: Item Fit and Non-Convergence

IRT Statistics

Item Fit

- IRT statistic obtained after item calibration
- Measures how well the student responses to each item fit the test data (by comparing parameter estimation prediction relative to the observed data)
- Item is flagged when the observed data pattern differs from the predicted probability of responding to the item.
- There is no specific criterion value for the fit flag: criterion is dependent on the number of students taking the item
- Typically, not a serious flag by itself.

Item Non-Convergence

- Item parameters cannot be estimated, and the item is not eligible for future use

IRT Statistics on Item Cards

Fit Statistics

Outfit t	Infit t	Outfit MnSq	Infit MnSq	Chi-sq	Deg Free	Item Fit	Fit
						11.58	MISFIT

Outfit t	Infit t	Outfit MnSq	Infit MnSq	Chi-sq	Deg Free	Item Fit	Fit
						2.58	

Non-Convergent Items (no Item Fit or IRT Stats)

Outfit t	Infit t	Outfit MnSq	Infit MnSq	Chi-sq	Deg Free	Item Fit	Fit

Differential Item Functioning

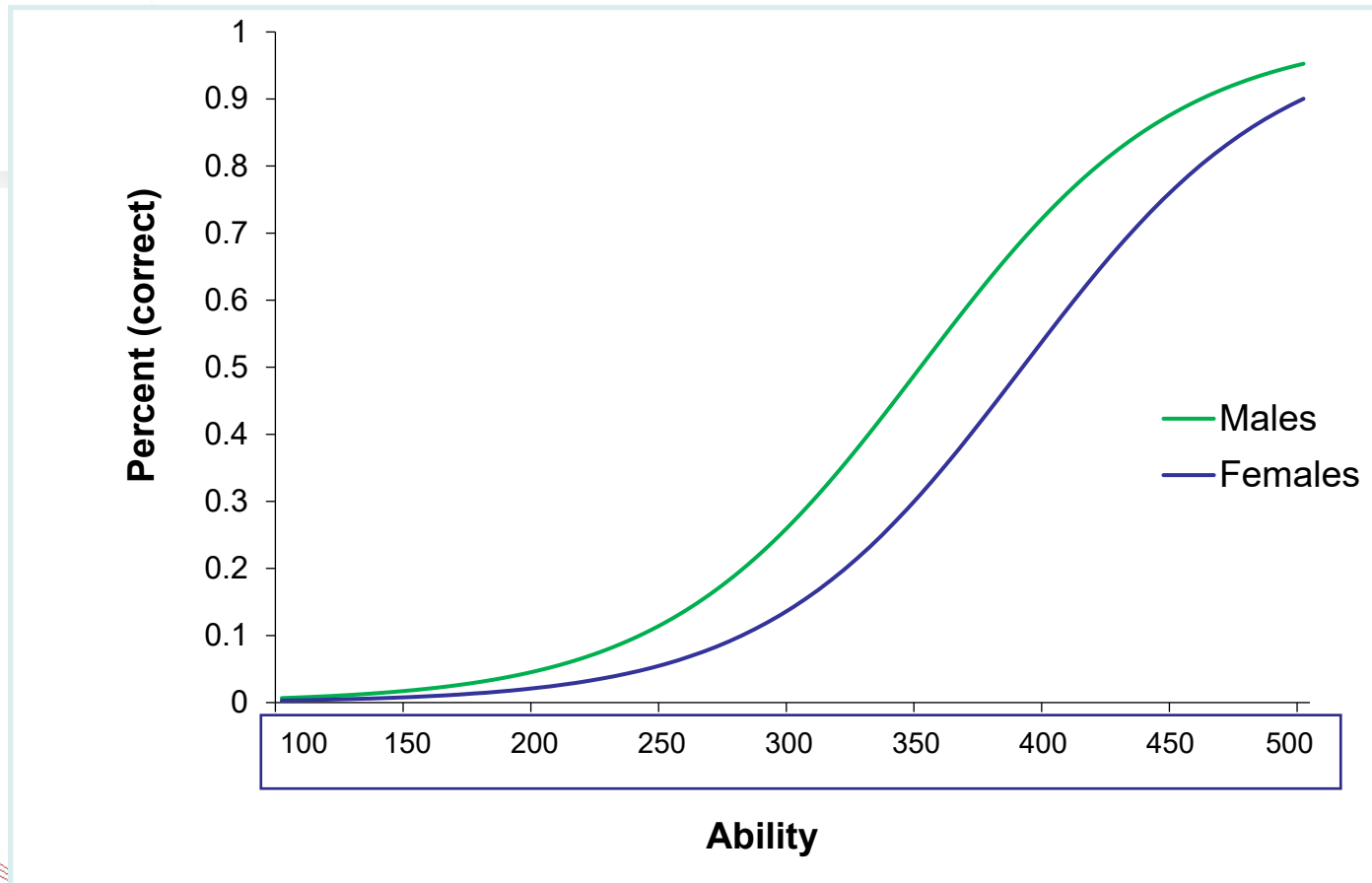
DIF

- Procedure used to identify items that function differently for particular groups of students (e.g., gender, ethnicity, and disability status, SES status, and LEP status).
- Hypothesis is that test takers with similar knowledge or ability should perform in similar ways on a test item.
- Items are flagged if they do not behave the same in different groups of students, after controlling for student ability.

Procedure

- Compares “focal” vs. “reference” groups.
- Reference groups: Males, Whites, students w/out disabilities, students not SES-disadvantaged, English proficient students, students not using accommodations.
- Focal groups: Females, non-White ethnic groups, students with disabilities, SES-disadvantaged students, LEP students, and students using accommodations

Visualizing DIF (Gender)



Differential Item Functioning

Guideline

- Each item is assigned a bias code of A, B, or C.
 - A – minor DIF (no DIF)
 - B – moderate DIF
 - C – Large DIF

DIF signs: “-” against Focal group; “+” favors Focal group.

- Only items with C (i.e., large) DIF require review. Items with C DIF may be acceptable if no potential bias causes the differential item functioning.

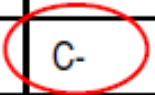
Content Consideration

- Is there anything in the content or format of the item that may interfere with, or advantage, one group of students over another based on:
 - Gender?
 - Ethnicity?
 - Disability status, SES status, LEP status, accommodation use?

DIF Statistics and Codes on Item Cards

DIF Analysis

Category	Bias Code	Num Value	N - Ref	N - Focal
ACC			0	0
DISAB			2270	136
ECODISAD	A	-0.01	1524	991
LEP			2173	101
MALEFEMALE	A	0.02	1231	1284
WHITEAMIN			1092	35
WHITEASIAN			1701	105
WHITEBLACK	C-	-0.12	1524	237
WHITEHISPANIC	A	0.00	1742	249
WHITEMULTI			1754	121



Reference Group/Focal Group

DIF: Summary

- All biased items should show DIF, but **Not** all items with DIF will be biased.
 - Smaller sample sizes of the focal groups may contribute to false positives.
 - DIF not computed if focal group $N < 200$ or reference group $N < 400$
 - You **must** be able to provide a reason for the bias to call the item biased.



Summary of Item Flags

- P-value less than 0.20 or higher than 0.90
- Item-total test correlation < 0.15
 - Negative or close to 0 item-total test correlation is a very serious flag, especially when combined with a positive correlation for a distractor for MC items
- Positive pt. biserial correlation for a distractor
 - Especially if pt. biserial for a distractor is higher than pt. biserial for the correct option
- Poor Fit
- Non-Convergence (kills the item)
- Large DIF (C +/-)
- Omit rates $> 3\%$ (not used in this data review)

Roles, Responsibilities, Questions

- DPI
 - Review Spring 2022 field test item data
 - Accept or reject items

- DRC
 - Facilitate Data Review
 - Answer DPI questions

- Questions?

Thank you!

APPENDIX C
SPRING 2023 ENGLISH LANGUAGE
ARTS OPERATIONAL TEST MAPS

List of Tables

Table C-1. English Language Arts, Grade 3 Test Map..... 3
Table C-2. English Language Arts, Grade 4 Test Map..... 4
Table C-3. English Language Arts, Grade 5 Test Map..... 5
Table C-4. English Language Arts, Grade 6 Test Map..... 6
Table C-5. English Language Arts, Grade 7 Test Map..... 7
Table C-6. English Language Arts, Grade 8 Test Map..... 8

Table C-1. English Language Arts, Grade 3 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
3	1	1	TDA	OP	4	3	3.W.2	Writing
3	2	2	MC	OP	1	2	3.W.8	Writing
3	2	3	MC	OP	1	2	3.W.8	Writing
3	2	4	MC	OP	1	1	3.L.1.b	Writing
3	2	5	MS	OP	2	2	3.L.2.d	Writing
3	2	6	MC	OP	1	1	3.W.1.c	Writing
3	2	7	MC	OP	1	2	3.W.8	Writing
3	2	8	MC	OP	1	2	3.L.1.i	Writing
3	2	9	TE	OP	1	3	3.W.8	Writing
3	2	10	EBSR	OP	2	2	3.W.2.b	Writing
3	2	11	MC	OP	1	2	3.W.1.d	Writing
3	2	12	TE	OP	2	2	3.W.8	Writing
3	2	13	MC	OP	1	1	3.L.1.d	Writing
3	2	14	MC	OP	1	2	3.L.1.g	Writing
3	3	15	MC	OP	1	2	3.SL.2	Listening
3	3	16	MC	OP	1	1	3.SL.2	Listening
3	3	17	MC	OP	1	1	3.SL.3	Listening
3	3	18	EBSR	OP	2	3	3.SL.2	Listening
3	3	19	MC	OP	1	1	3.SL.3	Listening
3	3	20	MC	OP	1	2	3.SL.2	Listening
3	4	21	MC	OP	1	2	3.RL.5	Reading
3	4	22	MC	OP	1	2	3.RL.1	Reading
3	4	23	MC	OP	1	2	3.RL.6	Reading
3	4	24	MC	OP	1	1	3.L.4.b	Reading
3	4	25	MC	OP	1	2	3.L.5.a	Reading
3	4	26	MS	OP	2	1	3.RL.1	Reading
3	4	27	TE	OP	2	3	3.RL.3	Reading
3	4	28	MC	OP	1	2	3.RL.4	Reading
3	4	29	MC	OP	1	3	3.RL.6	Reading
3	4	30	EBSR	OP	2	3	3.RL.2	Reading
3	4	31	MC	OP	1	2	3.RL.5	Reading
3	4	32	MC	OP	1	2	3.RL.6	Reading
3	4	33	EBSR	OP	2	2	3.RI.2	Reading
3	4	34	MC	OP	1	2	3.RI.1	Reading
3	4	35	MC	OP	1	2	3.RI.5	Reading
3	4	36	MC	OP	1	2	3.RI.1	Reading
3	4	37	MC	OP	1	2	3.RI.7	Reading
3	4	38	MC	OP	1	2	3.L.4.RI	Reading

Table C-2. English Language Arts, Grade 4 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
4	1	1	TDA	OP	4	3	4.W.9	Writing
4	2	2	MC	OP	1	2	4.W.3.e	Writing
4	2	3	MC	OP	1	2	4.W.1.c	Writing
4	2	4	EBSR	OP	2	3	4.W.1.b	Writing
4	2	5	MC	OP	1	1	4.L.3.a	Writing
4	2	6	MC	OP	1	2	4.W.8	Writing
4	2	7	TE	OP	2	2	4.W.8	Writing
4	2	8	MC	OP	1	2	4.W.8	Writing
4	2	9	MC	OP	1	3	4.W.8	Writing
4	2	10	MC	OP	1	2	4.W.8	Writing
4	2	11	MC	OP	1	1	4.L.2.c	Writing
4	2	12	TE	OP	2	2	4.L.1.b	Writing
4	2	13	MC	OP	1	1	4.L.2.a	Writing
4	2	14	MC	OP	1	2	4.L.3	Writing
4	3	15	MC	OP	1	2	4.SL.3	Listening
4	3	16	MC	OP	1	2	4.SL.3	Listening
4	3	17	EBSR	OP	2	3	4.SL.2	Listening
4	3	18	MC	OP	1	2	4.SL.3	Listening
4	3	19	MC	OP	1	1	4.SL.2	Listening
4	3	20	EBSR	OP	2	2	4.SL.3	Listening
4	4	21	MS	OP	2	2	4.RL.3	Reading
4	4	22	MC	OP	1	2	4.RL.5	Reading
4	4	23	MC	OP	1	2	4.RL.5	Reading
4	4	24	MC	OP	1	2	4.L.4	Reading
4	4	25	MC	OP	1	2	4.L.4.a	Reading
4	4	26	MC	OP	1	2	4.RI.5	Reading
4	4	27	EBSR	OP	2	3	4.RI.8	Reading
4	4	28	MC	OP	1	2	4.L.4.RI	Reading
4	4	29	MC	OP	1	2	4.RI.2	Reading
4	4	30	MC	OP	1	2	4.RL.2	Reading
4	4	31	MC	OP	1	2	4.RL.3	Reading
4	4	32	MC	OP	1	1	4.L.5.RL	Reading
4	4	33	TE	OP	2	2	4.RL.1	Reading
4	4	34	MC	OP	1	2	4.RL.6	Reading
4	4	35	MC	OP	1	1	4.RI.4	Reading
4	4	36	TE	OP	2	2	4.RI.8	Reading
4	4	37	MC	OP	1	2	4.L.5	Reading
4	4	38	EBSR	OP	2	3	4.RI.1	Reading
4	4	39	MC	OP	1	2	4.RI.9	Reading

Table C-3. English Language Arts, Grade 5 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
5	1	1	TDA	OP	4	3	5.W.9	Writing
5	2	2	MC	OP	1	2	5.L.2.c	Writing
5	2	3	MC	OP	1	3	5.W.3.d	Writing
5	2	4	MS	OP	2	2	5.L.1.a	Writing
5	2	5	MC	OP	1	1	5.W.2.d	Writing
5	2	6	MC	OP	1	2	5.W.8	Writing
5	2	7	TE	OP	2	2	5.W.1.c	Writing
5	2	8	MC	OP	1	2	5.L.2.e	Writing
5	2	9	MC	OP	1	2	5.L.2	Writing
5	2	10	MC	OP	1	2	5.L.1.e	Writing
5	2	11	TE	OP	2	2	5.W.8	Writing
5	2	12	MC	OP	1	2	5.L.2.b	Writing
5	2	13	TE	OP	2	3	5.W.8	Writing
5	3	14	MC	OP	1	1	5.SL.2	Listening
5	3	15	EBSR	OP	2	3	5.SL.3	Listening
5	3	16	MC	OP	1	2	5.SL.2	Listening
5	3	17	MC	OP	1	2	5.SL.3	Listening
5	3	18	MC	OP	1	2	5.SL.2	Listening
5	3	19	EBSR	OP	2	3	5.SL.3	Listening
5	4	20	MC	OP	1	1	5.RL.1	Reading
5	4	21	MC	OP	1	3	5.RL.5	Reading
5	4	22	EBSR	OP	2	3	5.RL.2	Reading
5	4	23	MC	OP	1	2	5.RL.6	Reading
5	4	24	MC	OP	1	1	5.RI.1	Reading
5	4	25	MC	OP	1	1	5.RI.4	Reading
5	4	26	MC	OP	1	2	5.RI.3	Reading
5	4	27	EBSR	OP	2	3	5.RI.8	Reading
5	4	28	MC	OP	1	2	5.RI.4	Reading
5	4	29	MC	OP	1	2	5.L.4.a	Reading
5	4	30	MC	OP	1	2	5.RI.2	Reading
5	4	31	MC	OP	1	3	5.RI.5	Reading
5	4	32	MC	OP	1	3	5.RI.9	Reading
5	4	33	MC	OP	1	2	5.RI.5	Reading
5	4	34	MC	OP	1	2	5.RI.8	Reading
5	4	35	MC	OP	1	2	5.RI.3	Reading
5	4	36	MC	OP	1	2	5.RI.2	Reading
5	4	37	MC	OP	1	2	5.RL.5	Reading
5	4	38	MC	OP	1	2	5.RL.6	Reading
5	4	39	EBSR	OP	2	3	5.RL.1	Reading
5	4	40	MC	OP	1	2	5.RL.4	Reading

Table C-4. English Language Arts, Grade 6 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
6	1	1	TDA	OP	4	3	6.W.9	Writing
6	2	2	MC	OP	1	2	6.W.1	Writing
6	2	3	MC	OP	1	1	6.W.3.b	Writing
6	2	4	MC	OP	1	1	6.L.2.a	Writing
6	2	5	MC	OP	1	2	6.L.3.b	Writing
6	2	6	EBSR	OP	2	3	6.W.8	Writing
6	2	7	TE	OP	2	2	6.W.8	Writing
6	2	8	MC	OP	1	1	6.L.1.a	Writing
6	2	9	EBSR	OP	2	3	6.W.1.a	Writing
6	2	10	TE	OP	2	2	6.L.1.d	Writing
6	2	11	MC	OP	1	2	6.W.8	Writing
6	2	12	MC	OP	1	2	6.L.1.c	Writing
6	2	13	TE	OP	1	2	6.W.8	Writing
6	3	14	MC	OP	1	2	6.SL.2	Listening
6	3	15	EBSR	OP	2	3	6.SL.3	Listening
6	3	16	MC	OP	1	2	6.SL.3	Listening
6	3	17	MC	OP	1	2	6.SL.3	Listening
6	3	18	MC	OP	1	2	6.SL.2	Listening
6	3	19	TE	OP	2	1	6.SL.2	Listening
6	4	20	MC	OP	1	2	6.RL.1	Reading
6	4	21	MC	OP	1	2	6.RL.4	Reading
6	4	22	MC	OP	1	2	6.RL.5	Reading
6	4	23	MS	OP	2	2	6.RL.2	Reading
6	4	24	EBSR	OP	2	3	6.RI.2	Reading
6	4	25	MC	OP	1	3	6.RI.8	Reading
6	4	26	MC	OP	1	2	6.RI.3	Reading
6	4	27	MC	OP	1	2	6.RI.6	Reading
6	4	28	MC	OP	1	2	6.RI.4	Reading
6	4	29	MC	OP	1	2	6.RI.1	Reading
6	4	30	MC	OP	1	2	6.RI.5	Reading
6	4	31	EBSR	OP	2	3	6.RI.6	Reading
6	4	32	TE	OP	2	2	6.RI.6	Reading
6	4	33	MC	OP	1	2	6.RI.4	Reading
6	4	34	MC	OP	1	2	6.RL.2	Reading
6	4	35	MC	OP	1	2	6.RL.4	Reading
6	4	36	MC	OP	1	2	6.RL.6	Reading
6	4	37	TE	OP	2	2	6.RL.3	Reading
6	4	38	MC	OP	1	2	6.RL.9	Reading

Table C-5. English Language Arts, Grade 7 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
7	1	1	TDA	OP	4	3	7.W.9	Writing
7	2	2	MC	OP	1	2	7.L.1.a	Writing
7	2	3	TE	OP	2	2	7.W.3.c	Writing
7	2	4	EBSR	OP	2	3	7.W.8	Writing
7	2	5	MC	OP	1	2	7.L.2	Writing
7	2	6	TE	OP	2	3	7.W.8	Writing
7	2	7	MC	OP	1	2	7.W.8	Writing
7	2	8	MC	OP	1	2	7.W.8	Writing
7	2	9	MC	OP	1	2	7.L.1a	Writing
7	2	10	MC	OP	1	1	7.W.1.d	Writing
7	2	11	MC	OP	1	3	7.L.1.b	Writing
7	2	12	TE	OP	1	1	7.L.2.b	Writing
7	2	13	EBSR	OP	2	3	7.W.1.a	Writing
7	3	14	EBSR	OP	2	3	7.SL.2	Listening
7	3	15	MC	OP	1	2	7.SL.3	Listening
7	3	16	MS	OP	2	2	7.SL.2	Listening
7	3	17	MC	OP	1	2	7.SL.2	Listening
7	3	18	EBSR	OP	2	3	7.SL.3	Listening
7	4	19	MC	OP	1	1	7.RI.1	Reading
7	4	20	MS	OP	2	3	7.RI.2	Reading
7	4	21	MC	OP	1	2	7.L.4.a	Reading
7	4	22	MC	OP	1	2	7.RI.5	Reading
7	4	23	EBSR	OP	2	3	7.RI.3	Reading
7	4	24	MC	OP	1	2	7.RI.6	Reading
7	4	25	MS	OP	2	2	7.RI.8	Reading
7	4	26	MC	OP	1	2	7.RI.3	Reading
7	4	27	MC	OP	1	2	7.L.4.a	Reading
7	4	28	MC	OP	1	2	7.L.4.a	Reading
7	4	29	MC	OP	1	2	7.RI.5	Reading
7	4	30	TE	OP	1	2	7.RL.1	Reading
7	4	31	MC	OP	1	2	7.RL.4	Reading
7	4	32	MC	OP	1	3	7.RL.9	Reading
7	4	33	MC	OP	1	2	7.RI.5	Reading
7	4	34	MC	OP	1	2	7.RL.4	Reading
7	4	35	TE	OP	2	2	7.RI.3	Reading
7	4	36	MC	OP	1	2	7.RL.2	Reading
7	4	37	TE	OP	2	2	7.RL.9	Reading

Table C-6. English Language Arts, Grade 8 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
8	1	1	TDA	OP	4	3	8.W.2	Writing
8	2	2	MC	OP	1	2	8.L.1.c	Writing
8	2	3	MC	OP	1	2	8.W.1.e	Writing
8	2	4	MC	OP	1	2	8.W.8	Writing
8	2	5	MC	OP	1	2	8.L.2.c	Writing
8	2	6	MC	OP	1	2	8.W.1.e	Writing
8	2	7	MC	OP	1	2	8.W.3.c	Writing
8	2	8	MC	OP	1	2	8.L.2.b	Writing
8	2	9	EBSR	OP	2	3	8.W.8	Writing
8	2	10	TE	OP	1	2	8.W.2.d	Writing
8	2	11	TE	OP	2	2	8.W.8	Writing
8	2	12	MC	OP	1	3	8.W.8	Writing
8	2	13	MS	OP	2	1	8.L.3	Writing
8	2	14	MC	OP	1	1	8.L.1.d	Writing
8	3	15	MC	OP	1	2	8.SL.3	Listening
8	3	16	MC	OP	1	2	8.SL.3	Listening
8	3	17	MC	OP	1	2	8.SL.2	Listening
8	3	18	MS	OP	2	2	8.SL.3	Listening
8	3	19	EBSR	OP	2	3	8.SL.3	Listening
8	3	20	MC	OP	1	2	8.SL.2	Listening
8	4	21	MC	OP	1	2	8.L.4.b	Reading
8	4	22	EBSR	OP	2	3	8.RI.2	Reading
8	4	23	MC	OP	1	2	8.RI.4	Reading
8	4	24	MC	OP	1	3	8.RI.6	Reading
8	4	25	MS	OP	2	2	8.RI.7	Reading
8	4	26	MC	OP	1	2	8.RI.3	Reading
8	4	27	MC	OP	1	2	8.RI.1	Reading
8	4	28	MC	OP	1	3	8.RI.5	Reading
8	4	29	EBSR	OP	2	2	8.L.4.RI	Reading
8	4	30	MC	OP	1	2	8.RI.5	Reading
8	4	31	MC	OP	1	2	8.RI.6	Reading
8	4	32	MC	OP	1	2	8.RL.1	Reading
8	4	33	MC	OP	1	2	8.RL.1	Reading
8	4	34	EBSR	OP	2	3	8.RL.3	Reading
8	4	35	MC	OP	1	2	8.RL.6	Reading
8	4	36	MS	OP	2	2	8.RL.1	Reading
8	4	37	MC	OP	1	2	8.RL.2	Reading
8	4	38	MC	OP	1	2	8.L.4.a	Reading
8	4	39	MC	OP	1	3	8.RL.6	Reading

APPENDIX D
SPRING 2023 MATHEMATICS
OPERATIONAL TEST MAPS

List of Tables

Table D-1. Mathematics, Grade 3 Test Map..... 3
Table D-2. Mathematics, Grade 4 Test Map..... 5
Table D-3. Mathematics, Grade 5 Test Map..... 7
Table D-4. Mathematics, Grade 6 Test Map..... 9
Table D-5. Mathematics, Grade 7 Test Map..... 11
Table D-6. Mathematics, Grade 8 Test Map..... 13

Table D-1. Mathematics, Grade 3 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
3	1	1	MC	OP	1	1	3.NBT.2	NBT
3	1	2	TE	OP	1	2	3.MD.3	MD
3	1	3	MC	OP	1	3	3.OA.6	OA
3	1	4	SA	OP	1	1	3.NF.1	NF
3	1	5	MC	OP	1	1	3.G.1	G
3	1	6	TE	OP	1	1	3.NBT.1	NBT
3	1	7	MC	OP	1	2	3.OA.2	OA
3	1	8	MC	OP	1	2	3.MD.4	MD
3	1	9	SA	OP	1	1	3.G.2	G
3	1	10	TE	OP	1	1	3.NF.3.d	NF
3	1	11	MC	OP	1	2	3.OA.8	OA
3	1	12	SA	OP	1	1	3.MD.5.b	MD
3	1	13	SA	OP	1	1	3.NBT.1	NBT
3	1	14	MC	OP	1	1	3.NF.2.a	NF
3	1	15	MC	OP	1	2	3.MD.7.d	MD
3	1	16	MC	OP	1	1	3.OA.1	OA
3	1	17	TE	OP	1	1	3.NBT.3	NBT
3	1	18	MC	OP	1	2	3.NF.3.b	NF
3	1	19	SA	OP	1	1	3.MD.6	MD
3	1	20	SA	OP	1	2	3.OA.7	OA
3	1	21	MC	OP	1	1	3.G.2	G
3	2	22	MC	OP	1	1	3.G.2	G
3	2	23	MC	OP	1	1	3.MD.5.a	MD
3	2	24	TE	OP	1	1	3.OA.5	OA
3	2	25	SA	OP	1	2	3.NBT.2	NBT
3	2	26	MC	OP	1	1	3.NF.3.c	NF
3	2	27	SA	OP	1	2	3.MD.7.b	MD
3	2	28	MC	OP	1	3	3.G.1	G
3	2	29	MC	OP	1	2	3.NF.1	NF
3	2	30	SA	OP	1	1	3.MD.8	MD
3	2	31	MC	OP	1	1	3.NBT.1	NBT
3	2	32	SA	OP	1	1	3.OA.3	OA
3	2	33	MC	OP	1	2	3.NF.3.a	NF
3	2	34	MC	OP	1	2	3.MD.1	MD
3	2	35	SA	OP	1	1	3.NBT.2	NBT
3	2	36	MC	OP	1	1	3.G.1	G
3	2	37	TE	OP	1	2	3.OA.2	OA
3	2	38	MC	OP	1	1	3.NF.2.b	NF

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
3	2	39	MC	OP	1	1	3.MD.2	MD
3	2	40	MC	OP	1	1	3.NBT.3	NBT
3	2	41	TE	OP	1	2	3.G.2	G
3	2	42	SA	OP	1	2	3.OA.9	OA

Domain Names: OA= Operations and Algebraic Thinking; NBT= Number and Operations in Base Ten; NF= Number and Operations – Fractions; MD= Measurement and Data; G=Geometry

Table D-2. Mathematics, Grade 4 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
4	1	1	TE	OP	1	2	4.MD.3	MD
4	1	2	MC	OP	1	2	4.NF.7	NF
4	1	3	SA	OP	1	1	4.OA.1	OA
4	1	4	MC	OP	1	1	4.NBT.2	NBT
4	1	5	TE	OP	1	2	4.MD.4	MD
4	1	6	MC	OP	1	2	4.NF.3.c	NF
4	1	7	TE	OP	1	2	4.G.3	G
4	1	8	MC	OP	1	2	4.MD.2	MD
4	1	9	MC	OP	1	2	4.NBT.5	NBT
4	1	10	MC	OP	1	2	4.OA.2	OA
4	1	11	SA	OP	1	2	4.MD.7	MD
4	1	12	MC	OP	1	1	4.G.1	G
4	1	13	TE	OP	1	1	4.NF.2	NF
4	1	14	MC	OP	1	1	4.NBT.6	NBT
4	1	15	MC	OP	1	1	4.OA.4	OA
4	1	16	MC	OP	1	2	4.NF.4.b	NF
4	1	17	MC	OP	1	1	4.G.2	G
4	1	18	SA	OP	1	1	4.NBT.3	NBT
4	1	19	TE	OP	1	2	4.OA.3	OA
4	1	20	MC	OP	1	1	4.MD.5.a	MD
4	1	21	SA	OP	1	1	4.G.1	G
4	1	22	MC	OP	1	1	4.NF.6	NF
4	1	23	MC	OP	1	3	4.OA.5	OA
4	2	24	MC	OP	1	2	4.NF.5	NF
4	2	25	TE	OP	1	1	4.NBT.2	NBT
4	2	26	MC	OP	1	2	4.OA.3	OA
4	2	27	SA	OP	1	1	4.MD.6	MD
4	2	28	MC	OP	1	3	4.NF.3.a	NF
4	2	29	TE	OP	1	1	4.NBT.4	NBT
4	2	30	MC	OP	1	3	4.MD.1	MD
4	2	31	MC	OP	1	2	4.OA.5	OA
4	2	32	MC	OP	1	1	4.G.1	G
4	2	33	SA	OP	1	2	4.NF.4.c	NF
4	2	34	MC	OP	1	2	4.NBT.1	NBT
4	2	35	TE	OP	1	2	4.OA.1	OA
4	2	36	MC	OP	1	1	4.MD.5.b	MD
4	2	37	SA	OP	1	1	4.NF.3.b	NF
4	2	38	MC	OP	1	2	4.NBT.6	NBT

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
4	2	39	MC	OP	1	2	4.MD.3	MD
4	2	40	MC	OP	1	2	4.OA.2	OA
4	2	41	SA	OP	1	1	4.G.2	G
4	2	42	MC	OP	1	3	4.NF.1	NF
4	2	43	MC	OP	1	2	4.MD.4	MD
4	2	44	SA	OP	1	1	4.OA.4	OA
4	2	45	MC	OP	1	1	4.G.3	G
4	2	46	TE	OP	1	1	4.NBT.3	NBT

Domain Names: OA= Operations and Algebraic Thinking; NBT= Number and Operations in Base Ten; NF= Number and Operations – Fractions; MD= Measurement and Data; G=Geometry

Table D-3. Mathematics, Grade 5 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
5	1	1	MC	OP	1	1	5.NBT.5	NBT
5	1	2	MC	OP	1	3	5.NF.5.b	NF
5	1	3	TE	OP	1	2	5.OA.2	OA
5	1	4	TE	OP	1	2	5.G.1	G
5	1	5	MC	OP	1	2	5.MD.2	MD
5	1	6	SA	OP	1	2	5.NF.7.b	NF
5	1	7	MC	OP	1	2	5.OA.3	OA
5	1	8	SA	OP	1	1	5.NBT.3.a	NBT
5	1	9	TE	OP	1	2	5.MD.1	MD
5	1	10	MC	OP	1	1	5.G.4	G
5	1	11	MC	OP	1	3	5.NBT.2	NBT
5	1	12	MC	OP	1	2	5.NF.6	NF
5	1	13	TE	OP	1	1	5.MD.3.b	MD
5	1	14	SA	OP	1	1	5.NBT.6	NBT
5	1	15	MC	OP	1	1	5.G.3	G
5	1	16	SA	OP	1	2	5.NF.4.a	NF
5	1	17	MC	OP	1	1	5.MD.4	MD
5	1	18	MC	OP	1	2	5.OA.3	OA
5	1	19	TE	OP	1	3	5.G.2	G
5	1	20	SA	OP	1	1	5.NBT.4	NBT
5	1	21	MC	OP	1	2	5.MD.5.a	MD
5	1	22	SA	OP	1	2	5.OA.1	OA
5	1	23	MC	OP	1	2	5.NF.2	NF
5	2	24	SA	OP	1	2	5.MD.5.b	MD
5	2	25	MC	OP	1	2	5.G.2	G
5	2	26	SA	OP	1	1	5.OA.1	OA
5	2	27	TE	OP	1	1	5.NBT.1	NBT
5	2	28	MC	OP	1	1	5.NF.4.b	NF
5	2	29	SA	OP	1	1	5.MD.3.a	MD
5	2	30	SA	OP	1	1	5.OA.2	OA
5	2	31	MC	OP	1	2	5.G.3	G
5	2	32	MC	OP	1	2	5.NF.7.a	NF
5	2	33	MC	OP	1	1	5.MD.3.b	MD
5	2	34	MC	OP	1	2	5.OA.3	OA
5	2	35	TE	OP	1	2	5.NBT.7	NBT
5	2	36	SA	OP	1	2	5.NF.1	NF
5	2	37	MC	OP	1	2	5.MD.2	MD
5	2	38	MC	OP	1	1	5.OA.1	OA
5	2	39	SA	OP	1	1	5.G.1	G
5	2	40	MC	OP	1	2	5.NBT.3.b	NBT

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
5	2	41	TE	OP	1	2	5.NF.3	NF
5	2	42	MC	OP	1	1	5.G.4	G
5	2	43	SA	OP	1	1	5.NBT.2	NBT
5	2	44	MC	OP	1	2	5.MD.5.c	MD
5	2	45	SA	OP	1	2	5.G.2	G
5	2	46	TE	OP	1	2	5.OA.2	OA

Domain Names: OA= Operations and Algebraic Thinking; NBT= Number and Operations in Base Ten; NF= Number and Operations – Fractions; MD= Measurement and Data; G=Geometry

Table D-4. Mathematics, Grade 6 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
6	1	1	MC	OP	1	3	6.RP.3.d	RP
6	1	2	MC	OP	1	1	6.NS.2	NS
6	1	3	MC	OP	1	1	6.EE.2.a	EE
6	1	4	TE	OP	1	2	6.RP.2	RP
6	1	5	MC	OP	1	3	6.EE.4	EE
6	1	6	SA	OP	1	1	6.RP.3.c	RP
6	1	7	MC	OP	1	1	6.EE.3	EE
6	1	8	TE	OP	1	2	6.RP.1	RP
6	1	9	MC	OP	1	2	6.EE.2.b	EE
6	1	10	MC	OP	1	2	6.NS.1	NS
6	1	11	MC	OP	1	1	6.EE.1	EE
6	1	12	MC	OP	1	2	6.RP.3.b	RP
6	1	13	MC	OP	1	2	6.NS.3	NS
6	1	14	MC	OP	1	2	6.RP.3.a	RP
6	1	15	TE	OP	1	2	6.NS.4	NS
6	1	16	SA	OP	1	1	6.RP.2	RP
6	2	17	TE	OP	1	1	6.SP.4	SP
6	2	18	SA	OP	1	2	6.G.3	G
6	2	19	MC	OP	1	1	6.NS.6.c	NS
6	2	20	MC	OP	1	1	6.SP.1	SP
6	2	21	MC	OP	1	2	6.G.1	G
6	2	22	SA	OP	1	2	6.SP.5.a	SP
6	2	23	TE	OP	1	2	6.EE.7	EE
6	2	24	MC	OP	1	3	6.NS.7.d	NS
6	2	25	MC	OP	1	3	6.SP.5.d	SP
6	2	26	SA	OP	1	1	6.G.2	G
6	2	27	TE	OP	1	2	6.NS.5	NS
6	2	28	MC	OP	1	2	6.SP.3	SP
6	2	29	SA	OP	1	2	6.EE.9	EE
6	2	30	SA	OP	1	1	6.NS.6.b	NS
6	2	31	MC	OP	1	2	6.SP.5.b	SP
6	2	32	MC	OP	1	2	6.G.4	G
6	2	33	MC	OP	1	1	6.NS.6.a	NS
6	2	34	SA	OP	1	1	6.EE.8	EE
6	2	35	MC	OP	1	1	6.SP.4	SP
6	2	36	TE	OP	1	2	6.G.3	G
6	2	37	MC	OP	1	2	6.EE.6	EE
6	2	38	SA	OP	1	1	6.NS.8	NS
6	2	39	MC	OP	1	2	6.SP.5.c	SP
6	2	40	TE	OP	1	2	6.G.1	G

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
6	2	41	SA	OP	1	2	6.EE.5	EE
6	2	42	MC	OP	1	2	6.SP.2	SP
6	2	43	MC	OP	1	2	6.G.4	G
6	2	44	TE	OP	1	1	6.NS.7.a	NS
6	2	45	MC	OP	1	1	6.SP.1	SP
6	2	46	TE	OP	1	2	6.EE.6	EE

Domain Names: G=Geometry; EE=Expressions and Equations; NS=The Number System; SP=Statistics and Probability; RP= Ratios and Proportional Relationships

Table D-5. Mathematics, Grade 7 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
7	1	1	SA	OP	1	1	7.NS.1.a	NS
7	1	2	MC	OP	1	2	7.EE.1	EE
7	1	3	SA	OP	1	1	7.NS.2.c	NS
7	1	4	MC	OP	1	2	7.EE.2	EE
7	1	5	TE	OP	1	2	7.NS.2.a	NS
7	1	6	MC	OP	1	2	7.EE.1	EE
7	1	7	SA	OP	1	2	7.NS.3	NS
7	1	8	MC	OP	1	2	7.NS.1.b	NS
7	1	9	SA	OP	1	2	7.EE.2	EE
7	1	10	MC	OP	1	3	7.NS.2.b	NS
7	1	11	MC	OP	1	1	7.NS.1.d	NS
7	2	12	MC	OP	1	2	7.SP.2	SP
7	2	13	TE	OP	1	2	7.G.2	G
7	2	14	MC	OP	1	1	7.RP.2.c	RP
7	2	15	MC	OP	1	1	7.SP.7.a	SP
7	2	16	TE	OP	1	1	7.G.3	G
7	2	17	SA	OP	1	2	7.SP.8.c	SP
7	2	18	TE	OP	1	2	7.EE.4.a	EE
7	2	19	MC	OP	1	2	7.G.1	G
7	2	20	SA	OP	1	1	7.RP.2.b	RP
7	2	21	MC	OP	1	3	7.SP.1	SP
7	2	22	SA	OP	1	2	7.EE.3	EE
7	2	23	MC	OP	1	2	7.G.6	G
7	2	24	TE	OP	1	2	7.RP.2.d	RP
7	2	25	MC	OP	1	1	7.SP.5	SP
7	2	26	MC	OP	1	2	7.G.5	G
7	2	27	MC	OP	1	2	7.EE.4.b	EE
7	2	28	TE	OP	1	2	7.SP.4	SP
7	2	29	MC	OP	1	2	7.G.3	G
7	2	30	MC	OP	1	2	7.RP.2.d	RP
7	2	31	MC	OP	1	2	7.SP.3	SP
7	2	32	MC	OP	1	1	7.G.4	G
7	2	33	SA	OP	1	1	7.RP.2.b	RP
7	2	34	MC	OP	1	2	7.SP.7.b	SP
7	2	35	MC	OP	1	2	7.G.6	G
7	2	36	MC	OP	1	2	7.RP.3	RP
7	2	37	SA	OP	1	2	7.EE.4.a	EE
7	2	38	TE	OP	1	1	7.SP.8.b	SP
7	2	39	SA	OP	1	2	7.G.5	G
7	2	40	MC	OP	1	1	7.RP.2.a	RP

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
7	2	41	SA	OP	1	1	7.SP.6	SP
7	2	42	TE	OP	1	3	7.EE.4.b	EE
7	2	43	MC	OP	1	2	7.G.2	G
7	2	44	SA	OP	1	1	7.RP.1	RP
7	2	45	MC	OP	1	2	7.SP.8.b	SP
7	2	46	MC	OP	1	2	7.EE.3	EE

Domain Names: G=Geometry; EE=Expressions and Equations; NS=The Number System; SP=Statistics and Probability; RP= Ratios and Proportional Relationships

Table D-6. Mathematics, Grade 8 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
8	1	1	SA	OP	1	1	8.NS.1	NS
8	1	2	MC	OP	1	1	8.NS.2	NS
8	1	3	TE	OP	1	2	8.EE.3	EE
8	1	4	MC	OP	1	1	8.NS.1	NS
8	1	5	MC	OP	1	2	8.EE.1	EE
8	1	6	SA	OP	1	1	8.NS.2	NS
8	1	7	MC	OP	1	2	8.EE.4	EE
8	1	8	MC	OP	1	1	8.NS.1	NS
8	1	9	SA	OP	1	1	8.NS.2	NS
8	1	10	MC	OP	1	1	8.EE.2	EE
8	1	11	MC	OP	1	3	8.NS.1	NS
8	1	12	MC	OP	1	1	8.EE.2	EE
8	1	13	MC	OP	1	2	8.NS.2	NS
8	2	14	SA	OP	1	2	8.G.4	G
8	2	15	MC	OP	1	2	8.F.3	F
8	2	16	MC	OP	1	2	8.SP.4	SP
8	2	17	MC	OP	1	2	8.G.8	G
8	2	18	MC	OP	1	2	8.F.1	F
8	2	19	SA	OP	1	2	8.EE.7.a	EE
8	2	20	TE	OP	1	2	8.G.1.c	G
8	2	21	MC	OP	1	2	8.F.2	F
8	2	22	MC	OP	1	1	8.SP.1	SP
8	2	23	MC	OP	1	2	8.F.4	F
8	2	24	MC	OP	1	1	8.G.1.b	G
8	2	25	MC	OP	1	1	8.SP.2	SP
8	2	26	MC	OP	1	2	8.F.5	F
8	2	27	SA	OP	1	1	8.EE.8.a	EE
8	2	28	TE	OP	1	2	8.G.6	G
8	2	29	SA	OP	1	2	8.F.2	F
8	2	30	SA	OP	1	2	8.SP.3	SP
8	2	31	MC	OP	1	2	8.EE.5	EE
8	2	32	TE	OP	1	2	8.F.3	F
8	2	33	MC	OP	1	2	8.G.7	G
8	2	34	MC	OP	1	2	8.SP.2	SP
8	2	35	TE	OP	1	2	8.F.4	F
8	2	36	SA	OP	1	2	8.G.5	G
8	2	37	MC	OP	1	2	8.SP.3	SP
8	2	38	MC	OP	1	3	8.EE.8.a	EE
8	2	39	MC	OP	1	2	8.G.2	G
8	2	40	TE	OP	1	2	8.F.1	F

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
8	2	41	MC	OP	1	2	8.SP.1	SP
8	2	42	TE	OP	1	2	8.G.3	G
8	2	43	MC	OP	1	1	8.F.5	F
8	2	44	MC	OP	1	3	8.EE.6	EE
8	2	45	MC	OP	1	1	8.G.1.a	G
8	2	46	TE	OP	1	2	8.SP.4	SP

Domain Names: G=Geometry; EE=Expressions and Equations; NS=The Number System; SP=Statistics and Probability; F= Functions

APPENDIX E
SPRING 2023 SCIENCE
OPERATIONAL TEST MAPS

List of Tables

Table E-1. Science, Grade 4 Test Map 3
Table E-2. Science, Grade 8 Test Map 5

Table E-1. Science, Grade 4 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain	Science and Engineering Practices Standard	Crosscutting Concepts Standard
4	1	1	TE	OP	1	3	SCI.PS3.B.4	Physical Science	SCI.SEP6.B.3-5	SCI.CC5.3-5
4	1	2	MC	OP	1	2	SCI.ESS2.A.4	Earth and Space Science	SCI.SEP3.A.3-5	
4	1	3	TE	OP	1	3	SCI.ESS3.A.4	Earth and Space Science	SCI.SEP8.A.3-5	
4	1	4	TE	OP	1	2	SCI.LS1.D.4	Life Science	SCI.SEP8.A.3-5	
4	1	5	MC	OP	1	2	SCI.LS1.A.4	Life Science	SCI.SEP7.A.3-5	
4	1	6	TE	OP	1	3	SCI.PS3.A.4	Physical Science	SCI.SEP6.A.3-5	SCI.CC5.3-5
4	1	7	TE	OP	1	2	SCI.PS3.A.4	Physical Science	SCI.SEP6.A.3-5	SCI.CC5.3-5
4	1	8	TE	OP	1	2	SCI.ETS1.A.3-5	Engineering	SCI.SEP6.B.3-5	SCI.CC5.3-5
4	1	9	TE	OP	1	3	SCI.LS1.D.4	Life Science	SCI.SEP2.A.3-5	SCI.CC4.3-5
4	1	10	MC	OP	1	2	SCI.ESS2.A.4	Earth and Space Science	SCI.SEP3.A.3-5	SCI.CC2.3-5
4	1	11	MC	OP	1	2	SCI.ETS1.B.3-5	Engineering	SCI.SEP3.A.3-5	
4	1	12	TE	OP	1	3	SCI.PS3.A.4	Physical Science	SCI.SEP6.A.3-5	SCI.CC5.3-5
4	1	13	TE	OP	1	2	SCI.ESS1.C.4	Earth and Space Science	SCI.SEP6.A.3-5	SCI.CC1.3-5
4	1	14	TE	OP	1	2	SCI.LS1.A.4	Life Science	SCI.SEP7.A.3-5	
4	1	15	MC	OP	1	3	SCI.ETS1.B.3-5	Engineering	SCI.SEP3.A.3-5	
4	2	16	TE	OP	1	2	SCI.ETS1.A.3-5	Engineering	SCI.SEP1.B.3-5	
4	2	17	EBSR	OP	1	3	SCI.ESS1.C.4	Earth and Space Science	SCI.SEP6.A.3-5	SCI.CC1.3-5
4	2	18	TE	OP	1	3	SCI.PS3.A.4	Physical Science	SCI.SEP6.A.3-5	SCI.CC5.3-5
4	2	19	EBSR	OP	1	3	SCI.ESS2.B.4	Earth and Space Science	SCI.SEP4.A.3-5	SCI.CC1.3-5
4	2	20	MC	OP	1	2	SCI.LS1.A.4	Life Science	SCI.SEP7.A.3-5	SCI.CC4.3-5
4	2	21	TE	OP	1	2	SCI.LS1.D.4	Life Science	SCI.SEP2.A.3-5	SCI.CC4.3-5
4	2	22	TE	OP	1	3	SCI.LS1.D.4	Life Science	SCI.SEP2.A.3-5	SCI.CC4.3-5
4	2	23	MC	OP	1	2	SCI.ETS1.A.3-5	Engineering	SCI.SEP1.B.3-5	
4	2	24	MC	OP	1	2	SCI.PS3.B.4	Physical Science	SCI.SEP6.B.3-5	

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain	Science and Engineering Practices Standard	Crosscutting Concepts Standard
4	2	25	TE	OP	1	2	SCI.PS4.C.4	Physical Science	SCI.SEP8.A.3-5	
4	3	26	MC	OP	1	2	SCI.PS4.A.4	Physical Science	SCI.SEP2.A.3-5	SCI.CC2.3-5
4	3	27	MC	OP	1	2	SCI.PS3.B.4	Physical Science	SCI.SEP8.A.3-5	SCI.CC5.3-5
4	3	28	MC	OP	1	2	SCI.ESS2.B.4	Earth and Space Science	SCI.SEP6.A.3-5	SCI.CC1.3-5
4	3	29	TE	OP	1	2	SCI.ESS2.A.4	Earth and Space Science	SCI.SEP4.A.3-5	SCI.CC2.3-5
4	3	30	TE	OP	1	3	SCI.ESS2.B.4	Earth and Space Science	SCI.SEP4.A.3-5	
4	3	31	TE	OP	1	3	SCI.ETS1.A.3-5	Engineering	SCI.SEP1.A.3-5	
4	3	32	TE	OP	1	2	SCI.ETS1.B.3-5	Engineering	SCI.SEP6.B.3-5	
4	3	33	MC	OP	1	2	SCI.LS1.A.4	Life Science	SCI.SEP6.A.3-5	SCI.CC6.3-5
4	3	34	TE	OP	1	3	SCI.PS3.A.4	Physical Science	SCI.SEP6.A.3-5	
4	3	35	TE	OP	1	2	SCI.LS1.A.4	Life Science		SCI.CC4.3-5
4	3	36	TE	OP	1	2	SCI.LS1.A.4	Life Science	SCI.SEP8.A.3-5	SCI.CC4.3-5
4	3	37	TE	OP	1	3	SCI.LS1.A.4	Life Science	SCI.SEP7.A.3-5	SCI.CC4.3-5
4	3	38	EBSR	OP	1	2	SCI.ETS1.B.3-5	Engineering	SCI.SEP6.B.3-5	
4	3	39	TE	OP	1	3	SCI.LS1.A.4	Life Science	SCI.SEP7.A.3-5	SCI.CC4.3-5
4	3	40	MC	OP	1	2	SCI.PS3.C.4	Physical Science		SCI.CC2.3-5

Table E-2. Science, Grade 8 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain	Science and Engineering Practices Standard	Crosscutting Concepts Standard
8	1	1	TE	OP	1	3	SCI.PS2.B.m	Physical Science		SCI.CC4.m
8	1	2	MS	OP	1	3	SCI.LS4.A.m	Life Science	SCI.SEP4.A.m	SCI.CC1.m
8	1	3	TE	OP	1	2	SCI.LS2.C.m	Life Science	SCI.SEP7.A.m	SCI.CC1.m
8	1	4	TE	OP	1	2	SCI.LS2.A.m	Life Science	SCI.SEP4.A.m	SCI.CC1.m
8	1	5	TE	OP	1	2	SCI.PS4.A.m	Physical Science	SCI.SEP2.A.m	SCI.CC6.m
8	1	6	TE	OP	1	2	SCI.ESS2.C.m	Earth and Space Science	SCI.SEP2.A.m	SCI.CC4.m
8	1	7	TE	OP	1	2	SCI.PS1.B.m	Physical Science	SCI.SEP7.A.m	
8	1	8	MC	OP	1	3	SCI.ETS1.B.m	Engineering	SCI.SEP7.A.m	
8	1	9	TE	OP	1	2	SCI.ESS1.B.m	Earth and Space Science	SCI.SEP2.A.m	SCI.CC1.m
8	1	10	MC	OP	1	2	SCI.ESS1.B.m	Earth and Space Science	SCI.SEP2.A.m	SCI.CC4.m
8	1	11	MC	OP	1	3	SCI.LS4.C.m	Life Science	SCI.SEP5.A.m	SCI.CC2.m
8	1	12	MC	OP	1	2	SCI.LS3.B.m	Life Science	SCI.SEP2.A.m	
8	1	13	TE	OP	1	3	SCI.LS4.B.m	Life Science	SCI.SEP8.A.m	SCI.CC2.m
8	1	14	TE	OP	1	3	SCI.ESS3.A.m	Earth and Space Science	SCI.SEP6.A.m	SCI.CC2.m
8	1	15	MC	OP	1	2	SCI.LS2.A.m	Life Science		SCI.CC5.m
8	2	16	MC	OP	1	2	SCI.PS4.B.m	Physical Science	SCI.SEP2.A.m	SCI.CC4.m
8	2	17	TE	OP	1	3	SCI.PS2.A.m	Physical Science	SCI.SEP3.A.m	SCI.CC7.m
8	2	18	TE	OP	1	2	SCI.PS1.A.m	Physical Science	SCI.SEP2.A.m	SCI.CC2.m
8	2	19	MC	OP	1	2	SCI.ETS1.A.m	Engineering		SCI.CC6.m
8	2	20	TE	OP	1	2	SCI.ETS1.A.m	Engineering	SCI.SEP6.A.m	
8	2	21	TE	OP	1	2	SCI.ETS1.A.m	Engineering	SCI.SEP1.A.m	
8	2	22	EBSR	OP	1	3	SCI.ESS1.C.m	Earth and Space Science	SCI.SEP6.A.m	SCI.CC3.m
8	2	23	TE	OP	1	2	SCI.ESS1.C.m	Earth and Space Science	SCI.SEP6.A.m	SCI.CC3.m

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain	Science and Engineering Practices Standard	Crosscutting Concepts Standard
8	2	24	TE	OP	1	3	SCI.ESS3.C.m	Earth and Space Science	SCI.SEP7.A.m	SCI.CC2.m
8	2	25	TE	OP	1	3	SCI.LS3.B.m	Life Science	SCI.SEP2.A.m	SCI.CC6.m
8	3	26	MC	OP	1	3	SCI.LS1.B.m	Life Science	SCI.SEP4.A.m	
8	3	27	TE	OP	1	2	SCI.LS2.C.m	Life Science		SCI.CC7.m
8	3	28	MC	OP	1	3	SCI.LS2.A.m	Life Science		SCI.CC7.m
8	3	29	TE	OP	1	2	SCI.ETS1.B.m	Engineering	SCI.SEP4.A.m	
8	3	30	TE	OP	1	2	SCI.PS1.B.m	Physical Science	SCI.SEP2.A.m	SCI.CC5.m
8	3	31	MC	OP	1	3	SCI.PS1.B.m	Physical Science	SCI.SEP6.B.m	
8	3	32	TE	OP	1	3	SCI.ETS1.A.m	Engineering	SCI.SEP1.A.m	
8	3	33	EBSR	OP	1	3	SCI.ETS1.A.m	Engineering	SCI.SEP1.A.m	
8	3	34	TE	OP	1	3	SCI.ESS2.D.m	Earth and Space Science	SCI.SEP4.A.m	SCI.CC1.m
8	3	35	MC	OP	1	2	SCI.PS4.C.m	Physical Science	SCI.SEP1.A.m	
8	3	36	TE	OP	1	3	SCI.PS2.A.m	Physical Science	SCI.SEP3.A.m	SCI.CC7.m
8	3	37	MC	OP	1	2	SCI.ETS1.B.m	Engineering	SCI.SEP4.A.m	
8	3	38	TE	OP	1	2	SCI.ESS3.D.m	Earth and Space Science	SCI.SEP1.A.m	SCI.CC7.m
8	3	39	TE	OP	1	3	SCI.PS3.A.m	Physical Science	SCI.SEP2.A.m	
8	3	40	EBSR	OP	1	2	SCI.PS3.B.m	Physical Science	SCI.SEP6.A.m	SCI.CC5.m

APPENDIX F
SPRING 2023 SOCIAL STUDIES
OPERATIONAL TEST MAPS

List of Tables

Table F-1. Social Studies, Grade 4 Test Map 3
Table F-2. Social Studies, Grade 8 Test Map 4
Table F-3. Social Studies, Grade 10 Test Map 5

Table F-1. Social Studies, Grade 4 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
4	1	1	TE	OP	1	2	BH.2.b.4	Behavioral Sciences
4	1	2	MC	OP	1	2	Geog.4.a.4	Geography
4	1	3	TE	OP	1	2	PS.3.a.4-5	Political Science
4	1	4	MC	OP	1	2	Hist.2.a.i	History
4	1	5	MC	OP	1	2	BH.4.a.i	Behavioral Sciences
4	1	6	MS	OP	1	2	BH.2.b.4	Behavioral Sciences
4	1	7	TE	OP	1	3	Geog.2.a.3	Geography
4	1	8	MC	OP	1	2	BH.1.b.4	Behavioral Sciences
4	1	9	MC	OP	1	2	Econ.1.b.4	Economics
4	1	10	MC	OP	1	2	PS.2.a.i	Political Science
4	1	11	MC	OP	1	2	Geog.1.a.4-5	Geography
4	1	12	MC	OP	1	3	BH.1.a.4	Behavioral Sciences
4	1	13	MS	OP	1	3	BH.2.a.4-5	Behavioral Sciences
4	1	14	MC	OP	1	2	Geog.2.a.3	Geography
4	1	15	MC	OP	1	3	Econ.3.a.4	Economics
4	1	16	MC	OP	1	2	Geog.4.a.4	Geography
4	1	17	MC	OP	1	3	Hist.1.b.i	History
4	1	18	TE	OP	1	3	Econ.2.c.3	Economics
4	1	19	MC	OP	1	3	Geog.3.b.4	Geography
4	1	20	MC	OP	1	2	Econ.2.a.3-4	Economics
4	2	21	TE	OP	1	2	PS.4.a.i	Political Science
4	2	22	MC	OP	1	2	Geog.1.b.i	Geography
4	2	23	MC	OP	1	3	Hist.4.d.i	History
4	2	24	TE	OP	1	3	Geog.3.b.4	Geography
4	2	25	MC	OP	1	2	Hist.4.c.i	History
4	2	26	TE	OP	1	2	Geog.2.d.4-5	Geography
4	2	27	TE	OP	1	2	PS.2.a.i	Political Science
4	2	28	MC	OP	1	2	Geog.1.b.i	Geography
4	2	29	TE	OP	1	2	BH.2.a.4-5	Behavioral Sciences
4	2	30	TE	OP	1	2	PS.1.a.i	Political Science
4	2	31	TE	OP	1	2	Econ.1.a.3	Economics
4	2	32	MC	OP	1	3	Hist.4.c.i	History
4	2	33	TE	OP	1	3	Hist.4.a.i	History
4	2	34	MC	OP	1	2	Geog.1.b.i	Geography
4	2	35	MC	OP	1	3	PS.4.a.i	Political Science
4	2	36	MC	OP	1	3	Hist.2.c.i	History
4	2	37	TE	OP	1	2	Hist.2.b.i	History
4	2	38	MC	OP	1	2	Econ.2.b.4-5	Economics
4	2	39	MC	OP	1	3	Hist.4.d.i	History
4	2	40	MC	OP	1	2	Econ.1.b.4	Economics

Table F-2. Social Studies, Grade 8 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
8	1	1	TE	OP	1	2	Geog.3.a.m	Geography
8	1	2	MC	OP	1	3	Geog.2.c.m	Geography
8	1	3	TE	OP	1	2	Econ.4.b.m	Economics
8	1	4	MS	OP	1	3	Hist.1.a.m	History
8	1	5	MC	OP	1	3	PS.3.c.m	Political Science
8	1	6	MC	OP	1	3	Hist.4.a.m	History
8	1	7	MC	OP	1	3	Geog.2.a.m	Geography
8	1	8	MC	OP	1	2	PS.1.a.m	Political Science
8	1	9	TE	OP	1	2	Geog.1.c.m	Geography
8	1	10	MC	OP	1	2	BH.3.a.m	Behavioral Sciences
8	1	11	MC	OP	1	2	Hist.2.c.m	History
8	1	12	MC	OP	1	2	PS.2.c.m	Political Science
8	1	13	MC	OP	1	2	PS.3.d.m	Political Science
8	1	14	MC	OP	1	3	Econ.2.b.m	Economics
8	1	15	MC	OP	1	2	Geog.5.b.m	Geography
8	1	16	MC	OP	1	2	Econ.2.c.m	Economics
8	1	17	MC	OP	1	3	PS.3.a.m	Political Science
8	1	18	MC	OP	1	2	Hist.2.a.m	History
8	1	19	TE	OP	1	3	Geog.4.a.m	Geography
8	1	20	MC	OP	1	3	Econ.4.a.m	Economics
8	2	21	MS	OP	1	2	PS.3.d.m	Political Science
8	2	22	TE	OP	1	2	Geog.3.b.m	Geography
8	2	23	MC	OP	1	3	Hist.3.c.m	History
8	2	24	MC	OP	1	3	BH.1.b.m	Behavioral Sciences
8	2	25	MC	OP	1	2	Hist.4.c.m	History
8	2	26	MC	OP	1	3	PS.3.a.m	Political Science
8	2	27	MC	OP	1	2	Geog.4.a.m	Geography
8	2	28	TE	OP	1	2	BH.4.a.m	Behavioral Sciences
8	2	29	MC	OP	1	2	Hist.3.c.m	History
8	2	30	MC	OP	1	2	BH.1.b.m	Behavioral Sciences
8	2	31	MC	OP	1	3	Econ.4.c.m	Economics
8	2	32	MC	OP	1	3	BH.3.a.m	Behavioral Sciences
8	2	33	MC	OP	1	3	Hist.2.c.m	History
8	2	34	MC	OP	1	2	BH.2.a.m	Behavioral Sciences
8	2	35	MS	OP	1	3	Hist.4.c.m	History
8	2	36	MC	OP	1	2	Econ.4.a.m	Economics
8	2	37	TE	OP	1	2	BH.1.a.m	Behavioral Sciences
8	2	38	MC	OP	1	2	Geog.2.b.m	Geography
8	2	39	MC	OP	1	3	Hist.2.b.m	History
8	2	40	TE	OP	1	2	Econ.3.b.m	Economics

Table F-3. Social Studies, Grade 10 Test Map

Grade	Session	Item Sequence	Item Type	Item Usage	Max Score Points	Depth of Knowledge	Standard	Domain
10	1	1	TE	OP	1	2	PS.2.b.h	Political Science
10	1	2	MC	OP	1	2	Geog.2.a.h	Geography
10	1	3	MC	OP	1	2	Hist.2.c.h	History
10	1	4	MC	OP	1	2	Geog.5.b.h	Geography
10	1	5	MC	OP	1	3	Hist.4.c.h	History
10	1	6	MC	OP	1	2	Hist.4.c.h	History
10	1	7	TE	OP	1	2	Geog.3.b.h	Geography
10	1	8	MS	OP	1	3	Hist.2.b.h	History
10	1	9	MC	OP	1	2	Econ.4.a.h	Economics
10	1	10	MC	OP	1	3	Hist.3.a.h	History
10	1	11	MC	OP	1	2	Hist.2.a.h	History
10	1	12	MC	OP	1	2	Econ.4.e.h	Economics
10	1	13	MC	OP	1	2	PS.1.b.h	Political Science
10	1	14	MC	OP	1	2	Econ.4.b.h	Economics
10	1	15	MC	OP	1	3	Hist.3.c.h	History
10	1	16	MC	OP	1	2	Econ.4.c.h	Economics
10	1	17	MC	OP	1	2	PS.2.a.h	Political Science
10	1	18	MC	OP	1	2	BH.2.b.h	Behavioral Sciences
10	1	19	MC	OP	1	2	Geog.3.b.h	Geography
10	1	20	MS	OP	1	2	BH.2.a.h	Behavioral Sciences
10	2	21	MC	OP	1	3	Hist.2.b.h	History
10	2	22	MS	OP	1	2	BH.1.a.h	Behavioral Sciences
10	2	23	TE	OP	1	3	BH.2.a.h	Behavioral Sciences
10	2	24	MC	OP	1	2	Econ.3.b.h	Economics
10	2	25	MC	OP	1	2	BH.1.a.h	Behavioral Sciences
10	2	26	MS	OP	1	2	Geog.1.a.h	Geography
10	2	27	MC	OP	1	2	Econ.3.b.h	Economics
10	2	28	MC	OP	1	2	Geog.4.a.h	Geography
10	2	29	MC	OP	1	2	PS.1.b.h	Political Science
10	2	30	MC	OP	1	3	BH.2.a.h	Behavioral Sciences
10	2	31	TE	OP	1	2	Geog.5.b.h	Geography
10	2	32	TE	OP	1	2	Econ.3.c.h	Economics
10	2	33	MC	OP	1	2	Geog.1.c.h	Geography
10	2	34	MC	OP	1	3	PS.3.b.h	Political Science
10	2	35	MC	OP	1	2	PS.2.a.h	Political Science
10	2	36	TE	OP	1	2	PS.2.a.h	Political Science
10	2	37	MC	OP	1	2	Geog.5.a.h	Geography
10	2	38	MC	OP	1	2	PS.3.d.h	Political Science
10	2	39	MC	OP	1	2	BH.1.b.h	Behavioral Sciences
10	2	40	MC	OP	1	2	Geog.2.c.h	Geography

APPENDIX G
SPRING 2023 TEST PARTICIPATION
RATES BY SUBGROUP

List of Tables

Table G-1. Test Participation Rates, Grade 3 3
Table G-2. Test Participation Rates, Grade 4 4
Table G-3. Test Participation Rates, Grade 5 5
Table G-4. Test Participation Rates, Grade 6 6
Table G-5. Test Participation Rates, Grade 7 7
Table G-6. Test Participation Rates, Grade 8 8
Table G-7. Test Participation Rates, Grade 10 9

Table G-1. Test Participation Rates, Grade 3

Group	Category	Enrolled	English Language Arts		Mathematics	
			Number Tested	Percent Tested	Number Tested	Percent Tested
State	All Students	61228	58497	95.54	58722	95.91
Gender	Male	31529	30027	95.24	30150	95.63
	Female	29692	28463	95.86	28565	96.20
	Non-Binary	7	7	100.00	7	100.00
Race/ Ethnicity	White	39463	37862	95.94	37895	96.03
	African American	6244	5948	95.26	5932	95.00
	Hispanic	8592	8088	94.13	8259	96.12
	Asian	2854	2743	96.11	2782	97.48
	American Indian	579	553	95.51	553	95.51
	Two or More	3496	3303	94.48	3301	94.42
Limited English Proficiency	No	55638	53346	95.88	53345	95.88
	Yes	5590	5151	92.15	5377	96.19
Disability Status	No	52053	50140	96.32	50365	96.76
	Yes	9175	8357	91.08	8357	91.08
Economically Disadvantaged	No	34771	33351	95.92	33466	96.25
	Yes	26457	25146	95.04	25256	95.46

Table G-2. Test Participation Rates, Grade 4

Group	Category	Enrolled	English Language Arts		Mathematics		Science		Social Studies	
			Number Tested	Percent Tested	Number Tested	Percent Tested	Number Tested	Percent Tested	Number Tested	Percent Tested
State	All Students	61689	58996	95.63	59165	95.91	59141	95.87	59131	95.85
Gender	Male	31449	29985	95.34	30064	95.60	30053	95.56	30047	95.54
	Female	30228	29001	95.94	29091	96.24	29078	96.20	29074	96.18
	Non-Binary	12	10	83.33	10	83.33	10	83.33	10	83.33
Race/Ethnicity	White	39918	38340	96.05	38345	96.06	38347	96.06	38340	96.05
	African American	6042	5734	94.90	5724	94.74	5717	94.62	5702	94.37
	Hispanic	8778	8286	94.40	8427	96.00	8413	95.84	8422	95.94
	Asian	2907	2803	96.42	2839	97.66	2838	97.63	2840	97.70
	American Indian	608	576	94.74	576	94.74	575	94.57	576	94.74
	Two or More	3436	3257	94.79	3254	94.70	3251	94.62	3251	94.62
Limited English Proficiency	No	56088	53787	95.90	53766	95.86	53745	95.82	53730	95.80
	Yes	5601	5209	93.00	5399	96.39	5396	96.34	5401	96.43
Disability Status	No	52610	50813	96.58	50989	96.92	50970	96.88	50965	96.87
	Yes	9079	8183	90.13	8176	90.05	8171	90.00	8166	89.94
Economically Disadvantaged	No	35494	34066	95.98	34147	96.20	34137	96.18	34138	96.18
	Yes	26195	24930	95.17	25018	95.51	25004	95.45	24993	95.41

Table G-3. Test Participation Rates, Grade 5

Group	Category	Enrolled	English Language Arts		Mathematics	
			Number Tested	Percent Tested	Number Tested	Percent Tested
State	All Students	62058	59386	95.69	59577	96.00
Gender	Male	31886	30408	95.36	30495	95.64
	Female	30158	28964	96.04	29068	96.39
	Non-Binary	14	14	100.00	14	100.00
Race/ Ethnicity	White	40201	38648	96.14	38666	96.18
	African American	5955	5652	94.91	5646	94.81
	Hispanic	9056	8561	94.53	8703	96.10
	Asian	2858	2756	96.43	2793	97.73
	American Indian	598	568	94.98	565	94.48
	Two or More	3390	3201	94.42	3204	94.51
Limited English Proficiency	No	56765	54535	96.07	54518	96.04
	Yes	5293	4851	91.65	5059	95.58
Disability Status	No	53353	51545	96.61	51743	96.98
	Yes	8705	7841	90.07	7834	89.99
Economically Disadvantaged	No	35656	34217	95.96	34309	96.22
	Yes	26402	25169	95.33	25268	95.70

Table G-4. Test Participation Rates, Grade 6

Group	Category	Enrolled	English Language Arts		Mathematics	
			Number Tested	Percent Tested	Number Tested	Percent Tested
State	All Students	62310	59412	95.35	59570	95.60
Gender	Male	31837	30307	95.19	30382	95.43
	Female	30458	29091	95.51	29174	95.78
	Non-Binary	15	14	93.33	14	93.33
Race/ Ethnicity	White	40620	39006	96.03	39012	96.04
	African American	6125	5703	93.11	5688	92.87
	Hispanic	9051	8527	94.21	8675	95.85
	Asian	2715	2626	96.72	2646	97.46
	American Indian	622	579	93.09	579	93.09
	Two or More	3177	2971	93.52	2970	93.48
Limited English Proficiency	No	57868	55353	95.65	55339	95.63
	Yes	4442	4059	91.38	4231	95.25
Disability Status	No	53805	51870	96.40	52031	96.70
	Yes	8505	7542	88.68	7539	88.64
Economically Disadvantaged	No	36265	34836	96.06	34911	96.27
	Yes	26045	24576	94.36	24659	94.68

Table G-5. Test Participation Rates, Grade 7

Group	Category	Enrolled	English Language Arts		Mathematics	
			Number Tested	Percent Tested	Number Tested	Percent Tested
State	All Students	63623	60413	94.95	60559	95.18
Gender	Male	32597	30933	94.90	31017	95.15
	Female	30995	29452	95.02	29514	95.22
	Non-Binary	31	28	90.32	28	90.32
Race/ Ethnicity	White	41683	39840	95.58	39852	95.61
	African American	6300	5842	92.73	5829	92.52
	Hispanic	9271	8729	94.15	8857	95.53
	Asian	2603	2509	96.39	2533	97.31
	American Indian	643	597	92.85	597	92.85
	Two or More	3123	2896	92.73	2891	92.57
Limited English Proficiency	No	59126	56284	95.19	56262	95.16
	Yes	4497	4129	91.82	4297	95.55
Disability Status	No	55121	52950	96.06	53103	96.34
	Yes	8502	7463	87.78	7456	87.70
Economically Disadvantaged	No	37200	35577	95.64	35639	95.80
	Yes	26423	24836	93.99	24920	94.31

Table G-6. Test Participation Rates, Grade 8

Group	Category	Enrolled	English Language Arts		Mathematics		Science		Social Studies	
			Number Tested	Percent Tested	Number Tested	Percent Tested	Number Tested	Percent Tested	Number Tested	Percent Tested
State	All Students	66060	62249	94.23	62360	94.40	62289	94.29	62261	94.25
Gender	Male	33938	32033	94.39	32094	94.57	32051	94.44	32021	94.35
	Female	32065	30165	94.07	30215	94.23	30187	94.14	30189	94.15
	Non-Binary	57	51	89.47	51	89.47	51	89.47	51	89.47
Race/Ethnicity	White	43159	40964	94.91	40964	94.91	40941	94.86	40943	94.87
	African American	6730	6171	91.69	6153	91.43	6136	91.17	6106	90.73
	Hispanic	9529	8917	93.58	9024	94.70	9005	94.50	9008	94.53
	Asian	2797	2687	96.07	2707	96.78	2705	96.71	2708	96.82
	American Indian	667	606	90.85	606	90.85	604	90.55	600	89.96
	Two or More	3178	2904	91.38	2906	91.44	2898	91.19	2896	91.13
Limited English Proficiency	No	61796	58365	94.45	58341	94.41	58277	94.31	58250	94.26
	Yes	4264	3884	91.09	4019	94.25	4012	94.09	4011	94.07
Disability Status	No	57324	54730	95.47	54849	95.68	54784	95.57	54777	95.56
	Yes	8736	7519	86.07	7511	85.98	7505	85.91	7484	85.67
Economically Disadvantaged	No	38944	37022	95.06	37072	95.19	37033	95.09	37058	95.16
	Yes	27116	25227	93.03	25288	93.26	25256	93.14	25203	92.95

Table G-7. Test Participation Rates, Grade 10

Group	Category	Enrolled	Social Studies	
			Number Tested	Percent Tested
State	All Students	69876	61819	88.47
Gender	Male	35825	31776	88.70
	Female	33982	29983	88.23
	Non-Binary	69	60	86.96
Race/ Ethnicity	White	46815	43110	92.09
	African American	6703	4602	68.66
	Hispanic	9854	8523	86.49
	Asian	2734	2510	91.81
	American Indian	704	526	74.72
	Two or More	3066	2548	83.11
Limited English Proficiency	No	65879	58453	88.73
	Yes	3997	3366	84.21
Disability Status	No	61296	55289	90.20
	Yes	8580	6530	76.11
Economically Disadvantaged	No	43440	39909	91.87
	Yes	26436	21910	82.88

APPENDIX H

CLASSICAL ITEM ANALYSIS RESULTS

List of Tables

Explanation of Data Columns in Tables H-1 through H-17	3
Table H-1. Item Statistics, English Language Arts Grade 3	4
Table H-2. Item Statistics, English Language Arts Grade 4	6
Table H-3. Item Statistics, English Language Arts Grade 5	8
Table H-4. Item Statistics, English Language Arts Grade 6	10
Table H-5. Item Statistics, English Language Arts Grade 7	12
Table H-6. Item Statistics, English Language Arts Grade 8	14
Table H-7. Item Statistics, Mathematics Grade 3	16
Table H-8. Item Statistics, Mathematics Grade 4	18
Table H-9. Item Statistics, Mathematics Grade 5	20
Table H-10. Item Statistics, Mathematics Grade 6	22
Table H-11. Item Statistics, Mathematics Grade 7	24
Table H-12. Item Statistics, Mathematics Grade 8	26
Table H-13. Item Statistics, Science Grade 4	28
Table H-14. Item Statistics, Science Grade 8	30
Table H-15. Item Statistics, Social Studies Grade 4	32
Table H-16. Item Statistics, Social Studies Grade 8	34
Table H-17. Item Statistics, Social Studies Grade 10	36

Explanation of Data Columns in Tables H-1 through H-17

Column Number	Data Description
1	Item Number
2	Item Type
3	Maximum Points
4	Number of Students
5	Item p-value
6	Item-Total Test Correlation
7	Proportion Omit
8	Proportion of Students at Score Point 0
9	Proportion of Students at Score Point 1 or Option 1
10	Proportion of Students at Score Point 2 or Option 2
11	Proportion of Students at Score Point 3 or Option 3
12	Proportion of Students at Score Point 4 or Option 4
13	Item-Total Test Correlation for Score Point 0
14	Item-Total Test Correlation for Score Point 1 or Option 1
15	Item-Total Test Correlation for Score Point 2 or Option 2
16	Item-Total Test Correlation for Score Point 3 or Option 3
17	Item-Total Test Correlation for Score Point 4 or Option 4

Table H-1. Item Statistics, English Language Arts Grade 3

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TDA	4	51304	0.34	0.39	0.01		0.61	0.25	0.03	0.00		-0.14	0.36	0.19	0.02
2	MC	1	58337	0.42	0.27	0.00		0.24	0.41	0.06	0.28		-0.18	0.27	-0.16	-0.04
3	MC	1	58276	0.79	0.43	0.00		0.10	0.79	0.07	0.04		-0.25	0.43	-0.23	-0.20
4	MC	1	58303	0.61	0.24	0.00		0.60	0.17	0.08	0.14		0.24	-0.03	-0.20	-0.14
5	MS	2	58296	0.30	0.12	0.00	0.44	0.52	0.03			-0.12	0.12	0.02		
6	MC	1	58318	0.67	0.23	0.00		0.14	0.67	0.06	0.14		-0.26	0.24	-0.12	0.03
7	MC	1	58306	0.52	0.41	0.00		0.52	0.21	0.17	0.10		0.41	-0.22	-0.16	-0.17
8	MC	1	58325	0.73	0.36	0.00		0.05	0.08	0.13	0.73		-0.27	-0.16	-0.16	0.36
9	TE	1	58203	0.46	0.16	0.00	0.54	0.45				-0.16	0.17			
10	EBSR	2	58345	0.56	0.54	0.00	0.34	0.20	0.46			-0.45	-0.13	0.54		
11	MC	1	58289	0.68	0.48	0.00		0.15	0.09	0.68	0.07		-0.27	-0.24	0.48	-0.21
12	TE	2	58284	0.50	0.50	0.00	0.25	0.50	0.24			-0.41	0.02	0.40		
13	MC	1	58308	0.44	0.36	0.00		0.18	0.22	0.44	0.15		-0.12	-0.19	0.36	-0.15
14	MC	1	58285	0.50	0.30	0.00		0.14	0.50	0.18	0.18		-0.22	0.30	-0.11	-0.08
15	MC	1	58344	0.58	0.33	0.00		0.58	0.27	0.03	0.12		0.33	-0.16	-0.17	-0.19
16	MC	1	58285	0.69	0.46	0.00		0.08	0.12	0.69	0.11		-0.24	-0.26	0.46	-0.20
17	MC	1	58296	0.59	0.25	0.00		0.17	0.13	0.11	0.59		-0.18	-0.04	-0.14	0.25
18	EBSR	2	58356	0.50	0.37	0.00	0.33	0.34	0.33			-0.30	-0.05	0.35		
19	MC	1	58286	0.56	0.38	0.00		0.23	0.05	0.16	0.56		-0.16	-0.21	-0.20	0.38
20	MC	1	58282	0.45	0.28	0.00		0.19	0.23	0.12	0.45		-0.22	0.03	-0.20	0.28
21	MC	1	58326	0.33	0.30	0.00		0.48	0.33	0.10	0.09		-0.12	0.30	-0.08	-0.21
22	MC	1	58264	0.64	0.39	0.00		0.64	0.09	0.08	0.18		0.39	-0.28	-0.21	-0.12
23	MC	1	58265	0.49	0.33	0.00		0.22	0.22	0.49	0.07		-0.09	-0.20	0.33	-0.17
24	MC	1	58247	0.85	0.45	0.00		0.04	0.06	0.85	0.05		-0.23	-0.25	0.45	-0.25
25	MC	1	58275	0.69	0.50	0.00		0.11	0.68	0.15	0.05		-0.25	0.50	-0.26	-0.27
26	MS	2	58293	0.58	0.44	0.00	0.18	0.48	0.34			-0.29	-0.17	0.41		
27	TE	2	57817	0.70	0.48	0.01	0.18	0.24	0.57			-0.35	-0.21	0.48		
28	MC	1	58277	0.58	0.42	0.00		0.58	0.20	0.11	0.11		0.42	-0.26	-0.16	-0.16
29	MC	1	58256	0.36	0.36	0.00		0.22	0.36	0.19	0.23		-0.17	0.36	-0.19	-0.05
30	EBSR	2	58322	0.52	0.52	0.00	0.40	0.15	0.45			-0.43	-0.16	0.54		

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	58257	0.74	0.46	0.00		0.74	0.07	0.07	0.11		0.46	-0.24	-0.25	-0.23
32	MC	1	58285	0.52	0.29	0.00		0.13	0.17	0.18	0.52		-0.23	-0.08	-0.10	0.29
33	EBSR	2	58308	0.48	0.57	0.00	0.46	0.11	0.43			-0.50	-0.11	0.58		
34	MC	1	58233	0.60	0.52	0.00		0.60	0.17	0.15	0.08		0.53	-0.20	-0.28	-0.29
35	MC	1	58216	0.45	0.39	0.00		0.19	0.13	0.44	0.24		-0.06	-0.19	0.39	-0.25
36	MC	1	58214	0.48	0.32	0.00		0.15	0.12	0.24	0.48		-0.12	-0.15	-0.16	0.32
37	MC	1	58249	0.44	0.33	0.00		0.15	0.44	0.08	0.32		-0.11	0.33	-0.22	-0.14
38	MC	1	58260	0.43	0.35	0.00		0.17	0.43	0.23	0.17		-0.18	0.35	-0.06	-0.20

Table H-2. Item Statistics, English Language Arts Grade 4

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TDA	4	51931	0.40	0.53	0.00		0.49	0.27	0.13	0.00		-0.29	0.31	0.39	0.05
2	MC	1	58863	0.78	0.28	0.00		0.04	0.15	0.77	0.03		-0.20	-0.14	0.28	-0.15
3	MC	1	58840	0.47	0.36	0.00		0.29	0.47	0.10	0.15		-0.16	0.36	-0.23	-0.11
4	EBSR	2	58866	0.61	0.54	0.00	0.19	0.39	0.42			-0.43	-0.13	0.48		
5	MC	1	58841	0.82	0.45	0.00		0.04	0.09	0.05	0.82		-0.25	-0.23	-0.27	0.45
6	MC	1	58832	0.60	0.33	0.00		0.26	0.09	0.05	0.59		-0.09	-0.25	-0.23	0.33
7	TE	2	58823	0.50	0.48	0.00	0.21	0.58	0.21			-0.38	0.00	0.38		
8	MC	1	58818	0.53	0.43	0.00		0.14	0.24	0.53	0.09		-0.23	-0.20	0.43	-0.17
9	MC	1	58841	0.59	0.37	0.00		0.23	0.05	0.13	0.59		-0.08	-0.24	-0.28	0.37
10	MC	1	58844	0.66	0.46	0.00		0.18	0.10	0.66	0.06		-0.23	-0.26	0.47	-0.23
11	MC	1	58829	0.43	0.25	0.00		0.13	0.43	0.18	0.26		-0.15	0.26	-0.11	-0.07
12	TE	2	58400	0.70	0.38	0.01	0.05	0.50	0.44			-0.15	-0.29	0.38		
13	MC	1	58824	0.63	0.48	0.00		0.63	0.12	0.12	0.13		0.48	-0.26	-0.21	-0.23
14	MC	1	58824	0.31	0.21	0.00		0.25	0.31	0.23	0.21		-0.12	0.21	-0.06	-0.06
15	MC	1	58857	0.65	0.32	0.00		0.11	0.65	0.11	0.13		-0.18	0.32	-0.15	-0.15
16	MC	1	58804	0.54	0.39	0.00		0.12	0.17	0.54	0.17		-0.29	-0.12	0.39	-0.14
17	EBSR	2	58867	0.36	0.33	0.00	0.52	0.24	0.24			-0.27	0.00	0.32		
18	MC	1	58837	0.56	0.30	0.00		0.11	0.18	0.56	0.14		-0.11	-0.10	0.30	-0.21
19	MC	1	58848	0.69	0.38	0.00		0.03	0.24	0.05	0.69		-0.21	-0.21	-0.23	0.38
20	EBSR	2	58858	0.30	0.34	0.00	0.58	0.23	0.19			-0.30	0.07	0.30		
21	MS	2	58830	0.46	0.34	0.00	0.24	0.60	0.16			-0.29	0.07	0.24		
22	MC	1	58774	0.44	0.21	0.00		0.05	0.44	0.35	0.16		-0.18	0.21	-0.01	-0.16
23	MC	1	58750	0.58	0.25	0.00		0.17	0.11	0.58	0.14		-0.09	-0.17	0.25	-0.11
24	MC	1	58738	0.62	0.26	0.00		0.04	0.17	0.62	0.17		-0.21	-0.19	0.26	-0.03
25	MC	1	58786	0.79	0.50	0.00		0.78	0.10	0.08	0.03		0.50	-0.29	-0.29	-0.22
26	MC	1	58718	0.56	0.45	0.00		0.56	0.13	0.14	0.17		0.45	-0.23	-0.25	-0.16
27	EBSR	2	58837	0.32	0.44	0.00	0.58	0.20	0.22			-0.34	-0.05	0.45		
28	MC	1	58805	0.56	0.36	0.00		0.56	0.19	0.14	0.10		0.36	-0.04	-0.24	-0.26
29	MC	1	58765	0.71	0.50	0.00		0.11	0.07	0.11	0.70		-0.20	-0.27	-0.30	0.50
30	MC	1	58767	0.57	0.41	0.00		0.57	0.21	0.12	0.10		0.41	-0.21	-0.23	-0.14

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	58740	0.69	0.38	0.00		0.18	0.69	0.08	0.05		-0.14	0.38	-0.26	-0.23
32	MC	1	58774	0.61	0.45	0.00		0.12	0.17	0.10	0.61		-0.22	-0.19	-0.24	0.45
33	TE	2	58791	0.50	0.48	0.00	0.21	0.58	0.21			-0.38	0.01	0.37		
34	MC	1	58770	0.64	0.47	0.00		0.06	0.13	0.64	0.16		-0.22	-0.19	0.47	-0.28
35	MC	1	58730	0.68	0.51	0.00		0.09	0.14	0.09	0.68		-0.24	-0.24	-0.30	0.51
36	TE	2	58514	0.70	0.52	0.01	0.1	0.40	0.50			-0.30	-0.33	0.51		
37	MC	1	58750	0.43	0.37	0.00		0.13	0.43	0.31	0.12		-0.16	0.37	-0.18	-0.14
38	EBSR	2	58807	0.38	0.40	0.00	0.45	0.33	0.21			-0.35	0.08	0.34		
39	MC	1	58782	0.48	0.41	0.00		0.09	0.47	0.23	0.20		-0.18	0.41	-0.21	-0.17

Table H-3. Item Statistics, English Language Arts Grade 5

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TDA	4	56546	0.38	0.50	0.00		0.54	0.33	0.08	0.00		-0.38	0.33	0.31	0.04
2	MC	1	59268	0.44	0.19	0.00		0.21	0.16	0.19	0.44		-0.12	-0.02	-0.09	0.19
3	MC	1	59241	0.28	0.21	0.00		0.25	0.28	0.14	0.33		-0.12	0.21	-0.15	0.02
4	MS	2	59244	0.52	0.42	0.00	0.22	0.50	0.27			-0.35	-0.01	0.34		
5	MC	1	59254	0.38	0.27	0.00		0.38	0.11	0.27	0.24		0.27	-0.17	-0.05	-0.13
6	MC	1	59225	0.60	0.41	0.00		0.25	0.07	0.60	0.08		-0.18	-0.27	0.41	-0.20
7	TE	2	59233	0.73	0.35	0.00	0.04	0.44	0.51			-0.24	-0.21	0.31		
8	MC	1	59251	0.27	0.16	0.00		0.13	0.09	0.27	0.50		-0.04	-0.14	0.16	-0.03
9	MC	1	59248	0.47	0.29	0.00		0.15	0.24	0.47	0.14		-0.23	-0.06	0.29	-0.11
10	MC	1	59237	0.78	0.33	0.00		0.77	0.15	0.03	0.04		0.33	-0.25	-0.19	-0.08
11	TE	2	59241	0.69	0.53	0.00	0.11	0.41	0.48			-0.35	-0.28	0.50		
12	MC	1	59250	0.70	0.43	0.00		0.10	0.07	0.70	0.14		-0.19	-0.25	0.44	-0.23
13	TE	2	59092	0.36	0.40	0.00	0.48	0.32	0.20			-0.29	-0.04	0.41		
14	MC	1	59252	0.76	0.29	0.00		0.76	0.03	0.14	0.07		0.29	-0.19	-0.18	-0.12
15	EBSR	2	59268	0.62	0.41	0.00	0.33	0.10	0.57			-0.36	-0.14	0.43		
16	MC	1	59220	0.60	0.20	0.00		0.60	0.23	0.05	0.11		0.20	-0.09	-0.20	-0.06
17	MC	1	59237	0.63	0.42	0.00		0.13	0.17	0.63	0.07		-0.30	-0.12	0.42	-0.23
18	MC	1	59247	0.59	0.37	0.00		0.11	0.18	0.12	0.59		-0.15	-0.19	-0.20	0.37
19	EBSR	2	59267	0.57	0.51	0.00	0.28	0.30	0.42			-0.44	-0.05	0.45		
20	MC	1	59217	0.79	0.38	0.00		0.09	0.79	0.08	0.04		-0.27	0.38	-0.17	-0.15
21	MC	1	59167	0.55	0.29	0.00		0.20	0.11	0.14	0.55		-0.09	-0.11	-0.20	0.29
22	EBSR	2	59232	0.36	0.40	0.00	0.57	0.13	0.30			-0.34	-0.06	0.42		
23	MC	1	59156	0.61	0.38	0.00		0.08	0.21	0.10	0.61		-0.20	-0.16	-0.23	0.38
24	MC	1	59093	0.51	0.51	0.00		0.17	0.13	0.18	0.51		-0.18	-0.20	-0.29	0.51
25	MC	1	59171	0.67	0.47	0.00		0.67	0.15	0.08	0.10		0.47	-0.30	-0.20	-0.20
26	MC	1	58971	0.45	0.21	0.00		0.22	0.45	0.14	0.18		0.04	0.22	-0.14	-0.19
27	EBSR	2	59210	0.30	0.17	0.00	0.59	0.22	0.19			-0.14	0.02	0.16		
28	MC	1	59152	0.78	0.48	0.00		0.07	0.09	0.78	0.06		-0.24	-0.29	0.48	-0.23
29	MC	1	59162	0.79	0.40	0.00		0.05	0.78	0.08	0.08		-0.18	0.40	-0.19	-0.26
30	MC	1	59126	0.36	0.22	0.00		0.23	0.18	0.22	0.36		-0.01	-0.15	-0.10	0.22

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	59124	0.38	0.32	0.00		0.38	0.21	0.24	0.17		0.32	-0.13	-0.11	-0.13
32	MC	1	59157	0.31	0.20	0.00		0.42	0.31	0.18	0.10		0.05	0.20	-0.18	-0.17
33	MC	1	58953	0.47	0.43	0.01		0.14	0.11	0.28	0.47		-0.15	-0.27	-0.17	0.43
34	MC	1	59100	0.50	0.37	0.00		0.12	0.50	0.18	0.21		-0.19	0.37	-0.23	-0.08
35	MC	1	59078	0.63	0.51	0.00		0.63	0.12	0.13	0.12		0.51	-0.25	-0.26	-0.23
36	MC	1	59127	0.49	0.40	0.00		0.16	0.49	0.19	0.15		-0.12	0.40	-0.22	-0.18
37	MC	1	59100	0.59	0.51	0.00		0.14	0.59	0.10	0.16		-0.19	0.51	-0.27	-0.27
38	MC	1	59144	0.60	0.45	0.00		0.18	0.11	0.10	0.60		-0.11	-0.27	-0.29	0.45
39	EBSR	2	59190	0.51	0.60	0.00	0.43	0.12	0.45			-0.50	-0.20	0.63		
40	MC	1	59055	0.58	0.25	0.00		0.11	0.07	0.58	0.24		-0.15	-0.25	0.26	-0.03

Table H-4. Item Statistics, English Language Arts Grade 6

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TDA	4	56568	0.42	0.52	0.00		0.46	0.36	0.12	0.01		-0.40	0.30	0.32	0.14
2	MC	1	59264	0.67	0.25	0.00		0.09	0.67	0.12	0.11		-0.15	0.25	-0.16	-0.06
3	MC	1	59219	0.43	0.35	0.00		0.25	0.19	0.14	0.43		-0.20	-0.05	-0.19	0.35
4	MC	1	59230	0.48	0.20	0.00		0.47	0.13	0.28	0.11		0.21	-0.08	-0.04	-0.18
5	MC	1	59206	0.81	0.41	0.00		0.07	0.81	0.07	0.04		-0.23	0.41	-0.22	-0.21
6	EBSR	2	59276	0.66	0.49	0.00	0.28	0.12	0.60			-0.40	-0.22	0.51		
7	TE	2	59198	0.28	0.29	0.00	0.54	0.36	0.09			-0.24	0.11	0.23		
8	MC	1	59219	0.36	0.25	0.00		0.04	0.14	0.36	0.46		-0.19	-0.20	0.25	-0.03
9	EBSR	2	59260	0.23	0.21	0.00	0.71	0.12	0.16			-0.14	-0.11	0.27		
10	TE	2	59220	0.65	0.46	0.00	0.13	0.45	0.43			-0.37	-0.13	0.38		
11	MC	1	59196	0.54	0.35	0.00		0.04	0.26	0.54	0.16		-0.25	-0.12	0.35	-0.19
12	MC	1	59232	0.42	0.30	0.00		0.12	0.21	0.25	0.42		-0.05	-0.12	-0.18	0.30
13	TE	1	59139	0.40	0.27	0.00	0.60	0.40				-0.26	0.27			
14	MC	1	59235	0.65	0.33	0.00		0.65	0.11	0.17	0.07		0.33	-0.21	-0.15	-0.14
15	EBSR	2	59257	0.31	0.33	0.00	0.65	0.08	0.28			-0.25	-0.18	0.38		
16	MC	1	59223	0.49	0.23	0.00		0.25	0.14	0.49	0.12		-0.03	-0.14	0.23	-0.17
17	MC	1	59228	0.66	0.47	0.00		0.66	0.18	0.08	0.07		0.47	-0.22	-0.31	-0.21
18	MC	1	59198	0.51	0.33	0.00		0.06	0.51	0.27	0.15		-0.15	0.33	-0.20	-0.11
19	TE	2	59212	0.58	0.45	0.00	0.18	0.48	0.35			-0.36	-0.08	0.37		
20	MC	1	59242	0.65	0.44	0.00		0.22	0.65	0.09	0.04		-0.27	0.44	-0.23	-0.17
21	MC	1	59172	0.70	0.36	0.00		0.70	0.19	0.06	0.05		0.36	-0.17	-0.24	-0.18
22	MC	1	59138	0.68	0.45	0.00		0.06	0.09	0.68	0.16		-0.22	-0.31	0.45	-0.18
23	MS	2	59122	0.65	0.50	0.00	0.12	0.46	0.42			-0.39	-0.15	0.42		
24	EBSR	2	59229	0.63	0.55	0.00	0.28	0.18	0.54			-0.50	-0.08	0.51		
25	MC	1	59146	0.71	0.46	0.00		0.08	0.11	0.10	0.71		-0.17	-0.25	-0.28	0.46
26	MC	1	59137	0.44	0.29	0.00		0.30	0.12	0.44	0.14		0.06	-0.30	0.29	-0.21
27	MC	1	59122	0.70	0.52	0.00		0.12	0.10	0.08	0.70		-0.20	-0.27	-0.34	0.52
28	MC	1	59155	0.53	0.28	0.00		0.53	0.20	0.12	0.14		0.29	-0.22	-0.10	-0.05
29	MC	1	59093	0.76	0.52	0.00		0.76	0.07	0.12	0.04		0.52	-0.28	-0.31	-0.23
30	MC	1	59116	0.51	0.38	0.00		0.22	0.16	0.51	0.10		-0.10	-0.17	0.39	-0.28

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	EBSR	2	59208	0.60	0.50	0.00	0.30	0.21	0.49			-0.43	-0.09	0.47		
32	TE	2	59141	0.69	0.57	0.00	0.14	0.34	0.52			-0.44	-0.22	0.51		
33	MC	1	59146	0.68	0.47	0.00		0.18	0.09	0.68	0.05		-0.22	-0.31	0.47	-0.21
34	MC	1	59107	0.54	0.46	0.00		0.54	0.12	0.22	0.11		0.46	-0.30	-0.22	-0.11
35	MC	1	59133	0.82	0.51	0.00		0.08	0.81	0.04	0.06		-0.29	0.51	-0.28	-0.24
36	MC	1	59131	0.45	0.45	0.00		0.16	0.23	0.16	0.45		-0.15	-0.20	-0.22	0.45
37	TE	2	58951	0.47	0.31	0.01	0.21	0.64	0.15			-0.25	0.04	0.24		
38	MC	1	59165	0.48	0.44	0.00		0.07	0.26	0.19	0.47		-0.23	-0.10	-0.28	0.44

Table H-5. Item Statistics, English Language Arts Grade 7

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TDA	4	56284	0.47	0.53	0.00		0.34	0.38	0.17	0.04		-0.37	0.17	0.34	0.22
2	MC	1	60222	0.43	0.34	0.00		0.27	0.17	0.13	0.43		-0.20	-0.06	-0.16	0.34
3	TE	2	60202	0.49	0.40	0.00	0.28	0.46	0.26			-0.32	-0.01	0.34		
4	EBSR	2	60240	0.58	0.48	0.00	0.31	0.22	0.47			-0.42	-0.07	0.45		
5	MC	1	60232	0.80	0.34	0.00		0.80	0.07	0.08	0.05		0.34	-0.22	-0.26	-0.05
6	TE	2	60197	0.45	0.25	0.00	0.26	0.59	0.15			-0.19	0.03	0.19		
7	MC	1	60154	0.51	0.20	0.00		0.25	0.12	0.50	0.12		-0.01	-0.22	0.20	-0.06
8	MC	1	60172	0.56	0.46	0.00		0.56	0.26	0.10	0.08		0.46	-0.26	-0.24	-0.16
9	MC	1	60197	0.68	0.36	0.00		0.17	0.08	0.68	0.06		-0.09	-0.24	0.36	-0.27
10	MC	1	60200	0.33	0.26	0.00		0.13	0.33	0.40	0.14		-0.11	0.26	-0.04	-0.20
11	MC	1	60159	0.41	0.29	0.00		0.41	0.24	0.21	0.13		0.29	-0.12	-0.13	-0.11
12	TE	1	60233	0.83	0.32	0.00	0.17	0.83				-0.32	0.32			
13	EBSR	2	60213	0.47	0.24	0.00	0.41	0.24	0.35			-0.13	-0.19	0.30		
14	EBSR	2	60237	0.55	0.47	0.00	0.27	0.36	0.37			-0.42	0.00	0.38		
15	MC	1	60184	0.76	0.49	0.00		0.11	0.08	0.05	0.76		-0.26	-0.26	-0.26	0.49
16	MS	2	60184	0.63	0.52	0.00	0.10	0.53	0.37			-0.38	-0.19	0.44		
17	MC	1	60186	0.79	0.44	0.00		0.07	0.79	0.08	0.05		-0.29	0.44	-0.20	-0.21
18	EBSR	2	60228	0.66	0.54	0.00	0.31	0.07	0.62			-0.47	-0.21	0.56		
19	MC	1	60193	0.48	0.39	0.00		0.48	0.14	0.34	0.05		0.39	-0.23	-0.21	-0.07
20	MS	2	60131	0.36	0.24	0.00	0.37	0.54	0.09			-0.19	0.07	0.20		
21	MC	1	60129	0.75	0.41	0.00		0.09	0.09	0.75	0.07		-0.24	-0.20	0.41	-0.20
22	MC	1	59990	0.41	0.21	0.00		0.23	0.41	0.12	0.24		0.06	0.21	-0.26	-0.09
23	EBSR	2	60192	0.52	0.53	0.00	0.42	0.13	0.45			-0.42	-0.23	0.58		
24	MC	1	60079	0.75	0.39	0.00		0.13	0.06	0.75	0.05		-0.11	-0.28	0.39	-0.26
25	MS	2	60078	0.61	0.59	0.00	0.25	0.28	0.47			-0.44	-0.22	0.58		
26	MC	1	60075	0.40	0.18	0.00		0.18	0.29	0.40	0.13		-0.15	0.01	0.18	-0.09
27	MC	1	60023	0.64	0.39	0.00		0.19	0.08	0.10	0.63		-0.14	-0.25	-0.21	0.39
28	MC	1	60098	0.78	0.49	0.00		0.09	0.78	0.08	0.04		-0.27	0.49	-0.26	-0.25
29	MC	1	60062	0.41	0.31	0.00		0.23	0.41	0.13	0.22		-0.10	0.31	-0.27	-0.03
30	TE	1	59798	0.34	0.43	0.01	0.65	0.34				-0.42	0.44			

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	60118	0.71	0.55	0.00		0.06	0.16	0.07	0.71		-0.25	-0.30	-0.30	0.55
32	MC	1	60060	0.56	0.38	0.00		0.15	0.21	0.56	0.09		-0.14	-0.20	0.38	-0.20
33	MC	1	59965	0.55	0.37	0.00		0.25	0.10	0.11	0.54		-0.07	-0.20	-0.29	0.37
34	MC	1	60073	0.73	0.43	0.00		0.06	0.72	0.12	0.09		-0.22	0.43	-0.22	-0.23
35	TE	2	59793	0.56	0.44	0.01	0.08	0.71	0.20			-0.30	-0.11	0.36		
36	MC	1	60097	0.34	0.27	0.00		0.34	0.24	0.16	0.25		0.27	-0.15	-0.23	0.06
37	TE	2	60104	0.45	0.23	0.00	0.26	0.58	0.15			-0.20	0.07	0.16		

Table H-6. Item Statistics, English Language Arts Grade 8

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TDA	4	57794	0.55	0.57	0.00		0.25	0.34	0.27	0.08		-0.38	0.03	0.36	0.30
2	MC	1	62075	0.30	0.22	0.00		0.30	0.15	0.51	0.04		0.22	-0.13	-0.06	-0.12
3	MC	1	62008	0.49	0.22	0.00		0.07	0.24	0.20	0.49		-0.22	-0.08	-0.04	0.22
4	MC	1	62006	0.61	0.39	0.00		0.11	0.16	0.12	0.61		-0.21	-0.16	-0.20	0.40
5	MC	1	62000	0.79	0.34	0.00		0.08	0.79	0.03	0.10		-0.16	0.34	-0.23	-0.18
6	MC	1	62036	0.82	0.42	0.00		0.10	0.06	0.82	0.03		-0.22	-0.28	0.42	-0.20
7	MC	1	62028	0.58	0.35	0.00		0.17	0.09	0.17	0.58		-0.20	-0.19	-0.12	0.36
8	MC	1	62002	0.66	0.36	0.00		0.11	0.09	0.66	0.14		-0.16	-0.20	0.36	-0.18
9	EBSR	2	62063	0.50	0.54	0.00	0.46	0.08	0.45			-0.50	-0.07	0.54		
10	TE	1	61925	0.29	0.27	0.00	0.70	0.29				-0.26	0.27			
11	TE	2	61965	0.62	0.50	0.00	0.14	0.48	0.38			-0.37	-0.16	0.43		
12	MC	1	62003	0.50	0.25	0.00		0.17	0.50	0.19	0.14		-0.10	0.25	-0.15	-0.08
13	MS	2	62033	0.53	0.36	0.00	0.24	0.46	0.30			-0.24	-0.10	0.34		
14	MC	1	62024	0.39	0.16	0.00		0.19	0.19	0.39	0.23		0.02	-0.16	0.16	-0.05
15	MC	1	62026	0.29	0.34	0.00		0.32	0.18	0.20	0.29		-0.01	-0.23	-0.15	0.34
16	MC	1	61982	0.35	0.34	0.00		0.35	0.15	0.39	0.11		0.34	-0.17	-0.11	-0.15
17	MC	1	61987	0.47	0.34	0.00		0.26	0.17	0.11	0.47		-0.07	-0.14	-0.26	0.34
18	MS	2	61996	0.63	0.48	0.00	0.10	0.53	0.37			-0.30	-0.24	0.44		
19	EBSR	2	62029	0.63	0.53	0.00	0.17	0.39	0.43			-0.47	-0.07	0.43		
20	MC	1	61966	0.65	0.34	0.00		0.65	0.10	0.11	0.14		0.34	-0.15	-0.23	-0.13
21	MC	1	62008	0.73	0.40	0.00		0.18	0.73	0.07	0.02		-0.25	0.40	-0.24	-0.14
22	EBSR	2	62010	0.55	0.45	0.00	0.31	0.29	0.41			-0.41	0.01	0.38		
23	MC	1	61870	0.87	0.50	0.00		0.03	0.04	0.06	0.86		-0.23	-0.27	-0.31	0.50
24	MC	1	61869	0.39	0.42	0.00		0.39	0.27	0.06	0.27		0.42	-0.19	-0.25	-0.13
25	MS	2	61875	0.62	0.50	0.00	0.14	0.48	0.37			-0.36	-0.18	0.44		
26	MC	1	61892	0.74	0.46	0.00		0.08	0.08	0.74	0.09		-0.19	-0.25	0.46	-0.27
27	MC	1	61732	0.56	0.28	0.01		0.16	0.12	0.16	0.55		-0.13	-0.16	-0.11	0.29
28	MC	1	61892	0.53	0.30	0.00		0.53	0.12	0.28	0.06		0.31	-0.26	-0.01	-0.26
29	EBSR	2	61972	0.61	0.62	0.00	0.25	0.27	0.48			-0.57	-0.04	0.53		
30	MC	1	61885	0.79	0.52	0.00		0.06	0.78	0.08	0.07		-0.22	0.52	-0.32	-0.27

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	61908	0.43	0.21	0.00		0.27	0.19	0.11	0.43		0.06	-0.16	-0.21	0.21
32	MC	1	61875	0.55	0.35	0.00		0.25	0.55	0.11	0.09		-0.07	0.35	-0.24	-0.23
33	MC	1	61933	0.63	0.48	0.00		0.08	0.15	0.63	0.14		-0.21	-0.31	0.48	-0.17
34	EBSR	2	61956	0.26	0.18	0.00	0.68	0.12	0.19			-0.07	-0.22	0.28		
35	MC	1	61887	0.66	0.45	0.00		0.04	0.14	0.17	0.66		-0.21	-0.17	-0.31	0.45
36	MS	2	61865	0.57	0.43	0.00	0.17	0.53	0.30			-0.35	-0.05	0.34		
37	MC	1	61882	0.44	0.40	0.00		0.09	0.21	0.26	0.44		-0.20	-0.11	-0.22	0.40
38	MC	1	61889	0.68	0.47	0.00		0.09	0.68	0.10	0.13		-0.22	0.47	-0.27	-0.22
39	MC	1	61907	0.47	0.43	0.00		0.46	0.15	0.26	0.13		0.43	-0.21	-0.19	-0.16

Table H-7. Item Statistics, Mathematics Grade 3

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	MC	1	58646	0.80	0.49	0.00		0.14	0.79	0.03	0.03		-0.35	0.49	-0.21	-0.22
2	TE	1	58388	0.64	0.58	0.01	0.36	0.64				-0.57	0.58			
3	MC	1	58615	0.45	0.42	0.00		0.18	0.14	0.23	0.45		-0.33	-0.13	-0.08	0.42
4	SA	1	58571	0.54	0.56	0.00	0.46	0.54				-0.56	0.57			
5	MC	1	58602	0.57	0.30	0.00		0.57	0.25	0.11	0.06		0.30	-0.17	-0.10	-0.18
6	TE	1	58588	0.29	0.55	0.00	0.71	0.29				-0.55	0.55			
7	MC	1	58588	0.52	0.48	0.00		0.16	0.25	0.06	0.52		-0.38	-0.11	-0.23	0.48
8	MC	1	58591	0.36	0.42	0.00		0.29	0.28	0.36	0.07		-0.20	-0.22	0.42	-0.05
9	SA	1	58526	0.45	0.59	0.00	0.55	0.45				-0.58	0.59			
10	TE	1	58560	0.66	0.40	0.00	0.34	0.66				-0.40	0.40			
11	MC	1	58571	0.71	0.44	0.00		0.71	0.12	0.09	0.09		0.44	-0.26	-0.21	-0.19
12	SA	1	58461	0.49	0.53	0.00	0.51	0.49				-0.52	0.53			
13	SA	1	58519	0.38	0.57	0.00	0.62	0.38				-0.56	0.57			
14	MC	1	58575	0.61	0.38	0.00		0.61	0.12	0.12	0.15		0.38	-0.17	-0.12	-0.25
15	MC	1	58525	0.45	0.31	0.00		0.26	0.16	0.44	0.14		-0.20	-0.16	0.31	-0.02
16	MC	1	58554	0.69	0.48	0.00		0.15	0.04	0.12	0.69		-0.16	-0.23	-0.36	0.48
17	TE	1	58499	0.36	0.41	0.00	0.64	0.36				-0.40	0.41			
18	MC	1	58541	0.33	0.38	0.00		0.06	0.33	0.35	0.25		-0.15	0.38	-0.12	-0.20
19	SA	1	58507	0.55	0.54	0.00	0.44	0.55				-0.54	0.55			
20	SA	1	58506	0.69	0.59	0.00	0.31	0.69				-0.58	0.59			
21	MC	1	58563	0.79	0.46	0.00		0.11	0.08	0.79	0.02		-0.39	-0.18	0.47	-0.13
22	MC	1	58609	0.71	0.40	0.00		0.10	0.71	0.05	0.15		-0.31	0.40	-0.22	-0.12
23	MC	1	58585	0.34	0.27	0.00		0.34	0.24	0.38	0.03		0.27	-0.28	0.05	-0.18
24	TE	1	58567	0.82	0.42	0.00	0.18	0.82				-0.42	0.42			
25	SA	1	58541	0.34	0.56	0.00	0.65	0.34				-0.55	0.56			
26	MC	1	58559	0.37	0.56	0.00		0.29	0.37	0.21	0.13		-0.29	0.56	-0.04	-0.35
27	SA	1	58520	0.54	0.47	0.00	0.46	0.54				-0.47	0.48			
28	MC	1	58561	0.34	0.24	0.00		0.26	0.25	0.15	0.34		-0.26	0.03	-0.02	0.24
29	MC	1	58568	0.88	0.37	0.00		0.06	0.03	0.88	0.04		-0.24	-0.19	0.37	-0.18
30	SA	1	58548	0.37	0.51	0.00	0.63	0.37				-0.51	0.51			

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	58545	0.68	0.50	0.00		0.07	0.09	0.17	0.68		-0.22	-0.34	-0.22	0.50
32	SA	1	58525	0.60	0.60	0.00	0.40	0.60				-0.59	0.60			
33	MC	1	58545	0.39	0.44	0.00		0.38	0.14	0.09	0.39		-0.21	-0.14	-0.21	0.44
34	MC	1	58566	0.59	0.51	0.00		0.21	0.12	0.08	0.59		-0.32	-0.20	-0.20	0.51
35	SA	1	58530	0.48	0.55	0.00	0.52	0.48				-0.54	0.55			
36	MC	1	58536	0.66	0.46	0.00		0.17	0.66	0.06	0.11		-0.33	0.46	-0.22	-0.13
37	TE	1	58539	0.45	0.61	0.00	0.55	0.45				-0.61	0.61			
38	MC	1	58533	0.64	0.50	0.00		0.64	0.10	0.10	0.16		0.50	-0.24	-0.19	-0.30
39	MC	1	58554	0.69	0.59	0.00		0.19	0.07	0.69	0.04		-0.44	-0.21	0.59	-0.21
40	MC	1	58556	0.61	0.48	0.00		0.06	0.19	0.61	0.15		-0.25	-0.32	0.48	-0.14
41	TE	1	57810	0.46	0.54	0.01	0.53	0.45				-0.51	0.54			
42	SA	1	58505	0.20	0.49	0.00	0.79	0.20				-0.49	0.49			

Table H-8. Item Statistics, Mathematics Grade 4

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TE	1	59090	0.34	0.51	0.00	0.66	0.34				-0.50	0.51			
2	MC	1	59062	0.70	0.44	0.00		0.70	0.10	0.08	0.11		0.44	-0.23	-0.23	-0.22
3	SA	1	59047	0.57	0.58	0.00	0.43	0.57				-0.58	0.58			
4	MC	1	59047	0.79	0.40	0.00		0.03	0.78	0.09	0.10		-0.12	0.40	-0.19	-0.31
5	TE	1	58643	0.72	0.50	0.01	0.28	0.71				-0.49	0.50			
6	MC	1	59060	0.39	0.48	0.00		0.39	0.09	0.35	0.17		0.48	-0.15	-0.13	-0.33
7	TE	1	58812	0.73	0.41	0.01	0.27	0.73				-0.40	0.41			
8	MC	1	58900	0.43	0.38	0.00		0.14	0.29	0.43	0.14		-0.15	-0.15	0.38	-0.19
9	MC	1	59037	0.64	0.47	0.00		0.07	0.17	0.64	0.12		-0.24	-0.27	0.47	-0.20
10	MC	1	59038	0.62	0.43	0.00		0.30	0.62	0.05	0.03		-0.29	0.43	-0.23	-0.15
11	SA	1	59000	0.36	0.63	0.00	0.64	0.36				-0.63	0.63			
12	MC	1	59050	0.45	0.23	0.00		0.45	0.37	0.05	0.13		0.23	-0.18	-0.15	0.01
13	TE	1	59003	0.54	0.63	0.00	0.46	0.53				-0.62	0.63			
14	MC	1	59037	0.42	0.49	0.00		0.13	0.42	0.12	0.32		0.01	0.49	-0.34	-0.29
15	MC	1	59050	0.65	0.35	0.00		0.10	0.10	0.16	0.65		-0.12	-0.12	-0.26	0.35
16	MC	1	58902	0.41	0.45	0.00		0.28	0.14	0.17	0.41		-0.19	-0.23	-0.15	0.45
17	MC	1	59033	0.38	0.31	0.00		0.26	0.22	0.15	0.38		-0.01	-0.22	-0.14	0.31
18	SA	1	59009	0.31	0.51	0.00	0.69	0.31				-0.50	0.51			
19	TE	1	59029	0.34	0.61	0.00	0.66	0.34				-0.60	0.61			
20	MC	1	59030	0.66	0.40	0.00		0.05	0.66	0.12	0.16		-0.17	0.40	-0.22	-0.22
21	SA	1	59044	0.41	0.49	0.00	0.59	0.41				-0.49	0.49			
22	MC	1	59052	0.49	0.52	0.00		0.12	0.49	0.09	0.30		-0.06	0.52	-0.29	-0.34
23	MC	1	59045	0.27	0.31	0.00		0.22	0.28	0.23	0.27		-0.22	-0.11	0.02	0.31
24	MC	1	59060	0.37	0.49	0.00		0.34	0.37	0.24	0.04		-0.16	0.49	-0.30	-0.15
25	TE	1	58907	0.70	0.38	0.00	0.30	0.70				-0.38	0.38			
26	MC	1	59044	0.27	0.38	0.00		0.30	0.40	0.27	0.03		0.02	-0.30	0.38	-0.18
27	SA	1	59015	0.52	0.53	0.00	0.48	0.52				-0.52	0.53			
28	MC	1	59007	0.38	0.48	0.00		0.39	0.13	0.10	0.37		-0.13	-0.21	-0.32	0.48
29	TE	1	59035	0.68	0.38	0.00	0.32	0.68				-0.38	0.38			
30	MC	1	59029	0.36	0.21	0.00		0.42	0.36	0.14	0.08		0.08	0.21	-0.25	-0.20

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	58890	0.53	0.48	0.00		0.12	0.53	0.16	0.19		-0.20	0.48	-0.22	-0.23
32	MC	1	59027	0.33	0.34	0.00		0.12	0.49	0.33	0.06		-0.09	-0.24	0.34	-0.04
33	SA	1	58987	0.43	0.54	0.00	0.57	0.43				-0.53	0.54			
34	MC	1	59003	0.33	0.41	0.00		0.14	0.35	0.18	0.33		-0.23	-0.14	-0.13	0.41
35	TE	1	59004	0.84	0.45	0.00	0.16	0.84				-0.45	0.45			
36	MC	1	59018	0.52	0.21	0.00		0.03	0.08	0.52	0.38		-0.09	-0.18	0.21	-0.08
37	SA	1	59023	0.78	0.46	0.00	0.22	0.78				-0.46	0.46			
38	MC	1	59018	0.38	0.57	0.00		0.38	0.19	0.25	0.18		0.57	-0.09	-0.23	-0.36
39	MC	1	58777	0.31	0.42	0.01		0.10	0.15	0.31	0.43		-0.06	-0.12	0.42	-0.26
40	MC	1	58999	0.57	0.54	0.00		0.57	0.08	0.10	0.24		0.54	-0.28	-0.34	-0.20
41	SA	1	59007	0.61	0.49	0.00	0.39	0.61				-0.49	0.49			
42	MC	1	58988	0.50	0.16	0.00		0.50	0.34	0.11	0.05		0.16	0.05	-0.22	-0.15
43	MC	1	58998	0.37	0.35	0.00		0.18	0.30	0.37	0.15		-0.18	-0.11	0.35	-0.14
44	SA	1	58993	0.53	0.59	0.00	0.47	0.52				-0.59	0.59			
45	MC	1	59019	0.86	0.37	0.00		0.04	0.07	0.04	0.85		-0.20	-0.25	-0.15	0.37
46	TE	1	59022	0.25	0.53	0.00	0.75	0.25				-0.52	0.53			

Table H-9. Item Statistics, Mathematics Grade 5

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	MC	1	59501	0.69	0.49	0.00		0.08	0.17	0.07	0.68		-0.28	-0.26	-0.21	0.49
2	MC	1	59477	0.54	0.23	0.00		0.54	0.11	0.13	0.22		0.23	-0.24	-0.15	0.03
3	TE	1	59485	0.51	0.54	0.00	0.49	0.51				-0.54	0.54			
4	TE	1	59469	0.52	0.53	0.00	0.48	0.52				-0.52	0.53			
5	MC	1	59478	0.34	0.20	0.00		0.21	0.34	0.24	0.21		-0.10	0.20	-0.04	-0.09
6	SA	1	59451	0.29	0.34	0.00	0.71	0.28				-0.34	0.34			
7	MC	1	59460	0.41	0.45	0.00		0.27	0.15	0.17	0.40		-0.10	-0.21	-0.27	0.45
8	SA	1	59354	0.36	0.45	0.00	0.64	0.36				-0.44	0.45			
9	TE	1	59470	0.22	0.53	0.00	0.78	0.22				-0.52	0.53			
10	MC	1	59462	0.30	0.27	0.00		0.12	0.32	0.26	0.30		-0.10	-0.15	-0.05	0.27
11	MC	1	59457	0.43	0.31	0.00		0.27	0.43	0.20	0.10		0.03	0.31	-0.28	-0.18
12	MC	1	59448	0.35	0.21	0.00		0.18	0.31	0.35	0.15		-0.05	-0.05	0.21	-0.16
13	TE	1	59486	0.44	0.44	0.00	0.56	0.44				-0.44	0.44			
14	SA	1	59376	0.52	0.55	0.00	0.47	0.52				-0.54	0.55			
15	MC	1	59481	0.55	0.30	0.00		0.13	0.13	0.18	0.55		-0.10	-0.14	-0.18	0.30
16	SA	1	59318	0.48	0.56	0.00	0.52	0.48				-0.55	0.56			
17	MC	1	59443	0.68	0.33	0.00		0.12	0.68	0.10	0.10		-0.18	0.33	-0.19	-0.12
18	MC	1	59431	0.42	0.43	0.00		0.31	0.19	0.42	0.08		-0.21	-0.20	0.43	-0.12
19	TE	1	59112	0.24	0.48	0.01	0.75	0.24				-0.45	0.48			
20	SA	1	59431	0.43	0.48	0.00	0.57	0.43				-0.48	0.48			
21	MC	1	59455	0.63	0.53	0.00		0.07	0.16	0.15	0.63		-0.23	-0.25	-0.30	0.53
22	SA	1	59394	0.60	0.58	0.00	0.40	0.59				-0.58	0.58			
23	MC	1	59460	0.35	0.43	0.00		0.26	0.17	0.22	0.35		-0.22	-0.30	0.01	0.43
24	SA	1	59440	0.59	0.46	0.00	0.41	0.58				-0.46	0.46			
25	MC	1	59444	0.65	0.27	0.00		0.02	0.31	0.65	0.02		-0.18	-0.17	0.27	-0.16
26	SA	1	59425	0.62	0.57	0.00	0.38	0.62				-0.57	0.57			
27	TE	1	59444	0.46	0.48	0.00	0.54	0.46				-0.47	0.48			
28	MC	1	59435	0.53	0.33	0.00		0.19	0.53	0.13	0.14		-0.21	0.33	-0.13	-0.10
29	SA	1	59364	0.37	0.46	0.00	0.63	0.37				-0.46	0.46			
30	SA	1	59426	0.45	0.60	0.00	0.55	0.45				-0.60	0.60			

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	59332	0.34	0.17	0.00		0.14	0.16	0.34	0.35		-0.18	-0.20	0.17	0.12
32	MC	1	59438	0.48	0.40	0.00		0.48	0.11	0.13	0.28		0.40	-0.24	-0.24	-0.11
33	MC	1	59417	0.63	0.45	0.00		0.13	0.63	0.13	0.11		-0.22	0.45	-0.30	-0.13
34	MC	1	59418	0.45	0.31	0.00		0.45	0.17	0.22	0.16		0.31	-0.07	-0.16	-0.17
35	TE	1	59309	0.62	0.52	0.00	0.38	0.61				-0.51	0.52			
36	SA	1	59417	0.43	0.58	0.00	0.57	0.43				-0.57	0.58			
37	MC	1	59411	0.37	0.56	0.00		0.11	0.38	0.37	0.13		-0.29	-0.22	0.56	-0.21
38	MC	1	59433	0.28	0.14	0.00		0.10	0.28	0.26	0.37		-0.03	0.14	-0.29	0.16
39	SA	1	59201	0.49	0.58	0.01	0.51	0.48				-0.57	0.58			
40	MC	1	59407	0.61	0.58	0.00		0.61	0.18	0.11	0.10		0.58	-0.35	-0.29	-0.19
41	TE	1	59428	0.41	0.35	0.00	0.59	0.41				-0.35	0.35			
42	MC	1	59410	0.48	0.28	0.00		0.23	0.15	0.13	0.48		-0.07	-0.13	-0.19	0.28
43	SA	1	59376	0.28	0.57	0.00	0.72	0.28				-0.56	0.57			
44	MC	1	59424	0.31	0.48	0.00		0.29	0.14	0.26	0.31		-0.31	-0.23	0.00	0.48
45	SA	1	59399	0.62	0.53	0.00	0.38	0.62				-0.52	0.53			
46	TE	1	59438	0.69	0.55	0.00	0.31	0.69				-0.55	0.56			

Table H-10. Item Statistics, Mathematics Grade 6

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	MC	1	59473	0.54	0.41	0.00		0.11	0.25	0.10	0.54		-0.24	-0.11	-0.27	0.42
2	MC	1	59430	0.58	0.49	0.00		0.57	0.14	0.22	0.07		0.49	-0.23	-0.26	-0.22
3	MC	1	59455	0.38	0.49	0.00		0.15	0.37	0.10	0.38		-0.27	-0.21	-0.12	0.49
4	TE	1	59453	0.34	0.56	0.00	0.66	0.34					-0.55	0.56		
5	MC	1	59440	0.33	0.11	0.00		0.13	0.33	0.33	0.21		-0.04	0.12	-0.11	0.03
6	SA	1	59365	0.22	0.56	0.00	0.78	0.22					-0.55	0.56		
7	MC	1	59433	0.53	0.47	0.00		0.23	0.13	0.53	0.10		-0.26	-0.23	0.47	-0.16
8	TE	1	59431	0.79	0.45	0.00	0.21	0.79					-0.44	0.45		
9	MC	1	59432	0.48	0.52	0.00		0.23	0.10	0.18	0.48		-0.28	-0.23	-0.19	0.52
10	MC	1	59419	0.35	0.34	0.00		0.35	0.31	0.20	0.15		0.34	-0.16	-0.13	-0.11
11	MC	1	59427	0.31	0.29	0.00		0.19	0.34	0.31	0.15		-0.09	-0.23	0.29	0.02
12	MC	1	59440	0.76	0.40	0.00		0.76	0.10	0.07	0.07		0.41	-0.27	-0.26	-0.10
13	MC	1	59414	0.61	0.55	0.00		0.61	0.18	0.11	0.09		0.55	-0.25	-0.27	-0.29
14	MC	1	59421	0.52	0.51	0.00		0.23	0.14	0.11	0.52		-0.30	-0.21	-0.18	0.51
15	TE	1	59382	0.39	0.65	0.00	0.61	0.39					-0.65	0.65		
16	SA	1	59355	0.58	0.60	0.00	0.42	0.58					-0.59	0.60		
17	TE	1	59020	0.80	0.34	0.01	0.20	0.79					-0.33	0.35		
18	SA	1	59309	0.63	0.52	0.00	0.37	0.63					-0.51	0.52		
19	MC	1	59410	0.65	0.37	0.00		0.07	0.64	0.10	0.19		-0.11	0.37	-0.28	-0.16
20	MC	1	59403	0.39	0.29	0.00		0.20	0.39	0.18	0.22		-0.10	0.29	-0.16	-0.10
21	MC	1	59427	0.36	0.49	0.00		0.17	0.36	0.31	0.16		-0.30	0.49	-0.07	-0.25
22	SA	1	59372	0.60	0.57	0.00	0.40	0.60					-0.57	0.58		
23	TE	1	59394	0.54	0.52	0.00	0.46	0.54					-0.52	0.53		
24	MC	1	59257	0.41	0.19	0.00		0.15	0.16	0.27	0.41		-0.19	-0.27	0.17	0.20
25	MC	1	59258	0.37	0.24	0.00		0.16	0.27	0.37	0.20		-0.08	-0.15	0.25	-0.04
26	SA	1	59078	0.22	0.56	0.01	0.78	0.22					-0.54	0.56		
27	TE	1	59390	0.64	0.51	0.00	0.36	0.64					-0.51	0.51		
28	MC	1	59319	0.42	0.31	0.00		0.42	0.18	0.20	0.20		0.31	-0.08	-0.17	-0.12
29	SA	1	59309	0.38	0.57	0.00	0.62	0.38					-0.56	0.57		
30	SA	1	59275	0.36	0.22	0.00	0.64	0.36					-0.21	0.22		

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	59384	0.57	0.35	0.00		0.17	0.20	0.57	0.05		-0.21	-0.17	0.35	-0.11
32	MC	1	59322	0.34	0.42	0.00		0.17	0.22	0.34	0.27		-0.21	-0.23	0.42	-0.05
33	MC	1	59388	0.56	0.50	0.00		0.56	0.32	0.05	0.07		0.50	-0.35	-0.25	-0.10
34	SA	1	59155	0.25	0.54	0.01	0.74	0.25				-0.52	0.54			
35	MC	1	59318	0.69	0.48	0.00		0.05	0.69	0.25	0.02		-0.24	0.48	-0.34	-0.16
36	TE	1	58806	0.37	0.52	0.01	0.62	0.37				-0.49	0.52			
37	MC	1	59345	0.63	0.39	0.00		0.15	0.12	0.63	0.10		-0.17	-0.27	0.39	-0.14
38	SA	1	59275	0.67	0.52	0.00	0.33	0.67				-0.52	0.53			
39	MC	1	59353	0.34	0.11	0.00		0.12	0.34	0.21	0.33		-0.11	0.11	-0.16	0.11
40	TE	1	58647	0.51	0.57	0.01	0.48	0.50				-0.55	0.58			
41	SA	1	59251	0.42	0.54	0.00	0.58	0.41				-0.54	0.54			
42	MC	1	59342	0.40	0.26	0.00		0.12	0.16	0.31	0.40		-0.16	-0.18	-0.01	0.26
43	MC	1	59361	0.34	0.21	0.00		0.21	0.29	0.34	0.15		-0.20	-0.05	0.21	0.02
44	TE	1	59290	0.38	0.58	0.00	0.62	0.38				-0.57	0.58			
45	MC	1	59310	0.46	0.41	0.00		0.13	0.20	0.20	0.46		-0.18	-0.18	-0.18	0.41
46	TE	1	59288	0.66	0.57	0.00	0.34	0.66				-0.56	0.57			

Table H-11. Item Statistics, Mathematics Grade 7

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	SA	1	60400	0.77	0.45	0.00	0.23	0.77				-0.45	0.45			
2	MC	1	60420	0.53	0.52	0.00		0.06	0.29	0.53	0.12		-0.18	-0.25	0.52	-0.30
3	SA	1	60303	0.26	0.48	0.00	0.73	0.26				-0.48	0.48			
4	MC	1	60428	0.49	0.48	0.00		0.05	0.39	0.07	0.49		-0.18	-0.33	-0.15	0.48
5	TE	1	60355	0.31	0.49	0.00	0.69	0.31				-0.48	0.49			
6	MC	1	60398	0.42	0.28	0.00		0.20	0.25	0.42	0.13		-0.08	-0.17	0.28	-0.10
7	SA	1	60349	0.26	0.49	0.00	0.74	0.26				-0.48	0.49			
8	MC	1	60419	0.62	0.50	0.00		0.62	0.16	0.13	0.08		0.50	-0.29	-0.28	-0.15
9	SA	1	60179	0.14	0.44	0.00	0.86	0.14				-0.42	0.44			
10	MC	1	60369	0.53	0.34	0.00		0.13	0.53	0.12	0.21		-0.17	0.34	-0.29	-0.04
11	MC	1	60413	0.27	0.35	0.00		0.27	0.25	0.15	0.33		0.35	-0.17	-0.01	-0.17
12	MC	1	60341	0.43	0.24	0.00		0.14	0.20	0.43	0.22		-0.16	-0.14	0.24	-0.01
13	TE	1	60149	0.36	0.45	0.00	0.63	0.36				-0.44	0.45			
14	MC	1	60301	0.65	0.53	0.00		0.65	0.15	0.13	0.07		0.53	-0.32	-0.27	-0.18
15	MC	1	60295	0.48	0.38	0.00		0.18	0.31	0.48	0.03		-0.29	-0.10	0.38	-0.17
16	TE	1	60312	0.47	0.24	0.00	0.53	0.47				-0.24	0.24			
17	SA	1	60094	0.27	0.44	0.01	0.72	0.27				-0.42	0.44			
18	TE	1	60295	0.72	0.23	0.00	0.28	0.72				-0.23	0.23			
19	MC	1	60202	0.48	0.53	0.00		0.14	0.18	0.21	0.47		-0.32	-0.32	-0.07	0.53
20	SA	1	60168	0.73	0.53	0.00	0.27	0.72				-0.53	0.53			
21	MC	1	60169	0.35	0.30	0.00		0.13	0.26	0.25	0.35		-0.08	-0.23	-0.03	0.30
22	SA	1	60156	0.53	0.55	0.00	0.47	0.52				-0.54	0.56			
23	MC	1	60206	0.34	0.24	0.00		0.20	0.34	0.29	0.16		-0.05	0.24	-0.13	-0.09
24	TE	1	60006	0.33	0.63	0.01	0.66	0.33				-0.61	0.63			
25	MC	1	60271	0.47	0.56	0.00		0.47	0.14	0.11	0.28		0.56	-0.22	-0.13	-0.35
26	MC	1	60244	0.44	0.26	0.00		0.05	0.45	0.43	0.06		-0.14	-0.13	0.26	-0.14
27	MC	1	60228	0.52	0.46	0.00		0.52	0.17	0.25	0.06		0.46	-0.29	-0.17	-0.18
28	TE	1	60261	0.46	0.26	0.00	0.53	0.46				-0.26	0.26			
29	MC	1	60253	0.57	0.23	0.00		0.06	0.20	0.17	0.57		-0.08	-0.14	-0.09	0.23
30	MC	1	60255	0.69	0.49	0.00		0.12	0.69	0.13	0.06		-0.18	0.50	-0.35	-0.22

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	60168	0.31	0.15	0.00		0.14	0.21	0.33	0.31		0.10	-0.13	-0.10	0.16
32	MC	1	60121	0.43	0.34	0.01		0.43	0.26	0.24	0.06		0.34	-0.17	-0.12	-0.15
33	SA	1	59918	0.51	0.54	0.01	0.49	0.50				-0.53	0.54			
34	MC	1	60173	0.36	0.34	0.00		0.12	0.36	0.26	0.25		-0.04	0.34	-0.13	-0.21
35	MC	1	60135	0.38	0.13	0.00		0.20	0.38	0.29	0.12		0.09	0.14	-0.16	-0.07
36	MC	1	60208	0.48	0.53	0.00		0.48	0.22	0.19	0.10		0.53	-0.29	-0.13	-0.28
37	SA	1	59999	0.37	0.67	0.01	0.62	0.37				-0.65	0.67			
38	TE	1	60171	0.65	0.54	0.00	0.35	0.65				-0.53	0.54			
39	SA	1	59979	0.26	0.51	0.01	0.73	0.26				-0.49	0.51			
40	MC	1	60203	0.50	0.18	0.00		0.19	0.20	0.50	0.10		0.12	-0.21	0.19	-0.17
41	SA	1	60069	0.28	0.54	0.01	0.72	0.28				-0.52	0.54			
42	TE	1	60008	0.45	0.27	0.01	0.55	0.44				-0.25	0.27			
43	MC	1	60105	0.26	0.20	0.01		0.12	0.32	0.29	0.26		-0.07	-0.08	-0.05	0.20
44	SA	1	60031	0.68	0.61	0.01	0.32	0.67				-0.60	0.61			
45	MC	1	60126	0.59	0.50	0.01		0.06	0.15	0.59	0.20		-0.19	-0.25	0.51	-0.28
46	MC	1	60166	0.51	0.41	0.00		0.19	0.20	0.51	0.09		-0.18	-0.23	0.41	-0.12

Table H-12. Item Statistics, Mathematics Grade 8

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	SA	1	62193	0.19	0.47	0.00	0.81	0.19				-0.47	0.47			
2	MC	1	62231	0.61	0.48	0.00		0.16	0.61	0.08	0.14		-0.20	0.48	-0.23	-0.27
3	TE	1	62228	0.58	0.46	0.00	0.42	0.58				-0.46	0.46			
4	MC	1	62202	0.45	0.24	0.00		0.45	0.16	0.25	0.14		0.24	-0.04	-0.16	-0.09
5	MC	1	62218	0.30	0.47	0.00		0.30	0.47	0.11	0.12		0.47	-0.29	0.02	-0.23
6	SA	1	62103	0.38	0.50	0.00	0.62	0.38				-0.49	0.50			
7	MC	1	62207	0.33	0.32	0.00		0.33	0.16	0.41	0.10		0.32	-0.12	-0.14	-0.12
8	MC	1	62202	0.34	0.23	0.00		0.25	0.11	0.30	0.34		-0.32	-0.16	0.18	0.23
9	SA	1	62157	0.44	0.63	0.00	0.56	0.44				-0.63	0.63			
10	MC	1	62230	0.52	0.48	0.00		0.52	0.14	0.21	0.13		0.48	-0.11	-0.30	-0.24
11	MC	1	62198	0.30	0.42	0.00		0.35	0.18	0.17	0.30		-0.08	-0.19	-0.22	0.42
12	MC	1	62213	0.46	0.48	0.00		0.12	0.27	0.46	0.14		-0.27	-0.11	0.48	-0.29
13	MC	1	62219	0.47	0.58	0.00		0.31	0.14	0.47	0.08		-0.47	-0.18	0.58	-0.02
14	SA	1	61898	0.36	0.51	0.01	0.64	0.36				-0.50	0.51			
15	MC	1	62091	0.61	0.37	0.00		0.09	0.61	0.10	0.20		-0.11	0.37	-0.23	-0.19
16	MC	1	62137	0.75	0.37	0.00		0.75	0.12	0.08	0.05		0.37	-0.16	-0.25	-0.19
17	MC	1	62072	0.50	0.33	0.00		0.50	0.17	0.28	0.05		0.33	-0.23	-0.13	-0.08
18	MC	1	62126	0.42	0.27	0.00		0.25	0.15	0.18	0.42		-0.05	-0.10	-0.20	0.27
19	SA	1	61898	0.32	0.55	0.01	0.68	0.32				-0.53	0.55			
20	TE	1	62126	0.18	0.38	0.00	0.82	0.18				-0.38	0.38			
21	MC	1	61998	0.52	0.33	0.00		0.10	0.51	0.26	0.11		-0.14	0.33	-0.14	-0.17
22	MC	1	62044	0.45	0.27	0.00		0.41	0.07	0.45	0.07		-0.08	-0.25	0.27	-0.12
23	MC	1	62026	0.29	0.48	0.00		0.20	0.15	0.36	0.28		-0.23	-0.25	-0.06	0.48
24	MC	1	62046	0.40	0.25	0.00		0.25	0.20	0.40	0.15		0.02	-0.20	0.25	-0.13
25	MC	1	62076	0.69	0.44	0.00		0.12	0.69	0.13	0.06		-0.24	0.44	-0.23	-0.20
26	MC	1	62081	0.53	0.37	0.00		0.06	0.53	0.23	0.19		-0.15	0.37	-0.23	-0.13
27	SA	1	61996	0.70	0.53	0.00	0.30	0.69				-0.52	0.53			
28	TE	1	62033	0.54	0.30	0.00	0.46	0.54				-0.29	0.30			
29	SA	1	61891	0.33	0.69	0.01	0.67	0.33				-0.68	0.69			
30	SA	1	61904	0.29	0.63	0.01	0.71	0.29				-0.62	0.63			

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	62085	0.57	0.32	0.00		0.13	0.56	0.25	0.05		-0.19	0.32	-0.12	-0.19
32	TE	1	61948	0.34	0.52	0.00	0.65	0.34				-0.51	0.53			
33	MC	1	62015	0.44	0.26	0.00		0.22	0.16	0.43	0.18		-0.08	-0.17	0.26	-0.09
34	MC	1	61992	0.49	0.53	0.00		0.15	0.21	0.16	0.48		-0.19	-0.28	-0.23	0.53
35	TE	1	61974	0.19	0.62	0.00	0.80	0.19				-0.60	0.62			
36	SA	1	61438	0.24	0.59	0.01	0.75	0.24				-0.56	0.59			
37	MC	1	61982	0.47	0.54	0.00		0.47	0.20	0.19	0.14		0.54	-0.26	-0.26	-0.17
38	MC	1	62024	0.42	0.29	0.00		0.10	0.22	0.42	0.26		-0.15	-0.27	0.29	0.04
39	MC	1	62043	0.42	0.42	0.00		0.42	0.19	0.23	0.16		0.42	-0.18	-0.20	-0.13
40	TE	1	61842	0.46	0.38	0.01	0.53	0.46				-0.37	0.38			
41	MC	1	62054	0.52	0.37	0.00		0.51	0.18	0.24	0.06		0.37	-0.29	-0.07	-0.16
42	TE	1	61563	0.22	0.50	0.01	0.77	0.22				-0.48	0.50			
43	MC	1	62029	0.69	0.43	0.00		0.08	0.11	0.69	0.12		-0.20	-0.23	0.43	-0.22
44	MC	1	61976	0.37	0.22	0.00		0.14	0.27	0.22	0.37		-0.06	0.02	-0.23	0.23
45	MC	1	62006	0.66	0.44	0.00		0.07	0.66	0.18	0.09		-0.23	0.44	-0.21	-0.23
46	TE	1	61869	0.70	0.47	0.01	0.30	0.69				-0.47	0.48			

Table H-13. Item Statistics, Science Grade 4

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TE	1	58965	0.78	0.38	0.00	0.22	0.77				-0.37	0.38			
2	MC	1	59027	0.60	0.45	0.00		0.19	0.17	0.60	0.05		-0.24	-0.24	0.45	-0.18
3	TE	1	58900	0.32	0.57	0.00	0.68	0.32				-0.56	0.57			
4	TE	1	58967	0.85	0.35	0.00	0.14	0.85				-0.35	0.35			
5	MC	1	59018	0.70	0.53	0.00		0.11	0.07	0.13	0.70		-0.25	-0.26	-0.30	0.53
6	TE	1	58990	0.57	0.46	0.00	0.43	0.57				-0.46	0.46			
7	TE	1	58995	0.54	0.31	0.00	0.46	0.54				-0.30	0.31			
8	TE	1	58444	0.55	0.41	0.01	0.44	0.54				-0.40	0.41			
9	TE	1	59008	0.46	0.59	0.00	0.53	0.46				-0.59	0.59			
10	MC	1	59019	0.67	0.40	0.00		0.16	0.67	0.08	0.08		-0.22	0.40	-0.24	-0.14
11	MC	1	59001	0.67	0.51	0.00		0.11	0.08	0.15	0.67		-0.20	-0.27	-0.30	0.51
12	TE	1	58929	0.47	0.45	0.00	0.52	0.47				-0.45	0.45			
13	TE	1	58949	0.43	0.60	0.00	0.57	0.43				-0.59	0.60			
14	TE	1	58976	0.42	0.40	0.00	0.58	0.42				-0.40	0.40			
15	MC	1	58997	0.55	0.49	0.00		0.11	0.11	0.23	0.55		-0.24	-0.27	-0.21	0.49
16	TE	1	59013	0.56	0.49	0.00	0.44	0.56				-0.49	0.49			
17	EBSR	1	59015	0.29	0.35	0.00	0.71	0.29				-0.34	0.35			
18	TE	1	58975	0.19	0.32	0.00	0.81	0.19				-0.31	0.32			
19	EBSR	1	58993	0.69	0.39	0.00	0.31	0.69				-0.39	0.39			
20	MC	1	58952	0.53	0.33	0.00		0.53	0.21	0.15	0.11		0.33	-0.14	-0.22	-0.08
21	TE	1	58981	0.53	0.51	0.00	0.47	0.53				-0.50	0.51			
22	TE	1	59005	0.73	0.33	0.00	0.27	0.72				-0.33	0.34			
23	MC	1	58949	0.50	0.56	0.00		0.13	0.19	0.18	0.49		-0.26	-0.29	-0.21	0.56
24	MC	1	58980	0.55	0.44	0.00		0.08	0.19	0.17	0.55		-0.19	-0.20	-0.23	0.44
25	TE	1	59000	0.21	0.38	0.00	0.79	0.21				-0.38	0.38			
26	MC	1	59010	0.65	0.35	0.00		0.13	0.10	0.11	0.65		-0.26	-0.20	-0.06	0.36
27	MC	1	58999	0.82	0.32	0.00		0.04	0.82	0.06	0.08		-0.22	0.32	-0.21	-0.12
28	MC	1	58948	0.40	0.37	0.00		0.14	0.28	0.18	0.40		-0.13	-0.09	-0.24	0.37
29	TE	1	58853	0.39	0.33	0.00	0.61	0.39				-0.32	0.33			
30	TE	1	58964	0.36	0.15	0.00	0.64	0.36				-0.15	0.16			

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	TE	1	58923	0.45	0.62	0.00	0.55	0.45				-0.61	0.62			
32	TE	1	58828	0.35	0.38	0.00	0.65	0.35				-0.38	0.39			
33	MC	1	58980	0.65	0.43	0.00		0.65	0.18	0.09	0.08		0.43	-0.28	-0.19	-0.15
34	TE	1	58981	0.68	0.47	0.00	0.32	0.68				-0.47	0.47			
35	TE	1	58993	0.89	0.37	0.00	0.11	0.89				-0.37	0.37			
36	TE	1	58933	0.56	0.36	0.00	0.43	0.56				-0.35	0.36			
37	TE	1	58939	0.36	0.30	0.00	0.64	0.36				-0.30	0.30			
38	EBSR	1	58954	0.38	0.43	0.00	0.62	0.38				-0.43	0.43			
39	TE	1	58969	0.51	0.39	0.00	0.48	0.51				-0.39	0.39			
40	MC	1	58956	0.53	0.54	0.00		0.14	0.20	0.13	0.53		-0.24	-0.24	-0.26	0.54

Table H-14. Item Statistics, Science Grade 8

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TE	1	61990	0.73	0.47	0.00	0.27	0.73				-0.46	0.47			
2	MS	1	62105	0.42	0.46	0.00	0.57	0.42				-0.45	0.46			
3	TE	1	62074	0.73	0.32	0.00	0.27	0.73				-0.32	0.32			
4	TE	1	61837	0.65	0.34	0.01	0.35	0.65				-0.33	0.34			
5	TE	1	61931	0.42	0.39	0.00	0.58	0.42				-0.39	0.40			
6	TE	1	62150	0.22	0.23	0.00	0.78	0.22				-0.22	0.23			
7	TE	1	62086	0.70	0.33	0.00	0.30	0.70				-0.33	0.34			
8	MC	1	62059	0.40	0.23	0.00		0.09	0.38	0.12	0.40		-0.16	0.01	-0.22	0.23
9	TE	1	62048	0.49	0.41	0.00	0.51	0.48				-0.40	0.41			
10	MC	1	62099	0.61	0.39	0.00		0.10	0.61	0.24	0.06		-0.07	0.39	-0.30	-0.17
11	MC	1	62100	0.53	0.41	0.00		0.53	0.17	0.19	0.10		0.41	-0.22	-0.21	-0.12
12	MC	1	62031	0.46	0.22	0.00		0.21	0.45	0.22	0.11		-0.10	0.22	-0.17	0.00
13	TE	1	62114	0.40	0.50	0.00	0.60	0.40				-0.50	0.50			
14	TE	1	62099	0.32	0.39	0.00	0.68	0.32				-0.39	0.39			
15	MC	1	62090	0.47	0.36	0.00		0.47	0.25	0.15	0.13		0.37	-0.18	-0.22	-0.06
16	MC	1	62072	0.48	0.15	0.00		0.28	0.08	0.48	0.16		0.01	-0.09	0.15	-0.15
17	TE	1	62042	0.58	0.38	0.00	0.42	0.58				-0.38	0.38			
18	TE	1	62064	0.48	0.33	0.00	0.52	0.48				-0.32	0.33			
19	MC	1	61951	0.52	0.42	0.00		0.18	0.10	0.52	0.19		-0.10	-0.22	0.42	-0.25
20	TE	1	62031	0.79	0.37	0.00	0.21	0.79				-0.36	0.37			
21	TE	1	62057	0.51	0.45	0.00	0.48	0.51				-0.45	0.45			
22	EBSR	1	62071	0.44	0.54	0.00	0.56	0.44				-0.54	0.54			
23	TE	1	62037	0.68	0.45	0.00	0.32	0.68				-0.45	0.45			
24	TE	1	62027	0.61	0.46	0.00	0.39	0.61				-0.46	0.46			
25	TE	1	62052	0.51	0.44	0.00	0.49	0.51				-0.44	0.44			
26	MC	1	62015	0.70	0.39	0.00		0.09	0.12	0.70	0.09		-0.19	-0.19	0.39	-0.21
27	TE	1	61920	0.53	0.37	0.00	0.46	0.53				-0.37	0.37			
28	MC	1	61948	0.47	0.41	0.00		0.12	0.47	0.20	0.20		-0.07	0.41	-0.31	-0.13
29	TE	1	61966	0.74	0.38	0.00	0.26	0.74				-0.38	0.38			
30	TE	1	61822	0.41	0.37	0.00	0.59	0.40				-0.37	0.38			

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	61962	0.49	0.45	0.00		0.16	0.21	0.49	0.14		-0.27	-0.28	0.45	-0.03
32	TE	1	61964	0.60	0.42	0.00	0.40	0.60				-0.42	0.42			
33	EBSR	1	62009	0.42	0.49	0.00	0.58	0.42				-0.49	0.49			
34	TE	1	61989	0.30	0.29	0.00	0.70	0.30				-0.29	0.30			
35	MC	1	61988	0.55	0.50	0.00		0.08	0.24	0.13	0.55		-0.16	-0.23	-0.31	0.50
36	TE	1	61980	0.38	0.41	0.00	0.62	0.38				-0.41	0.41			
37	MC	1	61957	0.60	0.56	0.00		0.60	0.12	0.21	0.06		0.56	-0.24	-0.36	-0.20
38	TE	1	61960	0.26	0.33	0.00	0.74	0.26				-0.33	0.33			
39	TE	1	61949	0.64	0.42	0.00	0.36	0.64				-0.42	0.42			
40	EBSR	1	61971	0.43	0.34	0.00	0.57	0.43				-0.34	0.34			

Table H-15. Item Statistics, Social Studies Grade 4

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TE	1	59058	0.55	0.47	0.00	0.45	0.54				-0.47	0.48			
2	MC	1	59042	0.72	0.41	0.00		0.06	0.16	0.06	0.72		-0.27	-0.19	-0.22	0.41
3	TE	1	59026	0.25	0.21	0.00	0.74	0.25				-0.21	0.21			
4	MC	1	59014	0.72	0.47	0.00		0.07	0.11	0.11	0.72		-0.22	-0.32	-0.18	0.47
5	MC	1	59021	0.72	0.44	0.00		0.72	0.05	0.12	0.11		0.45	-0.26	-0.20	-0.25
6	MS	1	59017	0.37	0.51	0.00	0.63	0.37				-0.50	0.51			
7	TE	1	59026	0.37	0.40	0.00	0.63	0.37				-0.40	0.40			
8	MC	1	59027	0.54	0.37	0.00		0.19	0.14	0.54	0.13		-0.21	-0.16	0.37	-0.14
9	MC	1	59019	0.73	0.49	0.00		0.06	0.07	0.14	0.73		-0.25	-0.26	-0.26	0.49
10	MC	1	59013	0.51	0.29	0.00		0.12	0.25	0.12	0.51		-0.28	0.04	-0.23	0.29
11	MC	1	59014	0.65	0.38	0.00		0.10	0.12	0.65	0.13		-0.19	-0.21	0.38	-0.15
12	MC	1	59022	0.80	0.49	0.00		0.06	0.07	0.80	0.07		-0.26	-0.30	0.49	-0.22
13	MS	1	59042	0.39	0.50	0.00	0.61	0.39				-0.50	0.50			
14	MC	1	59035	0.48	0.43	0.00		0.19	0.22	0.11	0.48		-0.16	-0.15	-0.28	0.43
15	MC	1	59017	0.55	0.44	0.00		0.55	0.24	0.11	0.11		0.44	-0.17	-0.24	-0.22
16	MC	1	59007	0.57	0.39	0.00		0.13	0.18	0.57	0.11		-0.24	-0.14	0.39	-0.18
17	MC	1	59019	0.34	0.29	0.00		0.34	0.28	0.19	0.19		0.29	-0.15	-0.17	0.00
18	TE	1	58927	0.43	0.34	0.00	0.57	0.43				-0.33	0.34			
19	MC	1	59027	0.65	0.41	0.00		0.16	0.65	0.12	0.08		-0.15	0.41	-0.31	-0.16
20	MC	1	59018	0.39	0.39	0.00		0.19	0.23	0.19	0.39		-0.08	-0.20	-0.19	0.39
21	TE	1	59020	0.54	0.47	0.00	0.46	0.54				-0.47	0.47			
22	MC	1	58988	0.74	0.48	0.00		0.06	0.74	0.10	0.10		-0.21	0.48	-0.30	-0.24
23	MC	1	58992	0.71	0.53	0.00		0.05	0.71	0.13	0.11		-0.19	0.53	-0.31	-0.29
24	TE	1	58996	0.52	0.53	0.00	0.48	0.52				-0.53	0.53			
25	MC	1	58985	0.61	0.50	0.00		0.13	0.61	0.14	0.13		-0.18	0.50	-0.33	-0.22
26	TE	1	59004	0.68	0.60	0.00	0.31	0.68				-0.60	0.60			
27	TE	1	58981	0.50	0.37	0.00	0.50	0.50				-0.37	0.37			
28	MC	1	59013	0.76	0.49	0.00		0.76	0.06	0.14	0.05		0.49	-0.26	-0.28	-0.24
29	TE	1	59009	0.55	0.54	0.00	0.45	0.55				-0.54	0.54			
30	TE	1	58984	0.54	0.36	0.00	0.46	0.54				-0.36	0.36			

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	TE	1	58990	0.44	0.39	0.00	0.56	0.44				-0.39	0.39			
32	MC	1	58962	0.62	0.42	0.00		0.07	0.19	0.62	0.12		-0.21	-0.18	0.42	-0.24
33	TE	1	59008	0.31	0.35	0.00	0.69	0.31				-0.35	0.35			
34	MC	1	59012	0.76	0.48	0.00		0.76	0.09	0.11	0.04		0.48	-0.29	-0.26	-0.20
35	MC	1	58969	0.45	0.33	0.00		0.23	0.18	0.45	0.14		-0.02	-0.23	0.33	-0.19
36	MC	1	58940	0.48	0.35	0.00		0.10	0.47	0.28	0.15		-0.15	0.35	-0.14	-0.19
37	TE	1	59002	0.38	0.41	0.00	0.62	0.38				-0.40	0.41			
38	MC	1	58971	0.38	0.39	0.00		0.38	0.19	0.22	0.21		0.39	-0.18	-0.17	-0.11
39	MC	1	58983	0.68	0.41	0.00		0.12	0.68	0.10	0.10		-0.13	0.41	-0.28	-0.21
40	MC	1	59004	0.57	0.39	0.00		0.15	0.57	0.20	0.07		-0.18	0.39	-0.18	-0.21

Table H-16. Item Statistics, Social Studies Grade 8

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TE	1	62174	0.69	0.36	0.00	0.31	0.69				-0.36	0.36			
2	MC	1	62133	0.44	0.35	0.00		0.17	0.44	0.14	0.25		-0.08	0.35	-0.24	-0.14
3	TE	1	62108	0.50	0.48	0.00	0.50	0.50				-0.48	0.48			
4	MS	1	62118	0.25	0.25	0.00	0.75	0.25				-0.25	0.25			
5	MC	1	62116	0.64	0.47	0.00		0.08	0.64	0.13	0.14		-0.18	0.48	-0.28	-0.24
6	MC	1	62079	0.79	0.49	0.00		0.05	0.79	0.10	0.06		-0.24	0.49	-0.32	-0.22
7	MC	1	62106	0.42	0.29	0.00		0.42	0.17	0.17	0.24		0.29	-0.17	-0.19	-0.02
8	MC	1	62123	0.80	0.53	0.00		0.05	0.10	0.06	0.80		-0.25	-0.32	-0.28	0.53
9	TE	1	62123	0.78	0.37	0.00	0.22	0.78				-0.37	0.37			
10	MC	1	62110	0.79	0.40	0.00		0.03	0.79	0.14	0.03		-0.18	0.40	-0.26	-0.21
11	MC	1	62102	0.51	0.40	0.00		0.50	0.15	0.22	0.13		0.40	-0.28	-0.17	-0.09
12	MC	1	62088	0.53	0.51	0.00		0.10	0.53	0.19	0.18		-0.08	0.51	-0.31	-0.28
13	MC	1	62096	0.38	0.41	0.00		0.22	0.11	0.29	0.38		-0.12	-0.25	-0.15	0.41
14	MC	1	62106	0.70	0.54	0.00		0.06	0.14	0.10	0.70		-0.22	-0.29	-0.32	0.54
15	MC	1	62087	0.74	0.46	0.00		0.11	0.08	0.74	0.07		-0.22	-0.30	0.46	-0.20
16	MC	1	62086	0.71	0.38	0.00		0.71	0.14	0.07	0.08		0.38	-0.17	-0.27	-0.17
17	MC	1	62097	0.47	0.45	0.00		0.47	0.17	0.19	0.17		0.45	-0.11	-0.30	-0.18
18	MC	1	62091	0.72	0.39	0.00		0.05	0.18	0.05	0.72		-0.26	-0.15	-0.28	0.40
19	TE	1	62097	0.52	0.61	0.00	0.47	0.52				-0.61	0.61			
20	MC	1	62100	0.51	0.45	0.00		0.51	0.16	0.25	0.08		0.45	-0.30	-0.16	-0.17
21	MS	1	62089	0.44	0.54	0.00	0.56	0.44				-0.54	0.54			
22	TE	1	62016	0.52	0.29	0.00	0.48	0.51				-0.29	0.29			
23	MC	1	62004	0.36	0.43	0.00		0.16	0.17	0.31	0.36		-0.14	-0.32	-0.07	0.43
24	MC	1	62016	0.70	0.50	0.00		0.07	0.09	0.14	0.70		-0.22	-0.29	-0.26	0.50
25	MC	1	62033	0.55	0.34	0.00		0.55	0.06	0.23	0.16		0.34	-0.21	-0.07	-0.23
26	MC	1	62013	0.66	0.57	0.00		0.66	0.09	0.13	0.12		0.57	-0.33	-0.32	-0.20
27	MC	1	62049	0.73	0.40	0.00		0.06	0.09	0.73	0.12		-0.22	-0.24	0.40	-0.17
28	TE	1	61992	0.36	0.32	0.00	0.64	0.36				-0.31	0.32			
29	MC	1	62002	0.66	0.46	0.00		0.11	0.19	0.66	0.04		-0.26	-0.27	0.46	-0.16
30	MC	1	62002	0.63	0.51	0.00		0.63	0.12	0.12	0.13		0.51	-0.23	-0.31	-0.21

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	MC	1	61997	0.47	0.39	0.00		0.17	0.47	0.22	0.14		-0.12	0.39	-0.22	-0.17
32	MC	1	61978	0.62	0.45	0.00		0.21	0.62	0.11	0.06		-0.17	0.45	-0.28	-0.25
33	MC	1	62013	0.57	0.35	0.00		0.14	0.21	0.57	0.07		-0.22	-0.15	0.35	-0.13
34	MC	1	62027	0.67	0.37	0.00		0.10	0.11	0.66	0.13		-0.22	-0.26	0.37	-0.08
35	MS	1	62017	0.52	0.52	0.00	0.47	0.52				-0.51	0.52			
36	MC	1	61982	0.60	0.40	0.00		0.21	0.59	0.12	0.07		-0.19	0.40	-0.24	-0.15
37	TE	1	61986	0.60	0.44	0.00	0.40	0.60				-0.44	0.44			
38	MC	1	61991	0.60	0.41	0.00		0.09	0.21	0.60	0.10		-0.19	-0.24	0.41	-0.17
39	MC	1	62024	0.67	0.42	0.00		0.08	0.16	0.66	0.10		-0.13	-0.23	0.42	-0.26
40	TE	1	61998	0.40	0.46	0.00	0.60	0.40				-0.45	0.46			

Table H-17. Item Statistics, Social Studies Grade 10

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
1	TE	1	61393	0.68	0.44	0.00	0.32	0.67				-0.43	0.44			
2	MC	1	61519	0.64	0.47	0.00		0.11	0.64	0.07	0.18		-0.16	0.47	-0.24	-0.30
3	MC	1	61514	0.72	0.50	0.00		0.07	0.08	0.14	0.72		-0.14	-0.26	-0.35	0.50
4	MC	1	61485	0.60	0.39	0.00		0.14	0.21	0.60	0.05		-0.19	-0.22	0.39	-0.15
5	MC	1	61380	0.41	0.29	0.00		0.14	0.41	0.19	0.27		-0.05	0.29	-0.32	0.00
6	MC	1	61403	0.47	0.41	0.00		0.47	0.11	0.24	0.18		0.41	-0.26	-0.20	-0.08
7	TE	1	61487	0.62	0.42	0.00	0.38	0.62				-0.42	0.42			
8	MS	1	61515	0.38	0.45	0.00	0.62	0.38				-0.45	0.45			
9	MC	1	61506	0.62	0.30	0.00		0.10	0.62	0.14	0.14		-0.19	0.30	-0.07	-0.18
10	MC	1	61417	0.69	0.47	0.00		0.10	0.69	0.16	0.05		-0.21	0.47	-0.33	-0.12
11	MC	1	61450	0.57	0.49	0.00		0.11	0.18	0.14	0.57		-0.16	-0.28	-0.24	0.49
12	MC	1	61465	0.57	0.44	0.00		0.13	0.57	0.20	0.10		-0.04	0.44	-0.30	-0.27
13	MC	1	61485	0.53	0.53	0.00		0.12	0.21	0.14	0.52		-0.23	-0.26	-0.25	0.54
14	MC	1	61470	0.38	0.34	0.00		0.38	0.14	0.29	0.19		0.34	-0.18	-0.06	-0.19
15	MC	1	61433	0.57	0.40	0.00		0.57	0.14	0.18	0.11		0.40	-0.16	-0.15	-0.27
16	MC	1	61413	0.35	0.38	0.00		0.16	0.18	0.31	0.35		-0.10	-0.21	-0.12	0.38
17	MC	1	61403	0.48	0.46	0.00		0.48	0.21	0.17	0.14		0.46	-0.21	-0.25	-0.13
18	MC	1	61391	0.32	0.26	0.00		0.32	0.17	0.43	0.08		0.26	-0.22	0.05	-0.23
19	MC	1	61447	0.50	0.28	0.00		0.13	0.27	0.50	0.10		-0.10	-0.08	0.28	-0.22
20	MS	1	61457	0.30	0.30	0.00	0.70	0.30				-0.29	0.30			
21	MC	1	60905	0.59	0.40	0.00		0.27	0.12	0.59	0.02		-0.25	-0.19	0.40	-0.15
22	MS	1	60892	0.45	0.53	0.00	0.55	0.45				-0.53	0.53			
23	TE	1	60829	0.51	0.44	0.00	0.49	0.51				-0.44	0.44			
24	MC	1	60814	0.65	0.40	0.00		0.07	0.19	0.65	0.08		-0.16	-0.25	0.40	-0.19
25	MC	1	60835	0.66	0.57	0.00		0.09	0.12	0.14	0.65		-0.29	-0.26	-0.30	0.57
26	MS	1	60837	0.27	0.32	0.00	0.73	0.27				-0.32	0.32			
27	MC	1	60803	0.34	0.21	0.00		0.22	0.34	0.21	0.23		-0.15	0.21	-0.14	0.06
28	MC	1	60758	0.61	0.55	0.00		0.06	0.20	0.13	0.61		-0.24	-0.27	-0.30	0.55
29	MC	1	60787	0.41	0.35	0.00		0.16	0.20	0.41	0.22		-0.13	-0.22	0.35	-0.08
30	MC	1	60793	0.56	0.37	0.00		0.19	0.12	0.56	0.12		-0.08	-0.19	0.37	-0.26

1	2	3	4	5	6	7	Proportion of Students					Item-Total Test Correlation				
							8	9	10	11	12	13	14	15	16	17
31	TE	1	60742	0.59	0.48	0.00	0.41	0.59				-0.48	0.49			
32	TE	1	60359	0.43	0.42	0.01	0.57	0.42				-0.40	0.42			
33	MC	1	60767	0.46	0.47	0.00		0.46	0.18	0.26	0.09		0.47	-0.26	-0.19	-0.17
34	MC	1	60782	0.55	0.27	0.00		0.19	0.12	0.55	0.14		0.06	-0.23	0.27	-0.23
35	MC	1	60752	0.63	0.44	0.00		0.08	0.14	0.14	0.63		-0.17	-0.19	-0.28	0.44
36	TE	1	60735	0.37	0.49	0.00	0.63	0.37				-0.48	0.49			
37	MC	1	60761	0.73	0.45	0.00		0.03	0.06	0.18	0.73		-0.15	-0.21	-0.32	0.45
38	MC	1	60743	0.44	0.22	0.00		0.18	0.23	0.44	0.14		-0.12	-0.12	0.22	-0.03
39	MC	1	60751	0.47	0.39	0.00		0.46	0.19	0.07	0.27		0.39	-0.18	-0.22	-0.15
40	MC	1	60760	0.59	0.49	0.00		0.08	0.59	0.15	0.18		-0.19	0.49	-0.33	-0.19

APPENDIX I
CONDITIONAL STANDARD ERROR OF
MEASUREMENT WITH CUT SCORES

List of Figures

Figure I-1. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 3.....	3
Figure I-2. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 4.....	4
Figure I-3. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 5.....	5
Figure I-4. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 6.....	6
Figure I-5. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 7.....	7
Figure I-6. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 8.....	8
Figure I-7. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 3.....	9
Figure I-8. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 4.....	10
Figure I-9. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 5.....	11
Figure I-10. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 6.....	12
Figure I-11. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 7.....	13
Figure I-12. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 8.....	14
Figure I-13. Conditional Standard Error of Measurement with Cut Scores, Science Grade 4.....	15
Figure I-14. Conditional Standard Error of Measurement with Cut Scores, Science Grade 8.....	16
Figure I-15. Conditional Standard Error of Measurement with Cut Scores, Social Studies Grade 4.....	17
Figure I-16. Conditional Standard Error of Measurement with Cut Scores, Social Studies Grade 8.....	18
Figure I-17. Conditional Standard Error of Measurement with Cut Scores, Social Studies Grade 10.....	19

Figure I-1. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 3

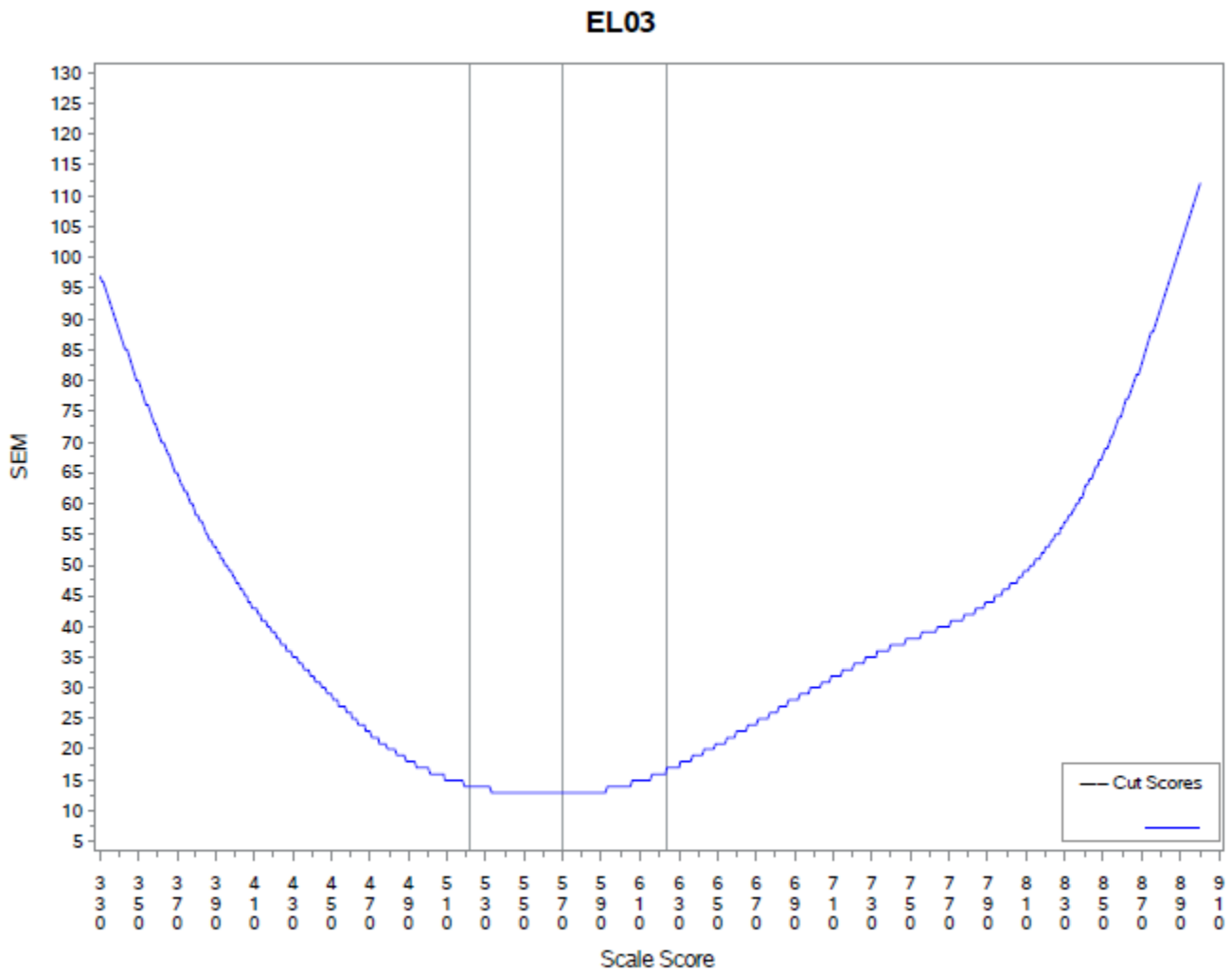


Figure I-2. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 4

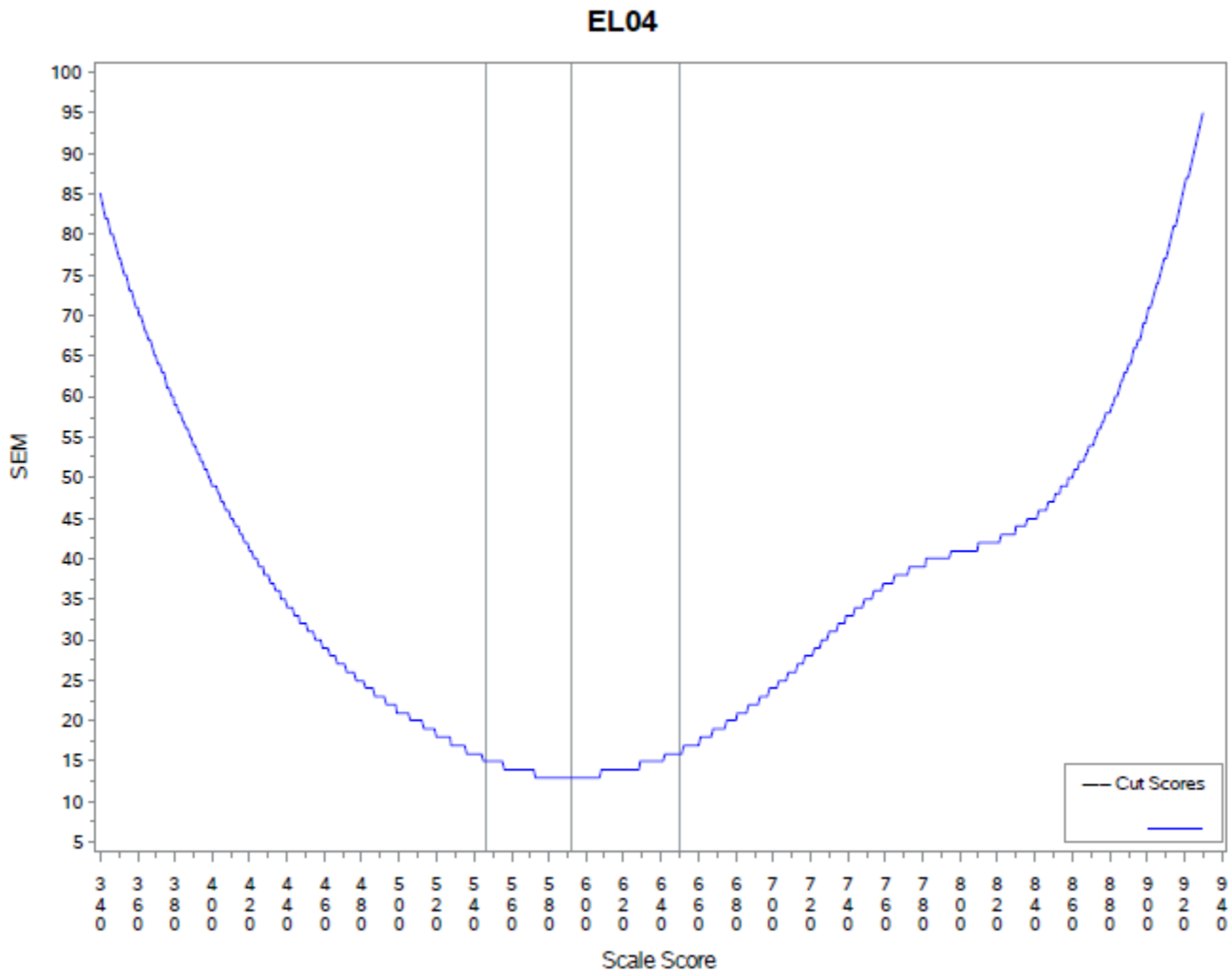


Figure I-3. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 5

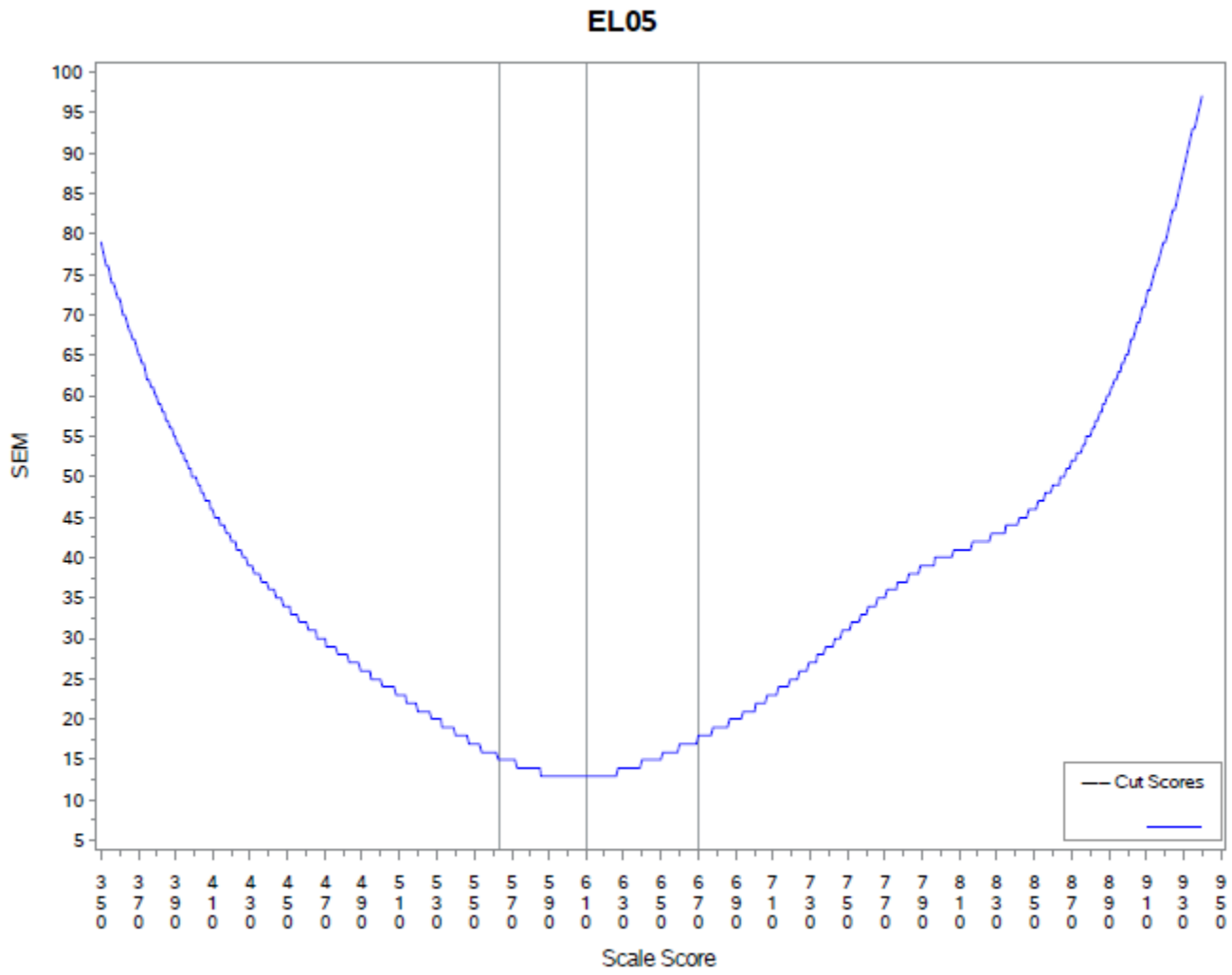


Figure I-4. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 6

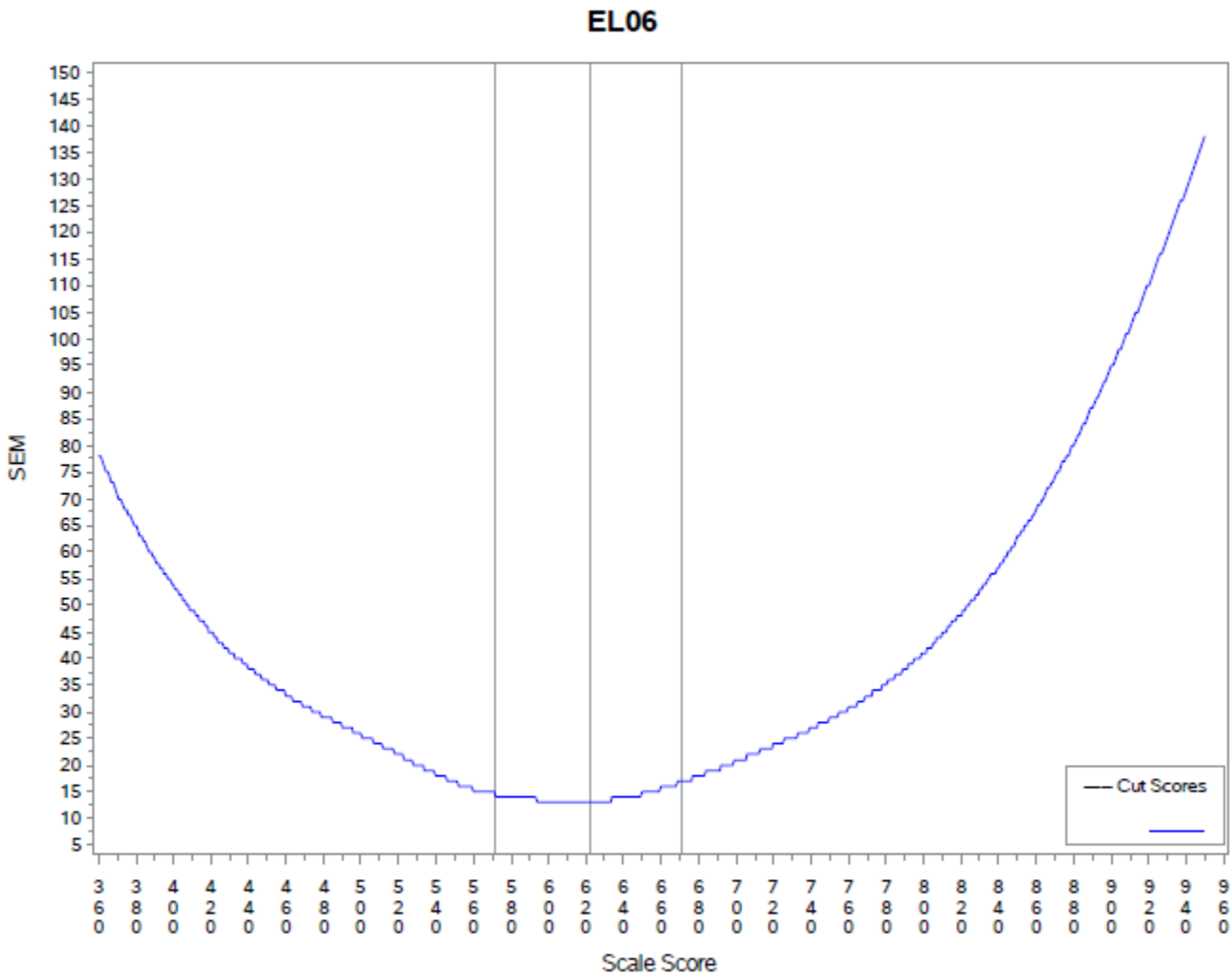


Figure I-5. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 7

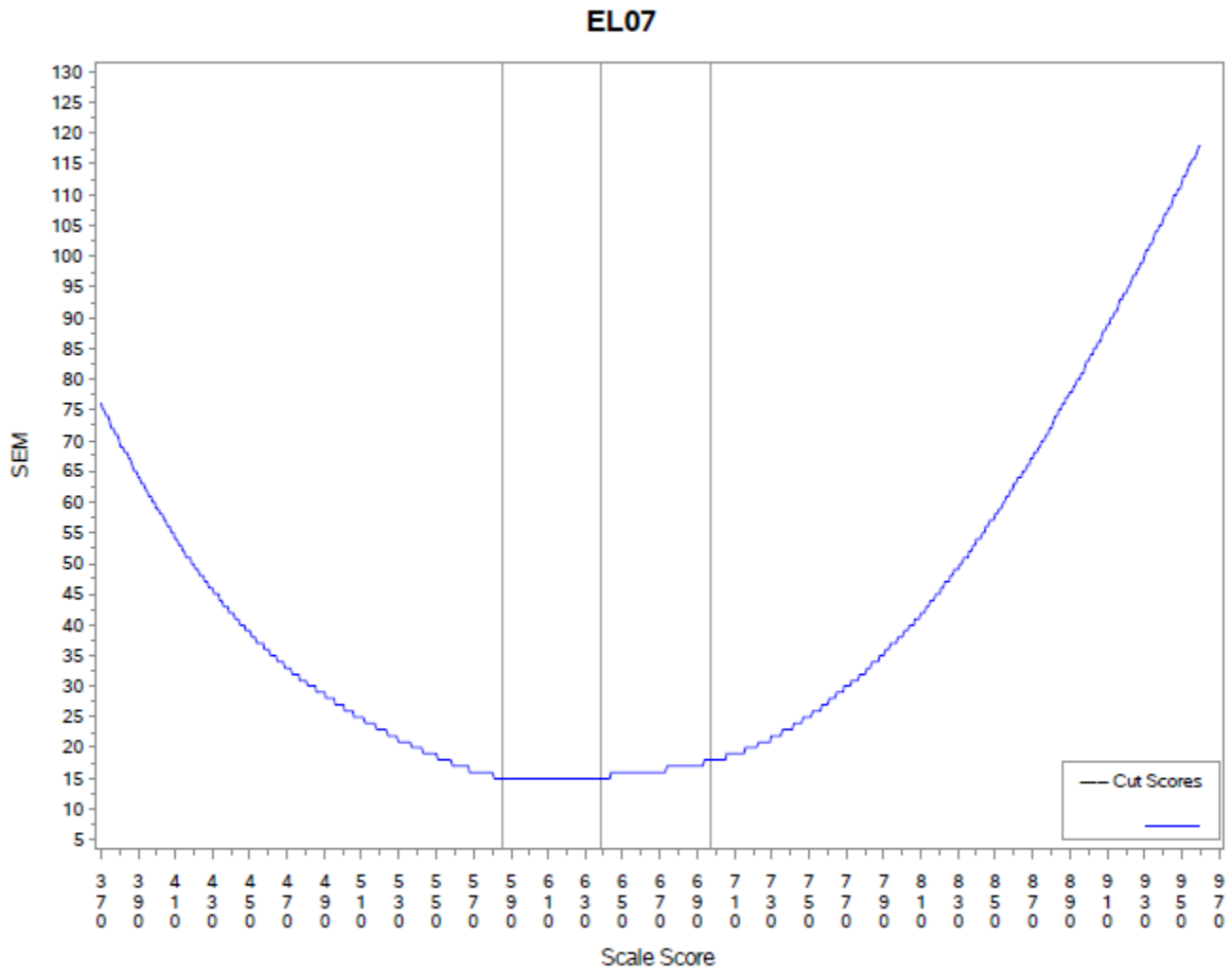


Figure I-6. Conditional Standard Error of Measurement with Cut Scores, English Language Arts Grade 8

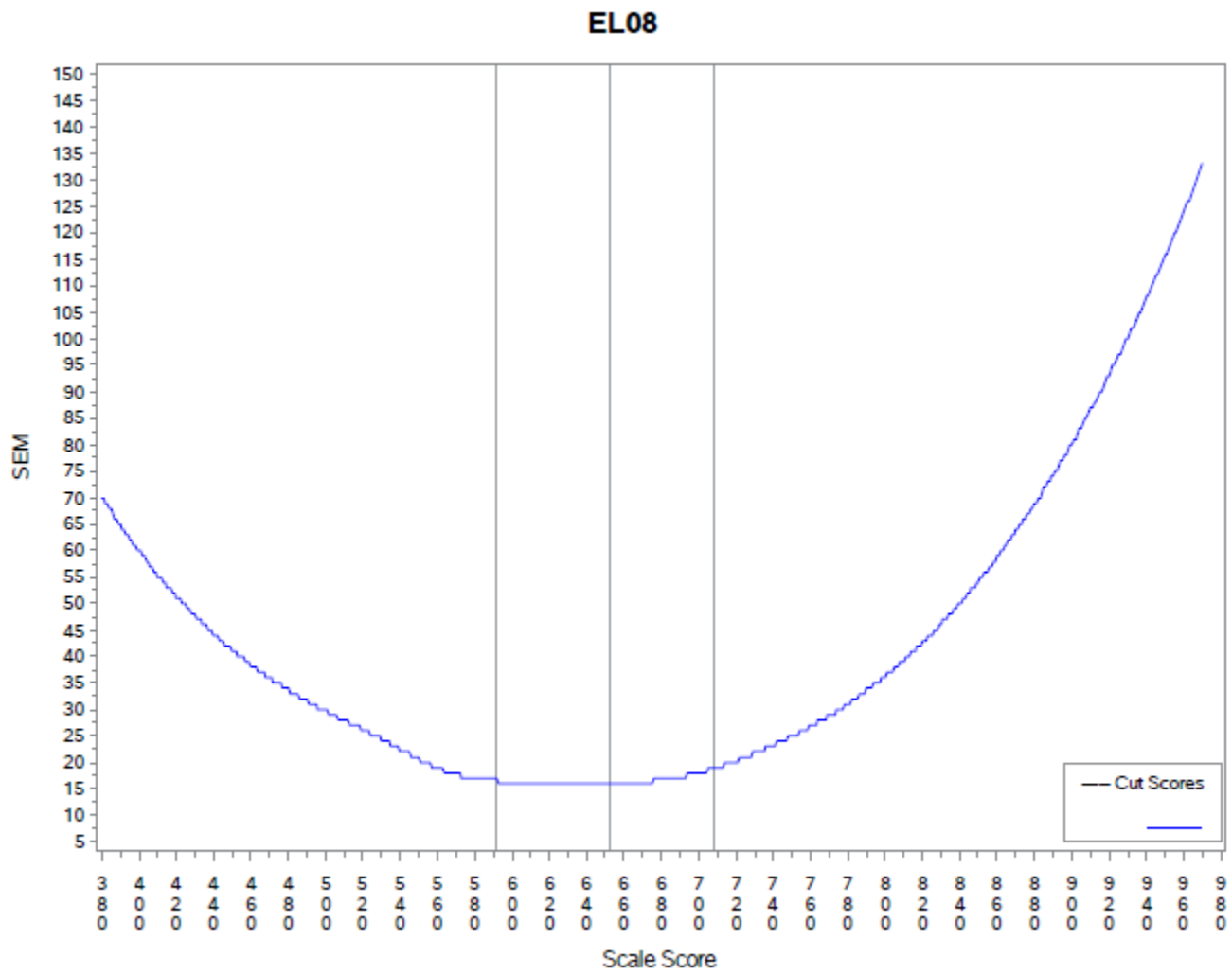


Figure I-7. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 3

MA03

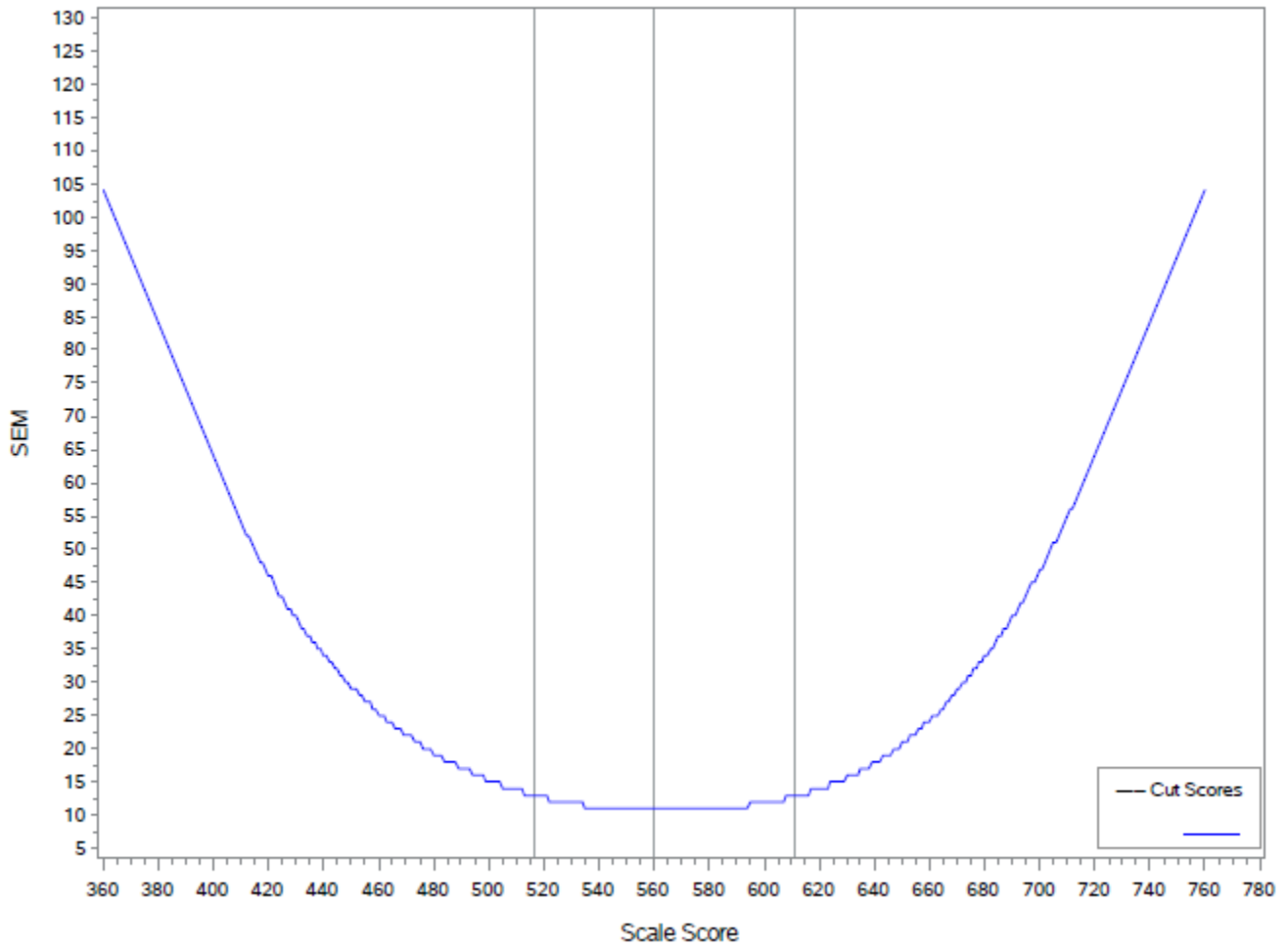


Figure I-8. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 4

MA04

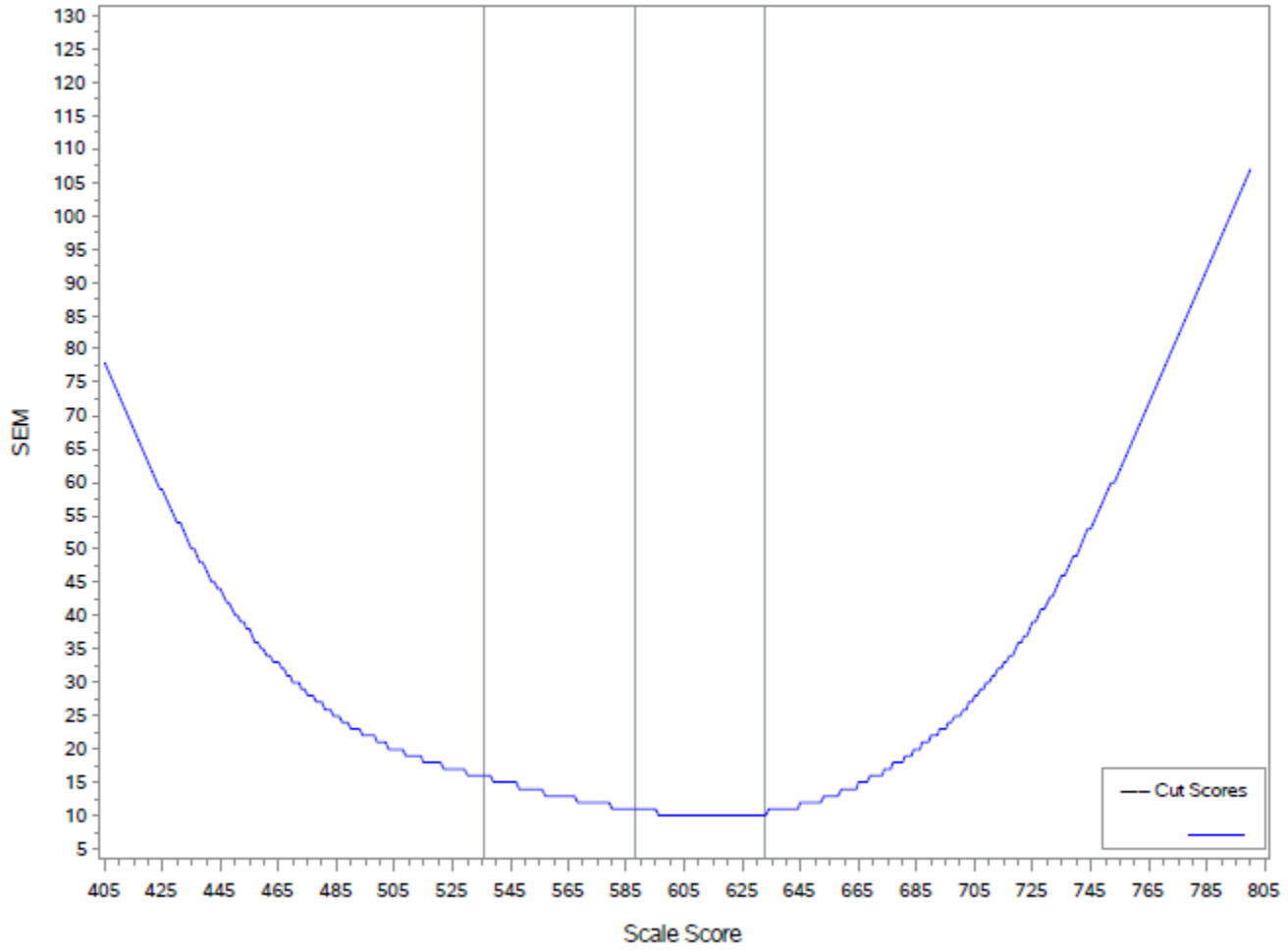


Figure I-9. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 5

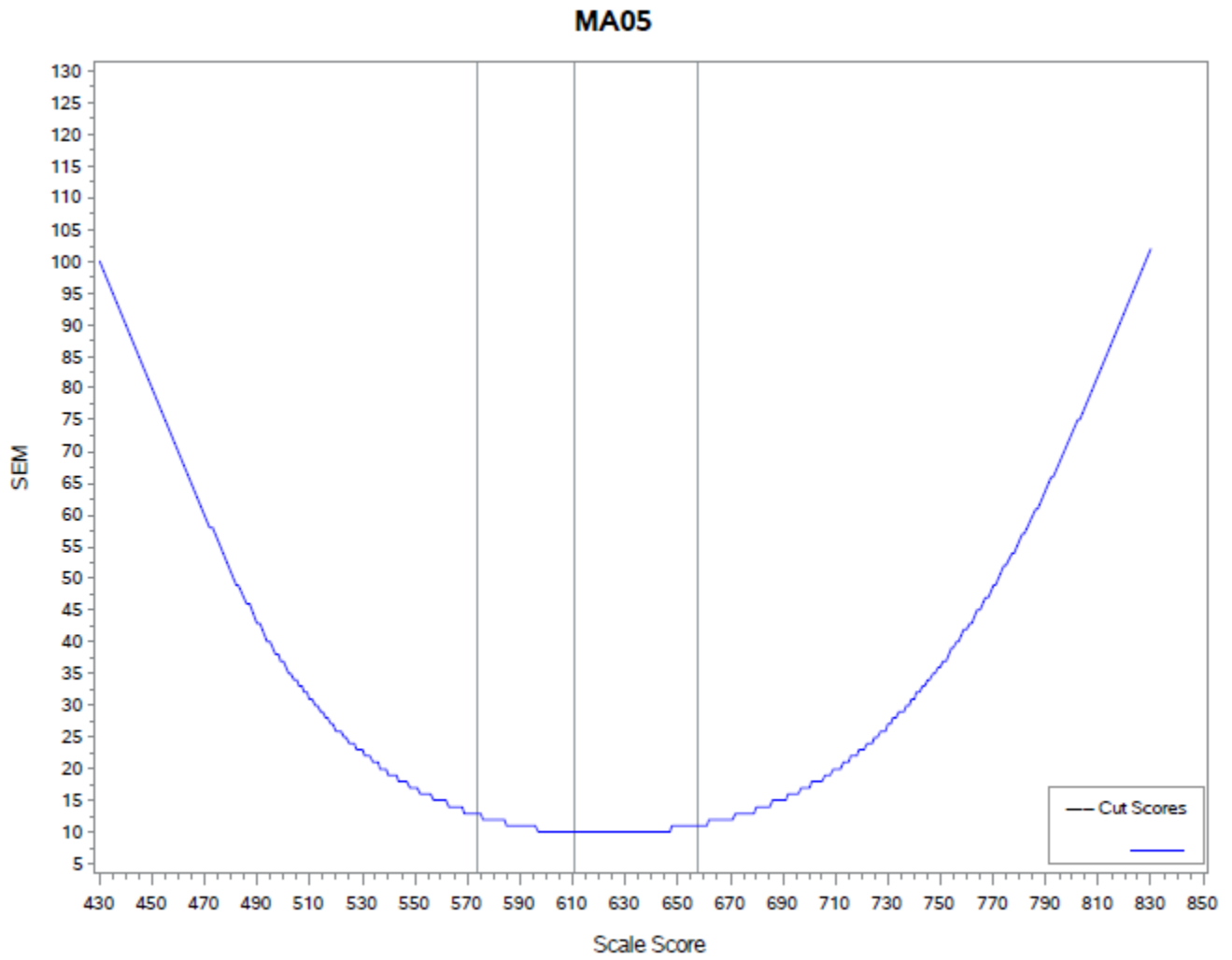


Figure I-10. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 6

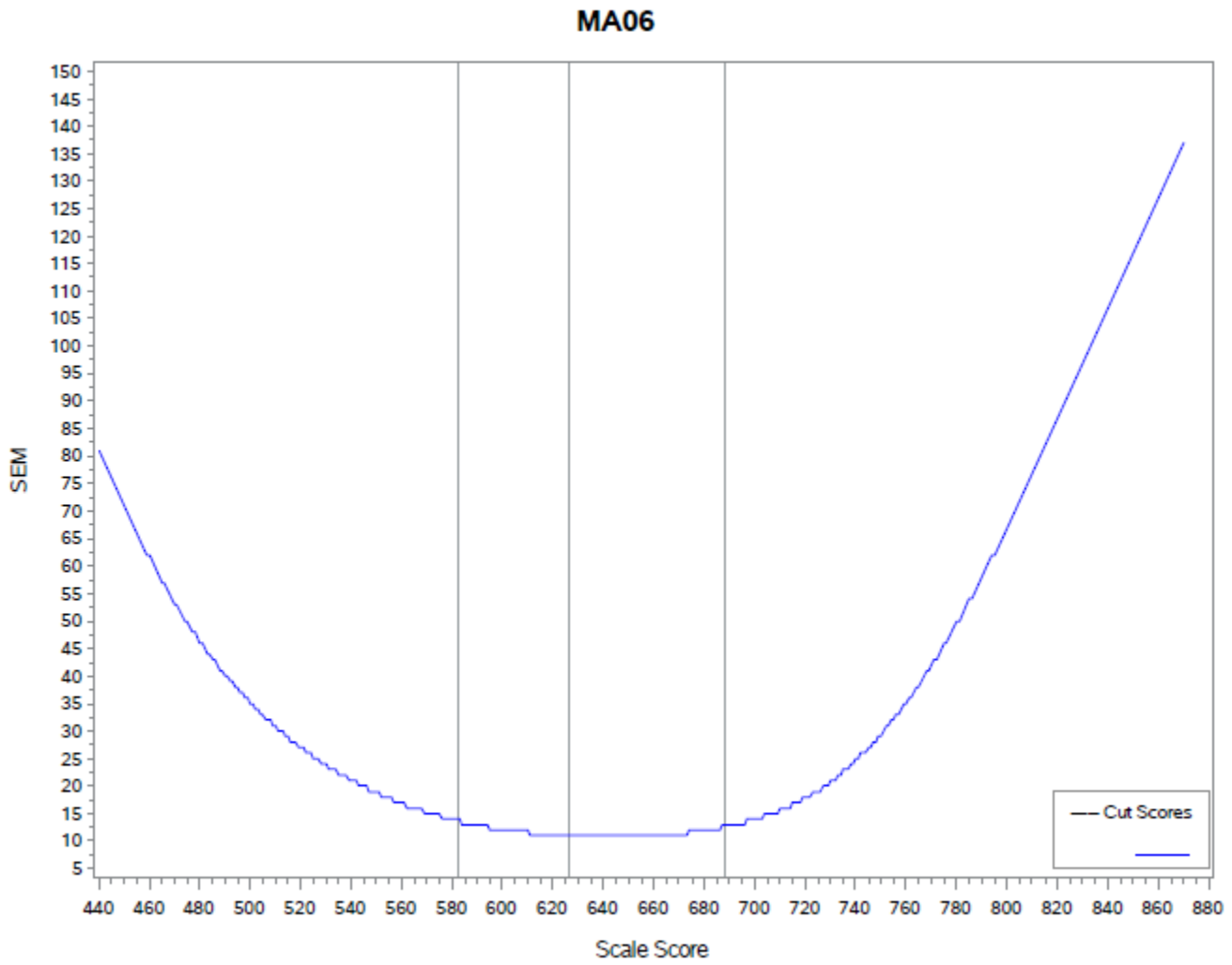


Figure I-11. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 7

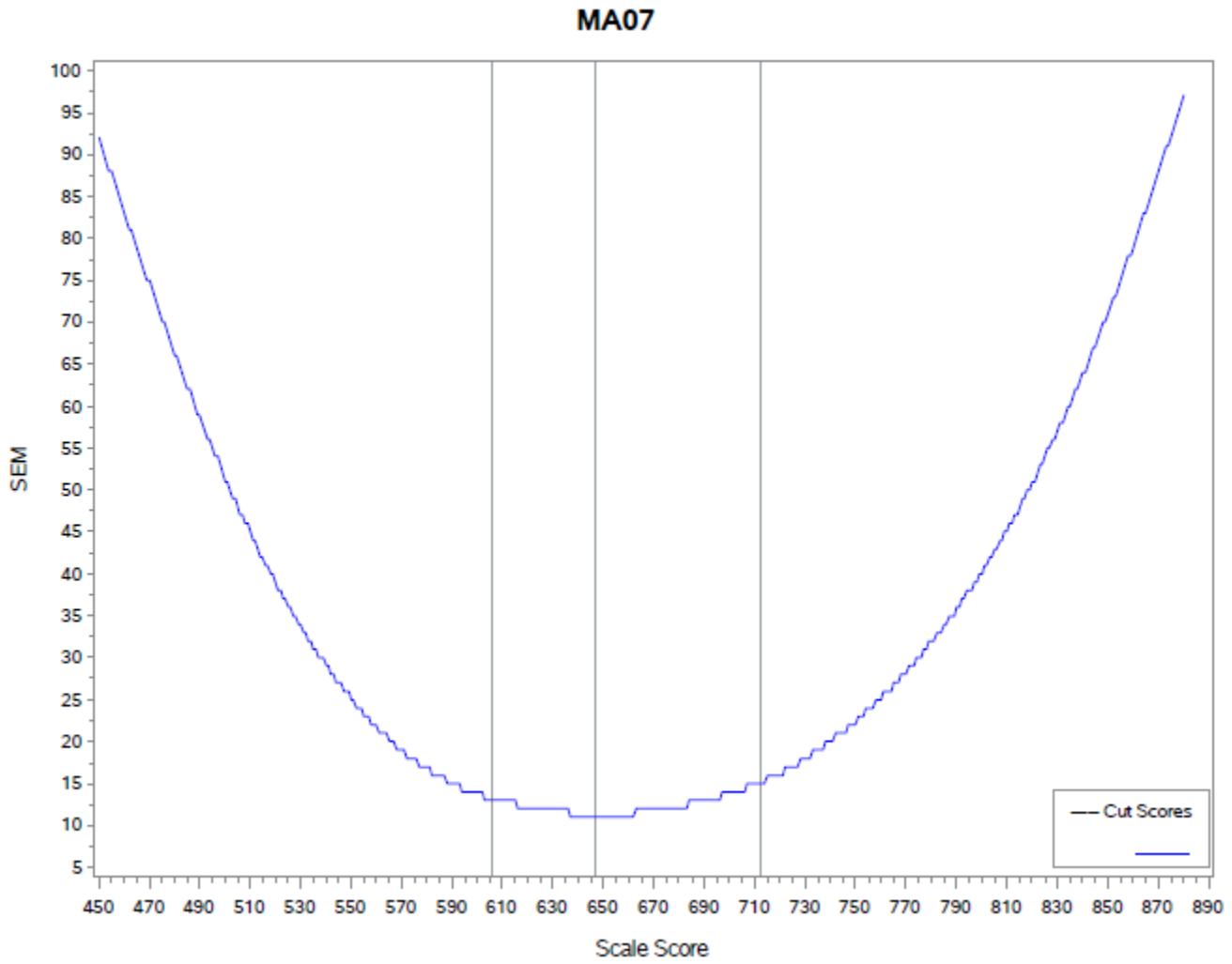


Figure I-12. Conditional Standard Error of Measurement with Cut Scores, Mathematics Grade 8

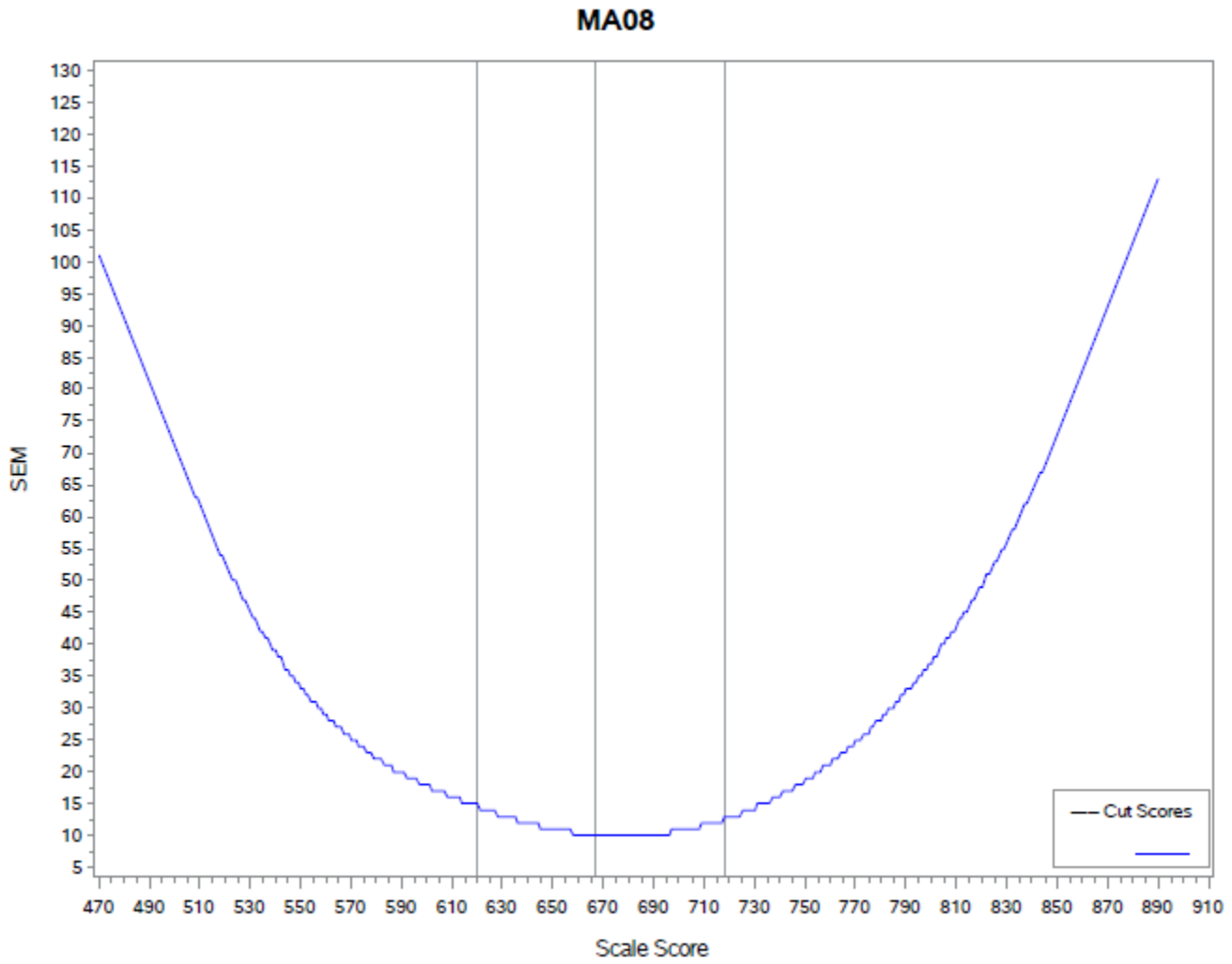


Figure I-13. Conditional Standard Error of Measurement with Cut Scores, Science Grade 4

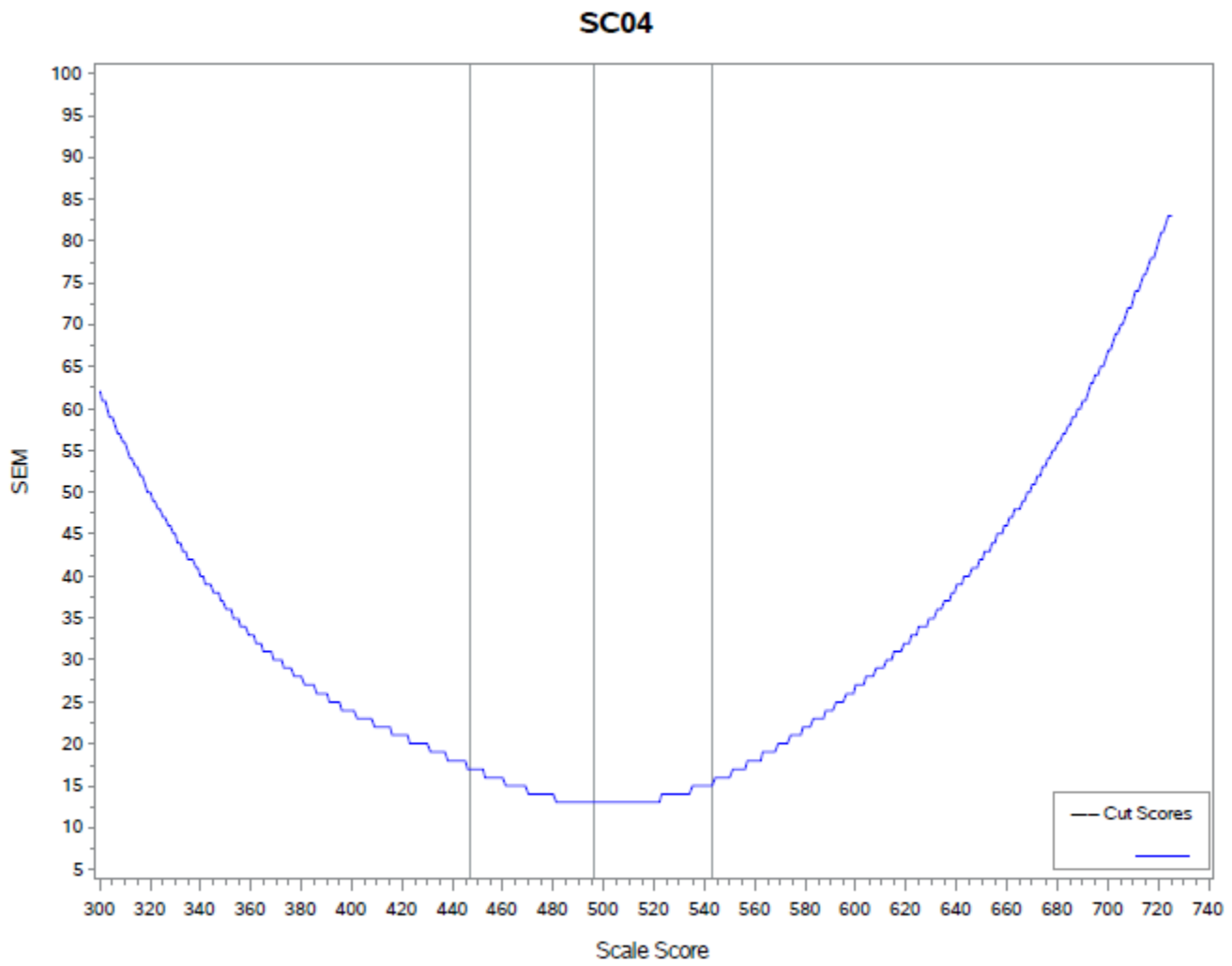


Figure I-14. Conditional Standard Error of Measurement with Cut Scores, Science Grade 8

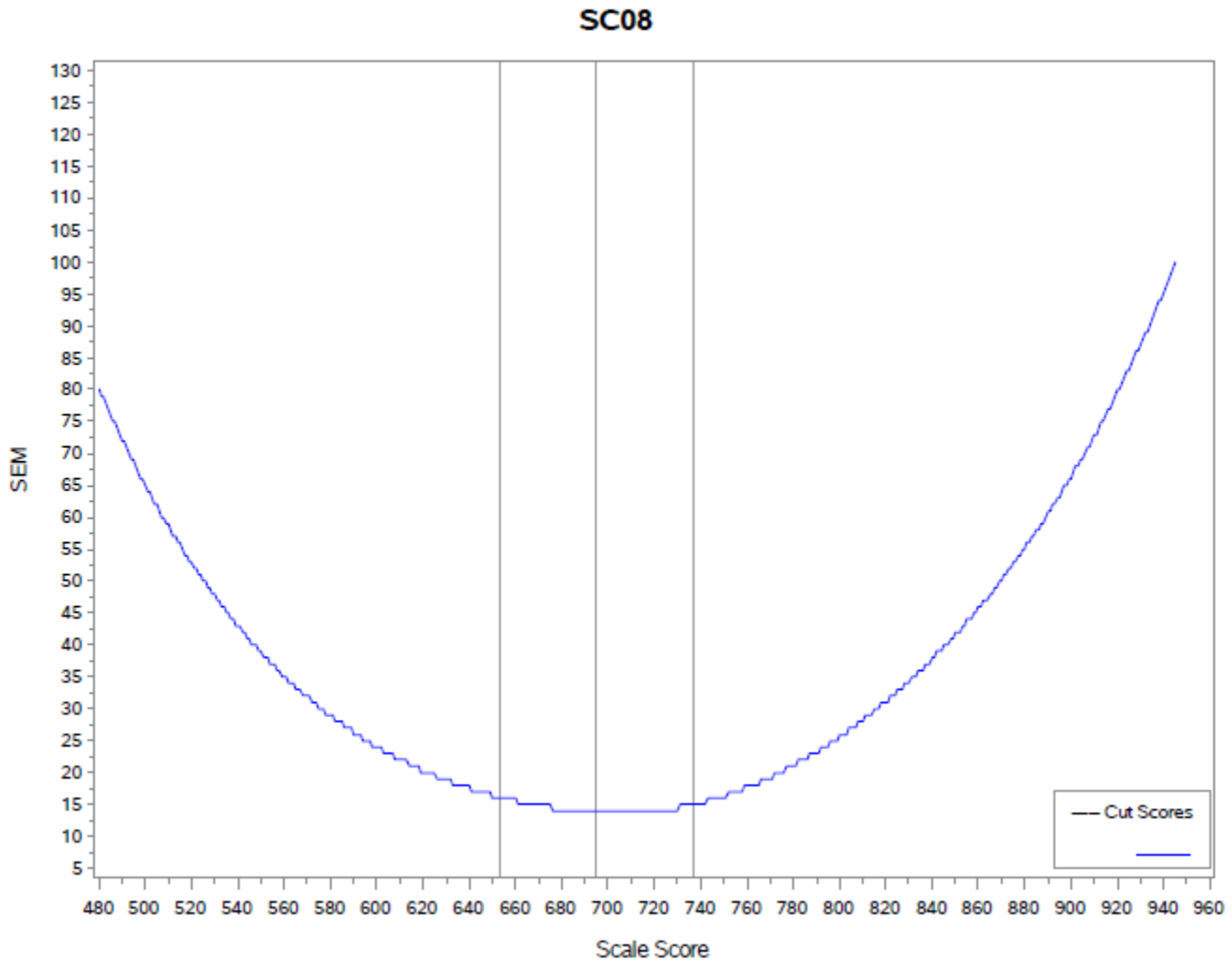


Figure I-15. Conditional Standard Error of Measurement with Cut Scores, Social Studies Grade 4

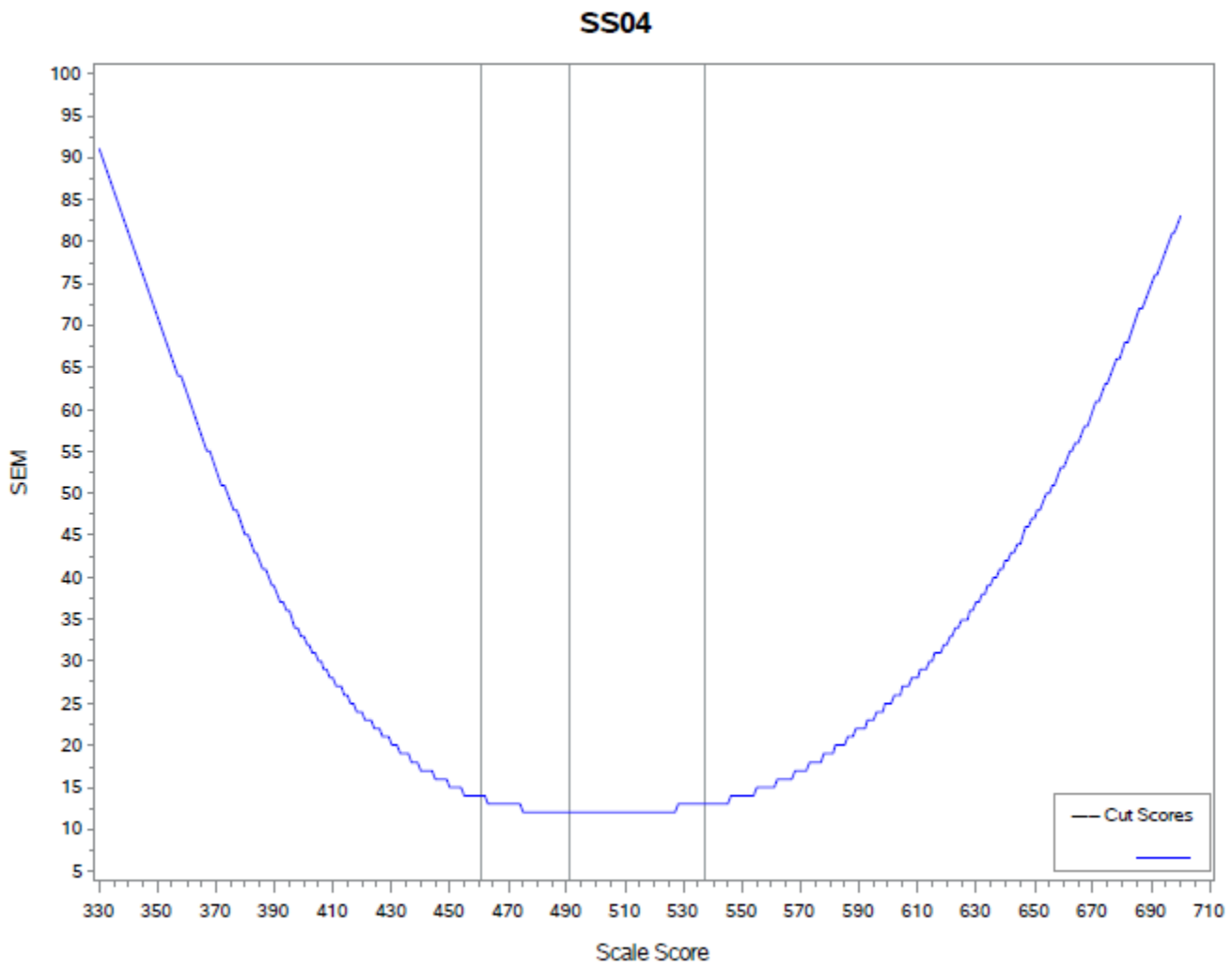


Figure I-16. Conditional Standard Error of Measurement with Cut Scores, Social Studies Grade 8

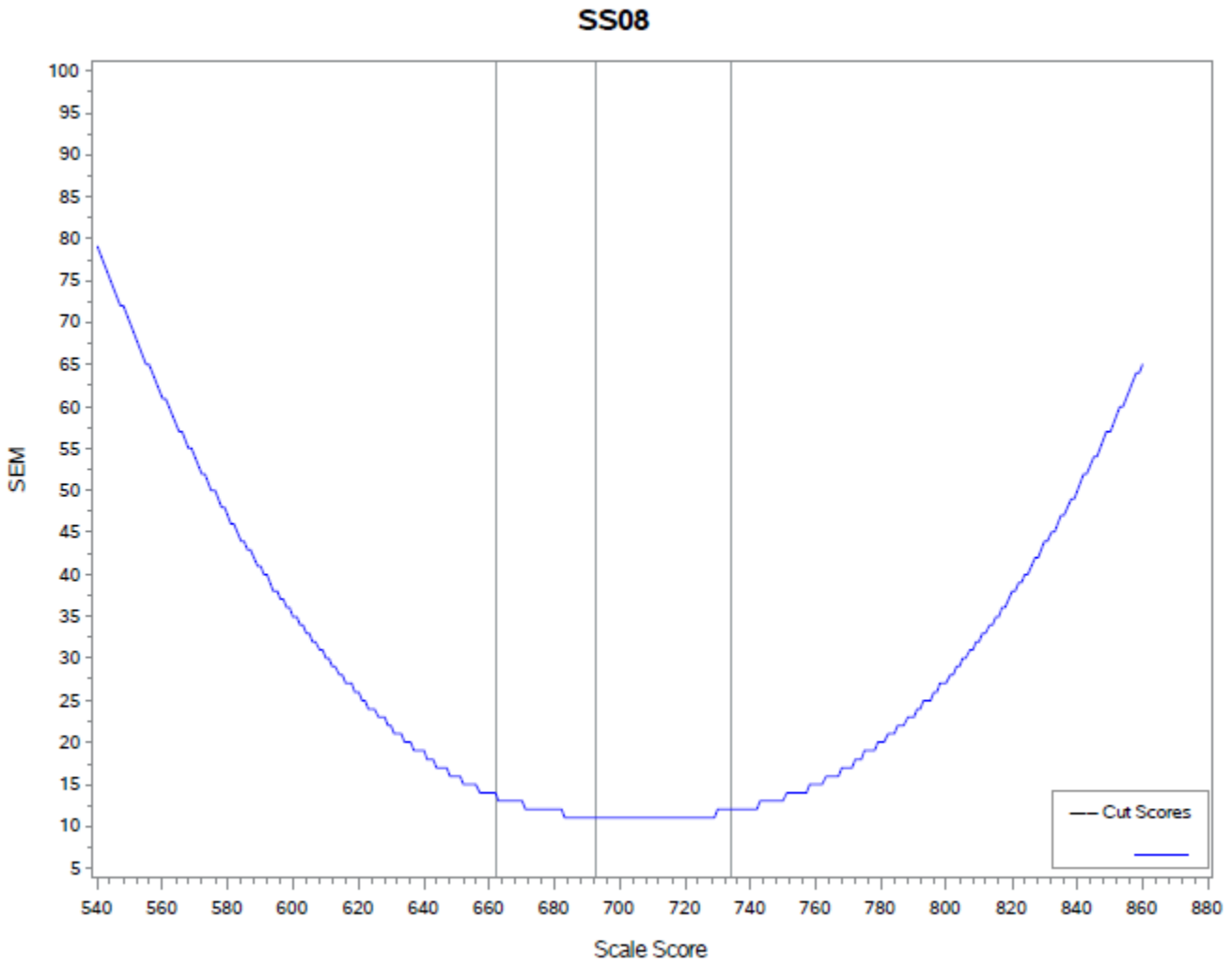
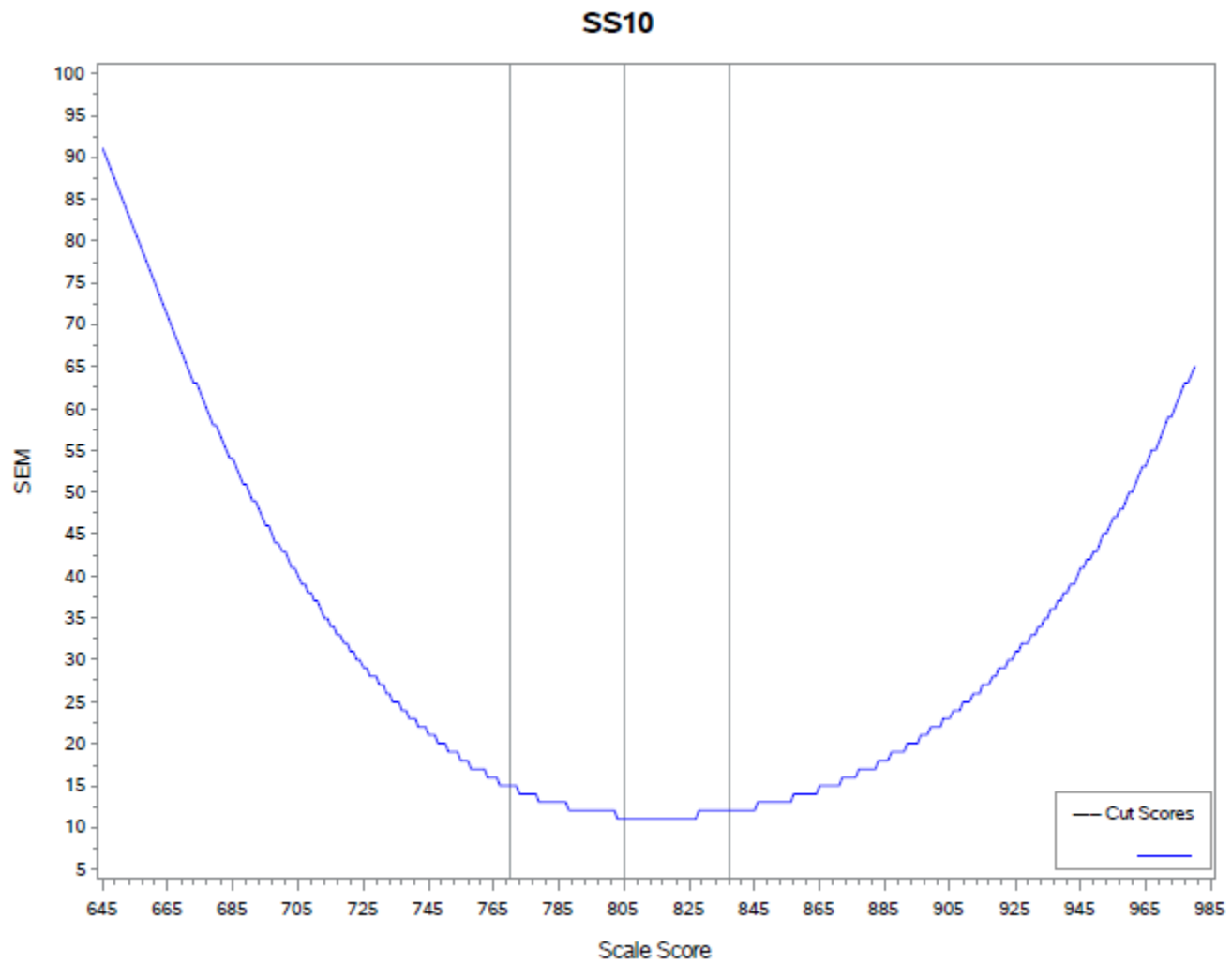


Figure I-17. Conditional Standard Error of Measurement with Cut Scores, Social Studies Grade 10



APPENDIX J
CLASSIFICATION CONSISTENCY AND
ACCURACY ANALYSIS BY SUBGROUP

List of Tables

Table J-1. Indices for Classification Consistency and Accuracy, ELA Grade 3	3
Table J-2. Indices for Classification Consistency and Accuracy, ELA Grade 4	5
Table J-3. Indices for Classification Consistency and Accuracy, ELA Grade 5	7
Table J-4. Indices for Classification Consistency and Accuracy, ELA Grade 6	9
Table J-5. Indices for Classification Consistency and Accuracy, ELA Grade 7	11
Table J-6. Indices for Classification Consistency and Accuracy, ELA Grade 8	13
Table J-7. Indices for Classification Consistency and Accuracy, Mathematics Grade 3	15
Table J-8. Indices for Classification Consistency and Accuracy, Mathematics Grade 4	17
Table J-9. Indices for Classification Consistency and Accuracy, Mathematics Grade 5	19
Table J-10. Indices for Classification Consistency and Accuracy, Mathematics Grade 6	21
Table J-11. Indices for Classification Consistency and Accuracy, Mathematics Grade 7	23
Table J-12. Indices for Classification Consistency and Accuracy, Mathematics Grade 8	25
Table J-13. Indices for Classification Consistency and Accuracy, Science Grade 4	27
Table J-14. Indices for Classification Consistency and Accuracy, Science Grade 8	28
Table J-15. Indices for Classification Consistency and Accuracy, Social Studies Grade 4	30
Table J-16. Indices for Classification Consistency and Accuracy, Social Studies Grade 8	31
Table J-17. Indices for Classification Consistency and Accuracy, Social Studies Grade 10	33

Table J-1. Indices for Classification Consistency and Accuracy, ELA Grade 3

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.90	0.90	0.94	0.73
		Probability of Chance	0.63	0.52	0.84	0.29
		Kappa (k)	0.73	0.78	0.61	0.63
		Classification Accuracy	0.93	0.92	0.96	0.81
	Male	Classification Consistency (P)	0.89	0.90	0.95	0.74
		Probability of Chance	0.59	0.55	0.88	0.30
		Kappa (k)	0.73	0.78	0.60	0.63
		Classification Accuracy	0.92	0.93	0.97	0.81
Race/Ethnicity	White	Classification Consistency (P)	0.91	0.88	0.93	0.72
		Probability of Chance	0.71	0.50	0.82	0.30
		Kappa (k)	0.70	0.76	0.60	0.60
		Classification Accuracy	0.94	0.91	0.95	0.80
	African American	Classification Consistency (P)	0.85	0.95	0.99	0.80
		Probability of Chance	0.52	0.82	0.98	0.45
		Kappa (k)	0.70	0.74	0.53	0.63
		Classification Accuracy	0.90	0.97	0.99	0.86
	Hispanic	Classification Consistency (P)	0.86	0.92	0.97	0.75
		Probability of Chance	0.52	0.66	0.94	0.34
		Kappa (k)	0.70	0.77	0.60	0.62
		Classification Accuracy	0.90	0.94	0.98	0.82
	Asian	Classification Consistency (P)	0.87	0.91	0.95	0.74
		Probability of Chance	0.57	0.57	0.86	0.30
		Kappa (k)	0.70	0.79	0.67	0.62
		Classification Accuracy	0.91	0.93	0.97	0.81
	American Indian	Classification Consistency (P)	0.85	0.92	0.98	0.75
		Probability of Chance	0.51	0.70	0.96	0.36
		Kappa (k)	0.69	0.74	0.54	0.61
		Classification Accuracy	0.89	0.94	0.99	0.82
Two or More	Classification Consistency (P)	0.88	0.90	0.95	0.73	
	Probability of Chance	0.60	0.55	0.88	0.30	
	Kappa (k)	0.71	0.78	0.60	0.62	
	Classification Accuracy	0.92	0.93	0.97	0.81	
Limited English Proficiency	Yes	Classification Consistency (P)	0.90	0.95	0.99	0.83
		Probability of Chance	0.50	0.75	0.93	0.41
		Kappa (k)	0.80	0.80	0.79	0.72
		Classification Accuracy	0.90	0.96	0.97	0.84
Disability Status	Yes	Classification Consistency (P)	0.86	0.94	0.98	0.77
		Probability of Chance	0.50	0.73	0.95	0.39
		Kappa (k)	0.72	0.77	0.58	0.63
		Classification Accuracy	0.90	0.95	0.99	0.84

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.86	0.92	0.98	0.75
		Probability of Chance	0.52	0.66	0.94	0.34
		Kappa (k)	0.71	0.76	0.56	0.63
		Classification Accuracy	0.90	0.94	0.98	0.82
Accommodation Use	Yes	Classification Consistency (P)	0.93	0.91	0.96	0.80
		Probability of Chance	0.50	0.60	0.91	0.36
		Kappa (k)	0.87	0.76	0.50	0.68
		Classification Accuracy	0.96	0.93	0.96	0.85

Table J-2. Indices for Classification Consistency and Accuracy, ELA Grade 4

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.91	0.90	0.93	0.73
		Probability of Chance	0.66	0.50	0.78	0.28
		Kappa (k)	0.72	0.79	0.68	0.63
		Classification Accuracy	0.94	0.93	0.95	0.81
	Male	Classification Consistency (P)	0.92	0.94	0.93	0.79
		Probability of Chance	0.58	0.50	0.72	0.26
		Kappa (k)	0.82	0.88	0.75	0.72
		Classification Accuracy	0.94	0.95	0.91	0.80
Race/Ethnicity	White	Classification Consistency (P)	0.92	0.89	0.92	0.73
		Probability of Chance	0.74	0.50	0.75	0.30
		Kappa (k)	0.69	0.77	0.67	0.61
		Classification Accuracy	0.94	0.92	0.94	0.81
	African American	Classification Consistency (P)	0.86	0.94	0.99	0.79
		Probability of Chance	0.51	0.77	0.96	0.43
		Kappa (k)	0.71	0.74	0.66	0.63
		Classification Accuracy	0.90	0.96	0.99	0.85
	Hispanic	Classification Consistency (P)	0.87	0.91	0.97	0.74
		Probability of Chance	0.53	0.60	0.91	0.31
		Kappa (k)	0.72	0.77	0.64	0.63
		Classification Accuracy	0.91	0.93	0.98	0.82
	Asian	Classification Consistency (P)	0.88	0.90	0.94	0.73
		Probability of Chance	0.61	0.52	0.81	0.28
		Kappa (k)	0.70	0.80	0.69	0.63
		Classification Accuracy	0.92	0.93	0.96	0.81
	American Indian	Classification Consistency (P)	0.86	0.89	0.97	0.73
		Probability of Chance	0.53	0.63	0.93	0.33
		Kappa (k)	0.71	0.72	0.57	0.60
		Classification Accuracy	0.91	0.93	0.98	0.81
Two or More	Classification Consistency (P)	0.89	0.89	0.94	0.73	
	Probability of Chance	0.61	0.52	0.82	0.28	
	Kappa (k)	0.73	0.77	0.68	0.62	
	Classification Accuracy	0.93	0.93	0.96	0.81	
Limited English Proficiency	Yes	Classification Consistency (P)	0.85	0.91	0.98	0.75
		Probability of Chance	0.51	0.69	0.97	0.36
		Kappa (k)	0.69	0.72	0.55	0.60
		Classification Accuracy	0.90	0.94	0.99	0.82
Disability Status	Yes	Classification Consistency (P)	0.86	0.93	0.98	0.77
		Probability of Chance	0.50	0.68	0.93	0.37
		Kappa (k)	0.72	0.78	0.67	0.63
		Classification Accuracy	0.90	0.95	0.99	0.84

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.87	0.91	0.97	0.74
		Probability of Chance	0.53	0.60	0.91	0.32
		Kappa (k)	0.72	0.77	0.63	0.62
		Classification Accuracy	0.91	0.93	0.98	0.82
Accommodation Use	Yes	Classification Consistency (P)	0.91	0.92	0.97	0.80
		Probability of Chance	0.50	0.60	0.91	0.35
		Kappa (k)	0.82	0.80	0.61	0.69
		Classification Accuracy	0.93	0.94	0.97	0.85

Table J-3. Indices for Classification Consistency and Accuracy, ELA Grade 5

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.90	0.90	0.94	0.73
		Probability of Chance	0.63	0.51	0.84	0.29
		Kappa (k)	0.73	0.79	0.62	0.62
		Classification Accuracy	0.93	0.92	0.96	0.81
	Male	Classification Consistency (P)	0.89	0.90	0.95	0.74
		Probability of Chance	0.57	0.54	0.89	0.30
		Kappa (k)	0.74	0.78	0.58	0.63
		Classification Accuracy	0.92	0.93	0.97	0.81
Race/Ethnicity	White	Classification Consistency (P)	0.91	0.89	0.93	0.72
		Probability of Chance	0.69	0.50	0.83	0.30
		Kappa (k)	0.70	0.77	0.60	0.61
		Classification Accuracy	0.93	0.91	0.95	0.80
	African American	Classification Consistency (P)	0.86	0.95	0.99	0.80
		Probability of Chance	0.53	0.81	0.98	0.46
		Kappa (k)	0.70	0.74	0.54	0.63
		Classification Accuracy	0.90	0.96	0.99	0.86
	Hispanic	Classification Consistency (P)	0.86	0.92	0.98	0.75
		Probability of Chance	0.51	0.65	0.94	0.34
		Kappa (k)	0.71	0.76	0.58	0.63
		Classification Accuracy	0.90	0.94	0.98	0.82
	Asian	Classification Consistency (P)	0.90	0.90	0.94	0.74
		Probability of Chance	0.59	0.53	0.84	0.29
		Kappa (k)	0.75	0.79	0.64	0.64
		Classification Accuracy	0.92	0.93	0.96	0.80
	American Indian	Classification Consistency (P)	0.86	0.91	0.98	0.76
		Probability of Chance	0.51	0.70	0.97	0.36
		Kappa (k)	0.71	0.71	0.52	0.62
		Classification Accuracy	0.90	0.94	0.99	0.82
Two or More	Classification Consistency (P)	0.88	0.90	0.95	0.74	
	Probability of Chance	0.57	0.55	0.88	0.30	
	Kappa (k)	0.72	0.79	0.59	0.63	
	Classification Accuracy	0.92	0.93	0.97	0.81	
Limited English Proficiency	Yes	Classification Consistency (P)	0.84	0.94	0.99	0.77
		Probability of Chance	0.51	0.81	0.99	0.43
		Kappa (k)	0.67	0.67	0.54	0.59
		Classification Accuracy	0.88	0.95	1.00	0.83
Disability Status	Yes	Classification Consistency (P)	0.87	0.95	0.99	0.80
		Probability of Chance	0.52	0.77	0.97	0.44
		Kappa (k)	0.72	0.77	0.52	0.64
		Classification Accuracy	0.91	0.96	0.99	0.86

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.86	0.92	0.98	0.76
		Probability of Chance	0.51	0.65	0.95	0.34
		Kappa (k)	0.72	0.76	0.55	0.63
		Classification Accuracy	0.90	0.94	0.98	0.82
Accommodation Use	Yes	Classification Consistency (P)	0.92	0.94	0.95	0.82
		Probability of Chance	0.51	0.66	0.88	0.39
		Kappa (k)	0.84	0.83	0.58	0.70
		Classification Accuracy	0.95	0.96	0.96	0.87

Table J-4. Indices for Classification Consistency and Accuracy, ELA Grade 6

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.92	0.88	0.91	0.72
		Probability of Chance	0.66	0.50	0.77	0.28
		Kappa (k)	0.77	0.77	0.60	0.61
		Classification Accuracy	0.94	0.92	0.94	0.79
	Male	Classification Consistency (P)	0.91	0.89	0.93	0.73
		Probability of Chance	0.58	0.54	0.84	0.29
		Kappa (k)	0.78	0.76	0.58	0.62
		Classification Accuracy	0.93	0.92	0.95	0.80
Race/Ethnicity	White	Classification Consistency (P)	0.93	0.87	0.90	0.70
		Probability of Chance	0.71	0.50	0.76	0.29
		Kappa (k)	0.75	0.75	0.59	0.59
		Classification Accuracy	0.94	0.91	0.93	0.78
	African American	Classification Consistency (P)	0.88	0.94	0.98	0.81
		Probability of Chance	0.51	0.78	0.96	0.44
		Kappa (k)	0.76	0.72	0.54	0.66
		Classification Accuracy	0.91	0.95	0.99	0.86
	Hispanic	Classification Consistency (P)	0.89	0.90	0.96	0.75
		Probability of Chance	0.53	0.62	0.90	0.32
		Kappa (k)	0.76	0.74	0.58	0.63
		Classification Accuracy	0.91	0.93	0.97	0.81
	Asian	Classification Consistency (P)	0.91	0.89	0.92	0.72
		Probability of Chance	0.63	0.52	0.79	0.28
		Kappa (k)	0.76	0.77	0.61	0.61
		Classification Accuracy	0.93	0.93	0.95	0.80
	American Indian	Classification Consistency (P)	0.88	0.90	0.97	0.76
		Probability of Chance	0.51	0.67	0.94	0.35
		Kappa (k)	0.76	0.69	0.50	0.62
		Classification Accuracy	0.91	0.93	0.98	0.83
Two or More	Classification Consistency (P)	0.91	0.90	0.93	0.74	
	Probability of Chance	0.58	0.54	0.83	0.29	
	Kappa (k)	0.79	0.77	0.59	0.63	
	Classification Accuracy	0.93	0.92	0.95	0.81	
Limited English Proficiency	Yes	Classification Consistency (P)	0.86	0.94	0.99	0.79
		Probability of Chance	0.51	0.81	0.98	0.43
		Kappa (k)	0.71	0.68	0.49	0.62
		Classification Accuracy	0.89	0.96	0.99	0.84
Disability Status	Yes	Classification Consistency (P)	0.88	0.95	0.98	0.80
		Probability of Chance	0.52	0.79	0.96	0.45
		Kappa (k)	0.74	0.75	0.55	0.65
		Classification Accuracy	0.91	0.96	0.99	0.86

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.89	0.90	0.96	0.75
		Probability of Chance	0.52	0.63	0.92	0.33
		Kappa (k)	0.76	0.74	0.53	0.63
		Classification Accuracy	0.91	0.93	0.97	0.82
Accommodation Use	Yes	Classification Consistency (P)	0.88	0.92	0.98	0.78
		Probability of Chance	0.50	0.67	0.95	0.36
		Kappa (k)	0.76	0.76	0.57	0.65
		Classification Accuracy	0.92	0.94	0.99	0.84

Table J-5. Indices for Classification Consistency and Accuracy, ELA Grade 7

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.90	0.85	0.93	0.69
		Probability of Chance	0.68	0.51	0.85	0.31
		Kappa (k)	0.69	0.70	0.55	0.55
		Classification Accuracy	0.92	0.89	0.94	0.76
	Male	Classification Consistency (P)	0.91	0.88	0.95	0.74
		Probability of Chance	0.60	0.53	0.86	0.29
		Kappa (k)	0.76	0.75	0.62	0.63
		Classification Accuracy	0.93	0.92	0.96	0.81
Race/Ethnicity	White	Classification Consistency (P)	0.92	0.87	0.92	0.72
		Probability of Chance	0.72	0.50	0.79	0.30
		Kappa (k)	0.73	0.74	0.64	0.60
		Classification Accuracy	0.94	0.91	0.94	0.79
	African American	Classification Consistency (P)	0.87	0.93	0.98	0.79
		Probability of Chance	0.50	0.74	0.96	0.41
		Kappa (k)	0.75	0.73	0.59	0.65
		Classification Accuracy	0.91	0.95	0.99	0.85
	Hispanic	Classification Consistency (P)	0.89	0.90	0.96	0.75
		Probability of Chance	0.54	0.60	0.91	0.32
		Kappa (k)	0.75	0.74	0.59	0.63
		Classification Accuracy	0.92	0.93	0.97	0.82
	Asian	Classification Consistency (P)	0.90	0.88	0.93	0.71
		Probability of Chance	0.66	0.51	0.79	0.28
		Kappa (k)	0.71	0.75	0.67	0.60
		Classification Accuracy	0.93	0.91	0.95	0.80
	American Indian	Classification Consistency (P)	0.87	0.90	0.98	0.75
		Probability of Chance	0.52	0.66	0.94	0.34
		Kappa (k)	0.72	0.71	0.63	0.61
		Classification Accuracy	0.90	0.93	0.98	0.81
Two or More	Classification Consistency (P)	0.91	0.89	0.94	0.74	
	Probability of Chance	0.60	0.53	0.83	0.28	
	Kappa (k)	0.77	0.77	0.64	0.63	
	Classification Accuracy	0.93	0.93	0.96	0.81	
Limited English Proficiency	Yes	Classification Consistency (P)	0.86	0.93	0.99	0.78
		Probability of Chance	0.50	0.79	0.98	0.42
		Kappa (k)	0.71	0.68	0.48	0.62
		Classification Accuracy	0.89	0.95	0.99	0.84
Disability Status	Yes	Classification Consistency (P)	0.88	0.95	0.99	0.81
		Probability of Chance	0.52	0.80	0.97	0.46
		Kappa (k)	0.74	0.73	0.58	0.65
		Classification Accuracy	0.91	0.96	0.99	0.86

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.88	0.90	0.97	0.75
		Probability of Chance	0.53	0.62	0.92	0.33
		Kappa (k)	0.75	0.73	0.59	0.63
		Classification Accuracy	0.91	0.93	0.98	0.82
Accommodation Use	Yes	Classification Consistency (P)	0.90	0.92	0.96	0.79
		Probability of Chance	0.50	0.66	0.88	0.35
		Kappa (k)	0.81	0.77	0.70	0.68
		Classification Accuracy	0.93	0.94	0.97	0.85

Table J-6. Indices for Classification Consistency and Accuracy, ELA Grade 8

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.94	0.91	0.95	0.79
		Probability of Chance	0.63	0.50	0.70	0.26
		Kappa (k)	0.84	0.81	0.82	0.72
		Classification Accuracy	0.96	0.94	0.92	0.81
	Male	Classification Consistency (P)	0.91	0.90	0.94	0.75
		Probability of Chance	0.57	0.56	0.84	0.29
		Kappa (k)	0.80	0.77	0.64	0.65
		Classification Accuracy	0.93	0.93	0.96	0.82
	Non-Binary	Classification Consistency (P)	0.95	0.87	0.90	0.73
		Probability of Chance	0.86	0.50	0.69	0.31
		Kappa (k)	0.66	0.75	0.69	0.61
		Classification Accuracy	0.97	0.92	0.93	0.81
Race/Ethnicity	White	Classification Consistency (P)	0.93	0.88	0.92	0.73
		Probability of Chance	0.68	0.51	0.77	0.28
		Kappa (k)	0.77	0.76	0.65	0.62
		Classification Accuracy	0.95	0.92	0.94	0.81
	African American	Classification Consistency (P)	0.89	0.94	0.98	0.81
		Probability of Chance	0.51	0.78	0.96	0.42
		Kappa (k)	0.78	0.72	0.57	0.68
		Classification Accuracy	0.92	0.96	0.99	0.87
	Hispanic	Classification Consistency (P)	0.90	0.91	0.96	0.77
		Probability of Chance	0.52	0.64	0.91	0.33
		Kappa (k)	0.79	0.74	0.60	0.66
		Classification Accuracy	0.92	0.94	0.97	0.83
	Asian	Classification Consistency (P)	0.91	0.89	0.93	0.73
		Probability of Chance	0.66	0.51	0.75	0.27
		Kappa (k)	0.75	0.76	0.71	0.62
		Classification Accuracy	0.93	0.92	0.95	0.80
	American Indian	Classification Consistency (P)	0.88	0.91	0.97	0.77
		Probability of Chance	0.51	0.69	0.94	0.35
		Kappa (k)	0.76	0.72	0.59	0.65
		Classification Accuracy	0.91	0.94	0.98	0.83
	Two or More	Classification Consistency (P)	0.91	0.90	0.94	0.76
		Probability of Chance	0.55	0.56	0.83	0.29
		Kappa (k)	0.79	0.78	0.67	0.66
		Classification Accuracy	0.93	0.93	0.96	0.83
Limited English Proficiency	Yes	Classification Consistency (P)	0.87	0.95	0.99	0.81
		Probability of Chance	0.51	0.85	0.98	0.45
		Kappa (k)	0.73	0.67	0.43	0.65
		Classification Accuracy	0.90	0.97	0.99	0.86

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Disability Status	Yes	Classification Consistency (P)	0.89	0.96	0.99	0.84
		Probability of Chance	0.55	0.84	0.97	0.50
		Kappa (k)	0.75	0.74	0.57	0.67
		Classification Accuracy	0.92	0.97	0.99	0.88
Economically Disadvantaged	Yes	Classification Consistency (P)	0.89	0.91	0.97	0.77
		Probability of Chance	0.52	0.65	0.91	0.34
		Kappa (k)	0.78	0.75	0.60	0.66
		Classification Accuracy	0.92	0.94	0.98	0.84
Accommodation Use	Yes	Classification Consistency (P)	0.90	0.93	0.97	0.80
		Probability of Chance	0.50	0.70	0.94	0.38
		Kappa (k)	0.80	0.76	0.51	0.68
		Classification Accuracy	0.93	0.95	0.98	0.86

Table J-7. Indices for Classification Consistency and Accuracy, Mathematics Grade 3

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.92	0.90	0.94	0.76
		Probability of Chance	0.63	0.51	0.79	0.28
		Kappa (k)	0.78	0.79	0.72	0.67
		Classification Accuracy	0.94	0.93	0.96	0.83
	Male	Classification Consistency (P)	0.93	0.90	0.93	0.77
		Probability of Chance	0.66	0.50	0.73	0.27
		Kappa (k)	0.80	0.80	0.74	0.68
		Classification Accuracy	0.95	0.93	0.95	0.84
Race/Ethnicity	White	Classification Consistency (P)	0.94	0.89	0.92	0.75
		Probability of Chance	0.77	0.51	0.70	0.29
		Kappa (k)	0.75	0.78	0.72	0.65
		Classification Accuracy	0.96	0.92	0.94	0.83
	African American	Classification Consistency (P)	0.89	0.94	0.99	0.82
		Probability of Chance	0.52	0.78	0.97	0.44
		Kappa (k)	0.77	0.75	0.65	0.68
		Classification Accuracy	0.92	0.96	0.99	0.87
	Hispanic	Classification Consistency (P)	0.89	0.91	0.97	0.77
		Probability of Chance	0.53	0.60	0.90	0.31
		Kappa (k)	0.77	0.77	0.71	0.66
		Classification Accuracy	0.92	0.94	0.98	0.84
	Asian	Classification Consistency (P)	0.92	0.90	0.95	0.77
		Probability of Chance	0.63	0.51	0.75	0.27
		Kappa (k)	0.79	0.81	0.79	0.69
		Classification Accuracy	0.94	0.93	0.96	0.84
	American Indian	Classification Consistency (P)	0.88	0.91	0.98	0.78
		Probability of Chance	0.52	0.65	0.93	0.34
		Kappa (k)	0.76	0.76	0.67	0.66
		Classification Accuracy	0.92	0.94	0.99	0.84
	Two or More	Classification Consistency (P)	0.91	0.90	0.95	0.77
		Probability of Chance	0.61	0.51	0.80	0.28
		Kappa (k)	0.78	0.80	0.75	0.68
		Classification Accuracy	0.94	0.93	0.97	0.84
Limited English Proficiency	Yes	Classification Consistency (P)	0.87	0.90	0.98	0.75
		Probability of Chance	0.52	0.64	0.96	0.34
		Kappa (k)	0.73	0.71	0.60	0.63
		Classification Accuracy	0.90	0.93	0.99	0.82
Disability Status	Yes	Classification Consistency (P)	0.90	0.93	0.97	0.80
		Probability of Chance	0.50	0.63	0.89	0.34
		Kappa (k)	0.80	0.81	0.73	0.70
		Classification Accuracy	0.93	0.95	0.98	0.86

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.90	0.91	0.97	0.77
		Probability of Chance	0.53	0.58	0.89	0.31
		Kappa (k)	0.78	0.78	0.70	0.67
		Classification Accuracy	0.93	0.94	0.98	0.84
Accommodation Use	Yes	Classification Consistency (P)	0.90	0.97	1.00	0.87
		Probability of Chance	0.61	0.90	0.99	0.58
		Kappa (k)	0.76	0.65	0.62	0.68
		Classification Accuracy	0.93	0.98	1.00	0.91

Table J-8. Indices for Classification Consistency and Accuracy, Mathematics Grade 4

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.90	0.91	0.95	0.75
		Probability of Chance	0.64	0.51	0.78	0.28
		Kappa (k)	0.72	0.81	0.74	0.66
		Classification Accuracy	0.93	0.94	0.96	0.83
	Male	Classification Consistency (P)	0.92	0.91	0.93	0.76
		Probability of Chance	0.67	0.50	0.71	0.26
		Kappa (k)	0.76	0.82	0.76	0.67
		Classification Accuracy	0.95	0.94	0.96	0.84
Race/Ethnicity	White	Classification Consistency (P)	0.93	0.90	0.92	0.75
		Probability of Chance	0.78	0.51	0.69	0.29
		Kappa (k)	0.68	0.79	0.74	0.65
		Classification Accuracy	0.95	0.93	0.95	0.83
	African American	Classification Consistency (P)	0.86	0.96	0.99	0.81
		Probability of Chance	0.52	0.80	0.97	0.45
		Kappa (k)	0.71	0.77	0.69	0.65
		Classification Accuracy	0.90	0.97	0.99	0.87
	Hispanic	Classification Consistency (P)	0.87	0.92	0.97	0.76
		Probability of Chance	0.53	0.62	0.89	0.32
		Kappa (k)	0.71	0.80	0.73	0.65
		Classification Accuracy	0.91	0.95	0.98	0.84
	Asian	Classification Consistency (P)	0.90	0.91	0.95	0.76
		Probability of Chance	0.66	0.51	0.73	0.27
		Kappa (k)	0.70	0.82	0.80	0.67
		Classification Accuracy	0.93	0.94	0.97	0.84
	American Indian	Classification Consistency (P)	0.86	0.92	0.97	0.76
		Probability of Chance	0.54	0.64	0.90	0.33
		Kappa (k)	0.70	0.79	0.73	0.64
		Classification Accuracy	0.90	0.95	0.98	0.83
	Two or More	Classification Consistency (P)	0.89	0.91	0.95	0.76
		Probability of Chance	0.61	0.53	0.79	0.28
		Kappa (k)	0.73	0.81	0.76	0.66
		Classification Accuracy	0.93	0.94	0.97	0.84
Limited English Proficiency	Yes	Classification Consistency (P)	0.85	0.92	0.98	0.75
		Probability of Chance	0.52	0.68	0.94	0.35
		Kappa (k)	0.69	0.76	0.71	0.62
		Classification Accuracy	0.90	0.95	0.99	0.83
Disability Status	Yes	Classification Consistency (P)	0.87	0.94	0.97	0.79
		Probability of Chance	0.50	0.66	0.89	0.35
		Kappa (k)	0.74	0.82	0.76	0.67
		Classification Accuracy	0.91	0.96	0.98	0.86

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.87	0.92	0.97	0.76
		Probability of Chance	0.54	0.60	0.89	0.31
		Kappa (k)	0.72	0.80	0.72	0.65
		Classification Accuracy	0.91	0.94	0.98	0.84
Accommodation Use	Yes	Classification Consistency (P)	0.84	0.97	1.00	0.81
		Probability of Chance	0.56	0.91	0.99	0.53
		Kappa (k)	0.64	0.71	0.55	0.59
		Classification Accuracy	0.89	0.98	1.00	0.87

Table J-9. Indices for Classification Consistency and Accuracy, Mathematics Grade 5

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.91	0.90	0.95	0.76
		Probability of Chance	0.61	0.51	0.81	0.28
		Kappa (k)	0.78	0.79	0.71	0.66
		Classification Accuracy	0.94	0.93	0.96	0.83
	Male	Classification Consistency (P)	0.93	0.90	0.93	0.77
		Probability of Chance	0.62	0.50	0.76	0.27
		Kappa (k)	0.81	0.81	0.73	0.68
		Classification Accuracy	0.95	0.93	0.95	0.83
Race/Ethnicity	White	Classification Consistency (P)	0.93	0.89	0.92	0.75
		Probability of Chance	0.73	0.51	0.73	0.29
		Kappa (k)	0.74	0.78	0.71	0.64
		Classification Accuracy	0.95	0.92	0.94	0.82
	African American	Classification Consistency (P)	0.89	0.95	0.99	0.84
		Probability of Chance	0.55	0.80	0.98	0.49
		Kappa (k)	0.76	0.74	0.64	0.68
		Classification Accuracy	0.92	0.96	0.99	0.88
	Hispanic	Classification Consistency (P)	0.89	0.91	0.97	0.78
		Probability of Chance	0.51	0.61	0.92	0.33
		Kappa (k)	0.78	0.77	0.70	0.67
		Classification Accuracy	0.92	0.94	0.98	0.84
	Asian	Classification Consistency (P)	0.91	0.91	0.94	0.76
		Probability of Chance	0.63	0.50	0.73	0.26
		Kappa (k)	0.76	0.82	0.77	0.67
		Classification Accuracy	0.94	0.93	0.95	0.82
	American Indian	Classification Consistency (P)	0.89	0.90	0.98	0.77
		Probability of Chance	0.51	0.62	0.94	0.33
		Kappa (k)	0.77	0.73	0.72	0.65
		Classification Accuracy	0.92	0.93	0.99	0.84
Two or More	Classification Consistency (P)	0.90	0.90	0.95	0.75	
	Probability of Chance	0.57	0.52	0.82	0.28	
	Kappa (k)	0.76	0.80	0.71	0.66	
	Classification Accuracy	0.93	0.93	0.96	0.82	
Limited English Proficiency	Yes	Classification Consistency (P)	0.87	0.93	0.99	0.78
		Probability of Chance	0.50	0.71	0.97	0.39
		Kappa (k)	0.74	0.74	0.61	0.65
		Classification Accuracy	0.91	0.95	0.99	0.85
Disability Status	Yes	Classification Consistency (P)	0.90	0.94	0.98	0.82
		Probability of Chance	0.51	0.68	0.93	0.40
		Kappa (k)	0.80	0.81	0.70	0.70
		Classification Accuracy	0.93	0.96	0.99	0.87

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.90	0.91	0.97	0.78
		Probability of Chance	0.51	0.59	0.91	0.32
		Kappa (k)	0.79	0.78	0.68	0.67
		Classification Accuracy	0.93	0.94	0.98	0.84
Accommodation Use	Yes	Classification Consistency (P)	0.89	0.97	1.00	0.86
		Probability of Chance	0.63	0.90	0.99	0.61
		Kappa (k)	0.70	0.70	0.65	0.63
		Classification Accuracy	0.92	0.98	1.00	0.90

Table J-10. Indices for Classification Consistency and Accuracy, Mathematics Grade 6

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.90	0.91	0.96	0.77
		Probability of Chance	0.59	0.52	0.87	0.30
		Kappa (k)	0.77	0.81	0.70	0.68
		Classification Accuracy	0.93	0.93	0.97	0.84
	Male	Classification Consistency (P)	0.91	0.91	0.96	0.78
		Probability of Chance	0.60	0.51	0.84	0.29
		Kappa (k)	0.78	0.82	0.71	0.69
		Classification Accuracy	0.94	0.94	0.97	0.84
Race/Ethnicity	White	Classification Consistency (P)	0.92	0.90	0.95	0.77
		Probability of Chance	0.71	0.50	0.82	0.31
		Kappa (k)	0.74	0.79	0.70	0.66
		Classification Accuracy	0.94	0.93	0.96	0.83
	African American	Classification Consistency (P)	0.88	0.96	1.00	0.84
		Probability of Chance	0.57	0.84	0.99	0.52
		Kappa (k)	0.72	0.76	0.67	0.66
		Classification Accuracy	0.92	0.97	1.00	0.89
	Hispanic	Classification Consistency (P)	0.87	0.93	0.99	0.78
		Probability of Chance	0.50	0.66	0.96	0.35
		Kappa (k)	0.74	0.78	0.69	0.66
		Classification Accuracy	0.91	0.95	0.99	0.85
	Asian	Classification Consistency (P)	0.90	0.91	0.96	0.77
		Probability of Chance	0.60	0.51	0.81	0.28
		Kappa (k)	0.76	0.82	0.77	0.69
		Classification Accuracy	0.93	0.94	0.97	0.84
	American Indian	Classification Consistency (P)	0.90	0.95	0.99	0.84
		Probability of Chance	0.50	0.67	0.93	0.38
		Kappa (k)	0.80	0.85	0.80	0.74
		Classification Accuracy	0.91	0.96	0.98	0.86
Two or More	Classification Consistency (P)	0.89	0.92	0.97	0.78	
	Probability of Chance	0.54	0.56	0.89	0.30	
	Kappa (k)	0.77	0.82	0.70	0.68	
	Classification Accuracy	0.93	0.94	0.98	0.84	
Limited English Proficiency	Yes	Classification Consistency (P)	0.86	0.96	1.00	0.81
		Probability of Chance	0.53	0.82	0.99	0.47
		Kappa (k)	0.70	0.76	0.71	0.64
		Classification Accuracy	0.90	0.97	1.00	0.87
Disability Status	Yes	Classification Consistency (P)	0.89	0.96	0.99	0.83
		Probability of Chance	0.53	0.78	0.97	0.47
		Kappa (k)	0.76	0.80	0.70	0.69
		Classification Accuracy	0.92	0.97	0.99	0.88

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.88	0.92	0.99	0.79
		Probability of Chance	0.50	0.65	0.96	0.35
		Kappa (k)	0.76	0.78	0.67	0.68
		Classification Accuracy	0.91	0.95	0.99	0.85
Accommodation Use	Yes	Classification Consistency (P)	0.92	0.99	1.00	0.91
		Probability of Chance	0.74	0.94	1.00	0.73
		Kappa (k)	0.70	0.79	0.74	0.65
		Classification Accuracy	0.93	0.99	1.00	0.92

Table J-11. Indices for Classification Consistency and Accuracy, Mathematics Grade 7

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.90	0.91	0.97	0.78
		Probability of Chance	0.53	0.55	0.91	0.31
		Kappa (k)	0.79	0.80	0.65	0.69
		Classification Accuracy	0.93	0.93	0.98	0.84
	Male	Classification Consistency (P)	0.91	0.92	0.96	0.79
		Probability of Chance	0.54	0.53	0.89	0.30
		Kappa (k)	0.81	0.82	0.66	0.71
		Classification Accuracy	0.94	0.94	0.98	0.85
Race/Ethnicity	White	Classification Consistency (P)	0.91	0.90	0.96	0.77
		Probability of Chance	0.62	0.51	0.87	0.30
		Kappa (k)	0.77	0.80	0.65	0.67
		Classification Accuracy	0.94	0.93	0.97	0.84
	African American	Classification Consistency (P)	0.91	0.97	1.00	0.87
		Probability of Chance	0.63	0.87	0.99	0.61
		Kappa (k)	0.75	0.78	0.62	0.68
		Classification Accuracy	0.93	0.98	1.00	0.91
	Hispanic	Classification Consistency (P)	0.89	0.94	0.99	0.82
		Probability of Chance	0.50	0.71	0.97	0.40
		Kappa (k)	0.78	0.80	0.60	0.70
		Classification Accuracy	0.92	0.96	0.99	0.87
	Asian	Classification Consistency (P)	0.90	0.93	0.96	0.78
		Probability of Chance	0.55	0.53	0.84	0.29
		Kappa (k)	0.77	0.85	0.72	0.69
		Classification Accuracy	0.93	0.95	0.97	0.85
	American Indian	Classification Consistency (P)	0.88	0.95	0.99	0.82
		Probability of Chance	0.52	0.74	0.99	0.43
		Kappa (k)	0.76	0.79	0.62	0.68
		Classification Accuracy	0.92	0.96	1.00	0.88
	Two or More	Classification Consistency (P)	0.90	0.92	0.97	0.80
		Probability of Chance	0.51	0.58	0.92	0.32
		Kappa (k)	0.79	0.82	0.70	0.70
		Classification Accuracy	0.93	0.94	0.98	0.85
Limited English Proficiency	Yes	Classification Consistency (P)	0.89	0.97	1.00	0.85
		Probability of Chance	0.58	0.85	0.99	0.55
		Kappa (k)	0.73	0.77	0.46	0.67
		Classification Accuracy	0.92	0.98	1.00	0.89
Disability Status	Yes	Classification Consistency (P)	0.91	0.97	0.99	0.88
		Probability of Chance	0.61	0.83	0.98	0.58
		Kappa (k)	0.78	0.81	0.64	0.71
		Classification Accuracy	0.94	0.98	1.00	0.91

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Economically Disadvantaged	Yes	Classification Consistency (P)	0.89	0.94	0.99	0.82
		Probability of Chance	0.50	0.70	0.97	0.40
		Kappa (k)	0.78	0.79	0.60	0.69
		Classification Accuracy	0.92	0.95	0.99	0.86
Accommodation Use	Yes	Classification Consistency (P)	0.93	0.99	1.00	0.92
		Probability of Chance	0.80	0.97	1.00	0.80
		Kappa (k)	0.66	0.66	0.58	0.61
		Classification Accuracy	0.95	0.99	1.00	0.95

Table J-12. Indices for Classification Consistency and Accuracy, Mathematics Grade 8

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.88	0.92	0.96	0.77
		Probability of Chance	0.55	0.57	0.88	0.30
		Kappa (k)	0.74	0.81	0.70	0.67
		Classification Accuracy	0.92	0.94	0.98	0.83
	Male	Classification Consistency (P)	0.89	0.93	0.96	0.78
		Probability of Chance	0.54	0.56	0.86	0.30
		Kappa (k)	0.77	0.83	0.72	0.69
		Classification Accuracy	0.92	0.95	0.97	0.84
	Non-Binary	Classification Consistency (P)	0.86	0.91	0.97	0.74
		Probability of Chance	0.68	0.55	0.83	0.32
		Kappa (k)	0.57	0.80	0.81	0.62
		Classification Accuracy	0.91	0.94	0.98	0.83
Race/Ethnicity	White	Classification Consistency (P)	0.90	0.91	0.95	0.76
		Probability of Chance	0.62	0.52	0.84	0.29
		Kappa (k)	0.73	0.81	0.71	0.66
		Classification Accuracy	0.93	0.93	0.97	0.83
	African American	Classification Consistency (P)	0.88	0.97	1.00	0.85
		Probability of Chance	0.60	0.90	0.99	0.57
		Kappa (k)	0.70	0.74	0.67	0.64
		Classification Accuracy	0.91	0.98	1.00	0.89
	Hispanic	Classification Consistency (P)	0.86	0.94	0.99	0.79
		Probability of Chance	0.50	0.74	0.96	0.40
		Kappa (k)	0.72	0.79	0.67	0.65
		Classification Accuracy	0.90	0.96	0.99	0.85
	Asian	Classification Consistency (P)	0.90	0.92	0.96	0.78
		Probability of Chance	0.58	0.53	0.79	0.28
		Kappa (k)	0.75	0.84	0.79	0.69
		Classification Accuracy	0.92	0.95	0.97	0.84
	American Indian	Classification Consistency (P)	0.86	0.94	0.99	0.79
		Probability of Chance	0.51	0.77	0.97	0.43
		Kappa (k)	0.72	0.75	0.72	0.64
		Classification Accuracy	0.90	0.96	0.99	0.86
	Two or More	Classification Consistency (P)	0.88	0.94	0.97	0.79
		Probability of Chance	0.51	0.62	0.89	0.33
		Kappa (k)	0.76	0.83	0.76	0.69
		Classification Accuracy	0.91	0.95	0.98	0.85
Limited English Proficiency	Yes	Classification Consistency (P)	0.86	0.97	1.00	0.83
		Probability of Chance	0.58	0.89	0.99	0.55
		Kappa (k)	0.68	0.73	0.74	0.63
		Classification Accuracy	0.90	0.98	1.00	0.87

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Disability Status	Yes	Classification Consistency (P)	0.89	0.97	0.99	0.85
		Probability of Chance	0.61	0.88	0.98	0.58
		Kappa (k)	0.71	0.79	0.73	0.65
		Classification Accuracy	0.92	0.98	1.00	0.90
Economically Disadvantaged	Yes	Classification Consistency (P)	0.86	0.95	0.98	0.80
		Probability of Chance	0.50	0.74	0.95	0.40
		Kappa (k)	0.73	0.80	0.66	0.66
		Classification Accuracy	0.90	0.96	0.99	0.85
Accommodation Use	Yes	Classification Consistency (P)	0.90	0.99	1.00	0.89
		Probability of Chance	0.75	0.97	1.00	0.75
		Kappa (k)	0.61	0.69	0.68	0.58
		Classification Accuracy	0.93	0.99	1.00	0.93

Table J-13. Indices for Classification Consistency and Accuracy, Science Grade 4

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.91	0.89	0.91	0.71
		Probability of Chance	0.69	0.50	0.69	0.26
		Kappa (k)	0.70	0.78	0.72	0.61
		Classification Accuracy	0.93	0.92	0.94	0.79
	Male	Classification Consistency (P)	0.91	0.90	0.91	0.72
		Probability of Chance	0.68	0.50	0.66	0.26
		Kappa (k)	0.71	0.80	0.73	0.62
		Classification Accuracy	0.94	0.93	0.93	0.79
Race/Ethnicity	White	Classification Consistency (P)	0.93	0.89	0.88	0.70
		Probability of Chance	0.80	0.53	0.61	0.28
		Kappa (k)	0.65	0.76	0.70	0.59
		Classification Accuracy	0.95	0.92	0.92	0.79
	African American	Classification Consistency (P)	0.84	0.94	0.98	0.75
		Probability of Chance	0.50	0.75	0.94	0.40
		Kappa (k)	0.67	0.74	0.66	0.59
		Classification Accuracy	0.89	0.95	0.99	0.82
	Hispanic	Classification Consistency (P)	0.86	0.90	0.95	0.71
		Probability of Chance	0.56	0.57	0.83	0.29
		Kappa (k)	0.69	0.77	0.70	0.59
		Classification Accuracy	0.90	0.93	0.96	0.79
	Asian	Classification Consistency (P)	0.89	0.90	0.92	0.71
		Probability of Chance	0.65	0.51	0.72	0.26
		Kappa (k)	0.69	0.80	0.72	0.61
		Classification Accuracy	0.92	0.93	0.94	0.80
	American Indian	Classification Consistency (P)	0.87	0.89	0.94	0.71
		Probability of Chance	0.60	0.55	0.84	0.29
		Kappa (k)	0.67	0.76	0.66	0.59
		Classification Accuracy	0.90	0.93	0.96	0.79
	Two or More	Classification Consistency (P)	0.89	0.90	0.91	0.70
		Probability of Chance	0.65	0.50	0.70	0.26
		Kappa (k)	0.67	0.80	0.71	0.60
		Classification Accuracy	0.92	0.93	0.94	0.79
Limited English Proficiency	Yes	Classification Consistency (P)	0.84	0.91	0.97	0.72
		Probability of Chance	0.53	0.65	0.92	0.33
		Kappa (k)	0.66	0.73	0.63	0.57
		Classification Accuracy	0.89	0.93	0.98	0.80
Disability Status	Yes	Classification Consistency (P)	0.86	0.91	0.96	0.73
		Probability of Chance	0.53	0.60	0.84	0.31
		Kappa (k)	0.69	0.78	0.73	0.60
		Classification Accuracy	0.90	0.94	0.97	0.80
Economically Disadvantaged	Yes	Classification Consistency (P)	0.87	0.90	0.95	0.71
		Probability of Chance	0.57	0.55	0.82	0.29
		Kappa (k)	0.69	0.77	0.70	0.60
		Classification Accuracy	0.91	0.93	0.96	0.79

Table J-14. Indices for Classification Consistency and Accuracy, Science Grade 8

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.90	0.88	0.91	0.69
		Probability of Chance	0.66	0.50	0.69	0.26
		Kappa (k)	0.69	0.76	0.72	0.59
		Classification Accuracy	0.93	0.91	0.93	0.77
	Male	Classification Consistency (P)	0.90	0.89	0.91	0.71
		Probability of Chance	0.64	0.50	0.65	0.25
		Kappa (k)	0.72	0.79	0.75	0.61
		Classification Accuracy	0.93	0.92	0.93	0.78
	Non-Binary	Classification Consistency (P)	0.95	0.86	0.89	0.70
		Probability of Chance	0.90	0.57	0.54	0.31
		Kappa (k)	0.51	0.67	0.75	0.57
		Classification Accuracy	0.97	0.89	0.92	0.78
Race/Ethnicity	White	Classification Consistency (P)	0.92	0.88	0.90	0.70
		Probability of Chance	0.74	0.51	0.61	0.26
		Kappa (k)	0.67	0.76	0.73	0.59
		Classification Accuracy	0.94	0.91	0.92	0.77
	African American	Classification Consistency (P)	0.84	0.92	0.98	0.74
		Probability of Chance	0.50	0.73	0.93	0.40
		Kappa (k)	0.67	0.71	0.67	0.57
		Classification Accuracy	0.88	0.95	0.98	0.81
	Hispanic	Classification Consistency (P)	0.86	0.89	0.95	0.71
		Probability of Chance	0.54	0.57	0.82	0.29
		Kappa (k)	0.70	0.75	0.70	0.58
		Classification Accuracy	0.90	0.92	0.96	0.78
	Asian	Classification Consistency (P)	0.90	0.88	0.91	0.69
		Probability of Chance	0.68	0.50	0.65	0.26
		Kappa (k)	0.67	0.77	0.74	0.58
		Classification Accuracy	0.93	0.92	0.93	0.78
	American Indian	Classification Consistency (P)	0.85	0.89	0.95	0.69
		Probability of Chance	0.54	0.58	0.84	0.30
		Kappa (k)	0.67	0.74	0.67	0.56
		Classification Accuracy	0.90	0.92	0.96	0.78
	Two or More	Classification Consistency (P)	0.88	0.89	0.92	0.69
		Probability of Chance	0.60	0.51	0.72	0.26
		Kappa (k)	0.71	0.77	0.71	0.58
		Classification Accuracy	0.92	0.92	0.94	0.78
Limited English Proficiency	Yes	Classification Consistency (P)	0.83	0.91	0.98	0.73
		Probability of Chance	0.50	0.75	0.96	0.39
		Kappa (k)	0.66	0.66	0.57	0.55
		Classification Accuracy	0.88	0.94	0.99	0.80

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Disability Status	Yes	Classification Consistency (P)	0.84	0.93	0.97	0.75
		Probability of Chance	0.50	0.70	0.89	0.38
		Kappa (k)	0.69	0.76	0.74	0.59
		Classification Accuracy	0.89	0.95	0.98	0.81
Economically Disadvantaged	Yes	Classification Consistency (P)	0.86	0.89	0.95	0.70
		Probability of Chance	0.54	0.56	0.81	0.29
		Kappa (k)	0.69	0.76	0.71	0.58
		Classification Accuracy	0.90	0.92	0.96	0.78

Table J-15. Indices for Classification Consistency and Accuracy, Social Studies Grade 4

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.91	0.89	0.90	0.71
		Probability of Chance	0.67	0.52	0.65	0.27
		Kappa (k)	0.73	0.78	0.71	0.60
		Classification Accuracy	0.94	0.93	0.93	0.79
	Male	Classification Consistency (P)	0.91	0.91	0.90	0.72
		Probability of Chance	0.65	0.52	0.62	0.26
		Kappa (k)	0.75	0.81	0.73	0.62
		Classification Accuracy	0.94	0.93	0.93	0.80
Race/Ethnicity	White	Classification Consistency (P)	0.93	0.90	0.88	0.71
		Probability of Chance	0.77	0.57	0.58	0.29
		Kappa (k)	0.70	0.76	0.71	0.59
		Classification Accuracy	0.95	0.93	0.91	0.79
	African American	Classification Consistency (P)	0.86	0.92	0.97	0.75
		Probability of Chance	0.50	0.65	0.91	0.39
		Kappa (k)	0.71	0.77	0.66	0.59
		Classification Accuracy	0.90	0.94	0.98	0.82
	Hispanic	Classification Consistency (P)	0.87	0.90	0.94	0.71
		Probability of Chance	0.55	0.52	0.79	0.28
		Kappa (k)	0.71	0.79	0.71	0.60
		Classification Accuracy	0.91	0.92	0.96	0.79
	Asian	Classification Consistency (P)	0.89	0.90	0.92	0.71
		Probability of Chance	0.62	0.50	0.68	0.26
		Kappa (k)	0.72	0.79	0.75	0.61
		Classification Accuracy	0.93	0.92	0.94	0.79
	American Indian	Classification Consistency (P)	0.90	0.91	0.95	0.76
		Probability of Chance	0.54	0.51	0.78	0.28
		Kappa (k)	0.78	0.81	0.77	0.67
		Classification Accuracy	0.93	0.93	0.96	0.82
	Two or More	Classification Consistency (P)	0.91	0.90	0.91	0.73
		Probability of Chance	0.64	0.51	0.67	0.26
		Kappa (k)	0.75	0.79	0.74	0.63
		Classification Accuracy	0.93	0.93	0.94	0.80
Limited English Proficiency	Yes	Classification Consistency (P)	0.85	0.89	0.96	0.71
		Probability of Chance	0.51	0.57	0.88	0.31
		Kappa (k)	0.69	0.75	0.63	0.57
		Classification Accuracy	0.90	0.92	0.97	0.79
Disability Status	Yes	Classification Consistency (P)	0.87	0.91	0.95	0.74
		Probability of Chance	0.51	0.56	0.82	0.31
		Kappa (k)	0.73	0.80	0.73	0.62
		Classification Accuracy	0.91	0.94	0.96	0.81
Economically Disadvantaged	Yes	Classification Consistency (P)	0.87	0.89	0.93	0.71
		Probability of Chance	0.55	0.52	0.79	0.28
		Kappa (k)	0.72	0.78	0.69	0.60
		Classification Accuracy	0.92	0.93	0.95	0.79

Table J-16. Indices for Classification Consistency and Accuracy, Social Studies Grade 8

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.92	0.90	0.90	0.72
		Probability of Chance	0.70	0.52	0.63	0.27
		Kappa (k)	0.73	0.78	0.73	0.62
		Classification Accuracy	0.95	0.92	0.93	0.79
	Male	Classification Consistency (P)	0.91	0.91	0.91	0.73
		Probability of Chance	0.65	0.51	0.62	0.26
		Kappa (k)	0.75	0.81	0.77	0.64
		Classification Accuracy	0.94	0.93	0.93	0.81
	Non-Binary	Classification Consistency (P)	0.98	0.88	0.90	0.76
		Probability of Chance	0.95	0.75	0.51	0.38
		Kappa (k)	0.50	0.51	0.79	0.61
		Classification Accuracy	0.98	0.92	0.92	0.83
Race/Ethnicity	White	Classification Consistency (P)	0.93	0.90	0.89	0.73
		Probability of Chance	0.76	0.56	0.57	0.28
		Kappa (k)	0.72	0.78	0.74	0.62
		Classification Accuracy	0.95	0.93	0.92	0.80
	African American	Classification Consistency (P)	0.86	0.90	0.96	0.74
		Probability of Chance	0.50	0.59	0.87	0.33
		Kappa (k)	0.73	0.76	0.71	0.61
		Classification Accuracy	0.91	0.93	0.97	0.81
	Hispanic	Classification Consistency (P)	0.88	0.89	0.93	0.71
		Probability of Chance	0.58	0.51	0.77	0.27
		Kappa (k)	0.72	0.78	0.72	0.61
		Classification Accuracy	0.92	0.92	0.95	0.79
	Asian	Classification Consistency (P)	0.93	0.89	0.90	0.72
		Probability of Chance	0.72	0.53	0.60	0.27
		Kappa (k)	0.74	0.77	0.75	0.62
		Classification Accuracy	0.95	0.92	0.93	0.80
	American Indian	Classification Consistency (P)	0.89	0.88	0.94	0.71
		Probability of Chance	0.57	0.51	0.80	0.28
		Kappa (k)	0.74	0.75	0.70	0.60
		Classification Accuracy	0.92	0.91	0.96	0.79
	Two or More	Classification Consistency (P)	0.91	0.90	0.91	0.73
		Probability of Chance	0.64	0.50	0.66	0.26
		Kappa (k)	0.76	0.80	0.74	0.63
		Classification Accuracy	0.94	0.93	0.94	0.82
Limited English Proficiency	Yes	Classification Consistency (P)	0.84	0.89	0.98	0.71
		Probability of Chance	0.50	0.63	0.93	0.34
		Kappa (k)	0.68	0.70	0.64	0.56
		Classification Accuracy	0.89	0.92	0.98	0.79

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Disability Status	Yes	Classification Consistency (P)	0.86	0.92	0.97	0.75
		Probability of Chance	0.50	0.64	0.88	0.36
		Kappa (k)	0.72	0.77	0.77	0.62
		Classification Accuracy	0.90	0.94	0.98	0.82
Economically Disadvantaged	Yes	Classification Consistency (P)	0.88	0.89	0.94	0.72
		Probability of Chance	0.57	0.51	0.77	0.27
		Kappa (k)	0.73	0.78	0.73	0.61
		Classification Accuracy	0.92	0.92	0.95	0.79

Table J-17. Indices for Classification Consistency and Accuracy, Social Studies Grade 10

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Gender	Female	Classification Consistency (P)	0.90	0.89	0.90	0.69
		Probability of Chance	0.62	0.50	0.64	0.25
		Kappa (k)	0.72	0.78	0.72	0.59
		Classification Accuracy	0.93	0.92	0.93	0.78
	Male	Classification Consistency (P)	0.89	0.90	0.91	0.72
		Probability of Chance	0.59	0.50	0.64	0.25
		Kappa (k)	0.74	0.81	0.76	0.62
		Classification Accuracy	0.93	0.93	0.94	0.80
	Non-Binary	Classification Consistency (P)	0.95	0.86	0.90	0.72
		Probability of Chance	0.79	0.54	0.52	0.29
		Kappa (k)	0.76	0.70	0.80	0.61
		Classification Accuracy	0.96	0.91	0.93	0.80
Race/Ethnicity	White	Classification Consistency (P)	0.91	0.89	0.89	0.69
		Probability of Chance	0.67	0.50	0.60	0.25
		Kappa (k)	0.72	0.77	0.72	0.58
		Classification Accuracy	0.94	0.92	0.93	0.78
	African American	Classification Consistency (P)	0.86	0.94	0.97	0.77
		Probability of Chance	0.52	0.71	0.88	0.42
		Kappa (k)	0.71	0.78	0.70	0.60
		Classification Accuracy	0.90	0.95	0.98	0.83
	Hispanic	Classification Consistency (P)	0.86	0.91	0.94	0.71
		Probability of Chance	0.51	0.57	0.79	0.30
		Kappa (k)	0.72	0.78	0.71	0.59
		Classification Accuracy	0.91	0.93	0.96	0.80
	Asian	Classification Consistency (P)	0.89	0.88	0.92	0.70
		Probability of Chance	0.63	0.50	0.64	0.25
		Kappa (k)	0.72	0.77	0.77	0.60
		Classification Accuracy	0.93	0.92	0.94	0.79
	American Indian	Classification Consistency (P)	0.84	0.90	0.94	0.69
		Probability of Chance	0.52	0.58	0.79	0.30
		Kappa (k)	0.67	0.77	0.71	0.56
		Classification Accuracy	0.89	0.93	0.96	0.78
	Two or More	Classification Consistency (P)	0.89	0.89	0.93	0.72
		Probability of Chance	0.57	0.51	0.67	0.26
		Kappa (k)	0.75	0.78	0.79	0.62
		Classification Accuracy	0.92	0.93	0.95	0.79
Limited English Proficiency	Yes	Classification Consistency (P)	0.84	0.94	0.98	0.76
		Probability of Chance	0.54	0.81	0.96	0.48
		Kappa (k)	0.65	0.68	0.56	0.54
		Classification Accuracy	0.88	0.96	0.99	0.83

Group	Category	Indices	Cut 1	Cut 2	Cut 3	All Cuts
Disability Status	Yes	Classification Consistency (P)	0.86	0.94	0.97	0.77
		Probability of Chance	0.53	0.74	0.89	0.45
		Kappa (k)	0.69	0.77	0.74	0.58
		Classification Accuracy	0.90	0.96	0.98	0.84
Economically Disadvantaged	Yes	Classification Consistency (P)	0.93	0.95	0.96	0.84
		Probability of Chance	0.50	0.57	0.72	0.35
		Kappa (k)	0.85	0.88	0.86	0.75
		Classification Accuracy	0.88	0.96	0.95	0.80

APPENDIX K
WISCONSIN STANDARD PERFORMANCE
INDEX SCORE COMPUTATION

Table of Contents

Standard Performance Index Computation.....	3
Estimation of the Prior Distribution of T_j	4
Check on Consistency and Adjustment of Weight Given to Prior Estimate	7
Possible Violations of the Assumptions	8
References.....	9

Technical details of the SPI estimation procedure described in this section are based on description of the SPI computation methodology included in the *TerraNova* 2nd Edition Technical Report (CTB/McGraw-Hill, 2000).

The Standard Performance Index (SPI) is an estimate of the true score (estimated proportion of total, or maximum, points possible) for a content standard based on the performance of a given student. Because most standards are measured by a relatively small number of items, a Bayesian procedure that takes into account the overall test performance is used to improve the reliability of the standard scores. Given a student's scale score on the test, item response theory (IRT) is used, via the 3-parameter logistic (3PL) model for MC items and the 2-parameter-partial credit (2PPC) model for CR items, to compute the estimated proportion of the maximum points obtained for that standard.

The estimated proportion of the maximum points obtained for the standard provides the initial (Bayesian prior) estimate of the student's mastery score. If this initial estimate is consistent with the student's observed proportion, as indicated by a chi-square test, the two scores are combined as a weighted average to obtain the SPI score (the estimated true score). The appropriate weight for the Bayesian prior estimate is computed as a function of the standard error (SE) of the scale score on which it is based: the smaller the SE, the larger the weight. If the prior estimate and the observed proportion differ significantly, the observed proportion of the maximum score is used without the prior estimate to compute the student's score on that objective.

Standard Performance Index Computation

The standard performance index (SPI) is an estimated true score (estimated proportion of total or maximum points obtained) based on the performance of a given examinee for the items in a given learning strand. Assume a k -item test is composed of j strands with a maximum possible raw score of n . Also assume that each item contributes to, at most, one strand, and the k_j items in strand j contribute a maximum of n_j points. Define X_j as the observed raw score on strand j . The true score is

$$T_j \equiv E(X_j/n_j).$$

It is assumed that there is information available about the examinee in addition to the strand score, and this information provides a prior distribution for T_j . This prior distribution of T_j for a given examinee is assumed to be $\beta(r_j, s_j)$:

$$g(T_j) = \frac{(r_j + s_j - 1)! T_j^{r_j - 1} (1 - T_j)^{s_j - 1}}{(r_j - 1)! (s_j - 1)!} \quad (1)$$

for $0 \leq T_j \leq 1; r_j, s_j > 0$. Estimates of r_j and s_j are derived from IRT (Lord, 1980).

It is assumed that X_j follows a binomial distribution, given T_j :

$$p(X_j = x_j | T_j) = \text{Binomial} \left(n_j, T_j = \sum_{i=1}^{k_j} T_i/n_j \right)$$

where

T_j is the expected value of the score for item i in strand j for a given θ .

Given these assumptions, the posterior distribution of T_j , given X_j , is

$$g(T_j | X_j = x_j) = \beta(p_j, q_j), \quad (2)$$

with

$$p_j = r_j + x_j \quad (3)$$

and

$$q_j = s_j + n_j - x_j. \quad (4)$$

The SPI is defined to be the mean of this posterior distribution:

$$\tilde{T}_j = \frac{p_j}{p_j + q_j}.$$

Following Novick and Jackson (1974, p. 119), a mastery band is created to be the $C\%$ central credibility interval for T_j . It is obtained by identifying the values that place $\frac{1}{2}(100 - C)\%$ of the $\beta(p_j, q_j)$ density in each tail of the distribution.

Estimation of the Prior Distribution of T_j

The k items in each test are scaled together using a generalized IRT model (3PL/2PPC) that fits a three-parameter logistic model (3PL) to the MC items and a generalized partial-credit model (2PPC) to the CR items (Yen, 1993).

The 3PL model is

$$P_i(\theta) = P(X_i = 1 | \theta) = c_i + \frac{1 - c_i}{1 + \exp[-1.7A_i(\theta - B_i)]}, \quad (5)$$

where

A_i is the discrimination, B_i is the location, and C_i is the guessing parameter for item i .

A generalization of Master's (1982) partial credit (2PPC) model was used for the CR items. The 2PPC model, the same as Muraki's (1992) "generalized partial credit model," has been shown to fit response data obtained from a wide variety of mixed-item type achievement tests (Fitzpatrick, Link, Yen, Burket, Ito, & Sykes, 1996). For a CR item with 1_i score levels, integer scores were assigned that ranged from 0 to $1_i - 1$:

$$P_{im}(\theta) = P(X_i = m - 1 | \theta) = \frac{\exp(z_{im})}{\sum_{g=1}^{1_i} \exp(z_{ig})}, \quad m = 1, \dots, 1_i \quad (6)$$

where

$$z_{ig} = \alpha_i(m - 1)\theta - \sum_{h=0}^{m-1} \gamma'_{ih} \quad (7)$$

and

$$\gamma_{i0} = 0.$$

Alpha (α_i) is the item discrimination, and gamma (γ_{ih}) is related to the difficulty of the item levels; the trace lines for adjacent score levels intersect at γ_{ih}/α_i .

Item parameters estimated from the national standardization sample are used to obtain SPI values. $T_{ij}(\theta)$ is the expected score for item i in strand j , and θ is the common trait value to which the items are scaled:

$$T_{ij}(\theta) = \sum_{m=1}^{1_i} (m-1)P_{ijm}(\theta),$$

where

1_i is the number of score levels in item i , including 0.

T_j , the expected proportion of maximum score for strand j , is

$$T_j = \frac{1}{n_j} \left[\sum_{i=1}^{k_j} T_{ij}(\theta) \right]. \quad (8)$$

The expected score for item i and estimated proportion-correct of maximum score for strand j are obtained by substituting the estimate of the trait ($\hat{\theta}$) the actual trait value.

The theoretical random variation in item response vectors and resulting ($\hat{\theta}$) values for a given examinee produces the distribution $g(\hat{T}_j | \hat{\theta})$ with mean $\mu(\hat{T}_j | \theta)$ and variance $\sigma^2(\hat{T}_j | \theta)$. This distribution is used to estimate a prior distribution of T_j . Given that T_j is assumed to be distributed as a beta distribution (equation 1), the mean $[\mu(\hat{T}_j | \theta)]$ and variance $[\sigma^2(\hat{T}_j | \theta)]$ of this distribution can be expressed in terms of its parameters, r_j and s_j .

Expressing the mean and variance of the prior distribution in terms of the parameters of the beta distribution (Novick & Jackson, 1974, p. 113) produces

$$\mu(\hat{T}_j | \theta) = \frac{r_j}{r_j + s_j} \quad (9)$$

and

$$\sigma^2(\hat{T}_j | \theta) = \frac{r_j s_j}{(r_j + s_j)^2 (r_j + s_j + 1)}. \quad (10)$$

Solving these equations for r_j and s_j produces

$$r_j = \mu(\hat{T}_j | \theta) n_j^* \quad (11)$$

and

$$s_j = [1 - \mu(\hat{T}_j | \theta)] n_j^* \quad (12)$$

where

$$n_j^* = \frac{\mu(\hat{T}_j | \theta) [1 - \mu(\hat{T}_j | \theta)]}{\sigma^2(\hat{T}_j | \theta)} - 1. \quad (13)$$

Using IRT, $\sigma^2(\hat{T}_j | \theta)$ can be expressed in terms of item parameters (Lord, 1983):

$$\mu(\hat{T}_j | \theta) \approx \frac{1}{n_j} \sum_{i=1}^{k_j} \hat{T}_{ij}(\theta) \quad (14)$$

Because T_j is a monotonic transformation of θ (Lord, 1980, p.71),

$$\sigma^2(\hat{T}_j | \theta) = \sigma^2(\hat{T}_j | T_j) \approx I(T_j, \hat{T}_j)^{-1} \quad (15)$$

where

$I(T_j, \hat{T}_j)$ is the information that \hat{T}_j contributes about T_j .

Given these results, Lord (1980, p. 79 and 85) produces

$$I(T_j, \hat{T}_j) = \frac{I(\theta, \hat{T}_j)}{(\partial T_j / \partial \theta)^2}, \quad (16)$$

and

$$I(\theta, \hat{T}_j) \approx I(\theta, \hat{\theta}). \quad (17)$$

Thus,

$$\sigma^2(\hat{T}_j | \theta) \approx \frac{\left[\frac{1}{n_j} \sum_{i=1}^{k_j} \hat{T}_{ij}(\theta) \right]^2}{I(\theta, \hat{\theta})}$$

and the parameters of the prior beta distribution for T_j can be expressed in terms of the parameters of the 3PL IRT and 2PPC models. Furthermore, the parameters of the posterior distribution of T_j also can be expressed in terms of the IRT parameters:

$$p_j = \hat{T}_j n_j^* + x_j, \quad (18)$$

and

$$q_j = [1 - \hat{T}_j] n_j^* + n_j - x_j \quad (19)$$

The SPI is

$$\tilde{T}_j = \frac{p_j}{p_j + q_j} \quad (20)$$

$$= \frac{\hat{T}_j n_j^* + x_j}{n_j^* + n_j} \quad (21)$$

The SPI can also be written in terms of the relative contribution of the prior estimate \hat{T}_j and the observed proportion of maximum raw (correct score) (OPM), x_j/n_j , as

$$\tilde{T}_j = w_j \hat{T}_j + (1 - w_j) [x_j/n_j] \quad (22)$$

w_j , a function of the mean and variance of the prior distribution, is the relative weight given to the prior estimate:

$$w_j = \frac{n_j^*}{n_j^* + n_j} \quad (23)$$

The term n_j^* may be interpreted as the contribution of the prior in terms of theoretical numbers of items.

Check on Consistency and Adjustment of Weight Given to Prior Estimate

The item responses are assumed to be described by $P_i(\hat{\theta})$ or $P_{im}(\hat{\theta})$, depending on the type of item. Even if the IRT model accurately described item performance over examinees, their item responses grouped by strand may be multidimensional. For example, a particular examinee may be able to perform difficult addition but not easy subtraction. Under these circumstances, it is not appropriate to pool the prior estimate, \hat{T}_j , with x_j/n_j . In calculating the SPI, the following statistic was used to identify examinees with unexpected performance on the strands in a test:

$$Q = \sum_{j=1}^J n_j \left(\frac{x_j}{n_j} - \hat{T}_j \right)^2 / \left(\hat{T}_j (1 - \hat{T}_j) \right) \quad (24)$$

If $Q \leq \chi^2(J, .10)$, the weight, w_j , is computed and the SPI is produced. If $Q > \chi^2(J, .10)$, n_j^* and subsequently w_j is set equal to 0 and the OPM is used as the estimate of strand performance.

As previously noted, the prior is estimated using an ability estimate based on responses to all the items (including the items of strand j) and hence is not independent of X_j . An adjustment for the overlapping information that requires minimal computation is to multiply the test information in equation 5 by the factor $(n - n_j)/n$. The application of this factor produces an “adjusted” SPI estimate that can be compared to the “unadjusted” estimate.

Possible Violations of the Assumptions

Even if the IRT model fits the test items, the responses for a given examinee, grouped by strand, may be multidimensional. In these cases, it would not be appropriate to pool the prior estimate, \hat{T}_j , with x_j/n_j . A chi-square fit statistic is used to evaluate the observed proportion of maximum raw score (OPM) relative to that predicted for the items in the strand on the basis of the student's overall trait estimate. If the chi-square is significant, the prior estimate is not used and the OPM obtained becomes the student's strand score.

If the items in the strand do not permit guessing, it is reasonable to assume \hat{T}_j , the expected proportion correct of the maximum score for a strand, will be greater or equal to zero. If correct guessing is possible, as it is with MC items, there will be a non-zero lower limit to \hat{T}_j , and a three-parameter beta distribution, in which \hat{T}_j is greater than or equal to this lower limit (Johnson & Kotz, 1979, p. 37), would be more appropriate. The use of the two-parameter beta distribution would tend to underestimate T_j among very low-performing examinees. While working with tests containing exclusively MC items, Yen found that there does not appear to be a practical importance to this underestimation (Yen, 1997). The impact of any such effect would be reduced as the proportion of CR items in the test increases. The size of this effect, nonetheless, was evaluated using simulations (Yen, Sykes, Ito, & Julian, 1997).

The SPI procedure assumes that $p(X_j T_j)$ is a binomial distribution. This assumption is appropriate only when all the items in a strand have the same Bernoulli item response function. Not only do real items differ in difficulty, but when there are mixed-item types, X_j is not the sum of n_j independent Bernoulli variables. It is instead the total raw score. In essence, the simplifying assumption has been made that each CR item with a maximum score of $1_j - 1$ is the sum of $1_j - 1$ independent Bernoulli variables. Thus, a complex compound distribution is theoretically more applicable than the binomial. Given the complexity of working with such a model, it appears valuable to determine if the simpler model described here is sufficiently accurate to be useful.

Finally, because the prior estimate of T_j , \hat{T}_j , is based on performance on the entire test, including strand j , the prior estimate is not independent of X_j . The smaller the ratio n_j/n , the less impact this dependence will have. The effect of the overlapping information would be to understate the width of the credibility interval. The extent to which the size of the credibility interval is too small was examined (Yen et al, 1997) by simulating strands that contained varying proportions of the total test points.

References

- CTB/McGraw-Hill. (2000). *TerraNova* 2nd Edition. Monterey, CA.
- Fitzpatrick, A. R., V. Link, W. M. Yen, G. Burket, K. Ito & R. Sykes (1996). Scaling performance assessments: A comparison between one-parameter and two-parameter partial credit models. *Journal of Educational Measurement*, 33, 291–314.
- Johnson, N. L. & S. Kotz (1970). *Distributions in statistics: Continuous univariate distributions* (Vol. 2). New York: John Wiley.
- Lord, F. M. (1980). *Application of item response theory to practical testing problems*. Hillsdale, NJ: Lawrence Erlbaum.
- Muraki, E. (1992). A generalized partial credit model: Application of an EM algorithm. *Applied Psychological Measurement*, 16, 159–176.
- Novick, M. R. & P. H. Jackson (1974). *Statistical methods for educational and psychological research*. New York: McGraw-Hill.
- Yen, W. M. (1997). The technical quality of performance assessments: Standard errors of percents of students reaching standards. *Educational Measurement: Issues and Practice*, 5–15.
- Yen, W. M. (1993). Scaling performance assessments: Strategies for managing local item dependence. *Journal of Educational Measurement*, 30, 187–213.
- Yen, W. M., R. C. Sykes, K. Ito & M. Julian (1997). *A Bayesian/IRT index of objective performance for tests with mixed-item types*. Paper presented at the annual meeting of the National Council on Measurement in Education, Chicago.