

Utah Aspire Plus 2022-2023 Technical Report



2023

Passed ADA Accessibility Sept. 2023

Table of Contents

1	Introduction	9
1.1	Background	9
1.2	Purpose of the Operational Tests.....	10
1.3	Prior Administrations	10
1.4	Spring 2023 Administration	11
1.5	Composition of the Operational Tests.....	11
1.6	Intended Population of the Operational Tests.....	12
1.7	Overview of the Technical Report.....	13
2	Test Development	14
2.1	Overview of the Utah Aspire Plus Assessments, Claims, and Blueprints.....	14
2.1.1	English Assessment Claims.....	15
2.1.2	Reading Assessment Claims	16
2.1.3	Mathematics Assessment Claims.....	17
2.1.4	Science Assessment Claims	18
2.2	Utah Aspire Plus Blueprints	19
2.3	Test Development Activities.....	20
2.3.1	Operational Forms Development	21
2.3.2	Statistical Guidelines	22
2.3.3	2023 Match to Test Blueprint	23
3	Operational Administration.....	27
3.1	Testing Window.....	27
3.2	Test Administration and Security Policies.....	27
3.2.1	Online Administration and Monitoring.....	28
3.3	Test Accommodations and Supports	29
3.4	Test Taking Irregularities and Security Breaches	31
3.4.1	Test Interruptions.....	32
3.4.2	Scoring of Interrupted Tests	32
3.4.3	Wrong Test Form/Accommodation.....	32
3.4.4	Extended Time Accommodation Issues.....	32
3.4.5	Test Invalidation	32
3.5	Test Taker Characteristics	32
3.6	Testing Time	34
4	Score Reporting.....	37
4.1	IRT Pattern Scoring	37
4.1.1	Quality Control of IRT Scoring	38
4.2	Appropriate Uses for Scores and Reports	38
4.3	Utah Aspire Plus Reporting Scale	39
4.4	Standard Setting	40

4.5	ACT Predicted Score Ranges	41
4.6	2022–2023 Utah Aspire Plus Performance Results.....	42
5	Classical Item Analyses	43
5.1	Item Analyses	43
5.1.1	p -Value and Item Mean Scores.....	43
5.1.2	Item-Test Score Correlations	43
5.1.3	Differential Item Functioning.....	44
5.2	Classical Item Summaries for Operational Administration.....	45
6	Reliability.....	46
6.1	Classical Definition of Reliability	46
6.2	Classical Test Theory Reliability Estimates.....	47
6.2.1	Cronbach’s Alpha.....	47
6.2.2	Standard Error of Measurement.....	47
6.3	IRT-Based Reliability.....	48
6.4	Reliability of Performance Level Categorization	48
6.4.1	Accuracy and Consistency	49
6.4.2	Calculating Accuracy.....	49
6.4.3	Calculating Consistency	50
6.4.4	Calculating Kappa	51
7	Field Test Calibration and Drift Analyses	52
7.1	IRT Overview.....	52
7.2	IRT Data Preparation	53
7.2.1	Student Inclusion/Exclusion Rules	53
7.2.2	Quality Control of the IRT Data Matrix Files.....	53
7.3	Description of the Calibration, Equating, and Scaling Process	54
7.3.1	IRT Models.....	54
7.3.2	IRTPRO Calibration Procedures and Convergence Criteria.....	55
7.3.3	Calibration Quality Control	55
7.3.4	Equating	56
7.3.5	Drift Analysis	57
7.4	Model Fit Evaluation Criteria.....	59
8	Quality Control	61
8.1	Online Assessment Delivery	61
8.1.1	Item Validation.....	61
8.1.2	Test Administration	62
8.1.3	Operational Monitoring	63
8.2	Production System Testing	63
8.2.1	Functional Testing	63
8.2.2	Integration Testing	64
8.2.3	Program Validation End-to-End Testing	64

8.2.4	Load Testing.....	64
8.2.5	Performance Monitoring	65
8.2.6	Regression Testing	65
8.2.7	User Acceptance Testing.....	65
8.3	Reporting.....	66
8.4	Quality Control of Psychometric Processes.....	66
9	Validity.....	67
9.1	Evidence Based on Test Content	68
9.2	Evidence Based on Cognitive Process.....	70
9.3	Evidence Based on Internal Structure.....	71
9.3.1	Reliability	73
9.4	Evidence Based on Different Student Populations.....	74
9.5	Summary	74
10	References.....	75
	Appendix A : Test-Level Reporting Categories and Standards by Item Type and DOK	77
	Appendix B : Student Testing Time	83
	Appendix C : Item Statistics Summaries	85
	Appendix D : Reliability and Standard Error by Subgroup.....	87
	Appendix E : Conditional Standard Error of Scale Scores	95
	Appendix F : Accuracy and Consistency.....	99
	Appendix G : Common Item Scatterplots for 2023 Anchor Items.....	107
	Appendix H : Scale Score Descriptive Statistics by Subgroup.....	115
	Appendix I : Scale Score Distributions for Overall Testing Population	123
	Appendix J : Performance Level Distributions.....	127
	Appendix K : Principal Components Scree Plot.....	135
	Appendix L : Subscore Correlations.....	137
	Appendix M : Item Drift.....	140

List of Tables

Table 2.1. Field Test Forms	20
Table 2.2. Utah Aspire Plus English Grade 9 Operational Test Blueprint Match.....	23
Table 2.3. Utah Aspire Plus English Grade 10 Operational Test Blueprint Match.....	23
Table 2.4. Utah Aspire Plus Reading Grade 9 Operational Test Blueprint Match	24
Table 2.5. Utah Aspire Plus Reading Grade 10 Operational Test Blueprint Match.....	24
Table 2.6. Utah Aspire Plus Mathematics Grade 9 Operational Test Blueprint Match.....	25
Table 2.7. Utah Aspire Plus Mathematics Grade 10 Operational Test Blueprint Match.....	25
Table 2.8. Utah Aspire Plus Science Grade 9 Operational Test Blueprint Match	26
Table 2.9. Utah Aspire Plus Science Grade 10 Operational Test Blueprint Match.....	26
Table 3.1. Spring 2023 Participation Rates for Utah Aspire Plus	33
Table 3.2. Student Testing Time for Spring 2023 Utah Aspire Plus: English and Reading.....	35
Table 3.3. Student Testing Time for Spring 2023 Utah Aspire Plus: Math and Science.....	36
Table 4.1. IRT Summary Parameter Estimates for Utah Aspire Plus Operational Items	37
Table 4.2. Utah Aspire Plus Scale Score Cuts by Grade and Subject.....	40
Table 5.1. Item 2x2 Contingency Table for the k th Score Level	44
Table 6.1. Example Accuracy Classification Table.....	49
Table 6.2. Example Accuracy Classification Table for Proficient Cut Point.....	50
Table 6.3. Example Consistency Classification Table.....	51
Table 7.1. 2023 Final Stocking and Lord Scaling Constants	56
Table 7.2. 2023 Items Showing Drift	58
Table 9.1. Model Fit Indices for Confirmatory Factor Analyses	72
Table A.1. Test-Level Reporting Categories and Standards for English Grade 9.....	77
Table A.2. Test-Level Reporting Categories and Standards for English Grade 10.....	78
Table A.3. Test-Level Reporting Categories and Standards for Reading Grade 9	79
Table A.4. Test-Level Reporting Categories and Standards for Reading Grade 10.....	80
Table A.5. Test-Level Reporting Categories and Standards for Mathematics Grade 9.....	81
Table A.6. Test-Level Reporting Categories and Standards for Mathematics Grade 10.....	82
Table C.1. Item Mean for One-Point Items	85
Table C.2. Item Mean for Two-Point Items	85
Table C.3. Item Total Correlation for One-Point Items	85
Table C.4. Item Total Correlation for Two-Point Items	86
Table C.5. Differential Item Functioning	86
Table D.1. English Grade 9 Test Reliability.....	87
Table D.2. English Grade 10 Test Reliability	88
Table D.3. Reading Grade 9 Test Reliability.....	89
Table D.4. Reading Grade 10 Test Reliability.....	90
Table D.5. Mathematics Grade 9 Test Reliability.....	91
Table D.6. Mathematics Grade 10 Test Reliability	92
Table D.7. Science Grade 9 Test Reliability.....	93
Table D.8. Science Grade 10 Test Reliability.....	94
Table F.1. Accuracy Classification for English Grade 9	99

Table F.2. Accuracy Classification at Proficient Cut Point for English Grade 9	99
Table F.3. Consistency Classification for English Grade 9	99
Table F.4. Accuracy Classification for English Grade 10	100
Table F.5. Accuracy Classification at Proficient Cut Point for English Grade 10	100
Table F.6. Consistency Classification for English Grade 10	100
Table F.7. Accuracy Classification for Reading Grade 9	101
Table F.8. Accuracy Classification at Proficient Cut Point for Reading Grade 9	101
Table F.9. Consistency Classification for Reading Grade 9	101
Table F.10. Accuracy Classification for Reading Grade 10	102
Table F.11. Accuracy Classification at Proficient Cut Point for Reading Grade 10	102
Table F.12. Consistency Classification for Reading Grade 10	102
Table F.13. Accuracy Classification for Mathematics Grade 9	103
Table F.14. Accuracy Classification at Proficient Cut Point for Mathematics Grade 9	103
Table F.15. Consistency Classification for Mathematics Grade 9	103
Table F.16. Accuracy Classification for Mathematics Grade 10	104
Table F.17. Accuracy Classification at Proficient Cut Point for Mathematics Grade 10	104
Table F.18. Consistency Classification for Mathematics Grade 10	104
Table F.19. Accuracy Classification for Science Grade 9	105
Table F.20. Accuracy Classification at Proficient Cut Point for Science Grade 9	105
Table F.21. Consistency Classification for Science Grade 9	105
Table F.22. Accuracy Classification for Science Grade 10	106
Table F.23. Accuracy Classification at Proficient Cut Point for Science Grade 10	106
Table F.24. Consistency Classification for Science Grade 10	106
Table H.1. English Grade 9 Scale Score Descriptive Statistics	115
Table H.2. English Grade 10 Scale Score Descriptive Statistics	116
Table H.3. Reading Grade 9 Scale Score Descriptive Statistics	117
Table H.4. Reading Grade 10 Scale Score Descriptive Statistics	118
Table H.5. Mathematics Grade 9 Scale Score Descriptive Statistics	119
Table H.6. Mathematics Grade 10 Scale Score Descriptive Statistics	120
Table H.7. Science Grade 9 Scale Score Descriptive Statistics	121
Table H.8. Science Grade 10 Scale Score Descriptive Statistics	122
Table J.1. English Grade 9 Performance Level Distribution	127
Table J.2. English Grade 10 Performance Level Distribution	128
Table J.3. Reading Grade 9 Performance Level Distribution	129
Table J.4. Reading Grade 10 Performance Level Distribution	130
Table J.5. Mathematics Grade 9 Performance Level Distribution	131
Table J.6. Mathematics Grade 10 Performance Level Distribution	132
Table J.7. Science Grade 9 Performance Level Distribution	133
Table J.8. Science Grade 10 Performance Level Distribution	134
Table L.1. English Correlations of Total Score and Subscores	137
Table L.2. Reading Correlations of Total Score and Subscores	137
Table L.3. Mathematics Correlations of Total Score and Subscores	138
Table L.4. Science Correlations of Total Score and Subscores	139

List of Figures

Figure B.1. English Grade 9 Student Testing Time	83
Figure B.2. English Grade 10 Student Testing Time.....	83
Figure B.3. Reading Grade 9 Student Testing Time	83
Figure B.4. Reading Grade 10 Student Testing Time	83
Figure B.5. Mathematics Grade 9 Student Testing Time.....	84
Figure B.6. Mathematics Grade 10 Student Testing Time.....	84
Figure B.7. Science Grade 9 Student Testing Time	84
Figure B.8. Science Grade 10 Student Testing Time	84
Figure E.1. English Grade 9 Conditional Standard Error of Scale Scores.....	95
Figure E.2. English Grade 10 Conditional Standard Error of Scale Scores.....	95
Figure E.3. Reading Grade 9 Conditional Standard Error of Scale Scores	96
Figure E.4. Reading Grade 10 Conditional Standard Error of Scale Scores	96
Figure E.5. Mathematics Grade 9 Conditional Standard Error of Scale Scores.....	97
Figure E.6. Mathematics Grade 10 Conditional Standard Error of Scale Scores.....	97
Figure E.7. Science Grade 9 Conditional Standard Error of Scale Scores	98
Figure E.8. Science Grade 10 Conditional Standard Error of Scale Scores	98
Figure G.1. English Grade 9 IRT B Parameters for Operational Items.....	107
Figure G.2. English Grade 10 IRT B Parameters for Operational Items	108
Figure G.3. Reading Grade 9 IRT B Parameters for Operational Items.....	109
Figure G.4. Reading Grade 10 IRT B Parameters for Operational Items.....	110
Figure G.5. Mathematics Grade 9 IRT B Parameters for Operational Items.....	111
Figure G.6. Mathematics Grade 10 IRT B Parameters for Operational Items	112
Figure G.7. Science Grade 9 IRT B Parameters for Operational Items.....	113
Figure G.8. Science Grade 10 IRT B Parameters for Operational Items.....	114
Figure I.1. English Grade 9 Scale Score Distribution.....	123
Figure I.2. English Grade 10 Scale Score Distribution.....	123
Figure I.3. Reading Grade 9 Scale Score Distribution	124
Figure I.4. Reading Grade 10 Scale Score Distribution	124
Figure I.5. Mathematics Grade 9 Scale Score Distribution.....	125
Figure I.6. Mathematics Grade 10 Scale Score Distribution.....	125
Figure I.7. Science Grade 9 Scale Score Distribution	126
Figure I.8. Science Grade 10 Scale Score Distribution	126
Figure K.1. English Grade 9 Principal Components Scree Plot	135
Figure K.2. English Grade 10 Principal Components Scree Plot	135
Figure K.3. Reading Grade 9 Principal Components Scree Plot.....	135
Figure K.4. Reading Grade 10 Principal Components Scree Plot.....	135
Figure K.5. Mathematics Grade 9 Principal Components Scree Plot	136
Figure K.6. Mathematics Grade 9 Principal Components Scree Plot	136
Figure K.7. Science Grade 9 Principal Components Scree Plot.....	136
Figure K.8. Science Grade 10 Principal Components Scree Plot.....	136
Figure M.1. English Grade 9 Item Drift.....	140

Figure M.2. English Grade 10 Item Drift.....	141
Figure M.3. Reading Grade 9 Item Drift.....	142
Figure M.4. Reading Grade 10 Item Drift.....	143
Figure M.5. Mathematics Grade 9 Item Drift.....	144
Figure M.6. Mathematics Grade 10 Item Drift.....	145
Figure M.7. Science Grade 9 Item Drift.....	146
Figure M.8. Science Grade 10 Item Drift.....	147

1 Introduction

1.1 Background

The Utah Aspire Plus summative assessments were created out of Utah Statute 53E-4-304 (https://le.utah.gov/xcode/Title53E/Chapter4/53E-4-S304.html?v=C53E-4-S304_2019051420190514). The statute requires the Utah State Board of Education (USBE) to administer assessments that are predictive of college readiness at grades 9 and 10 in addition to providing overall performance scores and proficiency indicators for English, reading, mathematics, and science. The Utah Aspire Plus assessments are a hybrid of ACT Aspire and Utah Core test items. These are computer-based, fixed-length tests intended to measure end-of-grade-level high school knowledge and skills for students in grades 9 and 10. Spring 2019 marked the first administration of the Utah Aspire Plus assessments and the creation of base reporting scales for each respective grade and subject assessment.

Prior to 2019, students were assessed on the core standards through the Utah Student Assessment of Growth and Excellence (SAGE) assessment program. The Utah Aspire Plus assessment program is an extension of the Utah SAGE, still intended to measure student performance in relation to the Utah Core Standards (<https://www.uen.org/core/>), but also intending to measure students' preparedness for meeting college readiness benchmarks. As such, the assessment content from Utah SAGE is used as one component of the Utah Aspire Plus assessments.

Additional content from ACT Aspire is used to provide predictions of performance on the ACT[®]. This content also aligns to the Utah Core Standards and is counted toward Utah Aspire Plus scores too. The ACT[®] is the primary college readiness assessment submitted to local universities in Utah. As such, the Utah Aspire Plus assessments incorporate test questions from the ACT Aspire assessments that are used not only to contribute to student overall scores but also to provide a predictive indicator of performance on the ACT[®]. Students receive predicted ACT[®] score ranges for each ACT[®] subtest (English, reading, mathematics, and science), as well as an overall predicted composite ACT[®] score range.

As required by the statute noted previously, the assessments also provide overall scores as indicators of end-of-grade-level expectations for 9th and 10th grade students and performance level indicators (*Below Proficient, Approaching Proficient, Proficient, and Highly Proficient*) for English, reading, mathematics, and science.

1.2 Purpose of the Operational Tests

The Utah Aspire Plus assessments are designed for several purposes. First, the tests are intended to measure the breadth and depth of the Utah Core Standards and measure across all levels of student performance. Second, the tests are created to provide awareness of individual achievement in relation to stated performance expectations. Third, performance on the tests is intended to provide evidence of whether students are on track for college and career readiness. Finally, the tests are used to evaluate growth between 9th and 10th grade.

1.3 Prior Administrations

As stated, the first operational administration was conducted in the spring of 2019 at grades 9 and 10 for English, reading, mathematics, and science. Data from that administration were used to establish the initial Utah Aspire reporting scales and the setting of performance levels. Technical details of these features and activities are presented in the *Utah Aspire Plus 2018–2019 Technical Report* (http://utah.pearsonaccessnext.com/resources/additional-services/UT1132740_UTPlusTechReportv4.3_WebTag.pdf).

Note that spring 2020 was intended to be the second operational administration of the Utah Aspire Plus tests. In spring of 2020, Senate Bill 3005, which included a waiver of the Utah Aspire Plus assessment requirements, was passed during the Utah Legislature’s 3rd Special Session of 2020 and signed into law on April 22, 2020. As a result, the spring testing of Utah Aspire Plus was cancelled. As a result, spring 2021 marked the second administration of the Utah Aspire Plus assessments. However, it should be noted that a waiver was sought and granted by the U.S. Department of Education (Department) to waive the accountability, school identification, and related reporting requirements for the 2020–2021 school year (<https://www.schools.utah.gov/file/829f7300-020d-456e-85ac-49e85ef0795a>).

Mathematics, reading, and English summative assessments for the Utah Aspire Plus administration were created in 2019 for use in spring 2020. Given the cancellation of testing in spring 2020, the tests were instead rolled over and administered in spring 2021. Spring 2021 also marked the initial administration of new science tests. The Utah Aspire Plus Science with Engineering Education Standards (SEEds) summative assessments were administered to Utah students in spring 2021. These assessments are composed of test units that are designed to measure multi-dimensional knowledge and skill interactions across different scientific phenomena within core disciplines.

The tests were administered as an operational field test, meaning that items used to provide scores for students were identified after the administration. That identification activity was akin to the standard test construction process involving Pearson and USBE content experts and psychometricians working to identify the best forms based on match to blueprint and statistical indices. After these forms were determined, they were then used to set performance standards in August of 2021.

1.4 Spring 2023 Administration

Spring of 2023 marked the fourth administration of the Utah Aspire Plus assessment for English, reading and math and the third administration of the science assessment (following the establishment of the base scale in 2021).

For the first time, the Utah Aspire Plus tests were pre-equated. This allowed scores to be delivered on-demand following administration.

1.5 Composition of the Operational Tests

Each operational Utah Aspire Plus test form was constructed to reflect the full test blueprint in terms of content, standards measured, and item types (<http://utah.pearsonaccessnext.com/additional-services/>). All blueprints were designed to measure knowledge and skills described in the Utah Core Standards (<https://www.uen.org/core/>). For science, the operational assessments were created to measure the new Science with Engineering Education Standards (SEEds). The standards were derived from several research-based sources such as A Framework for K–12 Science Education and the Next Generation Science Standards (NGSS).

The Utah Aspire Plus tests are composed of several different types of items to measure student performance. These include multiple choice, multiple select, evidence-based selected response, and technology enhanced (TE). Multiple-choice items present students with four or five responses, of which there is one correct answer. Multiple-select items require students to select two or three correct choices from several presented choices. Evidence-based selected response items have two parts: Part A is designed as an *identification* component, where Part B is designed to elicit an *evidence*-based component. Further, these types can be designed as two multiple-choice items, or a combination of multiple-choice and technology-enhanced (TE) items. Technology-enhanced (TE) items require specialized interactions within the online presentation for capturing student responses (e.g., drag and drop).

The Utah Aspire Plus English tests target language conventions and comprehension. Students should be able to demonstrate command of standard English grammar, usage, capitalization, punctuation, and spelling. In addition, students should be able to demonstrate vocabulary knowledge in comprehending complex texts.

The Utah Core Standards in Reading define expectations of comprehension skills, understanding tone and point of view of texts, and evaluating texts. On the Utah Aspire Plus Reading tests, students must demonstrate these skills with different types of text sources.

The assessment context for Utah Aspire Plus Mathematics is grounded in five conceptual categories from the Utah Core Standards: Number and Quantity, Algebra, Functions, Geometry, and Statistics and Probability. There are two general levels of math content for Utah Aspire Plus. The first level, referred to as Secondary Math I, extends the mathematics from the middle grades, particularly on linear and exponential relationships. The next level, Secondary Math II, focuses on quadratic relationships and comparing them to the linear and exponential relationships from Secondary Math I.

The primary emphasis of the new Utah Aspire Plus Science tests is on the multidimensional nature as expressed within the NGSS. Specific Science and Engineering Practices (SEP) and Cross-Cutting Concepts (CCC) are identified within four reporting targets (Gathering and Investigating, Developing Models, Using Mathematical Thinking, and Constructing Explanations). These are further represented within the Disciplinary Core Ideas (DCI) of Life Science, Physical Science, and Earth and Space Science.

1.6 Intended Population of the Operational Tests

The Utah Aspire Plus tests are designed for students completing their 9th and 10th grade courses in English Language Arts (ELA), mathematics, and science. The English and reading tests are designed to assess the skills that 9th and 10th grade ELA students should have by the end of those respective years. The mathematics tests are designed to assess the skills that 9th (Secondary Math I) and 10th grade (Secondary Math II) math students should have by the end of those respective years. The science tests are designed to assess the skills that 9th and 10th grade students taking biology, chemistry, Earth science, or physics should have by the end of instruction (regardless of the specific course).

1.7 Overview of the Technical Report

The intended audience of the report are those with a basic technical understanding of large-scale assessment systems and their uses. It assumes some technical knowledge of how score scales are developed and derived and how scores are intended to support valid interpretations of intended claims.

This report provides details of the maintenance of the Utah Aspire Plus testing system at grades 9 and 10 for mathematics, science, reading and English. In addition to a general overview that provides a frame of reference around key attributes of the assessments, the report provides details around development of items and test forms, the administration of operational tests, the maintenance of existing scales, and of scoring and reporting for all tests. Throughout the report, the narrative is intended to present an interpretive argument whereby the various claims of the assessment system are identified and described throughout the test development process from creation through administration and score reporting. Technical details are presented in the following chapters and address test design, development and implementation, test administration, test taker characteristics, classical item analyses, reliability analyses, item response theory (IRT) calibrations, equating, and scaling, quality control procedures, and evidence of validity.

2 Test Development

2.1 Overview of the Utah Aspire Plus Assessments, Claims, and Blueprints

The Utah Aspire Plus assessments are aligned to the Utah Core Standards and designed to measure the breadth and depth of the Utah Core Standards across all levels of student performance, to provide awareness of individual achievement in relation to stated performance expectations, and to provide evidence of whether students are on track for college and career readiness. Utah Aspire Plus content follows a rigorous development process that meets and often exceeds industry standards for best practices in assessment. Every item, written by Utah teachers, goes through an extensive review designed to ensure adherence to high quality and the principles of universal design.

This chapter describes the claims intended to support the purposes outlined in Chapter 1; the development of blueprints defining the components of the Utah Aspire Plus assessments that reflect the breadth of the Utah Core Standards across different levels of student understanding; and the development of tasks (items) intended to fulfill the respective blueprints and provide evidence of varying levels of performance reflective of each of the stated claims.

It should be noted that while both claims and sub claims are presented here for each subject, only the claims are reported on individual student reports (ISR). Sub claims currently only provide structure within the respective blueprints but are not reported at the individual student level.

2.1.1 English Assessment Claims

The Utah Aspire Plus English tests target language conventions and comprehension. Students should be able to demonstrate command of standard English grammar, usage, capitalization, punctuation, and spelling. In addition, students should be able to demonstrate vocabulary knowledge in comprehending complex texts.

The claim structure for the Utah Aspire Plus English tests is drawn from the Utah Core Standards and frames the design and development of the summative tests at grades 9 and 10.

Claims: The primary claims reflect the main goals for the use of the Utah Aspire Plus English tests. The first is that student performance reflects an indicator of career and college readiness as demonstrated through students' understanding of language conventions and comprehension as expected to have been attained by the end of each respective year as a prediction of performance on the ACT[®] English test. Second is that overall performance reflects students' understanding of language conventions and comprehension with respect to the breadth and depth of the Utah Core Standards and measures across all levels of student performance.

Sub Claims*: The sub claims further explicate what is measured on Utah Aspire Plus English tests and are grouped into the following categories:

- Production of Writing
- Knowledge of Language
- Conventions of Standard English

* It should be noted that sub claims are *not* reported on individual student reports but form an important structural element within the blueprints. They are included in this technical report for completeness.

2.1.2 Reading Assessment Claims

The Utah Aspire Plus Reading tests define expectations of comprehension skills, understanding tone and point of view of texts, and evaluating texts. On the Utah Aspire Plus Reading tests, students must demonstrate these skills with different types of text sources.

The claim structure for the Utah Aspire Plus Reading tests is drawn from the Utah Core Standards and frames the design and development of the summative tests at grades 9 and 10.

Claims: The primary claims reflect the main goals for the use of the Utah Aspire Plus Reading tests. The first is that student performance reflects an indicator of career and college readiness as demonstrated through students' ability to read and comprehend complex informational and literary texts as expected to have been attained by the end of each respective year as a prediction of performance on the ACT[®] Reading test. Second is that overall performance reflects students' understanding of reading and comprehending complex informational and literary texts with respect to the breadth and depth of the Utah Core Standards and measures across all levels of student performance.

Sub Claims*: The sub claims further explicate what is measured on Utah Aspire Plus Reading tests and are grouped into the following categories:

- Key Ideas
- Craft and Structure
- Integration of Knowledge and Ideas

* It should be noted that sub claims are *not* reported on individual student reports but form an important structural element within the blueprints. They are included in this technical report for completeness.

2.1.3 Mathematics Assessment Claims

The Utah Aspire Plus Mathematics tests are grounded in five conceptual categories from the Utah Core Standards: Number and Quantity, Algebra, Functions, Geometry, and Statistics and Probability. There are two levels of math content for Utah Aspire Plus that reflect expectations at grades 9 and 10, respectively. The first level (grade 9), referred to as Secondary Math I, extends the mathematics from the middle grades, particularly on linear and exponential relationships. The next level, Secondary Math II (grade 10), focuses on quadratic relationships and comparing them to the linear and exponential relationships from Secondary Math I.

The claim structure for the Utah Aspire Plus Math tests is drawn from the Utah Core Standards and frames the design and development of the summative tests at grades 9 and 10.

Claims: The primary claims reflect the main goals for the use of the Utah Aspire Plus Mathematics tests. The first is that student performance reflects an indicator of career and college readiness as demonstrated through students' ability to understand linear relationships, abstract and quantitative reasoning, and problem solving as expected to have been attained by the end of each respective year as a prediction of performance on the ACT[®] Math test. Second is that overall performance reflects students' understanding of linear relationships, abstract and quantitative reasoning, and problem solving with respect to the breadth and depth of the Utah Core Standards and measures across all levels of student performance.

Sub Claims*: The sub claims further explicate what is measured on Utah Aspire Plus Math tests and are grouped into the following categories:

* It should be noted that sub claims are *not* reported on individual student reports but form an important structural element within the blueprints. They are included in this technical report for completeness.

Math I (Grade 9)

- Algebra
- Functions
- Geometry
- Statistics and Probability

Math II (Grade 10)

- Number and Quantity
- Algebra
- Functions
- Geometry
- Statistics and Probability

2.1.4 Science Assessment Claims

The Utah Aspire Plus Science tests are developed around the Utah Core Standards for science as described in the Science with Engineering Education Standards (SEEds). These skills are applicable regardless of domain (Biology, Physics, Earth Science, and Chemistry). The claim structure for the Utah Aspire Plus Science tests is drawn from the Utah Core Standards as described in the SEEds and frames the design and development of the summative tests at grades 9 and 10.

Claims: The primary claims reflect the main goals for the use of the new Utah Aspire Plus Science tests. The first is that student performance reflects an indicator of career and college readiness as demonstrated through students' ability to understand and apply science as defined by the SEEds. Further, as expected to have been attained by the end of each respective year as a prediction of performance on the ACT[®] Science test. Second is that overall performance reflects students' understanding of science as defined by the SEEds with respect to the breadth and depth of the Utah Core Standards and measuring across all levels of student performance.

Sub Claims*: The sub claims further explicate what is measured on the new Utah Aspire Plus Science tests and are grouped into the following categories with respective SEP and CCC targets:

* It should be noted that sub claims are *not* reported on individual student reports but form an important structural element within the blueprints. They are included in this technical report for completeness.

- Gathering and Investigating
 - SEPs: Asking questions and defining problems; Obtaining, evaluating, and communicating information; Planning and carrying out investigations
 - CCCs: Patterns; Cause and effect; Systems and system models; Energy and matter; Structure and function; Stability and change Use Science Process and Thinking Skills
- Developing Models
 - SEPs: Developing and using models
 - CCCs: Patterns; Cause and effect; Scale, proportion and quantity; Systems and system models; Energy and matter; Stability and change
- Using Mathematical Thinking –
 - SEPs: Analyzing and interpreting data; Using mathematics and computational thinking
 - CCCs: Patterns; Cause and effect; Scale, proportion, and quantity; Systems and system models; Energy and matter; Stability and change
- Constructing Explanations –
 - SEPs: Constructing explanations and designing solutions; Engaging in argument from evidence
 - CCCs: Patterns; Cause and effect; Systems and system models; Energy and matter; Structure and function; Stability and change

These are expressed across the Life Science, Earth and Space Science, and Physical Science DCIs.

2.2 Utah Aspire Plus Blueprints

The Utah Aspire Plus tests are administered in English, reading, mathematics, and science in grades 9 and 10 and are described in Section 1.5. For the Utah Aspire Plus tests, the creation of test blueprints was driven by the intended purposes detailed previously in order to support the respective claim structures. The blueprints for Utah Aspire Plus are the distribution of item types across domains/reporting categories, level of cognitive demand, and the number of total points associated with each.

For the science tests, the SEEds blueprints assume a design in which one of the three DCIs will be assessed by two clusters and the other two DCIs with a single cluster. Coverage of the respective DCIs rotates across forms (either within a given year or across years) to ensure the standards are fully represented over time.

The 2023 Utah Aspire Plus blueprints can be found at:
<http://utah.pearsonaccessnext.com/additional-services/>.

2.3 Test Development Activities

Prior to the creation of Utah Aspire Plus, students were tested on the Utah Core Standards through the Utah Student Assessment of Growth and Excellence (SAGE). The Utah Aspire Plus assessments were built from existing Utah SAGE banked content combined with items from ACT Aspire to allow for predictions of students' preparedness for meeting college readiness. All available content for creation of the 2023 Utah Aspire Plus tests was based on the existing item banks described in the *Utah Aspire Plus 2018–2019 Technical Report* (available at http://utah.pearsonaccessnext.com/resources/additional-services/UT1132740_UTPlusTechReportv4.3_WebTag.pdf).

In prior years, forms were built with “linking sets” consisting of items that were previously delivered operationally. These served as common items to equate the Utah Aspire Plus tests to the base scales using a common item non-equivalent groups equating design (Kolen and Brennan, 2014). Since the 2023 tests were pre-equated, a linking set was not required to place the forms on the base scale. However, for the first year of pre-equating, prior practice was followed in building the tests, including the creation of a “linking set” of common items. For test development purposes, this meant selecting sets of items to ideally reflect a miniature version of the overall test (typically at least 20 percent) in content as well as statistical characteristics. For Mathematics, English, and Reading, the ACT Aspire forms that are used to source items are alternated each year. This helps limit exposure of the Aspire content that might otherwise negatively impact ACT predication score activities.

For 2023, there was one core operational form for regular online and text-to-speech forms. Mathematics and Science forms consisted of operational items and a small set of field-test items. The number of field test forms for 2023 by grade and subject is shown in Table 2.1.

Table 2.1. Field Test Forms

Subject	Grade	Number of FT Versions
Math	9	15
	10	15
Science	9	20
	10	19

In addition to the ONEN forms, there are several accommodated forms. These include:

- Non-screen reader (NREN)
- Screen reader (SREN)
- Spanish (ONSP)

The 2023 accommodated forms were a reuse of the 2021 forms. In Grade 10 math, one item from this form was replaced.

2.3.1 Operational Forms Development

The construction of test forms for the 2023 Utah Aspire Plus was a coordinated effort between experts from the Utah State Board of Education, Pearson, and ACT. This process required adhering to guidelines that promote fair and ethical testing practices. Using the content developed to measure the Utah Core Standards, specialists worked through an iterative process to evaluate the specific items, passages, and stimuli that best met the intended measurement targets and to support all stated claims.

The Utah Aspire Plus assessments measure students' mastery of the Utah Core Standards or the Utah Aspire Plus Science with Engineering Education Standards. These standards are used to drive Utah instruction as well as developing the Utah Aspire Plus tests. As stated earlier, the Utah Aspire Plus assessments are designed so that test scores can be linked to ACT scales to provide students with indicators of being prepared for meeting college readiness benchmark. In order to accomplish this, approximately 50% of the Utah Aspire Plus tests (less for mathematics) are composed of items from ACT Aspire. As noted, these items serve multiple purposes, which include being used to derive prediction scores between the Utah Aspire Plus scales and ACT scales.

The general test development process for Utah Aspire Plus was initiated with the selection of items from ACT Aspire. Items were selected based on match to blueprint, as well as statistical indicators of item quality and fairness provided from the SAGE and ACT Aspire banks, respectively. ACT Aspire items were positioned within each form in the same locations as originally administered within ACT Aspire forms to help facilitate the derivation of the predictive scores on Utah Aspire Plus.

Once the ACT Aspire items were selected, Pearson psychometrics selected sets of items common to 2019, 2021, or 2022 forms. In previous years, prior to the adoption of pre-equating, these sets were used to equate the base scales (2019 for English, reading and math; 2021 for science). In addition to selecting items to be as similar as possible to the overall blueprints, they were also targeted to the original base scale difficulties. This step was followed during test construction for spring 2023, although the set did not serve as a linking set due to the use of pre-equating.

The test construction procedure was an iterative process whereby the first proposed form was evaluated by each party (Pearson, USBE, and ACT) for content and psychometric quality, feedback provided, and revisions made until a best final version was approved by all. It should be noted that without new development of content, bank limitations meant an inability to strictly meet the new blueprint in all cases (see below). It also meant that there were also instances where items with poorer statistical indices were included to meet the blueprint. These were infrequent and, in all cases, deemed reasonable in supporting the intended claims without negative impact. Moving forward, newly developed content will fill gaps and address such limitations as the assessments mature.

2.3.2 Statistical Guidelines

While the initial Utah Aspire Plus tests were primarily driven by content considerations, statistical indices were available based on use within the SAGE and ACT Aspire Plus assessments. For creation of Utah Aspire Plus tests, some general guidelines were used to help support selection of a range of item difficulties and evaluate item quality to ensure the best overall test forms. These indices are described in detail further on in the report.

The guidelines for creation of the Utah Aspire Plus forms were as follows:

- **Target item difficulty range of between 0.30 and 0.85.** Based on p -values, where the percentage reflects the percentage of students correctly responding to the item. Items awarding more than one point used the item mean divided by the maximum points possible to place on the p -value metric.
- **Target threshold for item discrimination of 0.20 and above.** Where item discrimination is defined by item-total score correlations.
- **Extreme differential item functioning (DIF) indices should be avoided.** A standard flagging convention indicates differences of magnitude and classifies the most extreme cases of DIF as “C,” moderate DIF as “B,” and minor to no DIF as “A.” As such, items flagged “C” should be avoided and minimal use of items flagged “B” should be used and/or balanced within a form where possible.

More detailed description of the statistical indices reflecting item functioning for the Utah Aspire Plus tests appears later in this report, and distributional results by grade and subject test from the 2023 operational administration are presented in Appendix C. It should be noted that Appendix C reflects post hoc calculations, not what was available within the context of test construction. It should further be noted that while most items selected to appear on the initial Utah Aspire Plus forms were within the guidelines described here, there were instances in which bank limitations meant some items did fall outside the thresholds.

2.3.3 2023 Match to Test Blueprint

Tables 2.2 through 2.9 present the match between the final 2023 operational forms of Utah Aspire Plus and the test blueprints. English, reading, math, and science final forms reasonably matched all targets by item type, depth of knowledge, and reporting category (within 5 percent).

Table 2.2. Utah Aspire Plus English Grade 9 Operational Test Blueprint Match

	Number of Items	Minimum %	Maximum %	2023 Form
Item Type				
Multiple Choice	24-31	60%	89%	63%
Technology Enhanced	8-13	20%	37%	37%
Depth of Knowledge				
Level 1	15-20	38%	57%	52%
Level 2	8-16	20%	46%	22%
Level 3	12-15	30%	43%	26%
Reporting Categories				
Conventions of Standard English	20-30	50%	86%	65%
Knowledge of Language	4-10	10%	29%	13%
Production of Writing	7-12	18%	34%	22%

Table 2.3. Utah Aspire Plus English Grade 10 Operational Test Blueprint Match

	Number of Items	Minimum %	Maximum %	2023 Form
Item Type				
Multiple Choice	24-31	60%	89%	67%
Technology Enhanced	8-13	20%	37%	33%
Depth of Knowledge				
Level 1	15-20	38%	57%	47%
Level 2	8-16	20%	46%	26%
Level 3	12-15	30%	43%	28%
Reporting Categories				
Conventions of Standard English	20-30	50%	86%	56%
Knowledge of Language	4-10	10%	29%	21%
Production of Writing	7-12	18%	34%	23%

Table 2.4. Utah Aspire Plus Reading Grade 9 Operational Test Blueprint Match

	Number of Items	Minimum %	Maximum %	2023 Form
Item Type				
Evidence-Based Selected Response	3–6	9%	17%	11%
Multiple Choice	22–30	63%	86%	71%
Technology Enhanced	2–7	6%	20%	17%
Depth of Knowledge				
Level 1	4–7	11%	20%	17%
Level 2	14–20	40%	57%	43%
Level 3	12–15	34%	43%	40%
Reporting Categories				
Craft and Structure	12–16	34%	46%	40%
Integration of Knowledge and Ideas	3–7	9%	20%	17%
Key Ideas	12–18	34%	51%	43%

Table 2.5. Utah Aspire Plus Reading Grade 10 Operational Test Blueprint Match

	Number of Items	Minimum %	Maximum %	2023 Form
Item Type				
Evidence-Based Selected Response	3–6	9%	17%	14%
Multiple Choice	22–30	63%	86%	77%
Technology Enhanced	2–7	6%	20%	9%
Depth of Knowledge				
Level 1	4–7	11%	20%	20%
Level 2	14–20	40%	57%	46%
Level 3	12–15	34%	43%	34%
Reporting Categories				
Craft and Structure	12–16	34%	46%	37%
Integration of Knowledge and Ideas	3–7	9%	20%	11%
Key Ideas	12–18	34%	51%	51%

Table 2.6. Utah Aspire Plus Mathematics Grade 9 Operational Test Blueprint Match

	Number of Items	Minimum %	Maximum %	2023 Form
Item Type				
Multiple Choice	30–33	75%	83%	78%
Technology Enhanced	7–10	18%	25%	23%
Depth of Knowledge				
Level 1	8–12	20%	30%	30%
Level 2	15–20	38%	50%	48%
Level 3	9–13	23%	33%	23%
Reporting Categories				
Algebra	9–11	23%	28%	25%
Functions	10–12	25%	30%	28%
Geometry	9–11	23%	28%	25%
Statistics and Probability	7–9	18%	23%	23%

Table 2.7. Utah Aspire Plus Mathematics Grade 10 Operational Test Blueprint Match

	Number of Items	Minimum %	Maximum %	2023 Form
Item Type				
Multiple Choice	30–33	75%	83%	83%
Technology Enhanced	7–10	18%	25%	18%
Depth of Knowledge				
Level 1	8–12	20%	30%	33%
Level 2	15–20	38%	50%	48%
Level 3	9–13	23%	33%	20%
Reporting Categories				
Algebra	9–11	23%	28%	23%
Functions	10–12	25%	30%	28%
Geometry	11–13	28%	33%	30%
Number and Quantity	2–4	5%	10%	10%
Statistics and Probability	2–4	5%	10%	10%

Table 2.8. Utah Aspire Plus Science Grade 9 Operational Test Blueprint Match

	Number of Items	Minimum %	Maximum %	2023 Form
Item Type				
Multiple Choice/Select	18-21	78%	91%	87%
Technology Enhanced	3-6	13%	26%	13%
DCI				
Earth and Space	4-8	21%	34%	22%
Life Science	9-13	38%	55%	52%
Physical Science	4-8	21%	34%	26%
Reporting Categories				
Developing Models	4-8	14%	34%	30%
Gathering & Investigating	4-8	14%	34%	22%
Construct Explanations	5-9	17%	38%	22%
Mathematical Thinking	5-9	17%	38%	26%

Table 2.9. Utah Aspire Plus Science Grade 10 Operational Test Blueprint Match

	Number of Items	Minimum %	Maximum %	2023 Form
Item Type				
Multiple Choice/Select	18-21	78%	91%	87%
Technology Enhanced	3-6	13%	26%	13%
DCI				
Earth and Space	4-8	21%	34%	22%
Life	4-8	21%	34%	26%
Physical	9-13	38%	55%	52%
Reporting Categories				
Developing Models	4-8	14%	34%	26%
Gathering & Investigating	4-8	14%	34%	17%
Construct Explanations	5-9	17%	38%	35%
Mathematical Thinking	5-9	17%	38%	22%

For additional information on the 2023 operational forms, Appendix A contains a breakdown reporting categories and standards by item type and depth of knowledge (DOK), with the exception of science (which does not use DOK).

3 Operational Administration

3.1 Testing Window

The 2023 administration of the Utah Aspire Plus assessments was March 6–May 12, 2023. Utah Aspire Plus can be administered on a subject-by-subject basis or as a complete battery with all tests administered in one sitting. Each subject test, however, must be administered in one sitting. In other words, once a subject test is started, it must be completed within that sitting.

3.2 Test Administration and Security Policies

Comprehensive details of the Utah Aspire Plus test administration are detailed in the Test Administration Manual (TAM, http://utah.pearsonaccessnext.com/resources/training/SP23_Utah%20Summative%20TAM_Final%2011.pdf) as well as via the Utah Aspire Plus Resource Center (<http://utah.pearsonaccessnext.com/training/>). These resources cover all policies, procedures, specifications, training, instructions, security, accommodations, and oversight for every aspect of the Utah Aspire Plus test administration. These resources are further presented in a manner that addresses those responsible for carrying out the administration for all students as well as for educators and students to become familiar with the tests themselves (e.g., via practice tests and such) and for interpretation of test scores.

The Utah Aspire Plus tests are secure tests that follow the Utah Aspire Plus blueprints for each assessed subject area. All test items are secured items and may not be reviewed with students, discussed as a class, or reviewed during instructional conversations. Discussing, reviewing, recording, or transcribing test questions in any format is a violation of test security. All test security requirements of Utah Aspire Plus must be met. Personnel involved in test administration must complete testing ethics training. The Utah Standard Test Administration and Testing Ethics policy can be found here: <https://schools.utah.gov/file/2a1a1ecf-710e-439c-bd7d-318b0a9eb1c1>.

The LEA Assessment Director was responsible for ensuring that each student had an appropriate opportunity to demonstrate knowledge, skills, and abilities related to Utah Aspire Plus grade-based courses and assessments. This ensures that each student had a standardized (similar and fair) testing experience for a given assessment. Each LEA was responsible for determining school testing schedules. Subject tests did not have to be administered in any prescribed order. Subject tests could *not* be divided into multiple sessions. Once a subject test session began, the subject test had to be completed within that sitting.

It should be noted that the previous SAGE tests were untimed. To support the derivation of predictive scores on the ACT[®], the Utah Aspire Plus assessments follow the same fixed testing time conditions. For the 2022–2023 administration, the testing times were: 45 minutes for English, 75 minutes each for Reading and Mathematics, and 60 minutes for Science. It should be noted that students whose IEP, Section 504, or English Learner plan specified an accommodation for extended time were able to use extended time accommodations on Utah Aspire Plus as appropriate.

3.2.1 Online Administration and Monitoring

The Utah Aspire Plus tests are administered online via the Pearson test management and delivery systems. PearsonAccess^{next} is the web application used by test staff (i.e., test coordinators, room supervisors) to manage online testing and start and monitor tests. TestNav is the test delivery engine used by examinees to take the tests. TestNav provides advance warning of network issues that prevent sending student responses to the Pearson testing server. When the network is functioning normally, TestNav sends student responses to the Pearson testing server in real time, while the student is testing. If the student's device cannot connect to the Pearson servers, TestNav saves the response to an encrypted file and allows the student to continue testing. When the network connection is reestablished, the test proctor can upload a student's saved responses to Pearson's testing server, and then TestNav erases the encrypted response file from the student's device or local network. As part of test security, test administrators control individual student authorization by printing and distributing testing tickets with each student's identifying information and unique log-in credentials.

Pearson's operational monitoring practices and tools constantly verify that platforms remain available to users; that performance stays within acceptable limits; and that users do not encounter critical errors. Additionally, monitoring includes real-time security auditing and systems vulnerability monitoring throughout a given testing window.

3.3 Test Accommodations and Supports

The Utah Aspire Plus tests are provided to account for a range of accessibility features for all testers and accommodations for students with disabilities. Accommodations are determined by an EL, Individualized Education Program (IEP), or Section 504 team. Both federal and state laws require that all students be administered assessments intended to hold schools accountable for the academic performance of students. These laws include state statutes that regulate Utah's Accountability Systems. Additional laws include the 2015 reauthorization of ESEA, the Every Student Succeeds Act (ESSA), and the Individuals with Disabilities Education Improvement Act of 2004 (IDEA). All students are expected to participate in the state accountability system. This principle of full participation includes EL students, students with an Individualized Education Program (IEP), and students with a Section 504 plan.

For Utah Aspire Plus, accommodated test forms include Spanish-language forms and forms with assistive technology. These forms are modified reproductions of the original test forms. Modifications primarily involve incorporation of the accommodation with the intent of otherwise preserving the item content in its original form. Assistive technology within online test forms includes speech-to-text, magnification, and adaptive keyboard and mouse. Paper accommodations are also offered in the form of standard-print, large-print, and Braille reproductions.

For students requiring Braille, paper versions of the original forms are created, and student responses are transcribed into one of the assistive technology test formats. For items that are *not* able to be adopted as is, some modification must occur to create the accommodated parallel version. These are referred to as "sister" items and are created directly from the original item to preserve every aspect of the item as it is used in the original form, to include capture of student responses such that item characteristics are directly comparable. While this typically involves only a few items on a given assessment, the Spanish-language forms must be fully *transadapted*. This process is not only a matter of directly translating a test form's English text to Spanish, but also of adapting the content to account for the linguistic and cultural differences between speakers of the two different languages.

Creation of all transadapted and sister items for the Utah Aspire Plus assessments follow a similar process of creation and review as the original items, with an emphasis on fully matching to the original item in terms of content and function. That is, highly qualified item writers with extensive expert content experience are involved in the creation and review process of transadapted and/or sister item creation. Several reviews are held throughout the creative process involving Pearson and USBE content and psychometric experts to ensure match to source.

Testing accommodations and supports, including those mentioned above, are outlined in the TAM. (A complete list of accessibility and accommodation features for the Utah Aspire Plus assessments can be found in the accessibility and accommodations manual insert at http://utah.pearsonaccessnext.com/resources/training/UT1130483_SummSp21TAN_WebTag.pdf.)

Embedded and non-embedded supports are generally available to all students, whether through the online system or locally arranged. The list below provides the embedded and non-embedded supports provided within Utah Aspire Plus, as outlined in the TAM:

- Embedded:
 - In browser/app zoom
 - Answer eliminator
 - Calculator – Desmos graphing and Desmos scientific
 - Bookmarking items for review
 - Line reader mask
 - Color contrast
 - Answer masking
 - Highlighter
 - Keyboard navigation
 - Text-to-speech (English)
 - Directions reread (text-to-speech)
 - Text-to-speech (Spanish)
 - Personalized visual modification of remaining time
- Non-embedded:
 - Word to word dictionary
 - Scratch paper
 - Line reader
 - Supervised breaks within each day
 - Special seating/grouping
 - Location for movement
 - Separate/alternate location
 - Minimized distractions
 - Food or medication for individuals with medical needs
 - Administration and optimum time of day
 - Special lighting
 - Adaptive equipment/furniture
 - Wheelchair-accessible room

Testing accommodations require prior designation in a student's Individualized Education Program (IEP), 504, or English Learner (EL) plan. The list below provides the test accommodations, in addition to those supports previously mentioned.

- Assistive technology – screen reader
- Speech to text – assistive technology scribe
- Other assistive technology
- Spanish transadaptation
- Online test translation – other languages than Spanish or English
- Standard print
- Large print
- Braille plus tactile graphics
- Extra time
- Personalized auditory notification of remaining time
- Breaks: stop the clock
- Breaks: extending over multiple days
- Human scribe
- Home administration
- Human reader
- Signed exact English (directions only)
- Sign language interpretation
- Cued speech
- Alternate mouse pointer
- Zoom percentage
- Abacus

3.4 Test Taking Irregularities and Security Breaches

Test irregularities are non-standard situations that occur during test administration that affect one or more students. This includes students experiencing computer problems, experiencing a sudden illness, having to leave the room, or becoming unduly disturbed by the testing situation. Testing staff are trained to become familiar with the policy around unexpected/unforeseen circumstances prior to testing.

Some students may be unable to participate in regular testing schedules due to absence, technical difficulties, or other unforeseen circumstances. Opportunities for these students to complete each assessment were provided within the school's testing window. If there was an emergency that interrupted testing for an entire class or school, decisions about whether a test could be started again or not were to be made on a case-by-case basis by working with the Utah State Board of Education assessment team.

3.4.1 Test Interruptions

In the event that a student got sick, had to leave and could not return during the test, or for any other reason did not complete a test which had already begun, the test was to be concluded and submitted immediately. To maintain the security of the test questions, students were not allowed to restart or take a test over again.

3.4.2 Scoring of Interrupted Tests

If a student was interrupted and completed only part of a test before it was concluded and submitted, the student might not have received a score. A student must have attempted 85% of the questions to receive a score. If a student did not attempt at least 85% of the test questions, a score could not be generated, and no test score would be reported for that particular test. Overall composite scores would not be available for students who had missing subject test scores because the composite score is calculated using all four subject tests.

3.4.3 Wrong Test Form/Accommodation

If a student began a test using a test form or accommodation that they were not supposed to have, the teacher/proctor should have immediately stopped the test. In those instances, a new test assignment had to be created and a new test administration could proceed as normal from that point.

3.4.4 Extended Time Accommodation Issues

Extended time accommodations must be applied before preparing and starting sessions. In the event the accommodation is applied after the session has been prepared and started, students receive a time expired warning that has a link for "Proctor only." At that point a proctor can confirm the student should have extended time and is able to set the student up to continue testing as per their accommodation.

3.4.5 Test Invalidation

Tests could be invalidated when a student's performance was not deemed an accurate measure of their ability (e.g., the student cheated, used inappropriate materials, etc.). When a test is invalidated, the student is not given another opportunity to take the test. Invalidating a test has to be completed by the district testing administrator.

3.5 Test Taker Characteristics

Table 3.1 provides the participation rates for each Utah Aspire Plus test by subgroup. These are students that received a valid test score on a subject test. Cases that did not have a valid test score were excluded from being counted.

Table 3.1. Spring 2023 Participation Rates for Utah Aspire Plus

Students	Subgroup	English		Reading		Math		Science		
		Gr. 9	Gr.10	Gr. 9	Gr. 10	Gr. 9	Gr. 10	Gr. 9	Gr. 10	
All	Students Scored	46,511	43,766	46,680	43,536	45,163	42,908	46,592	43,280	
Sex	Female	47.78	47.27	47.70	47.20	47.41	47.17	47.68	47.19	
	Male	52.22	52.73	52.30	52.80	52.59	52.83	52.32	52.81	
Ethnicity	Hispanic or Latino									
	Ethnicity	19.43	19.29	19.66	19.36	19.34	19.35	19.61	19.45	
	Asian	1.75	1.75	1.75	1.75	1.76	1.78	1.75	1.76	
	Native Hawaiian or Other Pacific Islander	1.50	1.46	1.51	1.45	1.50	1.47	1.49	1.46	
	Black or African American	1.35	1.27	1.39	1.29	1.39	1.27	1.40	1.30	
	American Indian or Alaska Native	1.00	1.01	1.00	1.01	0.98	1.00	1.01	1.01	
	White	71.70	72.15	71.42	72.06	71.76	72.07	71.47	71.97	
	Other	3.27	3.07	3.27	3.08	3.27	3.07	3.27	3.05	
	Limited English Proficiency	No	91.11	92.51	90.93	92.42	91.00	92.42	90.99	92.40
		Yes	8.89	7.49	9.07	7.58	9.00	7.58	9.01	7.60
Economic Disadvantage	No	72.61	74.88	72.38	74.87	72.73	75.02	72.36	74.81	
	Yes	27.39	25.12	27.62	25.13	27.27	24.98	27.64	25.19	
Special Education	No	90.05	90.83	90.00	90.78	89.91	90.75	90.01	90.90	
	Yes	9.95	9.17	10.00	9.22	10.09	9.25	9.99	9.10	

3.6 Testing Time

One of the key questions in moving from an untimed to a timed test administration (from SAGE to Utah Aspire Plus) is gauging the extent to which the time allotted appears to be reasonable. As mentioned in Section 3.2, the operational testing times for the Utah Aspire Plus tests are: 45 minutes for English, 75 minutes for Reading, 75 minutes for Math, and 60 minutes for science. Students needing extra time fall into three categories: time and a half, double time, or triple time. After the spring 2023 test administration, student total testing time was analyzed for each test. Overall, students completed the assessments within the recommended testing times. Tables 3.2 and 3.3 provide breakdowns of student testing time across the full range of testing times. In other words, the percentile rankings are of the amount of time in minutes students took to complete the respective test. More specifically, with the Grade 9 English results for students testing using regular time (45 minutes), examination of the 95th percentile (P95) means that 95% of students finished the test in 43 minutes or less.

Additional information is presented in Appendix B, which provides a graphical display (box-and-whisker plot) of student testing time for each test. Box-and-whisker plots present the same information at each respective quartile, where the middle 50% of the given distribution is the box, and the whiskers represent the bottom 25% and top 25% of the distribution. Dots represent outliers and reflect very few overall cases. Most outliers for regular testers are still within the time allotment for the subject. For example, the outliers for grade 9 English for regular testers are all below the 75-minute time threshold. Based on these data and plots, the evidence suggests students in general had enough time to complete each respective test within the given allotments.

Table 3.2. Student Testing Time for Spring 2023 Utah Aspire Plus: English and Reading

Subject	Grade	Group	Testing Time (minutes)										
			N	Descriptive Statistics		Percentiles							
				Minimum	Maximum	Mean	St. Dev.	P50	P75	P80	P85	P90	P95
English	9	Regular Time	41,838	2	44	30	9	30	36	38	40	41	43
		Time and a Half	4,021	2	68	33	14	32	43	45	48	52	59
		Double Time	465	2	89	37	17	36	47	52	56	61	67
		Triple Time	161	4	134	35	22	33	46	48	54	59	68
	10	Regular Time	39,175	1	45	26	9	26	31	33	35	37	41
		Time and a Half	4,127	1	67	29	14	28	38	41	44	48	56
		Double Time	304	2	77	30	16	27	42	44	49	54	61
		Triple Time	129	4	102	32	17	28	35	38	44	54	65
Reading	9	Regular Time	41,941	1	75	42	15	43	53	55	58	62	67
		Time and a Half	4,085	1	134	41	21	39	53	57	62	68	77
		Double Time	470	2	149	39	23	36	52	58	63	69	83
		Triple Time	157	3	184	37	24	34	52	56	60	66	73
	10	Regular Time	38,922	1	74	37	15	37	47	49	52	56	61
		Time and a Half	4,150	1	112	36	20	33	47	51	56	62	74
		Double Time	307	1	128	36	23	33	48	52	58	67	85
		Triple Time	128	4	135	43	24	37	54	57	61	77	91

Table 3.3. Student Testing Time for Spring 2023 Utah Aspire Plus: Math and Science

Subject	Grade	Group	N	Testing Time (minutes)									
				Descriptive Statistics				Percentiles					
				Minimum	Maximum	Mean	St. Dev.	P50	P75	P80	P85	P90	P95
Math	9	Regular Time	40,502	1	74	51	17	53	64	66	69	71	73
		Time and a Half	4,016	2	140	48	23	46	62	66	71	77	92
		Double Time	463	2	201	50	29	44	65	72	78	86	104
		Triple Time	154	4	186	49	31	42	61	66	74	84	97
	10	Regular Time	38,362	1	75	44	18	46	58	61	64	67	71
		Time and a Half	4,082	1	112	39	23	37	53	58	62	69	81
		Double Time	308	2	149	44	29	38	62	66	74	82	101
		Triple Time	129	2	190	49	32	45	60	67	76	84	112
Science	9	Regular Time	41,872	1	80	32	12	32	40	42	44	48	52
		Time and a Half	4,070	1	89	30	16	27	39	43	47	52	59
		Double Time	469	1	110	31	19	27	42	45	50	55	65
		Triple Time	153	2	110	29	20	23	37	40	46	55	66
	10	Regular Time	38,770	0	66	26	13	26	35	37	39	43	48
		Time and a Half	4,049	1	89	24	16	21	33	36	40	46	57
		Double Time	303	1	100	26	20	21	36	41	45	51	62
		Triple Time	128	2	94	27	18	23	34	38	40	51	66

4 Score Reporting

4.1 IRT Pattern Scoring

Item parameters derived from previous IRT calibrations were used to estimate student ability (“theta”) scores by item response patterns. This is commonly referred to as pattern scoring. Pattern scoring takes advantage of the fact that items differ in their item characteristics and that an estimate of a student’s ability is based on their specific pattern of responses in combination with the item characteristics across all items. See Chapter 7 for more discussion of the IRT model and calibration methods.

The software package Operational Scoring: IRT Score Estimation (ISE V1.3.f; Chien & Shin, 2012) was used to perform the pattern scoring process and provide student scores on the IRT metric, using the student scored responses and the item response theory (IRT) item parameters for the operational items.

Two data-driven input files are required to execute the ISE software: a student response file and an item parameter file. The ISE algorithm combines the Newton-Raphson and Brute Force algorithms to generate the maximum likelihood estimated (MLE) of *theta* values. Specific configuration details include setting the upper- and lower-bound theta estimates, in this case +4 and -4, the number of iterations for the Newton-Raphson estimation method (30), the grid length interval for the Brute Force algorithm, the number of checking points for which the first derivatives are computed (120), and the number of decimal places for theta estimates (4).

IRT parameters for all 2023 Utah Aspire Plus operational items were used for estimating individual student scores for all forms. Table 4.1 presents the summary statistics for the IRT (*a*-, and *b*-) parameter estimates. The summary statistics shown include the total number of items, along with the mean, standard deviation (SD), minimum, and maximum.

Table 4.1. IRT Summary Parameter Estimates for Utah Aspire Plus Operational Items

Grade	Subject	No. of Items	Summary of <i>a</i> Estimates				Summary of <i>b</i> Estimates			
			Mean	SD	Min	Max	Mean	SD	Min	Max
9	English	46	0.86	0.40	0.21	1.77	-0.36	1.09	-2.70	1.71
	Reading	35	0.89	0.41	0.28	1.86	-0.12	0.87	-1.87	1.83
	Mathematics	40	0.97	0.35	0.33	1.76	0.25	0.81	-1.79	1.48
	Science	23	0.88	0.52	0.21	2.25	0.66	0.83	-0.67	2.22
10	English	43	0.92	0.35	0.35	1.77	-0.47	0.89	-1.84	2.06
	Reading	35	1.08	0.45	0.29	2.09	-0.06	0.66	-1.37	1.28
	Mathematics	40	1.04	0.27	0.48	1.51	0.43	0.64	-1.25	1.57
	Science	23	1.12	0.68	0.26	3.11	0.73	0.52	-0.01	2.31

4.1.1 Quality Control of IRT Scoring

Score tables used to estimate student scores on-demand were replicated independently through two parties internally. Additionally, a mock run of data was scored both using the on-demand process, and by two independent internal replicators. This scoring dry run was conducted at the overall test level as well as by reporting categories. Any differences were resolved and rerun until both parties' results were identical and deemed correct based on careful examination of output.

4.2 *Appropriate Uses for Scores and Reports*

As discussed, test forms constructed for Utah Aspire Plus cover a sampling of content as specified through test blueprints and reflective of the Utah Core Standards. The resulting scores reflect overall performance for each content area based on expectations of students' knowledge at the end of grades 9 and 10. It should be noted that while each test covers the standards, there is a limit to incorporating everything (e.g., given test time limits). Test scores should only be interpreted and used in the context from which they are obtained. In other words, Utah Aspire Plus test scores should be used to describe student achievement on the content assessed (i.e., grade level) and not used to generalize achievement beyond the test. In addition, academic placement decisions and promotions should not be based solely on these test scores but should include other indicators of achievement.

The Individual Student Report (ISR) communicates an individual student's test scores and interpretations of achievement based on those scores. The ISR provides the "snapshot" of achievement and explains the meaning of each piece of information provided, providing valuable information to students and parents. It is important that users of these reports do not extend the score information beyond the interpretations provided. A guide for understanding the ISR and its components can be found [online](#). For the Utah Aspire Plus tests, overall scale scores, performance level indicators, and predicted performance ranges for the ACT tests are provided. Note that no subscores are currently reported on student ISRs.

4.3 Utah Aspire Plus Reporting Scale

Commonly derived scores based on IRT are transformed to a reporting scale that is more consumable by users. The IRT metric being logit-based results in ability estimates typically ranging from -3.0 to 3.0 and to the second or third decimal. Interpreting differences across logits can be cumbersome. So scores are transformed to larger values without fractions. These are generally called scale scores. The purpose of scale scores is to facilitate interpretation and to report scores for all test-takers on a scale that remains consistent across multiple years or forms, even if the overall difficulty of the test varies slightly. Scale scores ensure that the test results mean the same thing regardless of which year the test was administered.

For the Utah Aspire Plus scales, the IRT metric uses a linear transformation to provide the final reporting scales as such:

$$SS = m\theta + b,$$

where m is the slope, and θ is the IRT person proficiency estimate obtained through pattern scoring. Using this equation, a scale score is transformed to the final reporting scale. The scale score metric for Utah Aspire Plus was chosen to range from 100 to 300, for each test and composite score. This range allows for the assessment to differ from the previous and remaining scales, and the slope chosen to spread final scores enough to contain each respective score distribution without floor or ceiling effects and to be dispersed enough to reasonably contain all transformed scores. The final transformation formula used for Utah Aspire Plus is:

$$SS = 25 \times \theta + 200$$

This transformation provides the following characteristics: 1) the mean of the scale is 200, 2) the standard deviation of the scale is 25, 3) the lowest operating scale score (LOSS) is 100, and 4) the highest operating scale score (HOSS) is 300. Composite scores were also created for Utah Aspire Plus. A composite score representing English Language Arts (ELA) is the average of a student's Reading and English scale scores, whereas a composite score representing Science, Technology, Engineering, and Mathematics (STEM) is the average of a student's Mathematics and Science scale scores.

4.4 Standard Setting

Descriptions of student performance are often used to help enhance the reporting of student scores beyond an overall reported score and references to other students or groups of students. Performance levels and descriptions of performance divide the test scores into meaningful categories and align to performance ranging from low to high. For Utah, these categories are called *Below Proficient*, *Approaching Proficient*, *Proficient*, and *Highly Proficient*. Performance level descriptions (PLDs) accompany these labels to describe typical performance of students within each group.

Standard settings were conducted in August of 2019 (for all subjects) and again for science in August of 2022 following the first administration of the new assessment based on the SEEds. PLDs are the core of all standard setting meetings. The PLDs for the Utah Aspire Plus assessments can be found [online](#).

Utah educators were convened to operationalize the PLDs through standard setting, a process of determining test score thresholds, or “cut points,” to divide the test scores into the four performance groups. Final scale score cuts for science, English, reading and mathematics are presented in Table 4.2.

Table 4.2. Utah Aspire Plus Scale Score Cuts by Grade and Subject

Grade	Subject	Scale Score Cut Points		
		Approaching Proficient	Proficient	Highly Proficient
9	English	165	202	242
	Reading	166	204	231
	Mathematics	172	206	233
	Science	187	211	237
10	English	161	200	245
	Reading	175	204	235
	Mathematics	181	210	236
	Science	187	210	240

4.5 ACT Predicted Score Ranges

As noted, one of the goals of the Utah Aspire Plus assessments is to be predictive of college readiness at grades 9 and 10, and the means of this is in terms of providing prediction score ranges of performance on the ACT for the four subject tests (English, math, reading, and science) and the Composite score (the average of the four subject tests). Predicted ranges of performance were determined originally between ACT Aspire scores and ACT scores, where for a given ACT Aspire score, there was a distribution of related ACT scores. The bounds of the range were denoted by the scores closest to the 25th and 75th percentiles of the ACT score distribution, conditional on ACT Aspire scores. For Utah Aspire Plus, an additional error term was added to account for error attributable to linking the Utah Aspire Plus scores.

Students can use the predicted scores together with the ACT College Readiness Benchmarks to monitor their preparedness to be college-ready by the end of high school. Utah students take the ACT® during their junior year of high school. Specific details from the original prediction score studies can be found in the *Utah Aspire Plus 2018–2019 Technical Report* (available at http://utah.pearsonaccessnext.com/resources/additional-services/UT1132740_UTPlusTechReportv4.3_WebTag.pdf).

In addition to relying on the relationship between the Utah Aspire Plus tests to the ACT Aspire scales for deriving the initial ACT prediction score ranges for the 2019 administration, the intention was to provide updated predictions based on longitudinal data as it becomes available. The updated ACT score ranges directly link the Utah Aspire Plus scores at grades 9 and 10 to ACT scores at grade 11. In spring 2020, the first longitudinal data was available for this purpose. The initial longitudinal Utah-to-ACT prediction studies were based on students who were in the 10th grade during the 2019 administration of the Utah Aspire Plus tests. The second longitudinal study was conducted in the spring of 2021. Details of this study can be found in Appendix J of the *Utah Aspire Plus 2020–2021 Technical Report* (available at http://utah.pearsonaccessnext.com/resources/additional-services/UT1140119_UTPlusTechReport2022_WebTag.pdf).

A third longitudinal study was conducted in 2022 to update the science grade 10 predictions. This study included students who were in 10th grade in 2021 and took the ACT as 11th grade students in spring 2022. Details of this study can be found in Appendix H of the *Utah Aspire Plus 2021–2022 Technical Report* (available at <https://schools.utah.gov/file/5da11ee6-32a3-4eeb-acdd-1b1e21348659>).

4.6 2022–2023 Utah Aspire Plus Performance Results

Descriptive statistics of the scale scores for each Utah Aspire Plus assessment are in Appendix H. The descriptive statistics are provided for the overall testing population, as well as by subgroups—sex, ethnicity, and special populations. Average scale scores as well as standard deviations, scores at the 25th, median, and 75th percentiles are also reported as well as skewness. Scale score distributions for each Utah Aspire Plus assessment are provided in Appendix I, for the overall testing population. Appendix J contains the performance level distributions of each Utah Aspire Plus. The tables contain the percentages of students being classified into each respective performance level.

While results can be compared directly to 2019, 2021, and 2022 performance within the same subject and grade, extra cautions should be taken with respect to interpretations beyond high-level due to impacts from the pandemic. These opportunity-to-learn (OTL) impacts are multi-faceted and differential across the state.

While comparatively, a similar number of students were tested in 2023 as compared to 2019, the percent of completed tests varied. In 2019 completion rates for registered testers was approximately 91–93%. In 2022, the completion rate ranged from 84–88%. In 2023, the completion rate ranged from 84–90%. Overall performance was similar in 2023 to 2022.

5 Classical Item Analyses

5.1 Item Analyses

In Chapter 2, statistical indices used in the test construction process were introduced. To build the initial test forms for Utah Aspire Plus, item statistics based on use within the SAGE and ACT Aspire tests served to guide test construction activities. As noted, while the best initial forms were created, there were instances in which not all statistical targets were fully met. This chapter describes in more detail those classical item statistics. Additionally, after the Utah Aspire Plus 2022–2023 operational administration, classical item statistics were also calculated. Results are presented in Appendix C.

5.1.1 p -Value and Item Mean Scores

Item difficulty offers an index of how easy or hard a given test question is to answer correctly or to earn a given score point for items scored according to a rubric. For dichotomously scored items (items scored correct or incorrect), item difficulty is indicated by its p -value, which is the proportion of test takers who answered that item correctly. The range for p -values is from 0 to 1.

For polytomously scored items (items scored according to a rubric with multiple points awarded), difficulty is indicated by the mean item score. Here the average ranges from 0 to the maximum total possible points for an item. To facilitate interpretation, the mean item values for polytomously scored items can also be expressed on the p -value metric as percentages of the maximum possible score.

5.1.2 Item-Test Score Correlations

Correlations between a given item score and total test score are used to evaluate how well items differentiate between “high” and “low” performing students. In general, the higher the correlation the better an item is at differentiating between high- and low-performing students. As this index is a correlation, it ranges from -1 to $+1$ (where $+/- 1$ reflects a perfect correlation and 0 reflects no correlation). When the correlation is negative, it means low-performing students on the test are answering the given question correctly more often than high-performing students, and this would be a reason to further investigate the item for potential flaws.

In addition to the correlation between item score and total test score, the same approach can be applied to each answer option of multiple-choice items. Although not provided in this report, this information is used within the context of data review and allows for further evaluation of the full functioning of multiple-choice items, as it focuses on the effective functioning of the options (distractors) which are other than the correct answer.

5.1.3 Differential Item Functioning

Differential item functioning (DIF) exists when an item functions differentially across identifiable subgroups (e.g., sex or ethnicity) where students are matched on ability (meaning comparisons are made between students of the same ability, so differences are not attributable to overall group performance differences). In this context, DIF may indicate an issue with fairness or that the item may be measuring something other than the intended construct (i.e., possible evidence of DIF). It is important, however, to recognize that item performance differences flagged for DIF might be related to actual differences in relevant knowledge or skills (item impact) or statistical Type I error. As a result, DIF statistics are used to identify potential biases. Subsequent reviews by content experts and bias/sensitivity committees are required to determine the source and meaning of performance differences.

There are multiple statistical procedures for analyzing DIF, one of which is based on the Mantel-Haenszel chi-square statistic (M-H χ^2) for multiple-choice items (Holland and Thayer, 1988). The chi-square statistic determines whether the odds of a correct response on an item is the same for both focal and reference groups, across all levels of proficiency. The Mantel-Haenszel odds ratio (α_{M-H}) is the odds of a correct response of the reference group divided by the odds of a correct response of the focal group. Data for these Mantel-Haenszel procedures are drawn from 2-by-2-by- k (score levels) contingency tables, for each item. As shown in Table 5.1 the number of focal and reference group members scoring in each possible item response is captured.

Table 5.1. Item 2x2 Contingency Table for the k^{th} Score Level

Group	Item Score		Total
	Correct (1)	Incorrect (0)	
Focal (f)	n_{f1k}	n_{f0k}	n_{fk}
Reference (r)	n_{r1k}	n_{r0k}	n_{rk}
Total (t)	n_{t1k}	n_{t0k}	n_{tk}

For classifications of DIF, the Mantel-Haenszel Delta DIF statistic (MHD: Dorans & Holland, 1993) is computed from the Mantel-Haenszel odds ratio and used in conjunction with M-H χ^2 to classify items into three categories distinguishing magnitudes of DIF: negligible DIF (A), moderate DIF (B), and large DIF (C). Classification is based on the following guidelines:

- M-H χ^2 not significantly different from 0 or |MHD| less than 1 results in a classification of A.
- M-H χ^2 significantly different from 0 and |MHD| at least 1 but less than 1.5 **or** M-H χ^2 not significantly different from 0 and |MHD| greater than 1 results in a classification of B.
- M-H χ^2 significantly different from 0 and |MHD| at least 1.5 results in a classification of C.

In addition to these classifications, notation of DIF includes a positive (+) sign, indicating that the item favors the focal group, or a negative (-) sign, indicating that the item favors the reference group. Items that are designated with “B” or “C” DIF classifications are recommended for review before continued use on assessments.

The standardized mean difference (SMD: Zwirk, Donoghue, and Grima, 1993) procedure is also used for detecting DIF, for items worth more than one point. SMD is a summary statistic used as an effect size estimate comparing the mean item score between the reference and focal groups (the two groups being compared). Although the numerical result of this statistical procedure is different from the M-H statistics, the classification of the results is the same—the results are classified into three categories indicating the magnitude of DIF with additional notation indicating the favored group.

5.2 Classical Item Summaries for Operational Administration

As noted, summaries of classical item statistics from the initial operational administration of Utah Aspire Plus are located in Appendix C. Examination of the distribution of items by difficulty across each test shows that items do vary in difficulty across each test, with most items between 0.30 and 0.75. There are items that did fall outside the guidelines outlined previously. Their inclusion was necessary to meet blueprints given limitations to the available item banks. The same can be said of the distributions of item-total correlations and DIF results, where there were items included in the tests that fell outside the guidelines but were ultimately included on final forms as the best available. Overall, even where items fell outside the guidelines, they were still useful. This was particularly true for the science assessments, where due to bank limitations and cluster design, some very difficult items and items with low discrimination were included on final operational forms to help hit blueprint targets.

6 Reliability

Estimation of reliability of a given assessment is critical in order to understand the precision of measurement for individual test scores. Test score reliability estimates are typically provided in both a classical as well as an item response theory (IRT) context. Classical reliability estimates such as standard error of measurement (SEM) or Cronbach's alpha are reliability measures of internal consistency. Where classical approaches are generally single indicators for a given assessment, IRT reliability reflects precision across the ability spectrum. There are a number of different approaches available to estimate reliability of test scores. For Utah Aspire Plus tests, both classical reliability and reliability within an item response theory framework were computed.

6.1 Classical Definition of Reliability

The basis of classical test theory is premised on the idea that a person's observed score is the sum of their true score (measured without error and not directly observable) plus error:

$$\text{Observed Score} = \text{True Score} + \text{Error}.$$

It provides a means of describing the quality of test scores through the interplay of these three elements. Arguably the most important descriptor is the concept of the reliability of test scores, where the reliability of observed scores is defined as follows:

$$\text{Reliability} = \frac{\sigma_T^2}{\sigma_O^2} = \frac{\sigma_T^2}{\sigma_T^2 + \sigma_E^2} = 1 - \frac{\sigma_E^2}{\sigma_O^2}$$

where σ_T^2 is the true score variance, σ_O^2 is the observed score variance, and σ_E^2 is the error variance. When there is no error, the reliability is the true score variance divided by true score variance, which is unity. However, as more error influences the measure, the error component in the denominator of the ratio increases and the reliability decreases.

6.2 Classical Test Theory Reliability Estimates

6.2.1 Cronbach's Alpha

Internal consistency methods use a single administration to estimate test score reliability. For state assessments where student testing time is at a premium, internal consistency procedures have a practical advantage over reliability estimation procedures requiring multiple tests. Probably the most frequently used internal consistency reliability estimate is the coefficient alpha (Cronbach, 1951). Coefficient alpha assumes that inter-item covariance constitutes true-score variance and the fact that the average true score variance of items is greater than or equal to the average inter-item covariance. The formula for the coefficient alpha is

$$\alpha = \left(\frac{N}{N-1} \right) \left(1 - \frac{\sum_{i=1}^N s_{Y_i}^2}{s_X^2} \right),$$

where N is the number of items on the test, $s_{Y_i}^2$ is the sample variance of the i^{th} item (or component), and s_X^2 is the observed score sample variance for the test.

Coefficient alpha reliability estimates are provided in Appendix D for the overall testing population as well as by sex, ethnicity, and other student breakout groups. In addition, they are also provided by each reporting category (though again it should be noted that currently, only overall scores are reported on individual student reports, and no subscores are reported).

6.2.2 Standard Error of Measurement

A reliability coefficient expresses test score consistency in terms of variance ratios. In contrast, the standard error of measurement (SEM) expresses score inconsistency (unreliability). The SEM is an estimate of how much error there is likely to be in an individual's observed score, or alternately, how much score variation would be expected if the individual were tested multiple times with equivalent forms of the test. The SEM is calculated using the following formula:

$$SEM = s_x \sqrt{1 - \rho_{XX'}},$$

where s_x is the standard deviation of the total test (standard deviation of the raw scores) and $\rho_{XX'}$ is a reliability estimate for the set of test scores. Test standard errors of measurement are provided in Appendix D and are presented on the Utah Aspire Plus scale score metric ($s_x = 25$).

6.3 IRT-Based Reliability

Where estimation of reliability is within a classical test theory frame, it should be noted that such measures are sample specific. Additionally, error estimates such as the SEM are group-level estimates that apply across test scores. And it is sometimes viewed as unrealistic that the size of errors would be unrelated to the “true scores” of examinees (identical for all).

For Utah Aspire Plus, student scores are derived within an item response theory framework (IRT) through pattern scoring based on the three-parameter logistic (3PL) and two-parameter logistic (2PL) measurement models (these are more thoroughly described later in this report). Under the IRT model, measurement precision is expressed as Conditional Standard Errors of Measurement (CSEM) and is equal to the inverse of the square root of the test information function across the ability continuum (see Hambleton and Swaminathan, 1985).

CSEMs depend upon both the unique set of items each student answers correctly and their estimated ability level (θ). Therefore, different students will likely have different CSEM values even if they have the same raw score and/or theta estimate. Each item contains a unique amount of information for a given ability level, which depends on each item’s discrimination, difficulty, and pseudo-guessing parameters.

The conditional standard errors for Utah Aspire Plus tests are provided in Appendix E, each including a line indicating the scale score cut score for Proficient. Ideally, the lowest value of conditional standard error of measurement occurs at the location of Proficient.

6.4 Reliability of Performance Level Categorization

Every test administration will result in some error in classifying examinees. The concept of the SEM provides a mechanism for explaining how measurement error can lead to classification errors when cut scores are used to classify students into different achievement levels. For example, some students may have a true achievement level greater than a cut score. However, due to random variations (measurement error), their observed test score may be below the cut score. As a result, the students may be classified as having a lower achievement level. As discussed in Section 6.2.2, a student’s true score is most likely to fall into a standard error band around their observed score. Thus, the classification of students into different achievement levels can be imperfect, especially for the borderline students whose true scores lie close to achievement-level cut scores.

For the Utah Aspire Plus assessment, the levels of achievement are *Below Proficient*, *Approaching Proficient*, *Proficient*, and *Highly Proficient*. A description and analysis of classification accuracy and consistency indices are provided below. All indices were calculated using the BB-CLASS software (Brennan, 2005).

6.4.1 Accuracy and Consistency

Accuracy refers to the extent to which achievement decisions based on test scores match those that would be made if the scores did not contain any measurement error, i.e., “true scores.” Since true scores are not available, an estimate of the true score distribution must be determined for classification accuracy to be estimated. Consistency, on the other hand, refers to the extent to which achievement classification decisions based on test scores match the decisions based on a second, parallel form of the same test. This index assumes that two parallel forms of the same test are administered to the same group of students. In Utah, however, this is impractical. Livingston and Lewis (1995) developed techniques to estimate both accuracy and consistency that overcome the constraints of true scores and multiple test forms on the same students. These procedures are used to generate accuracy and consistency indices on the Utah Aspire Plus assessments.

6.4.2 Calculating Accuracy

To calculate accuracy, a 4 x 4 contingency table is created for each subject area and grade. The $[x, y]$ entry of an accuracy table represents the estimated proportion of students whose true score fall into performance level x and whose observed scores fall into performance level y . Table 6.1 is an example of an accuracy table where the columns represent test-based student achievement, and the rows represent true achievement-level decisions. In this example, the total accuracy is approximately 75%, the sum of the diagonal (shaded) cells.

Table 6.1. Example Accuracy Classification Table

True Score	Observed Score				Total
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.117	0.034	0.000	0.001	0.152
Approaching Proficient	0.019	0.161	0.061	0.002	0.243
Proficient	0.000	0.034	0.294	0.061	0.389
Highly Proficient	0.000	0.000	0.036	0.179	0.215
Total	0.136	0.229	0.391	0.243	1.000

It is useful to consider decision accuracy based on a dichotomous classification of *Below Proficient* or *Approaching Proficient* versus *Proficient* or *Highly Proficient* because Utah uses *Proficient* and above as proficiency for accountability decision purposes as well as for an index tracking students' readiness to college and careers. To compute decision accuracy in this case, the table is dichotomized by combining cells associated with *Below Proficient* and *Approaching Proficient* and combining *Proficient* with *Highly Proficient*. The sum of the shaded cells in

Table 6.2 indicates classification accuracy around the Proficient cut point of approximately 90%. The percentage of examinees incorrectly classified as *Approaching Proficient* or lower, when their true score indicates *Proficient* or above, is approximately 3%.

Table 6.2. Example Accuracy Classification Table for Proficient Cut Point

True Score	Observed Score				Total
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.117	0.034	0.000	0.001	0.152
Approaching Proficient	0.019	0.161	0.061	0.002	0.243
Proficient	0.000	0.034	0.294	0.061	0.389
Highly Proficient	0.000	0.000	0.036	0.179	0.215
Total	0.136	0.229	0.391	0.243	1.000

6.4.3 Calculating Consistency

Consistency can be calculated in the same manner, via 4 x 4 contingency table, albeit with data indicating an estimate of the joint distribution of classifications on (hypothetically) two independent, parallel test forms. Table 6.3 shows sample statistics of consistency classification. Based on this sample data, the overall consistency is approximately 67%. The consistency at *Proficient* is 87%. The agreement rates are lower than those for accuracy because both classifications contain measurement error; whereas in the accuracy table, true score classification is assumed to be without error.

Table 6.3. Example Consistency Classification Table

First Form	Second Form				Total
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.111	0.043	0.009	0.001	0.164
Approaching Proficient	0.019	0.147	0.073	0.004	0.243
Proficient	0.006	0.038	0.252	0.075	0.371
Highly Proficient	0.000	0.002	0.056	0.163	0.221
Total	0.136	0.230	0.390	0.243	1.000

6.4.4 Calculating Kappa

Another way to express overall consistency is to use Cohen’s kappa (κ) coefficient (Cohen, 1960), which assesses the proportion of consistent classifications beyond chance. The coefficient is computed using

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

where P is the proportion of consistent classifications and P_c is the proportion of consistent classification by chance. Using Table 6.3, P is the sum of the shaded cells whereas P_c is

$$\sum_x C_x \cdot C_{x'}$$

where C_x is the proportion of students whose observed performance level would be x on the first form, and $C_{x'}$ is the proportion of students whose observed performance level would be x on the second form. Therefore, the kappa coefficient using the data from Table 6.3 is 0.548. Cohen suggested the Kappa result be interpreted as follows: values ≤ 0 as indicating no agreement and 0.01–0.20 as none to slight, 0.21–0.40 as fair, 0.41–0.60 as moderate, 0.61–0.80 as substantial, and 0.81–1.00 as almost perfect agreement. Estimates of classification accuracy and consistency indices—including kappa coefficients—for overall performance level classification and at the Proficient cut point are provided in Appendix F.

7 Field Test Calibration and Drift Analyses

7.1 IRT Overview

Item response theory (IRT) was used to create the base scales for the Utah Aspire Plus assessments. For the 2023 administration, the mathematics, reading, English, and science assessments were pre-equated for the first time. Item parameters were estimated either from prior operational post-equating, or field test calibration. See the *Utah Aspire Plus 2021–2022 Technical Report* (available at <https://schools.utah.gov/file/5da11ee6-32a3-4eeb-acdd-1b1e21348659>) and prior technical reports for details on these processes. Student scores were estimated using IRT and then transformed to the final Utah Aspire Plus scale score reporting metric. Scores were reported on-demand.

Following administration, a separate calibration and equating process was conducted. While these results did not affect student scores for Spring 2023, they served several purposes:

- Calibration of field test items
- Identification of items with parameter drift
- Update of bank parameters

Final parameters resulting from the calibration and equating processes were used to update parameters in the item bank for items in the following categories:

1. Item was field tested in 2023.
2. Item was used operationally for the first time in 2023 (prior parameters were from field test administration)
3. Item showed drift during the equating process, as described in Section 7.3.5.

In this section of the technical report, the following topics related to IRT calibration and equating are discussed:

- IRT Data Preparation
- Description of the Calibration Process
- Drift Analyses

7.2 IRT Data Preparation

7.2.1 Student Inclusion/Exclusion Rules

The data preparation for the IRT calibration process began with all Utah students who were administered the “base” forms (i.e., online, English-language forms).

The samples for item parameter estimation included the following:

- Students from the online, English language test forms,
- Students with the same grade battery of tests, and
- Students with a valid test score status for a subject test.

Students without a valid test score were excluded from calibration data.

7.2.2 Quality Control of the IRT Data Matrix Files

Student records in the calibration data files were ordered by ascending student identification number. In the case where field test forms are used, student records would first be sorted by form, then by student identification number. The array of item responses was presented in the order as administered in the test form, including items that are presented in field test slots.

The IRT data matrices were created independently by two Pearson psychometric staff. The matrices were checked for accuracy by comparing numbers of students (counts) and the item response arrays. Any discrepancy found was resolved. Final calibration data files matched perfectly.

7.3 Description of the Calibration, Equating, and Scaling Process

7.3.1 IRT Models

Multiple item types are used on Utah Aspire Plus assessments and require multiple measurement models. Traditional multiple-choice items, with one correct answer, are analyzed via the three-parameter logistic model (3PLM; Birnbaum, 1968), denoted as

$$p_i(\theta_j) = c_i + \frac{1-c_i}{1+e^{-Da_i(\theta_j-b_i)'}}$$

where $p_i(\theta_j)$ is the probability that student j would earn a score of 1 on item i , b_i is the difficulty parameter for item i , a_i is the slope (or discrimination) parameter for item i , c_i is the pseudo-chance (or guessing) parameter for item i , and D is the constant 1.7. Other selected response items worth one point (e.g., technology-enhanced items) are analyzed via the two-parameter logistic model (2PLM; Birnbaum, 1968), which is a reduced model from the 3PLM, where the pseudo-chance parameter, c , is assumed zero. Items worth two points were analyzed via the generalized partial credit model (GPCM; Muraki, 1992), denoted as

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

where $a_i(\theta_j - b_i + d_{i0}) \equiv 0$, $p_{im}(\theta_j)$ is the probability of an examinee with θ_j getting score m on item i , and M_i is the number of score categories of item i with possible item scores as consecutive integers from 0 to $M_i - 1$. In the GPCM, the d parameters define the “category intersections” (i.e., the θ value at which examinees have the same probability of scoring 0 and 1, 1 and 2).

7.3.2 IRTPRO Calibration Procedures and Convergence Criteria

The primary goal of the IRT calibration was to place the operational and field test items from a given test onto a common scale. The additional step of equating was also completed to place these parameters onto the original Utah Aspire Plus base scales.

Note that large enough samples are necessary to sufficiently estimate IRT parameters for a given test and across the respective models (generally for state summative tests similar to Utah Aspire Plus on order of 2,000). IRTPRO (Scientific Software International, Inc., 2017) was used to obtain the IRT parameter estimates using the measurement models described in Section 7.3.1. The software default estimation method, Bock-Aitkin (BAEM), was used for each calibration. The prior distributions for latent traits were set to a mean of zero and a standard deviation of one. The number of quadrature points used in the estimation was set to 49. For item parameters, a prior was placed on the lower asymptote (pseudo-chance) for the 3PLM: a normal distribution with a mean of -1.4 and a standard deviation of one. After calibration, convergence was checked.

To convert IRTPRO item parameters to the commonly used logistic parameter presentation, the a -parameter from the IRTPRO output needed to be converted since IRTPRO uses 1.0 for a scaling constant. The formula for this conversion is:

$$a_{new} = \frac{a_{irtpro}}{1.7}.$$

7.3.3 Calibration Quality Control

IRT calibrations were conducted independently by two Pearson psychometric staff using the same software program. All item parameters from both independent calibrations were compared. Item fit plots were generated as further analyses of reasonableness and support of decisions of items' future use.

7.3.4 Equating

A common item non-equivalent groups approach (Kolen and Brennan, 2014) was used for equating the 2023 forms to the base scales.

The Stocking and Lord (1983) test characteristic curve methodology was used to derive equating constants for each grade-subject test. The operational items were used as the common-item linking set. The banked IRT item parameter estimates for all of the Utah Aspire Plus operational items, and the respective item parameter estimates from the 2023 administration described in Section 7.3.2, were used to obtain transformation constants. This was conducted using the computer program STUIRT (Kim & Kolen, 2004).

Equating was carried out in conjunction with a drift analysis procedure, described in Section 7.3.5, which resulted in a final set of Stocking and Lord scaling constants. These constants were then applied to all 2023 calibrated items to obtain a set of parameters for the operational and field test items. Final Stocking and Lord scaling constants used for placing tests onto the Utah Aspire Plus base scales are presented in Table 7.1.

Table 7.1. 2023 Final Stocking and Lord Scaling Constants

Subject	Grade	Slope	Intercept
English	9	1.026	-0.160
	10	1.031	-0.148
Reading	9	1.034	-0.061
	10	0.989	-0.114
Math	9	1.041	-0.309
	10	1.024	-0.277
Science	9	1.047	0.070
	10	1.021	-0.160

Final parameters were then updated in the item bank for items in the following categories:

4. Item was field tested in 2023.
5. Item was used operationally for the first time in 2023 (prior parameters were from field test administration)
6. Item showed drift during the equating process, as described in Section 7.3.2.

7.3.5 Drift Analysis

A critical step in carrying out an equating is to evaluate the anchor items for stability in relation to its banked item characteristics. Items that deviate substantively in relation to the entire set of anchor items may be removed from contributing to the final equating solution. For Utah Aspire Plus, the item parameter stability check for the operational items was conducted using classical item analyses, scatter plots of item parameter estimates, and item-characteristic curve (ICC) comparison. For the ICC comparison, old and new ICCs were compared using the z-score approach based on D^2 (Wells, Hambleton, Kirkpatrick, & Meng, 2014) as outlined below:

1. Obtain the theoretically weighted estimated posterior theta distribution using 31 quadrature points (-5 to 5).
2. Compute the slope and intercept constants using Stocking and Lord in STUIRT with all operational items in the linking set.
3. Place the freely calibrated item parameter estimates onto the baseline scale by applying the constants obtained in Step 2.
4. For each operational item, calculate D^2 between the ICCs based on old (x) and new (y) parameters at each point in this theta distribution:

$$D_i^2 = \sum^k [P_{ix}(\theta_k) - P_{iy}(\theta_k)]^2 \cdot g(\theta_k)$$

where i = item, x = old form, y = new form, k = theta quadrature point, and g = theoretically weighted posterior theta distribution.

5. Flag items with a D^2 that is greater than the mean D^2 value, and whose distance from the mean D^2 value is greater than twice the standard deviation of the D^2 values.
6. Examine the impact of removing a flagged item on the content representativeness of the resulting anchor set. A flag alone is not the sole criteria for removing an item from the anchor set. It is important to also make sure that the remaining anchor set continues to be representative of the overall content and structure of the test.

Plots showing D^2 values following the initial equating are given in Appendix M. Counts of operational items showing drift are given in Table 7.2.

Table 7.2. 2023 Items Showing Drift

Subject	Grade	Number of items showing drift
English	9	1
	10	3
Reading	9	1
	10	1
Math	9	2
	10	1
Science	9	2
	10	1

Following removal of items for drift, the STUIRT equating process was repeated with the updated anchor set to obtain a final set of Stocking and Lord scaling constants, which were applied to the freely calibrated item parameters to obtain a final set of parameters. Parameters in the item bank were updated to these parameters for items showing drift, as well as for field test items and items which were operational in 2023 for the first time.

Scatterplots of the operational items can be found in Appendix G. Overall, item functioning of common items can be described as typical and stable. No more than three items in any of the common item sets were removed from final linking solutions. Scatterplots and correlations of IRT difficulty and discrimination parameters showed strong correlations.

7.4 Model Fit Evaluation Criteria

The Q_1 statistic (Yen, 1981) was used as an index of correspondence between observed and expected performance. To compute Q_1 , first the estimated item parameters and student response data (along with observed item scores) were used to estimate student ability ($\hat{\theta}$). Next, expected performance was computed for each item using students' ability estimates in combination with estimated item parameters. Differences between expected item performance and observed item performance were then compared at 10 intervals across the range of student achievement (with approximately the same number of students per interval). Q_1 was computed as a ratio involving expected and observed item performance. Q_1 is interpretable as a chi-squared (c^2) statistic, which can be compared to a critical chi-squared value to make a statistical inference about whether the data (observed item performance) were consistent with what might be observed if the IRT model was true (expected item performance). Q_1 is not directly comparable across different item types because items with different numbers of IRT parameters have different degrees of freedom (df). For that reason, a linear transformation (to a Z-score, Z_{Q_1}) was applied to Q_1 . This transformation also made item fit results easier to interpret and addressed the sensitivity of Q_1 to sample size.

To evaluate item fit, Yen's Q_1 statistic was calculated for all items. Q_1 is a fit statistic that compares observed and expected item performance. For dichotomous items, Q_1 was computed as

$$Q_{1i} = \sum_{j=1}^j \frac{N_{ij}(O_{ij}-E_{ij})^2}{E_{ij}(1-E_{ij})},$$

where N_{ij} was the number of examinees in interval (or group) j for item i , O_{ij} was the observed proportion of the students for the same cell, and E_{ij} was the expected proportions of the students for the same interval. The expected proportion was computed as

$$E_{ij} = \frac{1}{N_{ij}} \sum_{a \in j}^{N_{ij}} P_i(\hat{\theta}_a),$$

where $P_i(\hat{\theta}_a)$ was the item characteristic function for item i and students a . The summation is taken over students in interval j .

The generalization of Q_1 for items with multiple response categories is

$$Gen Q_{1i} = \sum_{j=1}^{10} \sum_{k=1}^{m_i} \frac{N_{ij}(O_{ikj}-E_{ikj})^2}{E_{ikj}},$$

where

$$E_{ikj} = \frac{1}{N_{ij}} \sum_{a \in j}^{N_{ij}} P_{ik}(\hat{\theta}_a).$$

Both Q_1 and generalized Q_1 results were transformed to ZQ_1 and were compared to a criterion $ZQ_{1,crit}$ to determine acceptable fit. The conversion formula was

$$ZQ_1 = \frac{Q_1 - df}{\sqrt{2df}},$$

and

$$ZQ_{1,crit} = \frac{N}{1500} * 4,$$

where df is the number of degrees of freedom. The number of degrees of freedom is equal to the number of independent cells less the number of independent item parameters. For example, the degrees of freedom for polytomous items equals $[10 \times (\text{number of score categories} - 1) - \text{number of independent item parameters}]$. For the GPCM, the number of independent item parameters equals 1 (for the α -parameter) plus the number of step values (e.g., for an item scored 0, 1, 2: there are 2 independent step values—the b parameter is simply the mean of the step values and is not, therefore, independent).

As all items were pre-equated, Q_1 statistics were calculated in previous administrations, along with item fit plots. All items included on previous forms showed adequate fit. Additionally, Q_1 and item fit plots were re-generated following the 2023 administration to assess pre-equating. Results were consistent with the drift analyses and did not suggest any concerns with model selection.

8 Quality Control

Quality control is a critically important element of every phase of the Utah Aspire Plus development, administration, and score reporting in ensuring the accuracy of student-, school- and district-level data. Pearson has developed and refined a set of quality procedures to help ensure that all USBE's testing requirements are met or exceeded. These quality control procedures are detailed in the paragraphs that follow. In general, Pearson's commitment to quality is incorporated in both task-specific quality standards applied to processing functions and services as well as a network of systems and procedures that coordinate quality steps across functions and services.

8.1 *Online Assessment Delivery*

8.1.1 Item Validation

Test items for Utah Aspire Plus are housed in Pearson's Automated Banking and Building for Interoperability (ABBI) platform. ABBI supports building and publishing online and paper-based tests and drives creation of those forms to both Pearson's paper and online publishing systems. Through ABBI, item scoring configuration is validated during initial item review (i.e., at the time of item writing) as well as during forms development.

8.1.2 Test Administration

PearsonAccess^{next} is Pearson's next-generation system for managing student data, paper, and online test administration, scoring, and reporting high-stakes assessments. This system provides comprehensive support for paper and online testing either through a single sign-on destination or by interfacing with other systems to provide a highly adaptable solution. TestNav delivers online tests. The core functionalities of TestNav include delivering tests to students, collecting student responses, and returning the responses to Pearson for scoring.

TestNav provides advance warning of network issues that prevent sending student responses to the Pearson testing server. When the network is functioning normally, TestNav sends student responses to the Pearson testing server in real time, while the student is testing. If the student's device cannot connect to the Pearson servers, TestNav saves the response to an encrypted file and allows the student to continue testing. When the network connection is reestablished, the test proctor can upload a student's saved responses to Pearson's testing server, and then TestNav erases the encrypted response file from the student's device or local network.

In the event of a non-network or non-Internet issue, such as a power outage or student device shutdown, student responses are saved to the encrypted file. When the student resumes testing, the system uploads the data in the file to the servers, and the student continues at the point in the test when the issue occurred.

As part of test security, test administrators control individual student authorization by printing and distributing testing tickets with each student's identifying information and unique log-in credentials. The student enters their log-in and password on the testing workstation to gain access to the test. To further secure the testing environment, a non allowed list capability sends notifications when unapproved applications are running when the test is started. Once all non allowed applications are shut down, TestNav starts in kiosk mode when a student signs into a secure test.

Kiosk mode locks down the testing computer or device, so the student cannot print, cut, or copy test content. Students cannot visit websites or access other installed applications not approved for use during the test.

8.1.3 Operational Monitoring

Pearson's operational monitoring practices and tools constantly verify that platforms remain available to users; that performance stays within acceptable limits; and that users do not encounter critical errors. The types of monitoring that Pearson performs to help keep testing on time and reduce the chance of interruptions include the following:

- Site Availability Monitoring – checking locations and providing alerts when response times or availability thresholds are crossed
- Synthetic User Monitoring – simulating key end-user actions (launching a test, logging into the administrative site, viewing reports, etc.) and running from several locations on the public internet
- End User Monitoring – analyzing page and click performance to verify that end users receive results in a reliable and timely manner
- Server Monitoring – collecting detailed metrics on server performance to gauge health
- Application Performance Monitoring – gathering detailed performance information about the health of Pearson's various assessment platforms
- Database Monitoring – using a variety of tools to watch performance in real time
- Event Monitoring and Real-Time Security Auditing – processing large volumes of machine-generated data in real time to look for trends, issues, or anomalies
- Systems Vulnerability Monitoring – monitoring multiple sources for newly identified vulnerabilities in systems and applications Pearson uses

8.2 *Production System Testing*

8.2.1 Functional Testing

Well before testing the entire system, Pearson engineers develop tests for each discrete software unit, and for small groups of related units. Debugging code is emphasized in the earliest stages of development, so during unit testing, each developer creates unique tests for code that has been written.

8.2.2 Integration Testing

Digital and traditional paper solutions require testing that is specific to its unique interactions and specifications. After testing each piece of component code, the behavior of the integrated parts is tested. In the first stage of integration testing, the testing is done at the base system level to verify and validate that the system components function together. The second stage of integration testing examines accuracy of the unique configuration to each administration specified in the contract.

Configuration requirements are the basis of our integration testing. This is documented, and test cases and results are maintained and verified prior to the final production scoring and reporting configuration, including item parameter files, keys, and cut scores.

8.2.3 Program Validation End-to-End Testing

After Product Testing approval, the Pearson Program Validation team uses a cross-system end-to-end approach to validate the user interface, scoring, data files, and reports. This testing confirms that all data are consistent with customer requirements by emulating the customer experience throughout the program lifecycle.

The Program Validation team coordinates test-material processing (distribution and data collection) with the same operational areas that process live material during production. Where appropriate, there is a Production Sample Verification process, which uses the first available student data as a final quality step before live production processing of materials to be distributed. An examination of the outputs verifies data are scored, aggregated, reported, and delivered accurately. After the Program Validation team approves, the delivery of code and configuration is moved to production.

8.2.4 Load Testing

To examine the system's expected performance during peak usage days, Pearson engineers will assemble the components and test the system under load conditions. During load testing, a period of peak production is modeled to identify any issues within the application that might be triggered by maximum activity. Load testing is performed several times per year so that the system can be scaled to meet anticipated customer demand in advance of when it is needed.

8.2.5 Performance Monitoring

Systems are constantly monitored for anomalous system behavior, with special care being taken during student testing cycles to provide the highest possible levels of availability and performance. Monitors watch for anomalous activity throughout the entire system, not just at the application or network layers. If suspicious activity shows up, the system triggers alerts to technical support staff for investigation and handling.

In addition to overall, system-wide monitoring for suspicious and anomalous system activity, systems are kept at current patch levels via a suite of tools to scan for vulnerabilities at the network, operating system, platform, and application layers.

8.2.6 Regression Testing

Core Regression Testing confirms that pre-existing functionality has not been adversely affected by changes introduced in a software update. The scope of regression testing is set up to match the changes that are being introduced into the systems by the implementation and testing teams. Regression testing is conducted for every release or patch that is created for our systems.

8.2.7 User Acceptance Testing

One of the testing steps includes the user acceptance test, which is performed by states. Pearson maintains a testing platform so that states can review system functionality prior to a production release.

The following steps are taken when designing the user acceptance testing plan:

1. Create release notes for all new or modified functionality.
2. Provide updated training and user documentation.
3. Review checklist and ask questions.
4. Provide user IDs and passwords to allow users to run tests on code along with associated documentation assisting users on the process and procedures.
5. Meet with users and share results to jointly establish appropriate action plans.

8.3 Reporting

From initial student data upload, through testing, data review, scoring, and reporting, Pearson completes multiple checks and confirms that all data are consistent with customer requirements. Quality Assurance (QA) tasks are part of the project schedule, which is built by working backwards from the reporting dates, to allow for QA work to flow effectively.

Solid requirements form the foundation of quality. USBE and Pearson collaborated to thoroughly and consistently document scoring and reporting requirements, so all involved have a clear understanding of desired results. Project management, product validation, reporting services, and Customer Data Quality (CDQ) teams also participated in requirements reviews to meet reporting requirements and provide accurate mockups.

All Utah Aspire Plus files go through a rigorous validation process as demonstrated by Pearson's comprehensive quality plan. The plan focuses on implementing test cases at the source of each activity, system, and process, thereby detecting defects at the earliest possible point. The impact, therefore, is minimized and resolution can be expedited. The mock data process has become a validation standard within Pearson. It demonstrates production readiness in advance of scoring and reporting actual student data.

CDQ uses industry-standard validation tools focusing on SAS, which allows Pearson the breadth and depth needed for large-scale, high-stakes assessment validation. Pearson's test plans and individual test cases target areas of historical risk (based on the knowledge of Utah Aspire Plus requirements and file layouts) to provide quality results.

8.4 Quality Control of Psychometric Processes

For all psychometric tasks, quality management is central to ensuring on-time and error-free results. Details of Pearson's quality and control procedures for all psychometric tasks conducted, to include test construction, calibration, equating, scaling, field test analysis, data review, item bank creation and management, standard setting, and technical reporting, can be found in the *Utah Aspire Plus 2018–2019 Technical Report* (available at http://utah.pearsonaccessnext.com/resources/additional-services/UT1132740_UTPlusTechReportv4.3_WebTag.pdf).

9 Validity

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

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

The purpose is not to validate the test itself but to validate interpretations of the test scores for specific uses. In that sense, then, test validation is not quantifiable but an ongoing process of evidence gathering beginning at initial conceptualization and continuing throughout the full cycle of an assessment. Every component of an assessment provides evidence in support of its validity, including design, content specifications, item development, and psychometric characteristics.

For the Utah Aspire Plus assessment, operational test development and administration provided the chance to collect initial validity evidence based on test content and internal structure of the tests. Validation is the process of collecting evidence to support inferences from assessment results. As noted, the Utah Aspire Plus assessments are designed to measure the breadth and depth of the Utah Core Standards across all levels of student performance, to provide awareness of individual achievement in relation to stated performance expectations, and to provide evidence of whether students are on track for college and career readiness. The Utah Core Standards define what students should know and be able to do by the end of each respective school year.

9.1 Evidence Based on Test Content

Content validity evidence addresses whether a given assessment adequately samples from the full given domain. Where the assessment is determined to be representative in terms of the standards and in the manner intended, it is said to have high content validity. For the Utah Aspire Plus assessments, they are designed to measure the Utah Core Standards broadly.

For the Utah Aspire Plus tests, design and blueprint specifications were developed in concert between USBE, Utah educators, and Pearson content experts well versed in the Utah Core Standards. As described in Chapter 2 of this report, item and stimulus development targets focused on the measurement of the Utah Core Standards (SAGE) and on providing predictive measures of college and career readiness (ACT Aspire). Blueprints reflect a policy definition of how the makeup of a given assessment is intended to reflect an appropriate sampling of the standards necessary to meet the underlying reporting claims reliably. USBE has published the Utah Aspire Plus blueprints publicly (<http://utah.pearsonaccessnext.com/additional-services/>).

As described in the respective SAGE and ACT Aspire technical manuals noted in Chapter 2, all items were developed to measure the breadth of the Utah Core Standards or related standards. All items were rigorously scrutinized during the various expert content reviews, from initial creation through data review. These expert reviews check for the appropriateness of test items as aligned to the given standard. They also check that items are measuring intended targets of measurement, are clear and concise, and are appropriately aligned to a depth of knowledge (DOK) level, as well as that vocabulary is appropriate for the given level, that the content is accurate and straightforward, and that supporting graphics or stimuli are necessary to answer the question. Further reviews check for cluing within the context of an item set or test form. Every item is also evaluated for fairness by bias and sensitivity committees who review the items for language, or content, that may be inappropriate or offensive to students, parents, or community members, or that contain stereotypical or biased references to sex, ethnicity, or culture. As noted, details of these procedures can be found in the respective technical manuals for SAGE and ACT Aspire referenced in Chapter 2 (see Volumes 2 and 4 of the 2016–2017 SAGE Technical Report and Chapter 2 of the ACT Aspire technical manual).

The process of developing the Utah Aspire Plus test design, development, and test construction is described, in Chapter 2 of this report, to include expert evaluation of the alignment of all content to the Utah Core Standards. As documented, USBE, Utah educators, Pearson, and the developers of the SAGE and ACT Aspire tests expended tremendous effort to ensure the Utah Aspire Plus tests are content-valid and support the intended claims detailed in this report. Additionally, evidence of the content coverage is presented in Appendix A.

Also described in Chapter 2, Utah educators created and recommended performance level descriptors for the Utah Aspire Plus tests, which provide a description of typical end-of-grade performance expectations for each level of achievement in relation to the Utah Core Standards. The PLDs are descriptions of the knowledge and skills demonstrated by students in each performance category. Higher scores translate to a greater level of knowledge and skills demonstrated. There is a link between the PLDs and the knowledge and skills required to meet proficiency according to the standards.

PLDs are used to relate performance on Utah Aspire Plus tests to the Utah Core Standards through the process of standard setting. As described, content experts and stakeholders participated in standard setting in August 2019 for mathematics, reading and English. In August 2022, similar meetings were conducted in support of the new Utah Aspire Plus SEEds science tests. These committees set the cut scores that delineate the four overall levels of achievement on the Utah Aspire Plus tests. Evidence of these activities is presented in the context of student performance on the Utah Aspire Plus tests described in Chapter 4.

9.2 Evidence Based on Cognitive Process

Content comprising the Utah Aspire Plus assessments is specified by standard as well as DOK levels. “Depth of knowledge” (DOK), or cognitive complexity, refers to the cognitive demand associated with interacting with a given item/task. *Levels* of cognitive demand generally focus on the type and level of thinking and reasoning required to answer a given question correctly or earn the most points. For Utah Aspire Plus content, Webb’s definitions of levels of cognitive demand (Webb, 2002) were used to define the DOK levels.

Evidence related to DOK for items developed to measure the Utah Core Standards is provided in volume 4 (Validity) of the SAGE 2016–2017 technical report. In Section 2.3.4 of that report, it is noted that *the alignment of items by DOK also represents a structural model that can be evaluated using confirmatory factor analysis*. Further, they present a confirmatory factor analytic approach to evaluating DOK, where each item is an indicator of a DOK-level first-order factor, and each DOK is in turn an indicator of subject area achievement. Further, in Section 2.4, they describe evidence related to cognitive processes for SAGE content as being “highly similar” to content from the Smarter Balanced assessments and proceed to cite several formal cognitive lab studies that evaluated several facets of items by type as well as across content area.

ACT Aspire content also targets DOK within their development. The content reflects expectations that students need to think, reason, and analyze at high levels of cognitive complexity to be college- and career-ready, and that items and tasks require sampling different levels of cognitive complexity with most targeted at upper levels. ACT’s definition of DOK is like Webb’s, assigned to reflect complexity of the cognitive process required, not the psychometric “difficulty” of the item.

Evidence of cognitive process is presented in Section 17.2.2 of their technical manual: https://actinc.my.salesforce.com/sfc/p/#300000000Wu5/a/4v0000005fHp/SLZ26Xzhfml8ibK_P_Ca5G94_T3HuveFbNgFmfcRaHoY. The pilot of the ACT Aspire CR items used think-aloud tasks, surveys, and interviews to provide evidence of cognitive process.

9.3 Evidence Based on Internal Structure

Internal structure evidence shows the degree to which items and test components conform to the construct on which the proposed test score interpretations are based (AERA, APA, and the NCME, 2014). For example, the Utah Aspire Plus tests report overall scale scores for individual students as well as performance level indicators and ACT prediction ranges for English, reading, math, and science at grades 9 and 10. Internal structure validity evidence identifies the degree to which the item relationships conform to the overall scores and individual subscales. It should be noted that, while information is provided in the appendices examining the Reporting Categories as structural elements of design, the focus of evidence is intended to support the primary claim of each subject test as being unidimensional in nature and supportive of reporting a single overall scale score reflective of the given grade/subject Utah Aspire Plus assessment.

While individual items may each measure multiple elements of the standards and dimensions, they are crafted without dependencies on other items. As such, the tests are designed to be unidimensional and to measure the overall Utah Core Standards primarily. Assuming this holds true, it is appropriate to apply a unidimensional IRT model for calibrating and scaling the Utah Aspire Plus assessments. The IRT model application assumes that the domain being measured by the test is essentially unidimensional. To test this assumption, a principal components analysis is performed.

A general rule of thumb suggests that a set of items may represent as many factors as there are eigenvalues greater than 1 in this analysis because there is one unit of information per item and the eigenvalues sum to the total number of items. However, a set of items may have multiple eigenvalues greater than 1 and still be sufficiently unidimensional for analysis within an IRT framework (Loehlin, 1987; Orlando, 2004). A scree plot is a convenient tool to examine results of factor analyses, as the resulting eigenvalues are plotted in order of magnitude. The scree plots for the principal component analyses for each subject and grade are provided in Appendix K.

In addition to the principal components analyses, confirmatory factor analyses were also conducted to test the model of one factor construct within the Utah Aspire Plus assessments. Indices of model fit are used to determine how well this model fits the data. McDonald and Ho (2002) define absolute fit indices as determining how well an *a priori* model fits the sample data. The chi-square statistic assesses the magnitude of discrepancy between the sample and fitted covariance matrices (Hu and Bentler, 1999). However, this statistic is sensitive to sample size and often rejects the model when large samples are used (Bentler and Bonnet, 1980).

Alternatives to the chi-square, the goodness-of-fit statistic (GFI: Jöresky and Sörbom, 1993), and adjusted goodness-of-fit (AGFI: Tabachnick and Fidell, 2007) are also sensitive to sample size, which has led to researchers reporting them along with other fit indices (Hooper, Coughlan, and Mullen, 2008).

The root mean square error of approximation (RMSEA), a comparative fit index, tells how well the model would fit the population covariance matrix (Byrne, 1998). This fit index favors parsimony since it is sensitive to the number of estimated parameters in the model. There have been a few suggestions of index threshold cut-offs of good fit. The most stringent criterion is 0.06, as suggested in Hu and Bentler (1999). In addition, a confidence interval can be constructed for RMSEA, with a lower limit close to 0 signifying a well-fitting model as well as an upper limit less than 0.08.

The root mean square residual (RMR) and standardized root mean square residual (SRMR) are the square root of the difference between the residuals of the sample covariance matrix and the hypothesized covariance model. The SRMR has a range of 0 to 1, with 0 indicating perfect fit. Byrne (1999) suggests well-fitting models having an SRMR less than 0.05. Hooper, Coughlan, and Mullen (2008) caution that SRMR will tend to be low with a high number of parameters and models with large sample sizes. Hu and Bentler (1999) suggested a two-index presentation when reporting model fit evaluation. One proposed combination is the RMSEA, with confidence interval, and the SRMR. The estimates of these indices are presented in Table 9.1.

Table 9.1. Model Fit Indices for Confirmatory Factor Analyses

Subject	Grade	SRMR	RMSEA	RMSEA 90% Lower CL	RMSEA 90% Upper CL
English	9	0.0314	0.0343	0.0341	0.0346
	10	0.0342	0.0382	0.0380	0.0385
Reading	9	0.0188	0.0224	0.0221	0.0227
	10	0.0251	0.0318	0.0315	0.0321
Mathematics	9	0.0215	0.0241	0.0238	0.0244
	10	0.0175	0.0206	0.0203	0.0209
Science	9	0.0225	0.0292	0.0287	0.0297
	10	0.0234	0.0292	0.0287	0.0298

Model-data fit based on the IRT model calibrations are also indicators of unidimensionality. To the extent that indicators of fit suggest data do not appropriately fit the model as applied may be the result of multidimensionality. Discussion of model fit is presented in Section 7.4 in terms of Q_1 indices. These statistics support the overall fit of Utah Aspire Plus items to the respective IRT models.

In addition to evidence of essential unidimensionality described here, it should be acknowledged that tests are not designed to be *strictly* unidimensional. It is common to observe what might be considered transient factors common to one or more test items in the face of a dominant overall factor. As discussed in Chapter 2, the Utah Aspire Plus blueprints were designed to reflect the Utah Core Standards partly around Reporting Categories. Correlations among the Utah Aspire Plus overall test scores and Reporting Categories offer additional evidence of the internal structure of the Utah Aspire Plus tests. These correlations quantify the strength of the relationships across structural elements of the assessments. Results of these analyses are presented in Appendix L.

9.3.1 Reliability

Additionally, the reliability analyses presented in Chapter 6 of this technical report provide information about the internal consistency of the Utah Aspire Plus tests. Internal consistency is typically measured by correlations among the items on a test and provides an indication of how much the items measure the same general construct.

9.4 Evidence Based on Different Student Populations

In addition, internal structure evidence should show that individual items are functioning similarly for different demographic subgroups within the population being measured. The Utah Aspire Plus tests are developed to assess the Utah Core Standards and are administered to all students irrespective of any particular demographic characteristic (as described in Chapter 2). Great care has been taken to ensure the items on the Utah Aspire Plus tests are fair and representative of the content domains expressed in the standards. Special attention is given to finding evidence that construct-irrelevant content has not been inadvertently included in the test, as such content could result in an unfair advantage for one group versus another.

This begins with item writers trained on how to avoid economic, regional, cultural, and ethnic biases when writing items. After items have been written, they are reviewed by a bias and sensitivity committee, which evaluates each item to identify language or content that might be inappropriate or offensive to students, parents, or other community members or that contain stereotypical or biased references to sex, ethnic, or cultural groups. The bias and sensitivity committee accepts, edits, or rejects each item for use prior to the items' administration.

Differential item functioning (DIF) analyses are conducted for the purpose of identifying items that are differentially difficult for different subpopulations of individuals. Section 5.1.3 details the methodology used to evaluate DIF for the Utah Aspire Plus items. Though DIF analyses flag items as being differentially difficult for one group as compared to another, it does not solely provide sufficient evidence for removing the item from use. Flagged items are re-examined post administration for any potentially overlooked biases attributable to the content of those items.

9.5 Summary

As noted, the process of validation involves accumulating relevant evidence to provide a sound scientific basis for stated score interpretations. Collection of validity evidence is an ongoing process and validity of interpretations are strengthened as positive evidence accrues. While this technical report reflects the continued administration of the Utah Aspire Plus assessments, sufficient evidence exists to support the primary claims detailed herein, including that test scores indicate the degree to which students achieved end-of-year expectations on the Utah Core Standards across subject tests in grades 9 and 10. Further, performance on the Utah Aspire Plus assessments could reasonably be linked to predictions of performance on the ACT college and career readiness benchmarks. These are supported by evidence of the content development processes that underpin the creation of assessments aligned to the Utah Core Standards and evidence that the internal structure aligns with the stated claims and is sound.

10 References

- ACT Aspire. (2017). *Summative Technical Manual*. Version 3. Iowa City, IA: ACT.
- American Educational Research Association, American Psychological Association, & the National Council on Measurement in Education. Joint Technical Committee. (2014). *Standards for Educational and Psychological Testing*. Washington, DC: American Educational Research Association.
- Bentler, P. M., & Bonnet, D. C. (1980). Significance tests and goodness of fit in the analysis of covariance structures. *Psychological Bulletin*, 88(3), 588–606.
- Birnbaum, A. (1968). Some latent trait models and their use in inferring an examinee's ability. In F. M. Lord and M. R. Novick (Eds.), *Statistical Theories of Mental Test Scores* (pp. 397–472). Reading, MA: Addison-Wesley.
- Byrne, B. M. (1998). *Structural Equation Modeling with LISREL, PRELIS and SIMPLIS: Basic Concepts, Applications and Programming*. Mahwah, NJ: Lawrence Erlbaum Associates.
- Chien, M. and Shin, D. (2012). *IRT Score Estimation Program, V1.3* [computer program]. Iowa City, IA: Pearson.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20, 37–47.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16, 297–334.
- Cronbach, L. J. (1988). Five perspectives on validity argument. In H. Wainer & H. Braun (Eds.), *Test Validity* (pp. 3–17). Hillsdale, NJ: Lawrence Erlbaum.
- Dorans, N. J., & Holland, P. W. (1993). DIF detection and description: Mantel-Haenszel and standardization. In P. W. Holland and H. Wainer (Eds.), *Differential Item Functioning* (pp. 35–66). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Holland, P. W., & Thayer, D. (1988). Differential item performance and the Mantel-Haenszel procedure. In H. Wainer and H. I. Braun (Eds.), *Test Validity* (pp. 129–145). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Hooper, D., Coughlan, J., & Mullen, M. R. (2008). Structural equation modelling: Guidelines for determining model fit. *Electronic Journal of Business Research Methods*, 6, 53–60.
- Hu, L. T., & Bentler, P. N. (1999). Cutoff criteria for fit indexes in covariance structure analysis: Conventional criteria versus new alternatives. *Structural Equation Modeling*, 6(1), 1–55.
- Jöresky, K., & Sörbom, D. (1993). *LISREL 8: Structural Equation Modeling with the SIMPLIS Command Language*. Chicago, IL: Scientific Software International Inc.

- Kane, M. T. (2006). Validation. In R. L. Brennan (Ed.), *Educational Measurement* (4th ed., pp. 17–64). Westport, CT: American Council on Education/Praeger.
- Kim, S. and Kolen, M. (2004). STUIRT [computer program]. Iowa City, IA: The University of Iowa.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32, 179–197.
- Loehlin, J. C. (1987). *Latent Variable Models*. Hillsdale, NJ: Lawrence Erlbaum Associates, Publishers.
- McDonald, R. P., & Ho, M.–H. R. (2002). Principles and practice in reporting statistical equation analyses. *Psychological Methods*, 7(1), 64–82.
- Muraki, E. (1992). A generalized partial credit model: Application of an EM algorithm. *Applied Psychological Measurement* 16, 159–176.
- 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>.
- Next Generation Science Standards (NGSS Lead States. 2013. *Next Generation Science Standards: For States, By States*. Washington, DC: The National Academies Press) <http://www.nextgenscience.org>
- Orlando, M. (2004, June). Critical issues to address when applying item response theory (IRT) models. Paper presented at the Drug Information Association, Bethesda, M. D.
- Scientific Software International, Inc. (2017). IRTPRO. Lincolnwood, IL: www.ssicentral.com.
- Tabachnick, B. G., & Fidell, L. S. (2007). *Using Multivariate Statistics* (5th ed.). New York: Allyn and Bacon.
- Wells, C. S., Hambleton, R. K., Kirkpatrick, R., & Meng, Y. (2014). An examination of two procedures for identifying consequential item parameter drift. *Applied Measurement in Education*, 27, 214–231.
- Yen, W. M. (1984). Effects of local item dependence on the fit and equating performance of the three-parameter logistic model. *Applied Psychological Measurement*, 8, 125–145.
- Yen, W. M. (1993). Scaling performance assessments: Strategies for managing local item dependence. *Journal of Educational Measurement*, 30, 187–213.
- Zwick, R., Donoghue, J. R., & Grima, A. (1993). Assessment of differential item functioning for performance tasks. *Journal of Educational Measurement*, 30, 233–251.

Appendix A: Test-Level Reporting Categories and Standards by Item Type and DOK

Table A.1. Test-Level Reporting Categories and Standards for English Grade 9

Grade	Reporting Category: Standard	Multiple Choice			Technology Enhanced		
		DOK 1	DOK 2	DOK 3	DOK 1	DOK 2	DOK 3
9	Conventions of Standard English: L.9-10.1	3	2	0	0	0	0
	Conventions of Standard English: L.9-10.1a	1	0	0	2	1	0
	Conventions of Standard English: L.9-10.1b	0	0	0	1	2	0
	Conventions of Standard English: L.9-10.2	2	0	0	1	0	0
	Conventions of Standard English: L.9-10.2a	1	0	0	2	0	0
	Conventions of Standard English: L.9-10.2b	2	0	0	1	0	0
	Conventions of Standard English: L.9-10.2c	1	0	0	2	0	0
	Conventions of Standard English: L.9-10.5b	1	1	0	1	1	0
	Conventions of Standard English: L.9-10.6	1	0	0	1	0	0
	Knowledge of Language: L.9-10.3	1	1	3	0	0	0
	Knowledge of Language: L.9-10.4b	1	0	0	0	0	0
	Production of Writing: W.9-10.4	5	0	0	1	0	0
	Production of Writing: W.9-10.5	3	0	0	1	0	0
	Total	46					

Table A.2. Test-Level Reporting Categories and Standards for English Grade 10

Grade	Reporting Category: Standard	Multiple Choice			Technology Enhanced		
		DOK 1	DOK 2	DOK 3	DOK 1	DOK 2	DOK 3
10	Conventions of Standard English: L.9-10.1	3	1	0	0	0	0
	Conventions of Standard English: L.9-10.1a	0	0	0	2	0	0
	Conventions of Standard English: L.9-10.1b	0	0	0	2	0	0
	Conventions of Standard English: L.9-10.2	2	0	0	0	0	0
	Conventions of Standard English: L.9-10.2a	1	0	0	3	0	0
	Conventions of Standard English: L.9-10.2b	1	0	0	1	0	0
	Conventions of Standard English: L.9-10.2c	2	0	0	1	0	0
	Conventions of Standard English: L.9-10.5a	0	0	0	1	0	0
	Conventions of Standard English: L.9-10.5b	0	0	0	1	2	0
	Conventions of Standard English: L.9-10.6	0	0	0	0	0	0
	Knowledge of Language: L.9-10.3	5	0	0	1	0	0
	Knowledge of Language: L.9-10.3a	0	0	0	1	0	0
	Knowledge of Language: L.9-10.4a	1					
	Knowledge of Language: L.9-10.4b	0					
	Knowledge of Language: L.9-10.4d	1					
	Production of Writing: W.9-10.4	2					
	Production of Writing: W.9-10.5	1					
	Total	43					

Table A.3. Test-Level Reporting Categories and Standards for Reading Grade 9

Reporting Category: Standard	Multiple Choice			Technology Enhanced			Evidence-Based Selected Response		
	DOK 1	DOK 2	DOK 3	DOK 1	DOK 2	DOK 3	DOK 1	DOK 2	DOK 3
Craft and Structure: RI.9-10.4	0	0	0	2	0	0	0	0	0
Craft and Structure: RI.9-10.5	1	1	0	0	0	0	0	0	0
Craft and Structure: RI.9-10.6	1	0	0	0	0	0	0	0	0
Craft and Structure: RL.9-10.4	2	1	0	0	0	0	1	0	0
Craft and Structure: RL.9-10.5	1	1	0	0	0	0	0	0	0
Craft and Structure: RL.9-10.6	1	0	0	0	0	0	0	0	0
Integration of Knowledge and Ideas: RI.9-10.7	1	0	0	0	0	0	0	0	0
Integration of Knowledge and Ideas: RI.9-10.9	0	0	0	0	0	0	2	0	0
Integration of Knowledge and Ideas: RL.9-10.7	0	0	0	0	0	0	1	0	0
Integration of Knowledge and Ideas: RL.9-10.9	0	0	0	2	0	0	0	0	0
Key Ideas: RI.9-10.1	1	2	0	0	0	0	0	0	0
Key Ideas: RI.9-10.2	0	0	0	1	0	0	0	0	0
Key Ideas: RI.9-10.5	1	0	0	0	0	0	0	0	0
Key Ideas: RL.9-10.1	2	1	0	0	0	0	0	0	0
Key Ideas: RL.9-10.2	1	0	0	0	0	0	0	0	0
Key Ideas: RL.9-10.3	0	0	0	1	0	0	0	0	0
Total	35								

Table A.4. Test-Level Reporting Categories and Standards for Reading Grade 10

Reporting Category: Standard	Multiple Choice			Technology Enhanced			Evidence-Based Selected Response		
	DOK 1	DOK 2	DOK 3	DOK 1	DOK 2	DOK 3	DOK 1	DOK 2	DOK 3
Craft and Structure: L.9-10.4a	1	0	0	0	0	0	0	0	0
Craft and Structure: L.9.10.6	1	0	0	0	0	0	0	0	0
Craft and Structure: RI.9-10.4	1	0	0	0	0	0	0	0	0
Craft and Structure: RI.9-10.5	1	1	0	0	0	0	1	0	0
Craft and Structure: RI.9-10.6	1	1	0	0	0	0	0	0	0
Craft and Structure: RL.9-10.4	1	1	0	0	0	0	1	0	0
Craft and Structure: RL.9-10.5	1	0	0	0	0	0	0	0	0
Craft and Structure: RL.9-10.6	1	0	0	0	0	0	0	0	0
Integration of Knowledge and Ideas: CCRA.R.5	1	0	0	0	0	0	0	0	0
Integration of Knowledge and Ideas: RI.9-10.7	0	0	0	0	0	0	1	0	0
Integration of Knowledge and Ideas: RL.9-10.7	1	1	0	0	0	0	0	0	0
Key Ideas: RI.9-10.1	0	0	0	1	0	0	0	0	0
Key Ideas: RI.9-10.2	1	2	0	0	0	0	1	0	0
Key Ideas: RL.9-10.1	1	1	0	0	0	0	1	0	0
Key Ideas: RL.9-10.2	2	0	0	0	0	0	0	0	0
Key Ideas: RL.9-10.3	0	0	0	2	0	0	0	0	0
Total	35								

Table A.5. Test-Level Reporting Categories and Standards for Mathematics Grade 9

Reporting Category: Standard	Multiple Choice			Technology Enhanced		
	DOK 1	DOK 2	DOK 3	DOK 1	DOK 2	DOK 3
Algebra: MI.A.CED.1	1	0	0	0	0	0
Algebra: MI.A.CED.2	1	0	0	0	0	0
Algebra: MI.A.CED.4	1	0	0	0	0	0
Algebra: MI.A.REI.10	1	0	0	0	0	0
Algebra: MI.A.REI.11	1	0	0	0	0	0
Algebra: MI.A.REI.12	1	0	0	0	0	0
Algebra: MI.A.REI.3	1	0	0	0	0	0
Algebra: MI.A.REI.6	0	0	0	1	0	0
Algebra: MI.A.SSE.1b	1	0	0	0	0	0
Functions: MI.F.BF.1a	1	0	0	0	0	0
Functions: MI.F.BF.2	1	0	0	0	0	0
Functions: MI.F.BF.3	1	0	0	0	0	0
Functions: MI.F.IF.2	1	0	0	0	0	0
Functions: MI.F.IF.4	0	0	0	1	0	0
Functions: MI.F.IF.7a	0	0	0	1	0	0
Functions: MI.F.IF.9	1	0	0	0	0	0
Functions: MI.F.LE.1b	0	0	0	1	0	0
Functions: MI.F.LE.3	0	0	0	1	0	0
Functions: MI.F.LE.5	1	0	0	0	0	0
Geometry: MI.G.CO.1	1	0	0	0	0	0
Geometry: MI.G.CO.3	1	0	0	0	0	0
Geometry: MI.G.CO.4	1	0	0	0	0	0
Geometry: MI.G.CO.5	0	0	0	1	0	0
Geometry: MI.G.CO.6	1	0	0	0	0	0
Geometry: MI.G.CO.7	1	1	0	0	0	0
Geometry: MI.G.CO.8	0	0	0	1	0	0
Geometry: MI.G.GPE.4	0	0	0	1	0	0
Geometry: MI.G.GPE.7	1	0	0	0	0	0
Statistics and Probability: MI.S.ID.1	1	0	0	0	0	0
Statistics and Probability: MI.S.ID.2	2	0	0	0	0	0
Statistics and Probability: MI.S.ID.3	1	1	0	0	0	0
Statistics and Probability: MI.S.ID.6	1	0	0	0	0	0
Statistics and Probability: MI.S.ID.6c	1	0	0	0	0	0
Statistics and Probability: MI.S.ID.7	0	0	0	1	0	0
Statistics and Probability: MI.S.ID.8	1	0	0	0	0	0
Total	40					

Table A.6. Test-Level Reporting Categories and Standards for Mathematics Grade 10

Reporting Category: Standard	Multiple Choice			Technology Enhanced		
	DOK 1	DOK 2	DOK 3	DOK 1	DOK 2	DOK 3
Algebra: MII.A.APR.1	1	0	0	0	0	0
Algebra: MII.A.CED.1	1	0	0	0	0	0
Algebra: MII.A.CED.4	1	0	0	0	0	0
Algebra: MII.A.REI.4b	0	0	0	1	0	0
Algebra: MII.A.REI.7	1	0	0	0	0	0
Algebra: MII.A.SSE.2	1	0	0	0	0	0
Algebra: MII.A.SSE.3a	1	0	0	0	0	0
Algebra: MII.A.SSE.3b	1	0	0	0	0	0
Functions: MII.F.BF.1a	0	0	0	1	0	0
Functions: MII.F.BF.3	1	0	0	0	0	0
Functions: MII.F.IF.4	1	0	0	0	0	0
Functions: MII.F.IF.7a	1	0	0	0	0	0
Functions: MII.F.IF.7b	2	0	0	0	0	0
Functions: MII.F.IF.8b	1	0	0	0	0	0
Functions: MII.F.IF.9	1	0	0	0	0	0
Functions: MII.F.LE.3	0	0	0	1	0	0
Functions: MII.F.TF.8	1	0	0	0	0	0
Geometry: MII.G.C.2	1	0	0	0	0	0
Geometry: MII.G.C.4	0	0	0	1	0	0
Geometry: MII.G.CO.10	1	0	0	0	0	0
Geometry: MII.G.CO.9	1	0	0	0	0	0
Geometry: MII.G.GMD.3	1	0	0	0	0	0
Geometry: MII.G.GPE.4	1	0	0	0	0	0
Geometry: MII.G.GPE.6	0	0	0	1	0	0
Geometry: MII.G.SRT.1a	0	0	0	1	0	0
Geometry: MII.G.SRT.2	1	1	0	0	0	0
Geometry: MII.G.SRT.3	0	0	0	1	0	0
Number and Quantity: MII.N.CN.2	1	0	0	0	0	0
Number and Quantity: MII.N.RN.2	3	0	0	0	0	0
Statistics and Probability: MII.S.CP.5	1	0	0	0	0	0
Statistics and Probability: MII.S.CP.6	1	0	0	0	0	0
Statistics and Probability: MII.S.ID.5	2	0	0	0	0	0
Total	40					

Appendix B: Student Testing Time

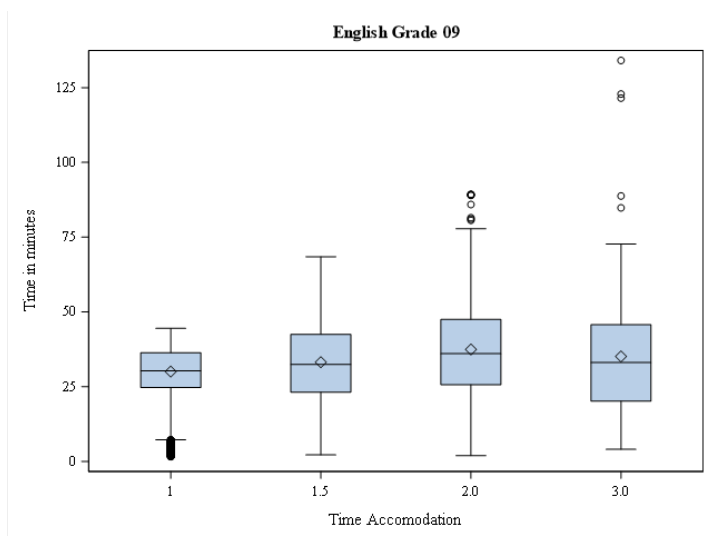


Figure B.1. English Grade 9 Student Testing Time

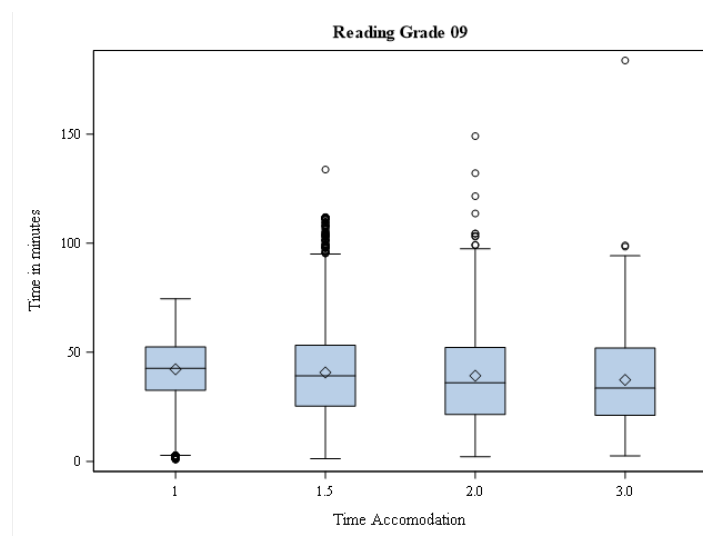


Figure B.3. Reading Grade 9 Student Testing Time

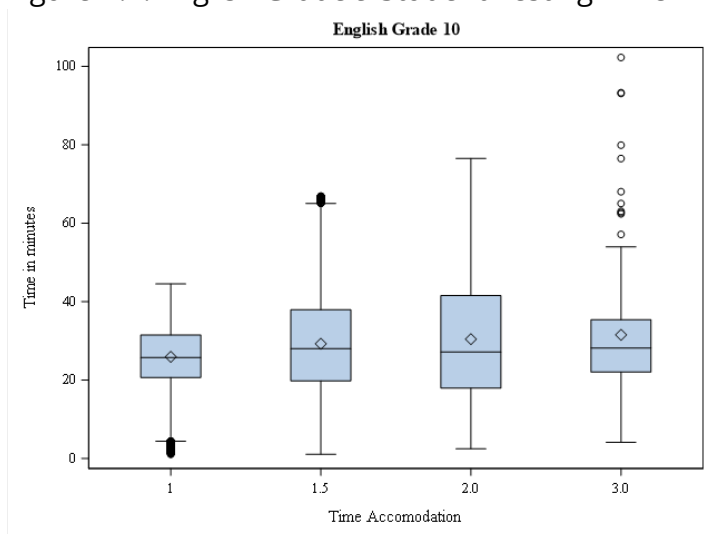


Figure B.2. English Grade 10 Student Testing Time

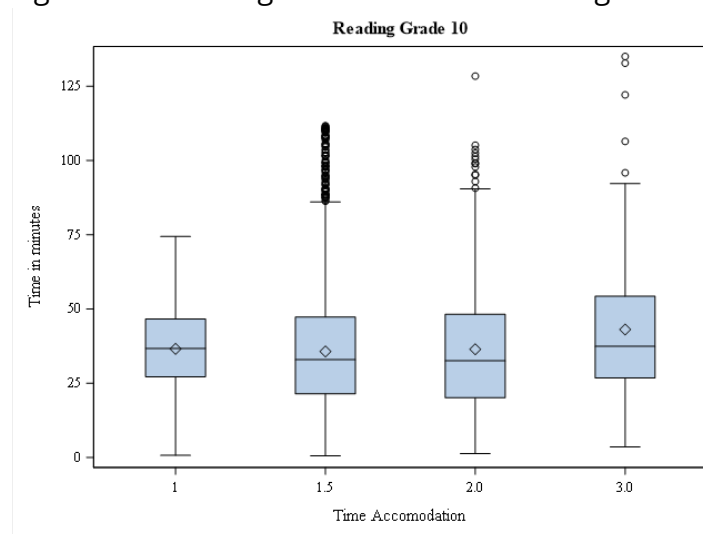


Figure B.4. Reading Grade 10 Student Testing Time

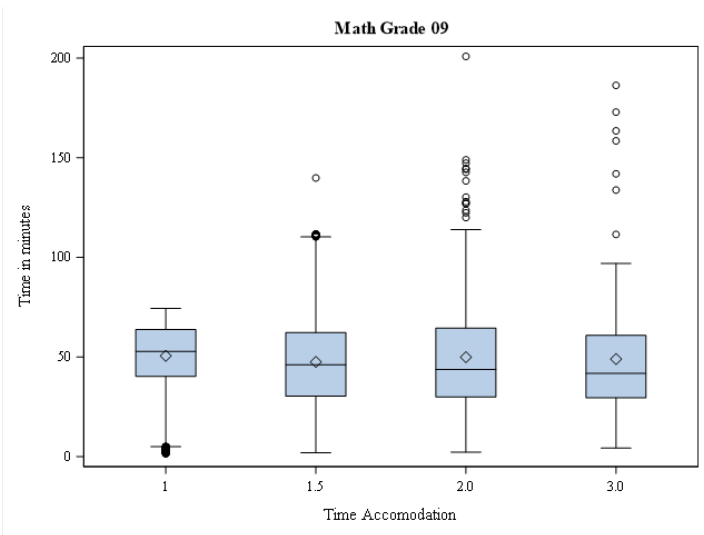


Figure B.5. Mathematics Grade 9 Student Testing Time

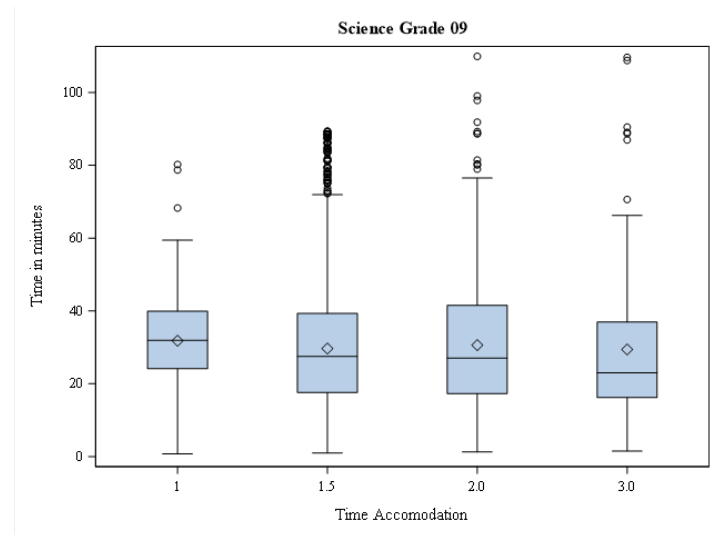


Figure B.7. Science Grade 9 Student Testing Time

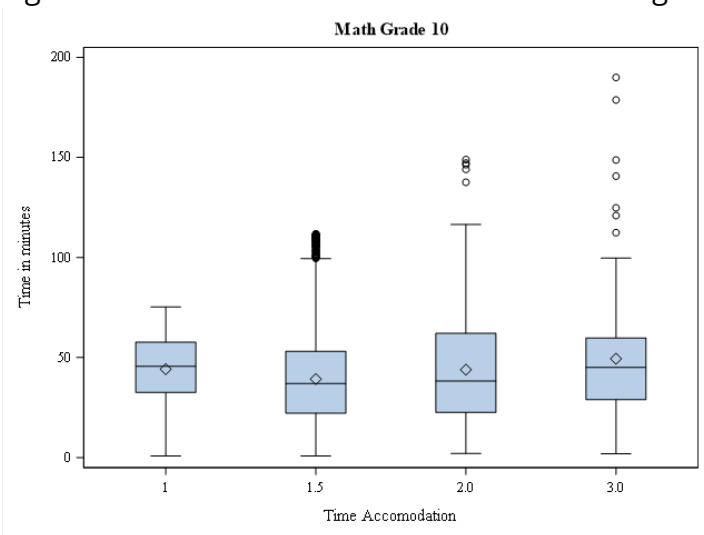


Figure B.6. Mathematics Grade 10 Student Testing Time

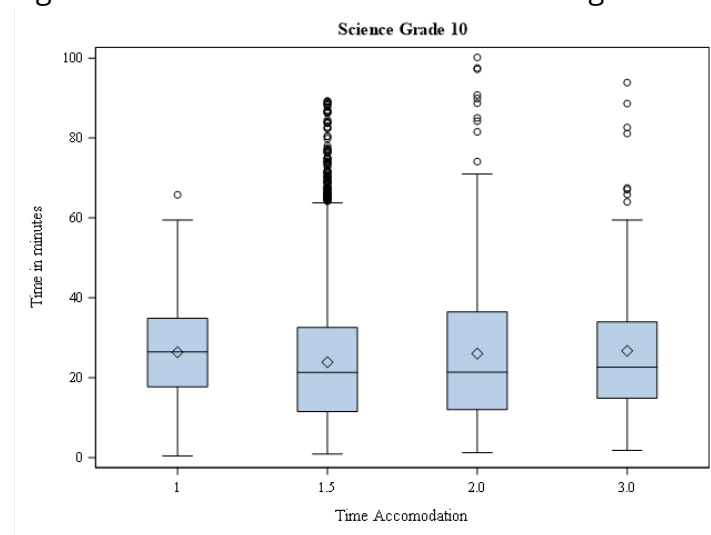


Figure B.8. Science Grade 10 Student Testing Time

Appendix C: Item Statistics Summaries

Table C.1. Item Mean for One-Point Items

Subject	Grade	N	$p < 0.30$	$0.30 \leq p < 0.55$	$0.55 \leq p < 0.75$	$0.75 \leq p < 0.95$	$p \geq 0.95$	Mean p
English	9	42	4	12	20	6	0	0.57
	10	36	0	11	15	10	0	0.64
Reading	9	27	1	6	16	4	0	0.60
	10	29	1	12	14	2	0	0.57
Mathematics	9	40	8	22	7	3	0	0.44
	10	40	4	27	9	0	0	0.46
Science	9	18	4	7	7	0	0	0.50
	10	18	1	14	3	0	0	0.45

Table C.2. Item Mean for Two-Point Items

Subject	Grade	N	Mean	Min	Max
English	9	4	1.20	0.79	1.51
	10	7	1.25	0.45	1.81
Reading	9	8	0.86	0.44	1.32
	10	6	1.05	0.78	1.37
Science	9	5	0.92	0.70	1.20
	10	5	0.66	0.50	0.93

Note: There were no 2-point mathematics items in Spring 2023.

Table C.3. Item Total Correlation for One-Point Items

Subject	Grade	N	$r < 0.20$	$0.20 \leq r < 0.40$	$0.40 \leq r < 0.60$	$0.60 \leq r < 0.80$	$r \geq 0.80$	Median ITC
English	9	42	3	21	1	0	0	0.40
	10	36	0	21	1	0	0	0.44
Reading	9	27	0	20	1	0	0	0.44
	10	29	0	22	1	0	0	0.44
Mathematics	9	40	1	21	0	0	0	0.41
	10	40	2	21	0	0	0	0.40
Science	9	18	2	7	0	0	0	0.37
	10	18	1	8	0	0	0	0.37

Note: ITC=Item Total Correlation

Table C.4. Item Total Correlation for Two-Point Items

Subject	Grade	N	Median r	Min r	Max r
English	9	4	0.49	0.46	0.66
	10	7	0.60	0.42	0.76
Reading	9	8	0.47	0.35	0.66
	10	6	0.58	0.44	0.66
Science	9	5	0.42	0.29	0.61
	10	5	0.41	0.28	0.53

Note: There were no 2-point mathematics items in Spring 2023.

Table C.5. Differential Item Functioning

Subject	Grade	Subgroups	DIF Categories				
			Negligible DIF	Moderate DIF		Substantial DIF	
				Focal	Reference	Focal	Reference
English	9	Male-Female	45	1	0	0	0
		White-Black	46	0	0	0	0
		White-Hispanic	46	0	0	0	0
	10	Male-Female	42	1	0	0	0
		White-Black	41	0	2	0	0
		White-Hispanic	43	0	0	0	0
Reading	9	Male-Female	33	0	0	0	0
		White-Black	35	0	0	0	0
		White-Hispanic	35	0	0	0	0
	10	Male-Female	34	0	0	0	0
		White-Black	35	0	0	0	0
		White-Hispanic	34	0	0	0	1
Mathematics	9	Male-Female	37	0	0	0	0
		White-Black	38	0	2	0	0
		White-Hispanic	40	0	0	0	0
	10	Male-Female	39	0	0	0	0
		White-Black	40	0	0	0	0
		White-Hispanic	40	0	0	0	0
Science	9	Male-Female	22	0	0	0	0
		White-Black	22	0	1	0	0
		White-Hispanic	23	0	0	0	0
	10	Male-Female	23	0	0	0	0
		White-Black	22	0	1	0	0
		White-Hispanic	23	0	0	0	0

Note: "Focal" indicates DIF in favor of Female, Black, or Hispanic students; "Reference" indicates DIF in favor of Male or White students.

Appendix D: Reliability and Standard Error by Subgroup

Table D.1. English Grade 9 Test Reliability

	Test Group	N	Alpha	SEM	Conventions of			
					Standard English	Knowledge of Language	Production of Writing	
All	Students							
	Tested	46,485	0.90	9.02	0.82	0.59	0.76	
Sex	Female	22,195	0.89	8.88	0.82	0.58	0.74	
	Male	24,264	0.90	9.10	0.82	0.60	0.77	
Ethnicity	Hispanic or Latino Ethnicity	9,032	0.88	9.37	0.78	0.53	0.74	
	Asian	815	0.90	9.10	0.84	0.58	0.76	
	Native Hawaiian or Other Pacific Islander	698	0.86	9.28	0.76	0.48	0.72	
	Black or African American	627	0.87	9.66	0.77	0.51	0.74	
	American Indian or Alaska Native	463	0.86	9.36	0.74	0.50	0.72	
	White	33,332	0.89	8.85	0.82	0.58	0.75	
	Other	1,518	0.89	8.85	0.81	0.57	0.75	
	Limited English Proficiency	No	42,350	0.89	8.88	0.82	0.58	0.75
		Yes	4,135	0.78	10.21	0.62	0.34	0.63
	Economic Disadvantage	No	33,751	0.89	8.88	0.82	0.58	0.75
	Yes	12,734	0.89	9.30	0.80	0.56	0.75	
Special Education	No	41,881	0.89	8.91	0.81	0.58	0.75	
	Yes	4,604	0.83	9.83	0.71	0.44	0.67	

Table D.2. English Grade 10 Test Reliability

	Test Group	N	Alpha	SEM	Conventions of Standard English	Knowledge of Language	Production of Writing	
All	Students Tested	43,735	0.91	8.56	0.84	0.72	0.70	
Sex	Female	20,651	0.90	8.54	0.83	0.69	0.69	
	Male	23,042	0.91	8.57	0.85	0.74	0.71	
Ethnicity	Hispanic or Latino Ethnicity	8,434	0.89	8.21	0.81	0.70	0.60	
	Asian	767	0.92	8.60	0.86	0.72	0.71	
	Native Hawaiian or Other Pacific Islander	639	0.87	7.99	0.78	0.64	0.54	
	Black or African American	556	0.90	8.48	0.83	0.73	0.65	
	American Indian or Alaska Native	443	0.88	7.85	0.80	0.66	0.60	
	White	31,554	0.91	8.59	0.84	0.71	0.70	
	Other	1,342	0.90	8.63	0.83	0.70	0.68	
	Limited English Proficiency	No	40,460	0.91	8.57	0.84	0.70	0.70
	Yes	3,275	0.82	8.21	0.72	0.64	0.32	
Economic Disadvantage	No	32,757	0.91	8.58	0.84	0.71	0.70	
Yes	10,978	0.90	8.37	0.83	0.71	0.65		
Special Education	No	39,750	0.90	8.59	0.83	0.70	0.69	
Yes	3,985	0.86	8.28	0.77	0.69	0.50		

Table D.3. Reading Grade 9 Test Reliability

	Test Group	N	Alpha	SEM	Craft and Structure	Integration of Knowledge and Ideas	Key Ideas
All	Students Tested	46,653	0.89	9.24	0.76	0.56	0.80
Sex	Female	22,241	0.89	9.16	0.74	0.54	0.79
	Male	24,387	0.90	9.28	0.77	0.57	0.81
Ethnicity	Hispanic or Latino						
	Ethnicity	9,172	0.87	9.60	0.71	0.51	0.77
	Asian	815	0.90	9.20	0.77	0.57	0.82
	Native Hawaiian or Other Pacific Islander	706	0.85	9.81	0.64	0.48	0.75
	Black or African American	650	0.87	9.61	0.69	0.52	0.77
	American Indian or Alaska Native	467	0.85	9.39	0.68	0.41	0.74
	White	33,319	0.89	9.12	0.74	0.54	0.80
	Other	1,524	0.89	9.31	0.75	0.54	0.79
	Limited English Proficiency	No	42,423	0.89	9.15	0.74	0.54
	Yes	4,230	0.77	10.62	0.53	0.35	0.63
Economic Disadvantage	No	33,763	0.89	9.17	0.75	0.55	0.80
	Yes	12,890	0.88	9.43	0.73	0.52	0.78
Special Education	No	42,010	0.89	9.17	0.74	0.54	0.79
	Yes	4,643	0.82	10.12	0.63	0.42	0.69

Table D.4. Reading Grade 10 Test Reliability

	Test Group	N	Alpha	SEM	Craft and Structure	Integration of Knowledge and Ideas	Key Ideas
All	Students Tested	43,507	0.91	8.22	0.79	0.48	0.84
Sex	Female	20,517	0.90	8.11	0.78	0.46	0.83
	Male	22,952	0.91	8.30	0.80	0.48	0.84
Ethnicity	Hispanic or Latino						
	Ethnicity	8,422	0.88	8.43	0.73	0.44	0.80
	Asian	761	0.91	8.52	0.81	0.49	0.84
	Native Hawaiian or Other Pacific						
	Islander	630	0.87	8.44	0.66	0.46	0.80
	Black or African American	560	0.90	8.31	0.77	0.53	0.83
	American Indian or Alaska Native	439	0.86	8.09	0.68	0.38	0.78
	White	31,352	0.91	8.17	0.79	0.46	0.84
	Other	1,343	0.91	8.28	0.79	0.49	0.84
Limited English Proficiency	No	40,213	0.91	8.18	0.79	0.46	0.84
	Yes	3,294	0.76	9.65	0.45	0.29	0.67
Economic Disadvantage	No	32,583	0.91	8.19	0.79	0.47	0.84
	Yes	10,924	0.90	8.33	0.76	0.45	0.82
Special Education	No	39,521	0.91	8.18	0.79	0.46	0.83
	Yes	3,986	0.85	8.93	0.66	0.37	0.76

Table D.5. Mathematics Grade 9 Test Reliability

	Test Group	N	Alpha	SEM Algebra	Functions	Geometry	Number and Quantity	
All	Students Tested	45,135	0.90	9.26	0.72	0.73	0.70	0.62
Sex	Female	21,384	0.89	9.08	0.69	0.70	0.67	0.58
	Male	23,727	0.91	9.37	0.74	0.75	0.73	0.65
Ethnicity	Hispanic or Latino							
	Ethnicity	8,732	0.85	10.77	0.60	0.60	0.61	0.51
	Asian	793	0.92	8.93	0.77	0.77	0.75	0.62
	Native Hawaiian or Other Pacific Islander	676	0.81	10.90	0.57	0.52	0.58	0.44
	Black or African American	628	0.82	11.69	0.54	0.59	0.56	0.49
	American Indian or Alaska Native	443	0.83	11.48	0.57	0.52	0.63	0.47
	White	32,391	0.90	8.89	0.71	0.73	0.69	0.61
	Other	1,472	0.91	9.15	0.74	0.76	0.71	0.64
Limited English Proficiency	No	41,072	0.90	9.04	0.71	0.73	0.69	0.61
	Yes	4,063	0.70	13.39	0.38	0.33	0.42	0.31
Economic Disadvantage	No	32,821	0.90	8.90	0.72	0.73	0.69	0.61
	Yes	12,314	0.87	10.36	0.64	0.66	0.65	0.56
Special Education	No	40,604	0.90	8.97	0.71	0.73	0.69	0.61
	Yes	4,531	0.78	12.42	0.50	0.49	0.54	0.38

Table D.6. Mathematics Grade 10 Test Reliability

	Test Group	N	Alpha	SEM	Algebra	Functions	Geometry	Number and Quantity	Statistics and Probability	
All	Students Tested	42,881	0.89	12.00	0.64	0.65	0.78	0.49	0.32	
	Female	20,210	0.87	11.95	0.60	0.60	0.75	0.47	0.29	
Sex	Male	22,639	0.90	12.01	0.67	0.69	0.79	0.52	0.35	
	Hispanic or Latino Ethnicity	8,296	0.82	15.40	0.52	0.47	0.67	0.39	0.17	
Ethnicity	Asian	762	0.92	11.01	0.70	0.74	0.80	0.58	0.42	
	Native Hawaiian or Other Pacific Islander	632	0.77	16.72	0.43	0.41	0.66	0.42	0.05	
	Black or African American	543	0.82	15.92	0.52	0.51	0.70	0.45	0.21	
	American Indian or Alaska Native	427	0.77	16.84	0.45	0.30	0.66	0.35	-0.04	
	White	30,905	0.89	11.27	0.63	0.66	0.77	0.48	0.34	
	Other	1,316	0.89	11.85	0.64	0.66	0.77	0.50	0.31	
	Limited English Proficiency	No	39,634	0.89	11.65	0.63	0.66	0.77	0.49	0.33
	Yes	3,247	0.63	20.80	0.34	0.18	0.47	0.26	0.04	
	Economic Disadvantage	No	32,177	0.89	11.44	0.64	0.66	0.77	0.50	0.34
Yes	10,704	0.85	14.23	0.56	0.53	0.72	0.42	0.23		
Special Education	No	38,938	0.89	11.53	0.63	0.65	0.77	0.48	0.33	
	Yes	3,943	0.73	18.77	0.41	0.30	0.56	0.29	0.08	

Table D.7. Science Grade 9 Test Reliability

	Test Group	N	Alpha	SEM	Construct Explanations	Developing Models	Gathering and Investigating	Using Mathematical Thinking
All	Students Tested	46,564	0.82	13.96	0.49	0.54	0.58	0.55
	Female	22,189	0.79	13.96	0.43	0.51	0.55	0.48
Sex	Male	24,350	0.84	13.95	0.53	0.56	0.61	0.59
	Hispanic or Latino							
Ethnicity	Ethnicity	9,132	0.74	15.77	0.38	0.42	0.50	0.39
	Asian	816	0.83	13.68	0.52	0.59	0.61	0.59
Ethnicity	Native Hawaiian or Other Pacific Islander	694	0.65	15.75	0.27	0.29	0.44	0.28
	Black or African American	652	0.65	16.96	0.30	0.24	0.45	0.28
Ethnicity	American Indian or Alaska Native	470	0.71	16.92	0.46	0.37	0.45	0.30
	White	33,279	0.81	13.49	0.48	0.53	0.58	0.55
Limited English Proficiency	Other	1,521	0.82	14.02	0.46	0.55	0.59	0.57
	No	42,370	0.81	13.72	0.48	0.53	0.58	0.55
Economic Disadvantage	Yes	4,194	0.52	18.59	0.24	0.17	0.31	0.18
	No	33,692	0.81	13.65	0.48	0.54	0.58	0.55
Special Education	Yes	12,872	0.78	14.96	0.44	0.46	0.54	0.46
	No	41,934	0.81	13.73	0.48	0.53	0.58	0.55
Special Education	Yes	4,630	0.69	16.69	0.37	0.37	0.39	0.31

Table D.8. Science Grade 10 Test Reliability

	Test Group	N	Alpha	SEM	Construct Explanations	Developing Models	Gathering and Investigating	Using Mathematical Thinking
All	Students Tested	43,250	0.82	15.03	0.61	0.50	0.54	0.58
	Female	20,391	0.78	15.74	0.55	0.43	0.49	0.56
Sex	Male	22,819	0.84	14.55	0.65	0.55	0.58	0.60
	Hispanic or Latino							
Ethnicity	Ethnicity	8,412	0.69	18.65	0.41	0.33	0.41	0.45
	Asian	763	0.85	14.58	0.68	0.58	0.55	0.64
Ethnicity	Native Hawaiian or Other Pacific Islander	631	0.61	20.37	0.23	0.33	0.30	0.41
	Black or African American	563	0.67	19.40	0.37	0.27	0.38	0.38
Ethnicity	American Indian or Alaska Native	435	0.62	20.14	0.28	0.22	0.38	0.44
	White	31,127	0.82	14.34	0.63	0.51	0.54	0.60
Ethnicity	Other	1,319	0.82	15.54	0.61	0.53	0.57	0.57
	Limited English Proficiency							
English Proficiency	No	39,966	0.82	14.71	0.62	0.50	0.54	0.59
	Yes	3,284	0.30	25.78	0.04	0.12	0.14	0.24
Economic Disadvantage	No	32,366	0.82	14.49	0.63	0.51	0.54	0.60
	Yes	10,884	0.75	17.27	0.49	0.40	0.48	0.50
Special Education	No	39,338	0.82	14.73	0.62	0.50	0.54	0.59
	Yes	3,912	0.59	20.68	0.28	0.28	0.37	0.33

Appendix E: Conditional Standard Error of Scale Scores

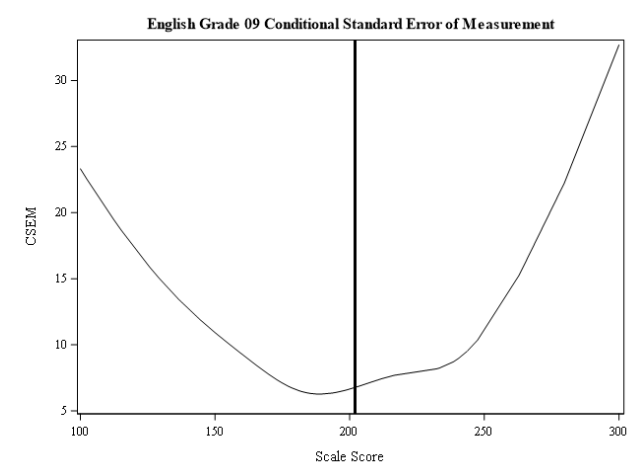


Figure E.1. English Grade 9 Conditional Standard Error of Scale Scores

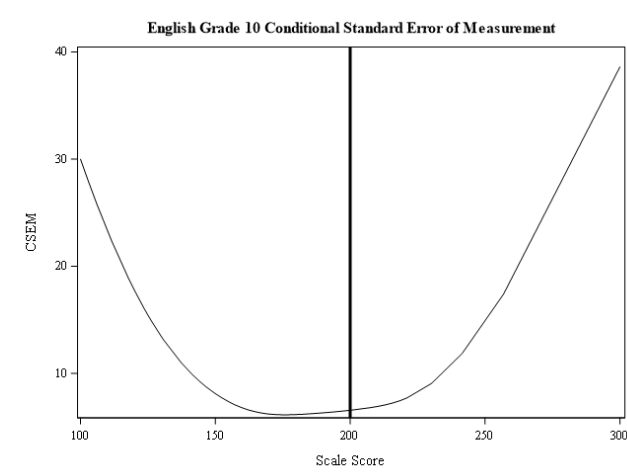


Figure E.2. English Grade 10 Conditional Standard Error of Scale Scores

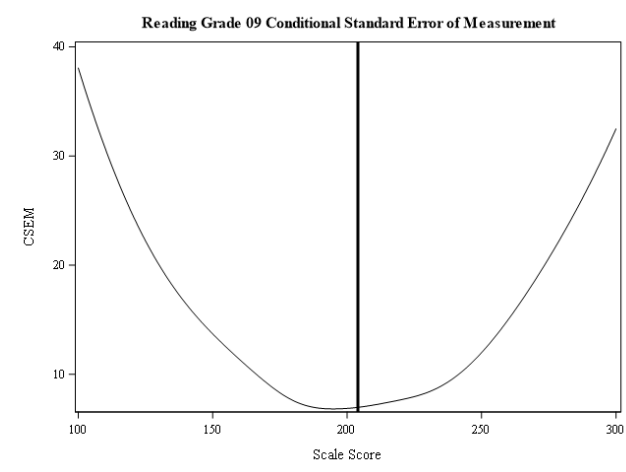


Figure E.3. Reading Grade 9 Conditional Standard Error of Scale Scores

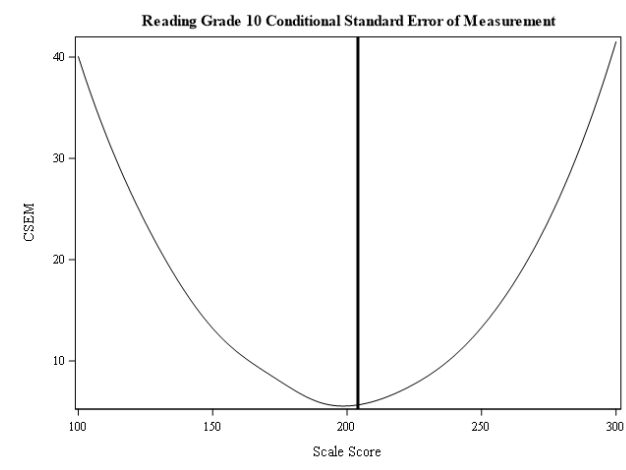


Figure E.4. Reading Grade 10 Conditional Standard Error of Scale Scores

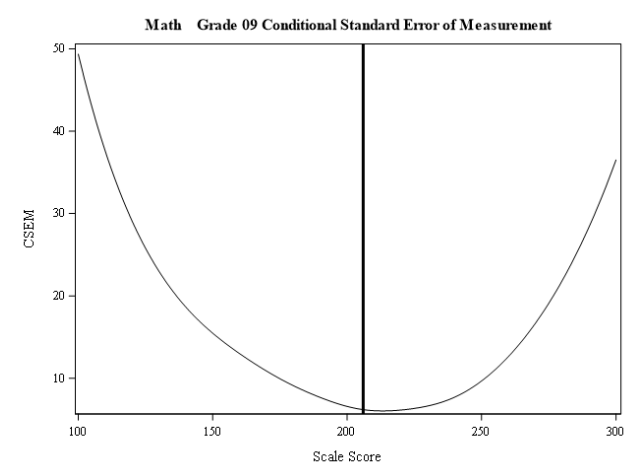


Figure E.5. Mathematics Grade 9 Conditional Standard Error of Scale Scores

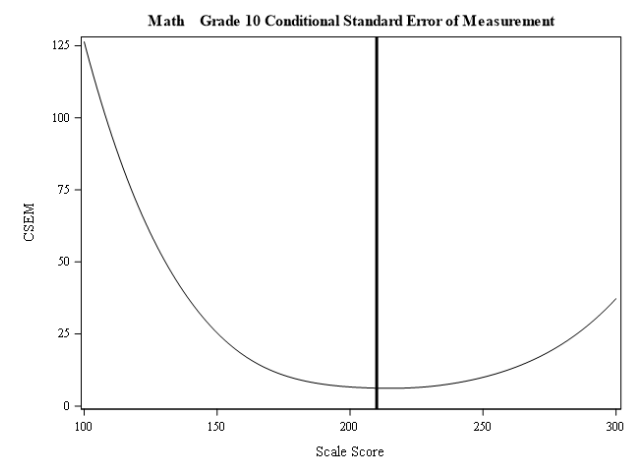


Figure E.6. Mathematics Grade 10 Conditional Standard Error of Scale Scores

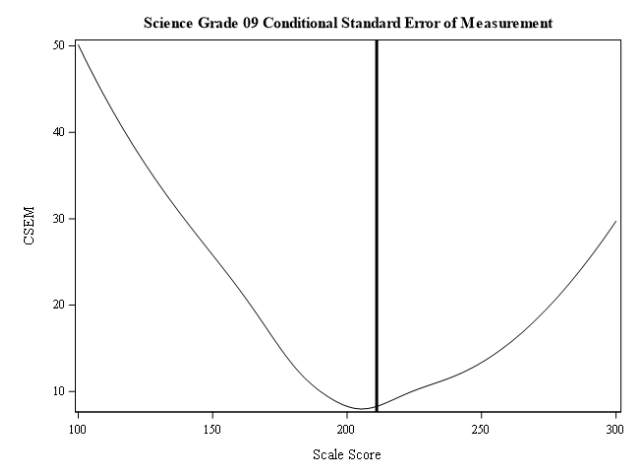


Figure E.7. Science Grade 9 Conditional Standard Error of Scale Scores

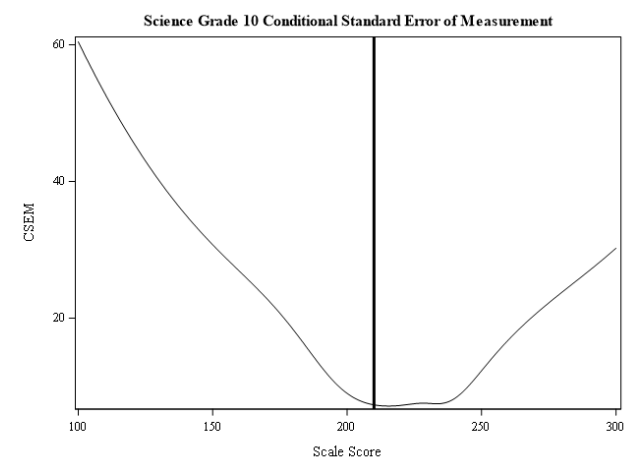


Figure E.8. Science Grade 10 Conditional Standard Error of Scale Scores

Appendix F: Accuracy and Consistency

Table F.1. Accuracy Classification for English Grade 9

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.109	0.023	0.000	0.000	81.85
Approaching Proficient	0.035	0.362	0.056	0.000	
Proficient	0.000	0.045	0.322	0.014	
Highly Proficient	0.000	0.000	0.009	0.025	

Table F.2. Accuracy Classification at Proficient Cut Point for English Grade 9

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.109	0.023	0.000	0.000	89.86
Approaching Proficient	0.035	0.362	0.056	0.000	
Proficient	0.000	0.045	0.322	0.014	
Highly Proficient	0.000	0.000	0.009	0.025	

Table F.3. Consistency Classification for English Grade 9

First Form	Alternate Form				Consistency %	Kappa
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient		
Below Proficient	0.104	0.041	0.000	0.000	74.21	0.600
Approaching Proficient	0.040	0.322	0.075	0.000		
Proficient	0.000	0.066	0.292	0.014		
Highly Proficient	0.000	0.000	0.021	0.024		

Table F.4. Accuracy Classification for English Grade 10

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.072	0.016	0.000	0.000	83.86
Approaching Proficient	0.026	0.368	0.049	0.000	
Proficient	0.000	0.047	0.368	0.015	
Highly Proficient	0.000	0.000	0.008	0.031	

Table F.5. Accuracy Classification at Proficient Cut Point for English Grade 10

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.072	0.016	0.000	0.000	90.43
Approaching Proficient	0.026	0.368	0.049	0.000	
Proficient	0.000	0.047	0.368	0.015	
Highly Proficient	0.000	0.000	0.008	0.031	

Table F.6. Consistency Classification for English Grade 10

First Form	Alternate Form				Consistency %	Kappa
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient		
Below Proficient	0.069	0.031	0.000	0.000	77.10	0.632
Approaching Proficient	0.029	0.333	0.068	0.000		
Proficient	0.000	0.066	0.340	0.016		
Highly Proficient	0.000	0.000	0.018	0.030		

Table F.7. Accuracy Classification for Reading Grade 9

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.092	0.021	0.000	0.000	78.28
Approaching Proficient	0.034	0.365	0.060	0.000	
Proficient	0.000	0.048	0.242	0.032	
Highly Proficient	0.000	0.000	0.023	0.085	

Table F.8. Accuracy Classification at Proficient Cut Point for Reading Grade 9

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.092	0.021	0.000	0.000	89.18
Approaching Proficient	0.034	0.365	0.060	0.000	
Proficient	0.000	0.048	0.242	0.032	
Highly Proficient	0.000	0.000	0.023	0.085	

Table F.9. Consistency Classification for Reading Grade 9

First Form	Alternate Form				Consistency %	Kappa
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient		
Below Proficient	0.088	0.039	0.000	0.000	69.46	0.550
Approaching Proficient	0.038	0.324	0.079	0.001		
Proficient	0.000	0.069	0.203	0.034		
Highly Proficient	0.000	0.001	0.043	0.081		

Table F.10. Accuracy Classification for Reading Grade 10

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.156	0.030	0.000	0.000	79.90
Approaching Proficient	0.037	0.308	0.053	0.000	
Proficient	0.000	0.045	0.281	0.021	
Highly Proficient	0.000	0.000	0.015	0.055	

Table F.11. Accuracy Classification at Proficient Cut Point for Reading Grade 10

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.156	0.030	0.000	0.000	90.18
Approaching Proficient	0.037	0.308	0.053	0.000	
Proficient	0.000	0.045	0.281	0.021	
Highly Proficient	0.000	0.000	0.015	0.055	

Table F.12. Consistency Classification for Reading Grade 10

First Form	Alternate Form				Consistency %	Kappa
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient		
Below Proficient	0.148	0.050	0.000	0.000	71.52	0.588
Approaching Proficient	0.044	0.267	0.071	0.000		
Proficient	0.000	0.066	0.248	0.023		
Highly Proficient	0.000	0.000	0.030	0.052		

Table F.13. Accuracy Classification for Mathematics Grade 9

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.191	0.038	0.000	0.000	78.70
Approaching Proficient	0.041	0.360	0.057	0.000	
Proficient	0.000	0.042	0.194	0.018	
Highly Proficient	0.000	0.000	0.016	0.043	

Table F.14. Accuracy Classification at Proficient Cut Point for Mathematics Grade 9

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.191	0.038	0.000	0.000	90.04
Approaching Proficient	0.041	0.360	0.057	0.000	
Proficient	0.000	0.042	0.194	0.018	
Highly Proficient	0.000	0.000	0.016	0.043	

Table F.15. Consistency Classification for Mathematics Grade 9

First Form	Alternate Form				Consistency %	Kappa
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient		
Below Proficient	0.181	0.062	0.000	0.000	69.83	0.557
Approaching Proficient	0.051	0.312	0.072	0.001		
Proficient	0.000	0.065	0.164	0.019		
Highly Proficient	0.000	0.001	0.031	0.041		

Table F.16. Accuracy Classification for Mathematics Grade 10

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.304	0.067	0.001	0.000	74.18
Approaching Proficient	0.042	0.259	0.062	0.000	
Proficient	0.000	0.051	0.144	0.014	
Highly Proficient	0.000	0.000	0.021	0.035	

Table F.17. Accuracy Classification at Proficient Cut Point for Mathematics Grade 10

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.304	0.067	0.001	0.000	88.55
Approaching Proficient	0.042	0.259	0.062	0.000	
Proficient	0.000	0.051	0.144	0.014	
Highly Proficient	0.000	0.000	0.021	0.035	

Table F.18. Consistency Classification for Mathematics Grade 10

First Form	Alternate Form				Consistency %	Kappa
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient		
Below Proficient	0.289	0.093	0.005	0.000	64.41	0.485
Approaching Proficient	0.053	0.206	0.069	0.002		
Proficient	0.003	0.072	0.115	0.014		
Highly Proficient	0.000	0.005	0.038	0.033		

Table F.19. Accuracy Classification for Science Grade 9

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.240	0.063	0.003	0.000	66.56
Approaching Proficient	0.058	0.190	0.075	0.002	
Proficient	0.003	0.064	0.167	0.038	
Highly Proficient	0.000	0.001	0.027	0.069	

Table F.20. Accuracy Classification at Proficient Cut Point for Science Grade 9

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.240	0.063	0.003	0.000	85.15
Approaching Proficient	0.058	0.190	0.075	0.002	
Proficient	0.003	0.064	0.167	0.038	
Highly Proficient	0.000	0.001	0.027	0.069	

Table F.21. Consistency Classification for Science Grade 9

First Form	Alternate Form				Consistency %	Kappa
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient		
Below Proficient	0.225	0.087	0.016	0.000	56.03	0.395
Approaching Proficient	0.064	0.140	0.077	0.006		
Proficient	0.012	0.080	0.128	0.035		
Highly Proficient	0.000	0.010	0.051	0.067		

Table F.22. Accuracy Classification for Science Grade 10

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.312	0.078	0.007	0.000	67.49
Approaching Proficient	0.054	0.167	0.077	0.001	
Proficient	0.004	0.062	0.165	0.022	
Highly Proficient	0.000	0.001	0.021	0.031	

Table F.23. Accuracy Classification at Proficient Cut Point for Science Grade 10

True Score	Observed Score				Accuracy %
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
Below Proficient	0.312	0.078	0.007	0.000	84.91
Approaching Proficient	0.054	0.167	0.077	0.001	
Proficient	0.004	0.062	0.165	0.022	
Highly Proficient	0.000	0.001	0.021	0.031	

Table F.24. Consistency Classification for Science Grade 10

First Form	Alternate Form				Consistency %	Kappa
	Below Proficient	Approaching Proficient	Proficient	Highly Proficient		
Below Proficient	0.293	0.099	0.023	0.000	57.38	0.388
Approaching Proficient	0.062	0.121	0.074	0.003		
Proficient	0.015	0.079	0.130	0.019		
Highly Proficient	0.000	0.008	0.043	0.031		

Appendix G: Common Item Scatterplots for 2023 Anchor Items

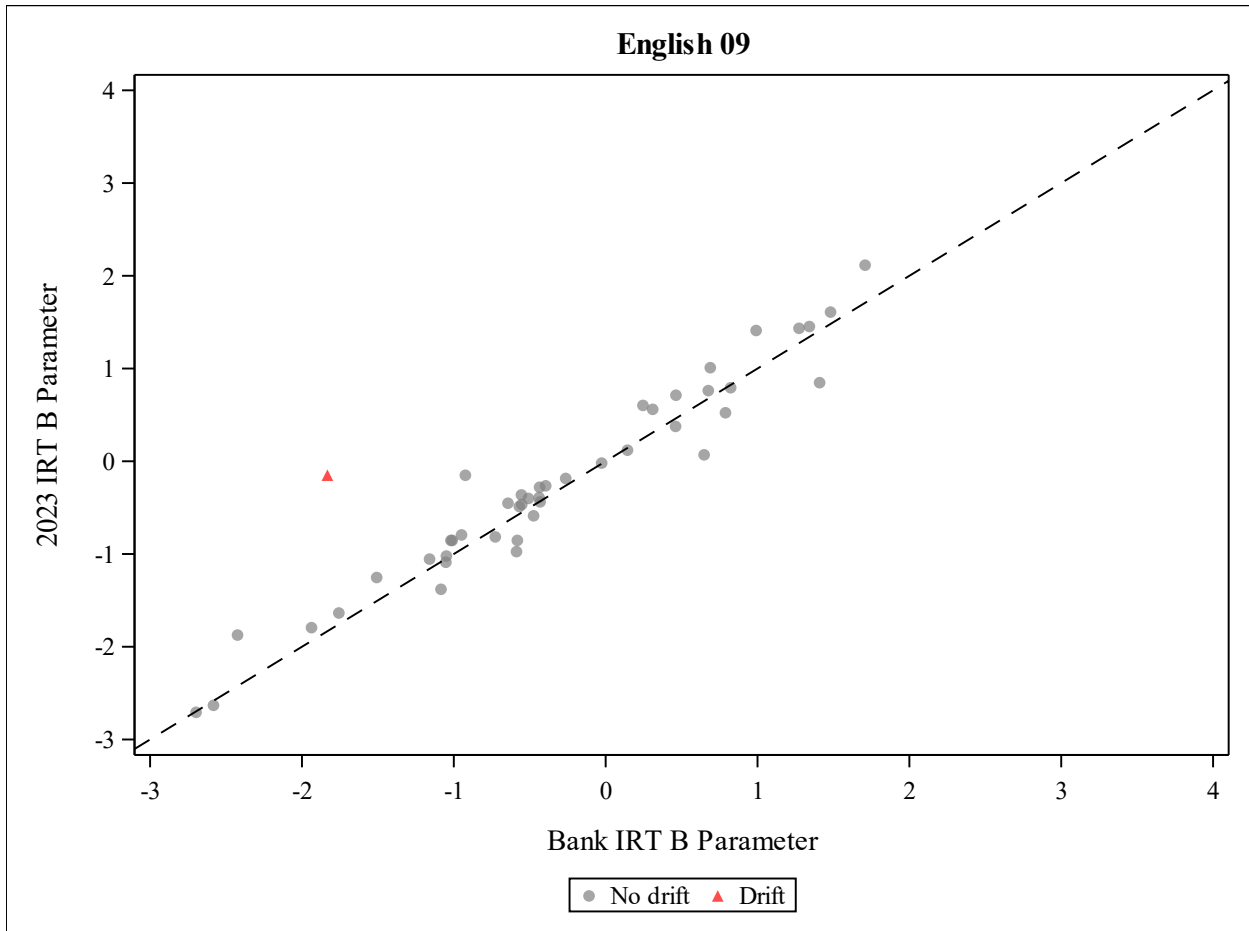


Figure G.1. English Grade 9 IRT B Parameters for Operational Items

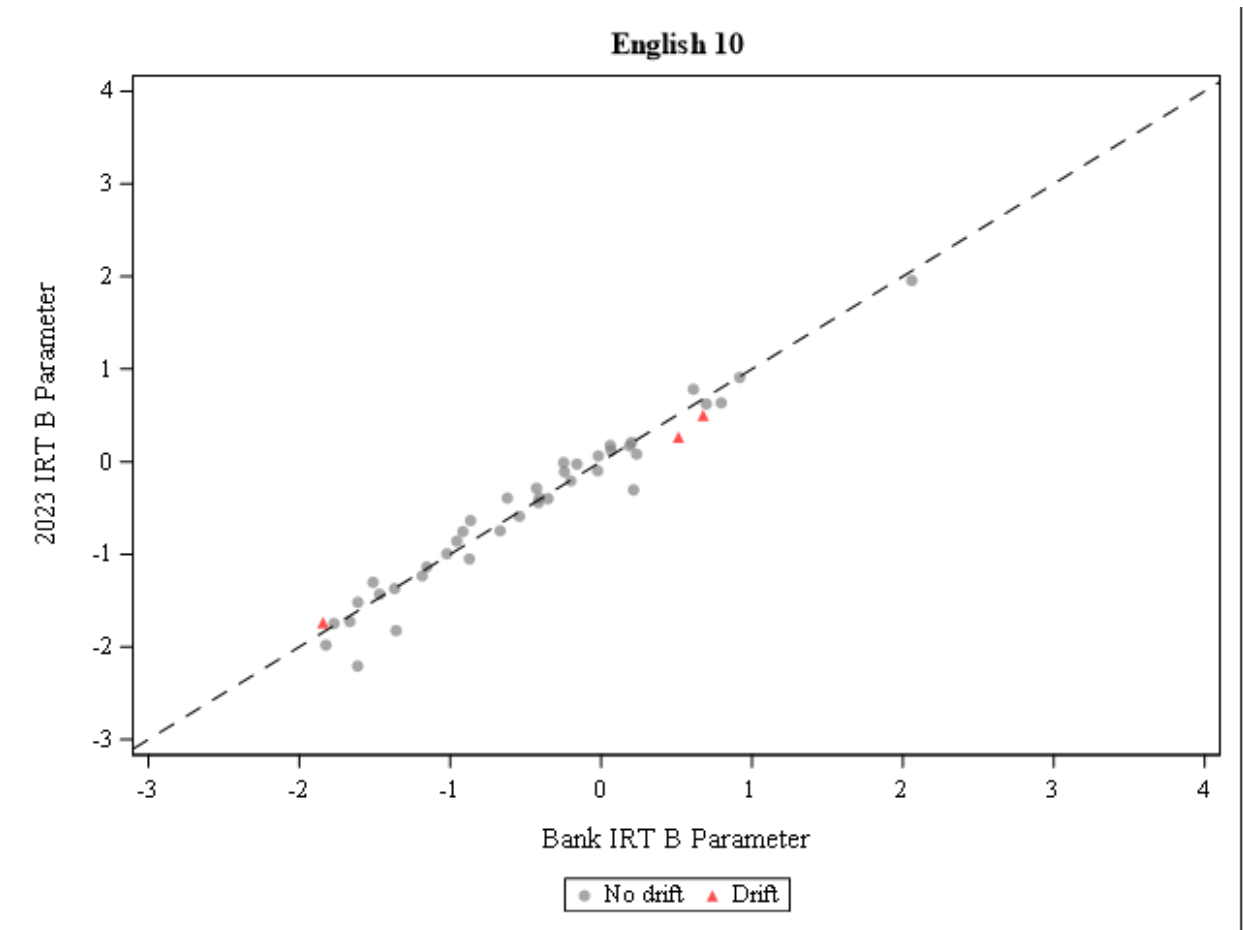


Figure G.2. English Grade 10 IRT B Parameters for Operational Items

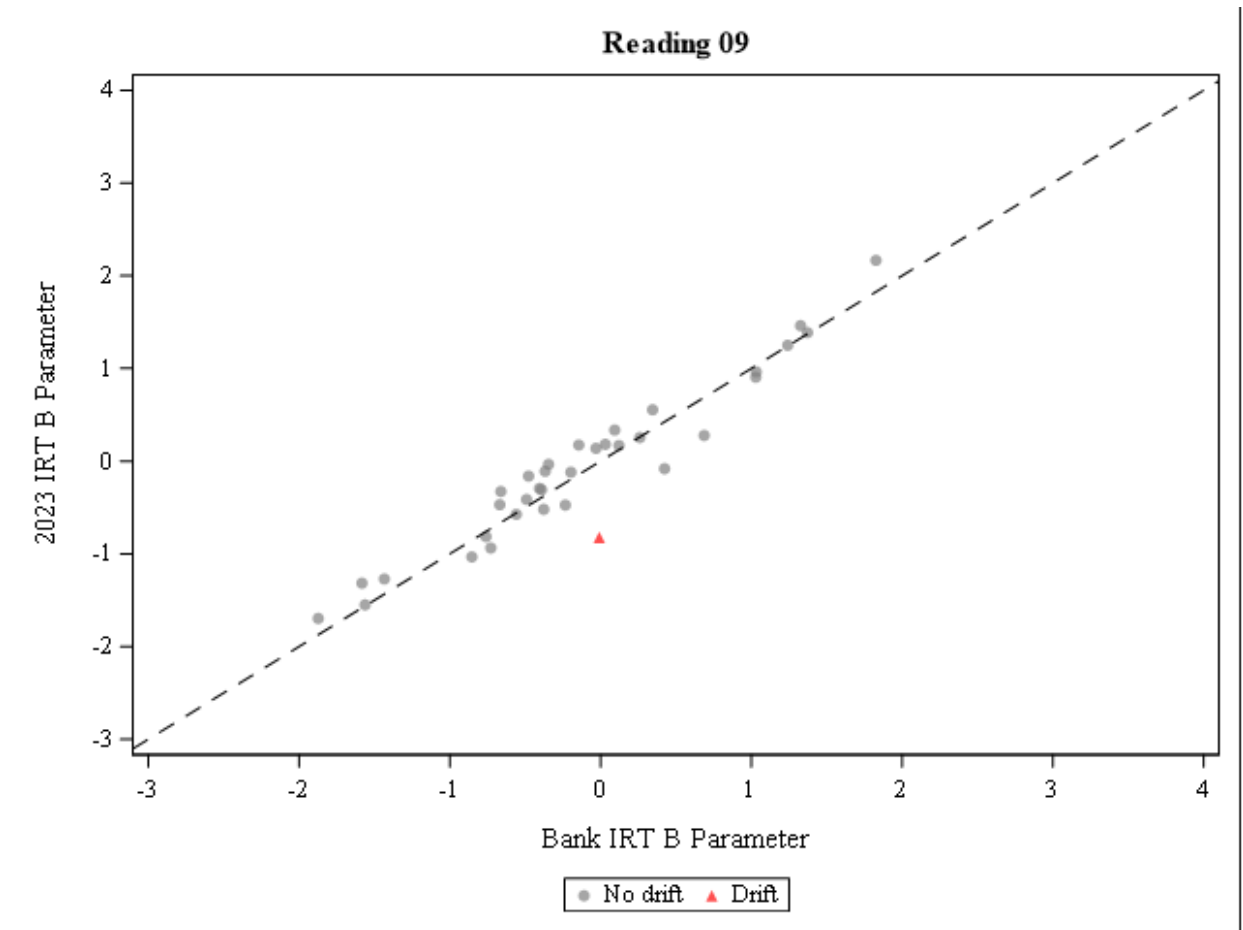


Figure G.3. Reading Grade 9 IRT B Parameters for Operational Items

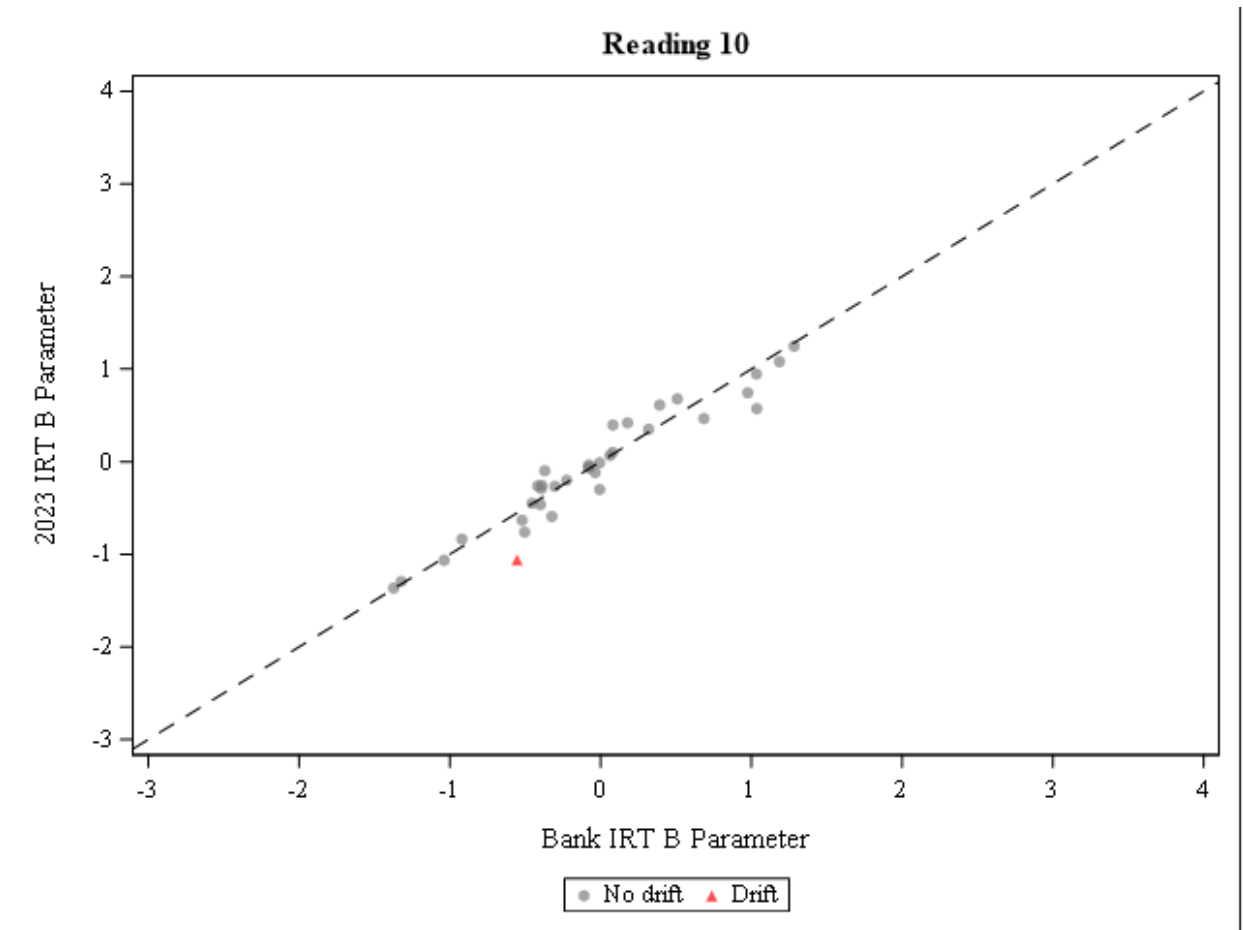


Figure G.4. Reading Grade 10 IRT B Parameters for Operational Items

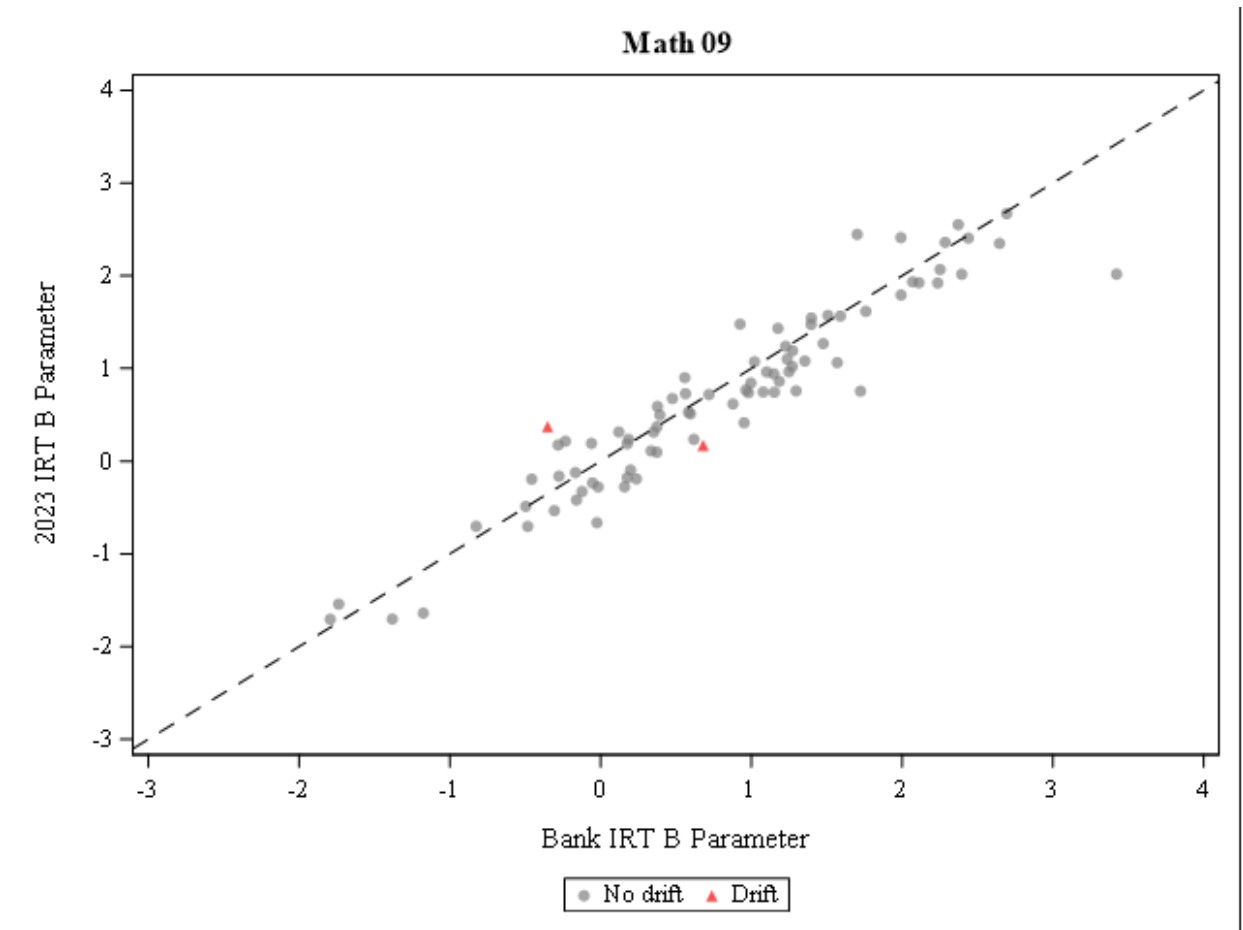


Figure G.5. Mathematics Grade 9 IRT B Parameters for Operational Items

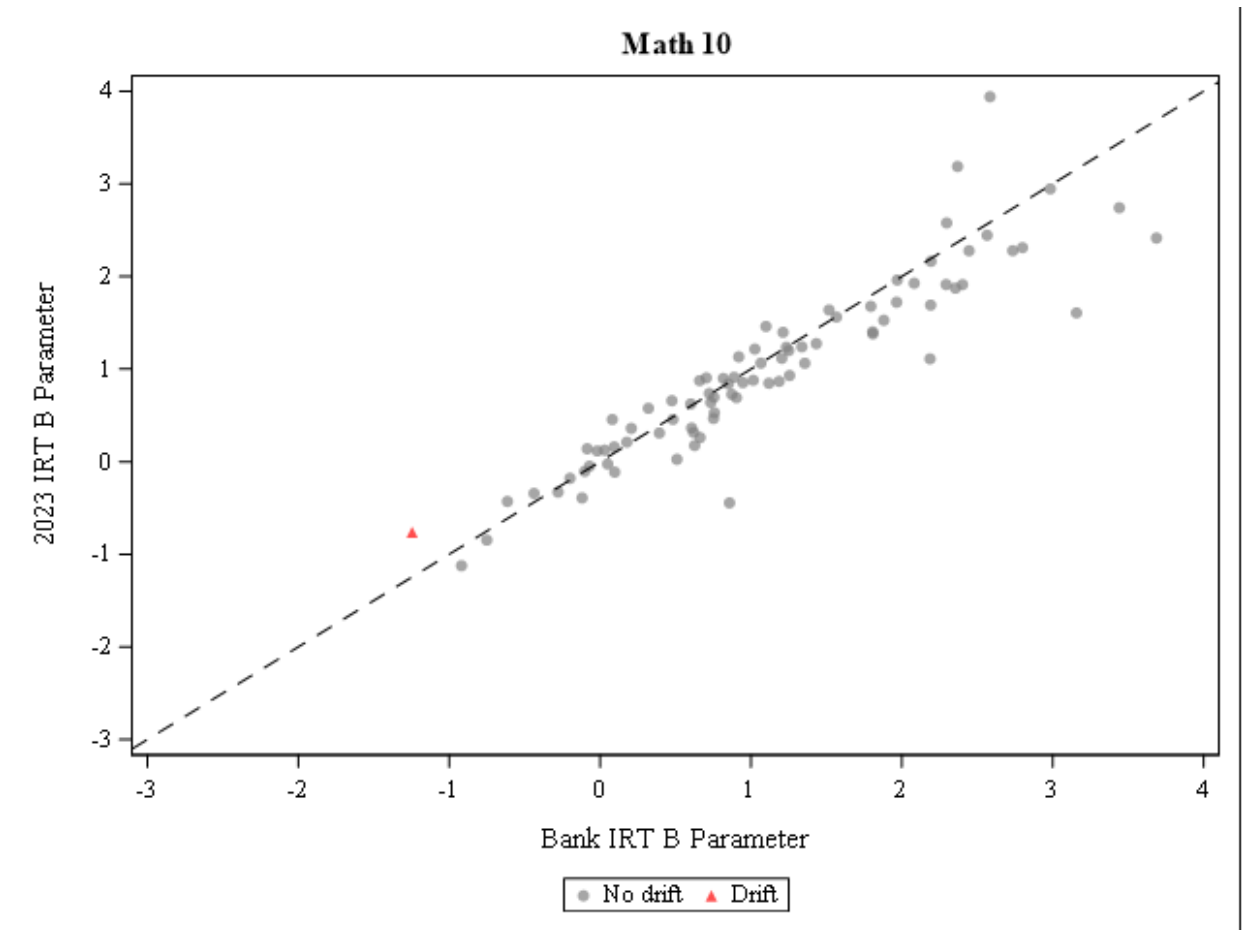


Figure G.6. Mathematics Grade 10 IRT B Parameters for Operational Items

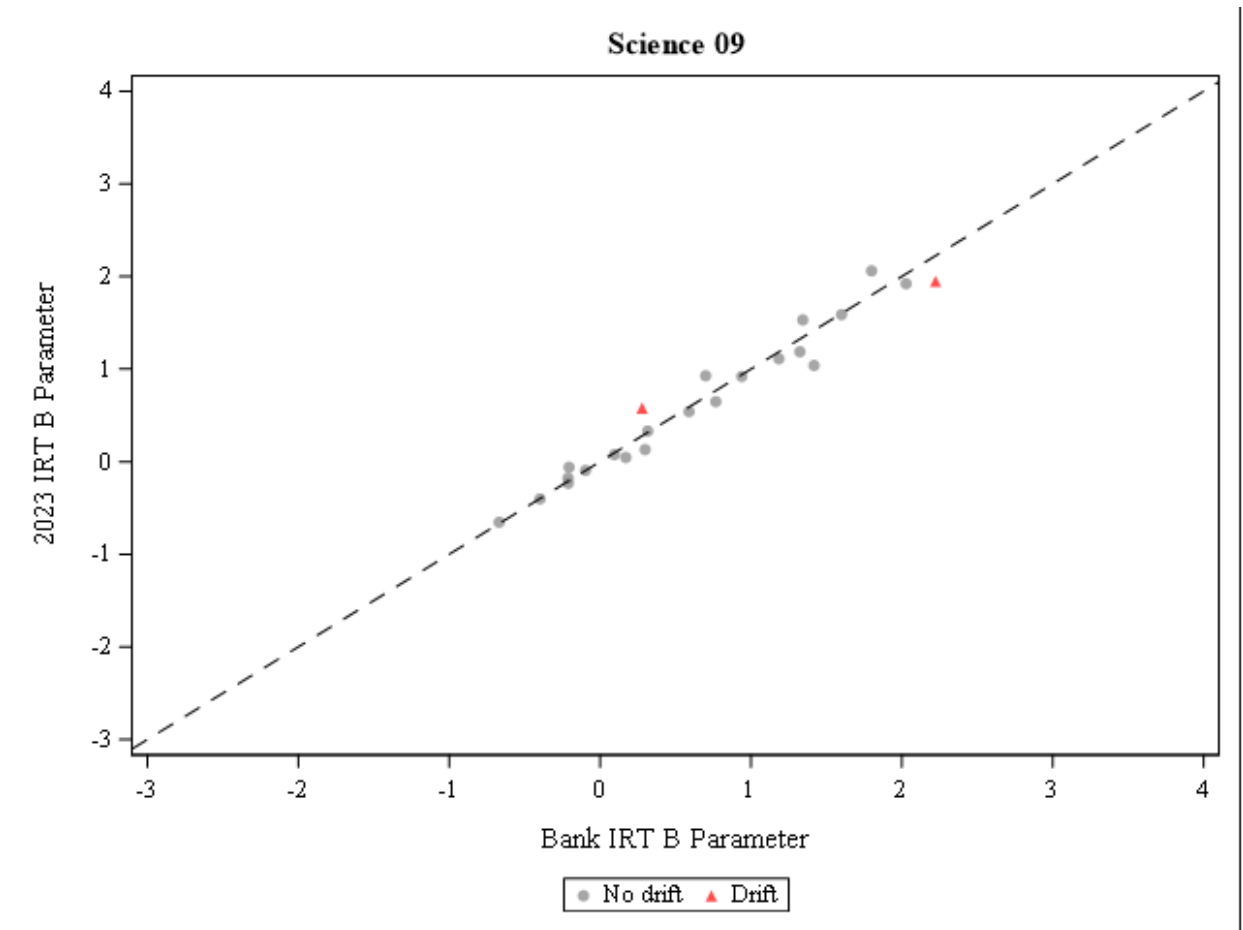


Figure G.7. Science Grade 9 IRT B Parameters for Operational Items

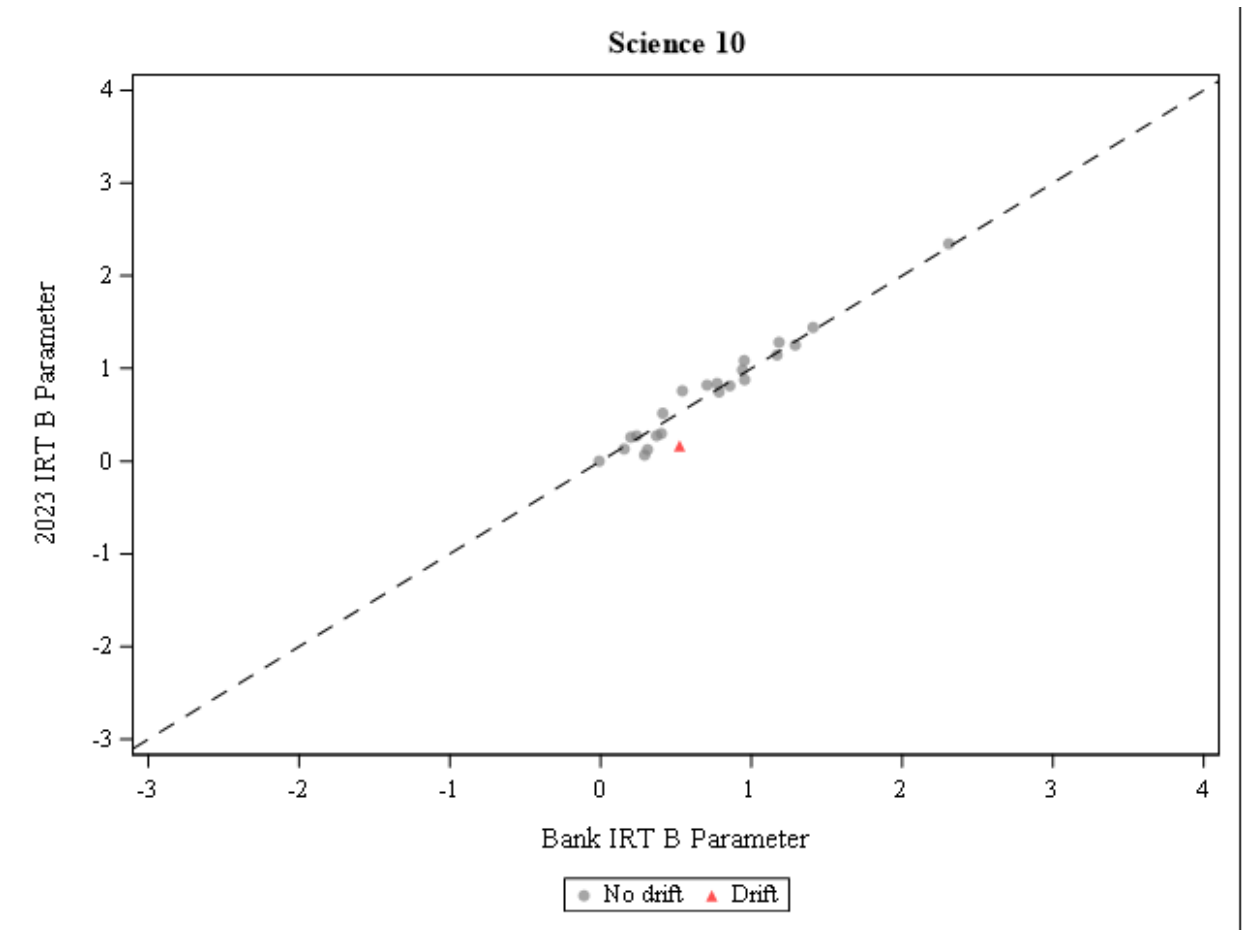


Figure G.8. Science Grade 10 IRT B Parameters for Operational Items

Appendix H: Scale Score Descriptive Statistics by Subgroup

Table H.1. English Grade 9 Scale Score Descriptive Statistics

	Test Group	N	Mean	SD	P25	Median	P75	Skew
All	Students Scored	46,511	195	28.33	177	197	214	-0.14
Sex	Female	22,210	199	27.24	182	200	217	-0.08
	Male	24,275	191	28.79	172	193	211	-0.15
Ethnicity	Hispanic or Latino Ethnicity	9,035	180	26.51	162	180	197	0.02
	Asian	816	201	29.47	183	201	219	-0.03
	Native Hawaiian or Other Pacific Islander	699	179	24.74	162	181	195	-0.18
	Black or African American	627	176	26.98	156	176	194	0.01
	American Indian or Alaska Native	463	175	24.74	159	175	192	0.05
	White	33,350	200	27.09	184	201.5	217	-0.19
	Other	1,521	197	26.76	180	197	215	-0.10
	Limited English Proficiency	No	42,375	198	27.19	181	199	216
	Yes	4,136	165	21.62	151	166	179	-0.18
Economic Disadvantage	No	33,773	199	27.29	183	201	217	-0.17
	Yes	12,738	183	27.51	164	184	202	-0.01
Special Education	No	41,883	198	27.12	181	199	216	-0.15
	Yes	4,628	168	24.07	152	167	183	0.2

Table H.2. English Grade 10 Scale Score Descriptive Statistics

	Test Group	N	Mean	SD	P25	Median	P75	Skew	
All	Students Scored	43,766	197	28.51	179	197	216	0.18	
Sex	Female	20,668	201	27.23	183	201	219	0.24	
	Male	23,056	194	29.14	174	194	213	0.20	
Ethnicity	Hispanic or Latino Ethnicity	8,441	183	24.62	167	182	199	0.22	
	Asian	767	200	29.77	181	199	219	0.20	
	Native Hawaiian or Other Pacific Islander	639	184	21.78	169	184	199	0.04	
	Black or African American	557	182	27.27	163	179	198	0.44	
	American Indian or Alaska Native	443	181	22.36	166	180	195	0.26	
	White	31,577	202	28.17	184	202	219	0.12	
	Other	1,342	198	27.37	181	198	215	0.17	
	Limited English Proficiency	No	40,486	200	27.83	182	200	217	0.18
	Yes	3,280	168	19.22	155	169	181	-0.05	
Economic Disadvantage	No	32,770	201	28.31	183	201	219	0.13	
	Yes	10,996	187	26.31	169	186	203	0.32	
Special Education	No	39,754	200	27.69	182	200	217	0.19	
	Yes	4,012	171	22.34	156	170	184	0.42	

Table H.3. Reading Grade 9 Scale Score Descriptive Statistics

	Test Group	N	Mean	SD	P25	Median	P75	Skew	
All	Students Scored	46,680	198	28.13	179	199	218	-0.18	
Sex	Female	22,256	201	27.03	183	202	220	-0.15	
	Male	24,399	196	28.86	176	197	216	-0.17	
Ethnicity	Hispanic or Latino Ethnicity	9,176	184	26.73	167	184	202	-0.02	
	Asian	816	204	28.94	187	206	224	-0.17	
	Native Hawaiian or Other Pacific Islander	707	183	25.18	168	184	200	-0.28	
	Black or African American	650	180	26.53	162	179	198	-0.02	
	American Indian or Alaska Native	467	181	24.51	166	181	195	0.19	
	White	33,337	203	27.06	185	204	221	-0.24	
	Other	1,527	200	27.61	181	201	219	-0.09	
	Limited English Proficiency	No	42,448	201	27.14	183	202	219	-0.20
	Yes	4,232	170	22.17	157	171	184	-0.30	
Economic Disadvantage	No	33,786	202	27.34	185	204	221	-0.24	
	Yes	12,894	187	27.13	169	187	206	-0.03	
Special Education	No	42,012	201	27.07	184	202	219	-0.20	
	Yes	4,668	172	24.15	158	172	187	0.16	

Table H.4. Reading Grade 10 Scale Score Descriptive Statistics

	Test Group	N	Mean	SD	P25	Median	P75	Skew
All	Students Scored	43,536	198	27.28	180	199	215	-0.11
Sex	Female	20,532	200	26.21	184	201	216.5	-0.06
	Male	22,966	196	28.05	177	197	214	-0.13
Ethnicity	Hispanic or Latino Ethnicity	8,429	184	24.63	168	185	200	0.00
	Asian	761	200	28.62	183	200	219	0.00
	Native Hawaiian or Other Pacific Islander	630	182	23.30	167	183	199	-0.31
	Black or African American	561	183	26.67	165	181	200	0.07
	American Indian or Alaska Native	439	184	21.78	171	183	199	-0.04
	White	31,373	202	26.68	186	203	219	-0.19
	Other	1,343	199	27.38	182	200	217	-0.19
	Limited English Proficiency	No	40,237	200	26.66	184	201	217
	Yes	3,299	171	19.75	160	172	184	-0.52
Economic Disadvantage	No	32,595	201	26.98	185	202	218	-0.17
	Yes	10,941	188	25.78	171	188	205	0.01
Special Education	No	39,524	200	26.60	184	201	217	-0.13
	Yes	4,012	175	22.98	161	174	189	-0.03

Table H.5. Mathematics Grade 9 Scale Score Descriptive Statistics

	Test Group	N	Mean	SD	P25	Median	P75	Skew
All	Students Scored	45,163	191	29.34	173	194	211	-0.40
Sex	Female	21,400	191	26.91	175	194	209	-0.46
	Male	23,739	191	31.36	172	194	213	-0.35
Ethnicity	Hispanic or Latino Ethnicity	8,736	175	27.42	158	176	193	-0.27
	Asian	794	200	30.94	182	201	218	-0.07
	Native Hawaiian or Other Pacific Islander	677	174	25.29	160	175	192	-0.47
	Black or African American	628	168	27.81	152	167	188	-0.27
	American Indian or Alaska Native	443	171	27.69	155	172	189	-0.16
	White	32,410	197	27.69	181	199	215	-0.50
	Other	1,475	193	29.98	174	194	212	-0.33
	Limited English Proficiency	No	41,098	194	28.27	177	196	213
	Yes	4,065	163	24.24	149	164	179	-0.45
Economic Disadvantage	No	32,845	196	28.02	180	199	215	-0.46
	Yes	12,318	178	28.68	160	179	197	-0.26
Special Education	No	40,606	194	27.95	178	196	213	-0.41
	Yes	4,557	164	26.70	148	164	180	-0.16

Table H.6. Mathematics Grade 10 Scale Score Descriptive Statistics

	Test Group	N	Mean	SD	P25	Median	P75	Skew	
All	Students Scored	42,908	188	36.04	172	193	212	-0.87	
Sex	Female	20,224	188	33.41	174	193	210	-1.02	
	Male	22,652	187	38.24	170	193	213	-0.76	
Ethnicity	Hispanic or Latino Ethnicity	8,303	169	35.86	156	175	193	-0.67	
	Asian	762	196	37.91	177	200	221	-0.79	
	Native Hawaiian or Other Pacific Islander	632	168	35.00	156	175	191	-0.71	
	Black or African American	544	163	38.01	149.5	170	189	-0.43	
	American Indian or Alaska Native	427	168	35.03	157	176	191	-0.81	
	White	30,924	194	33.81	179	199	215	-1.01	
	Other	1,316	188	35.73	172	191	211	-0.84	
	Limited English Proficiency	No	39,656	190	34.83	175	195	213	-0.94
	Yes	3,252	155	34.19	124	165	179	-0.53	
Economic Disadvantage	No	32,188	192	34.63	177	197	215	-0.97	
	Yes	10,720	174	36.47	160	179	198	-0.69	
Special Education	No	38,941	191	34.30	176	196	213	-0.95	
	Yes	3,967	154	35.80	115	164	179	-0.33	

Table H.7. Science Grade 9 Scale Score Descriptive Statistics

	Test Group	N	Mean	SD	P25	Median	P75	Skew
All	Students Scored	46,592	200	32.47	182	202	220	-0.29
Sex	Female	22,205	199	30.17	183	202	218	-0.54
	Male	24,362	202	34.38	181	202	223	-0.15
Ethnicity	Hispanic or Latino Ethnicity	9,136	184	30.75	167	186	203	-0.36
	Asian	817	207	33.16	188	209	227	-0.22
	Native Hawaiian or Other Pacific Islander	695	183	26.45	167	185	200	-0.61
	Black or African American	652	176	28.47	161	178	194	-0.42
	American Indian or Alaska Native	470	181	31.44	166	184	199	-0.31
	White	33,298	206	31.07	188	207	224	-0.30
	Other	1,524	202	32.63	183.5	203	222	-0.15
	Limited English Proficiency	No	42,396	203	31.56	185	205	222
	Yes	4,196	171	26.91	157	175	190	-0.60
Economic Disadvantage	No	33,716	205	31.64	187	207	224	-0.31
	Yes	12,876	188	31.58	171	190	208	-0.27
Special Education	No	41,936	203	31.64	185	205	222	-0.31
	Yes	4,656	176	29.78	160	178	194	-0.18

Table H.8. Science Grade 10 Scale Score Descriptive Statistics

	Test Group	N	Mean	SD	P25	Median	P75	Skew
All	Students Scored	43,280	192	35.07	174	197	215	-0.59
Sex	Female	20,407	191	33.58	173	196	213	-0.75
	Male	22,833	194	36.29	174	198	217	-0.50
Ethnicity	Hispanic or Latino Ethnicity	8,419	177	33.28	160	182	200	-0.58
	Asian	763	196	37.34	176	202	220	-0.52
	Native Hawaiian or Other Pacific Islander	631	174	32.49	160	180	197	-0.76
	Black or African American	564	174	33.58	154.5	178	196	-0.36
	American Indian or Alaska Native	435	176	32.64	158	182	198	-0.74
	White	31,149	197	34.05	180	202	219	-0.67
	Other	1,319	191	36.68	172	197	215	-0.53
	Limited English Proficiency	No	39,991	194	34.47	177	199	217
	Yes	3,289	165	30.79	147	171	188	-0.68
Economic Disadvantage	No	32,378	196	34.43	178	201	218	-0.64
	Yes	10,902	181	34.51	163	185	204	-0.54
Special Education	No	39,341	195	34.48	177	199	217	-0.64
	Yes	3,939	169	32.44	150	174	191	-0.47

Appendix I: Scale Score Distributions for Overall Testing Population

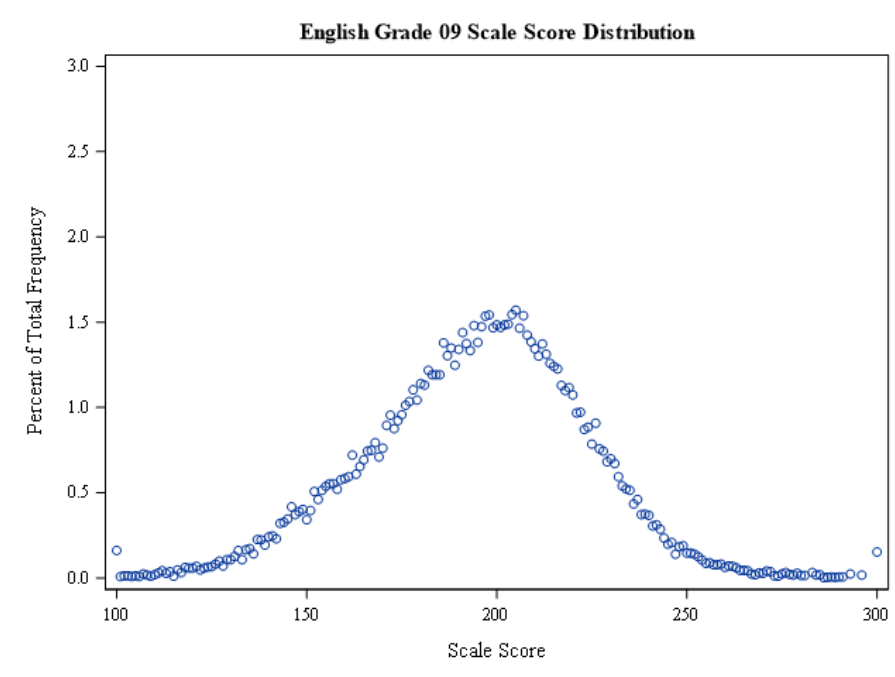


Figure I.1. English Grade 9 Scale Score Distribution

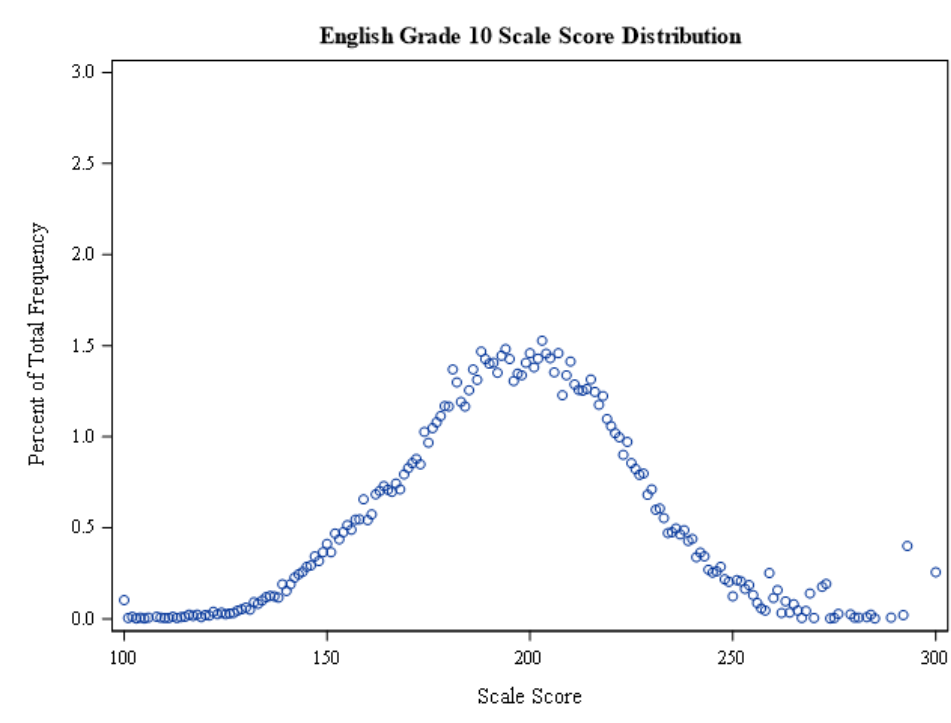


Figure I.2. English Grade 10 Scale Score Distribution

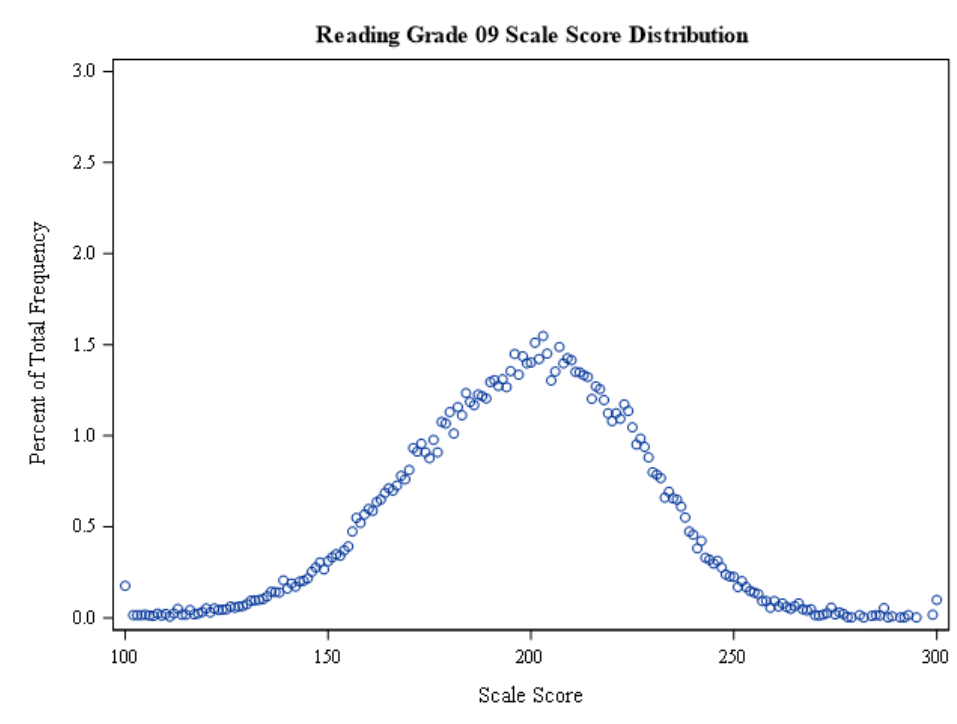


Figure I.3. Reading Grade 9 Scale Score Distribution

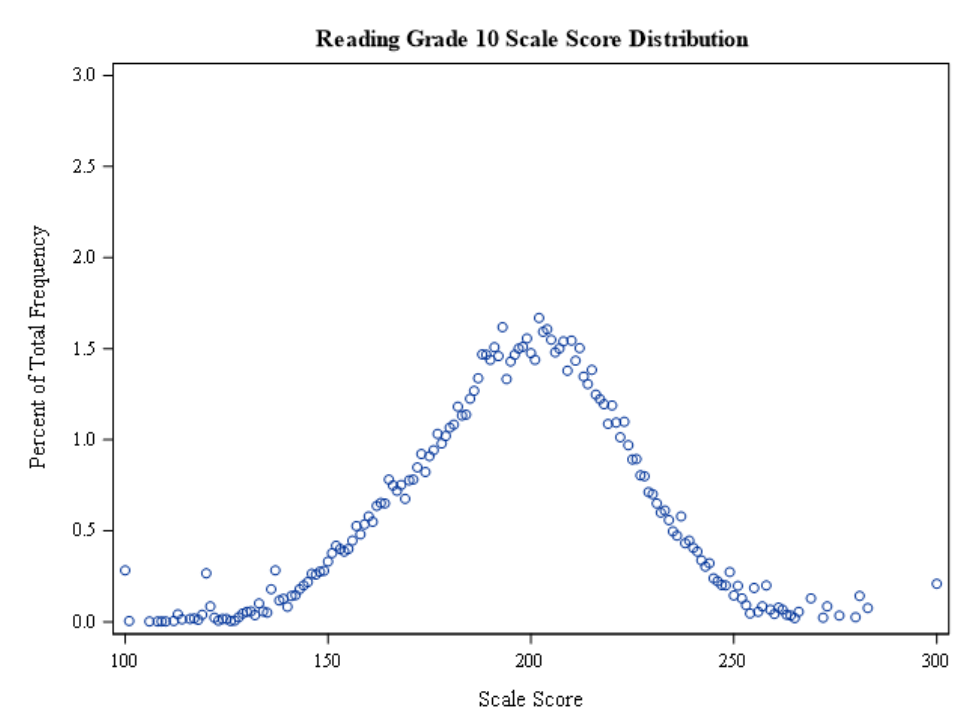


Figure I.4. Reading Grade 10 Scale Score Distribution

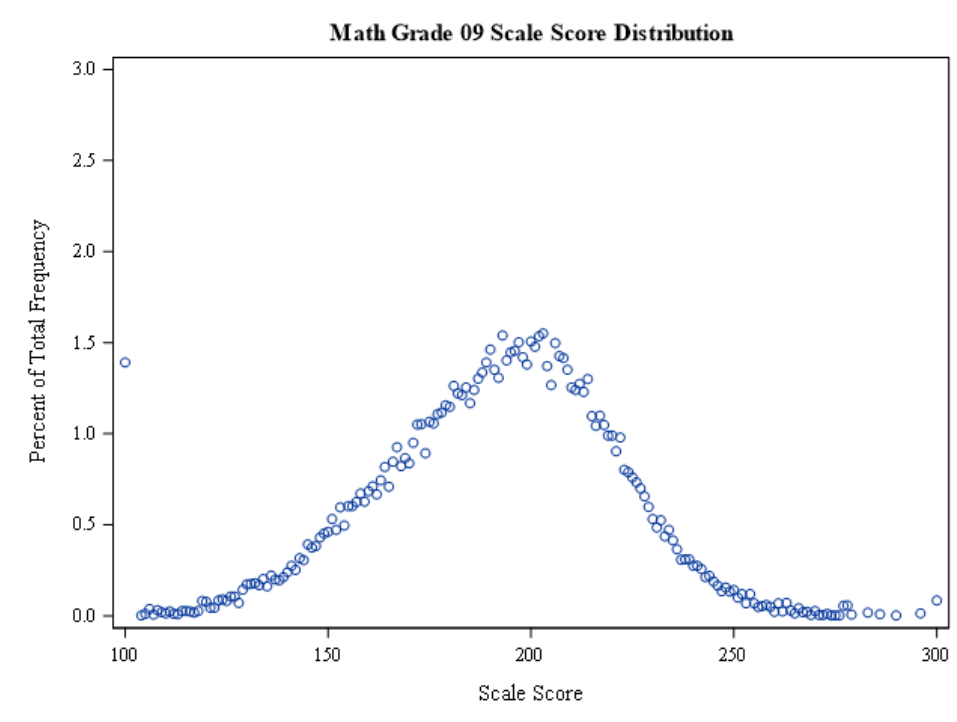


Figure I.5. Mathematics Grade 9 Scale Score Distribution

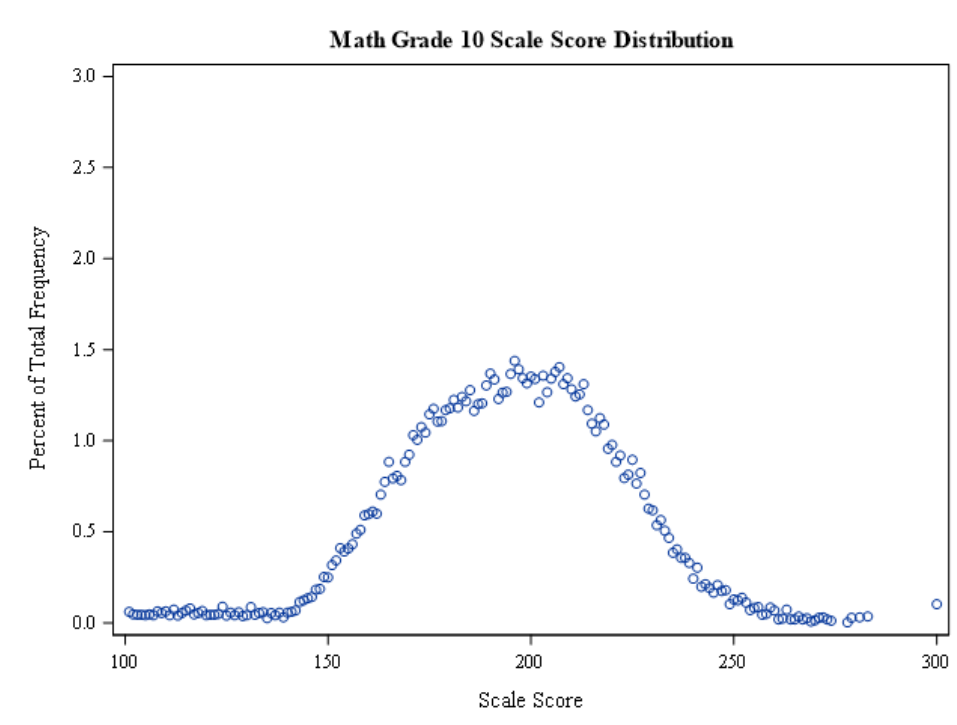


Figure I.6. Mathematics Grade 10 Scale Score Distribution

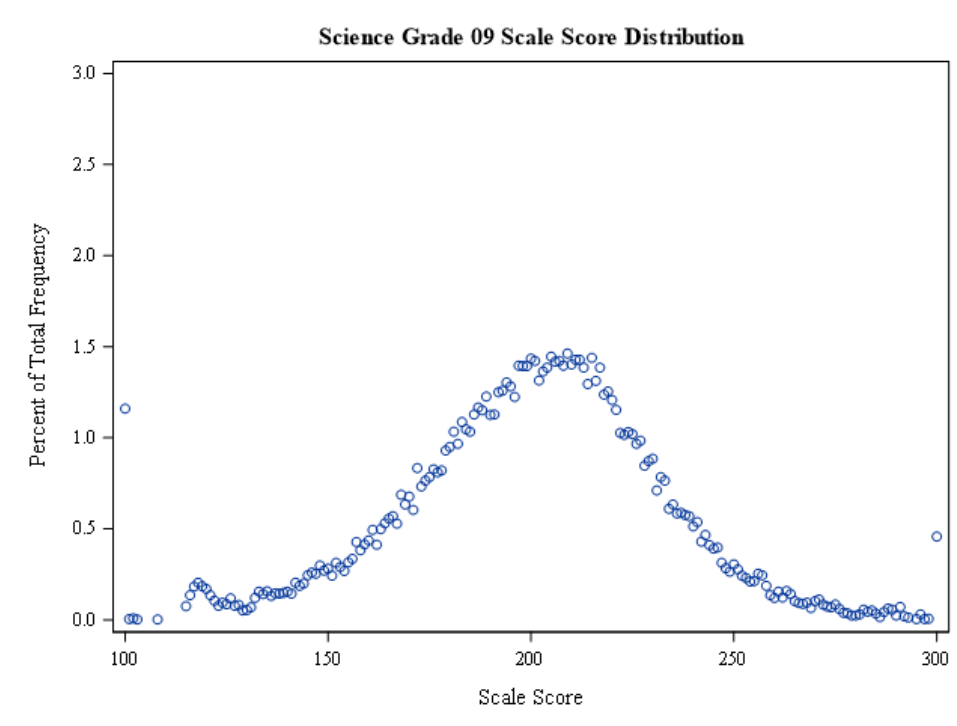


Figure I.7. Science Grade 9 Scale Score Distribution

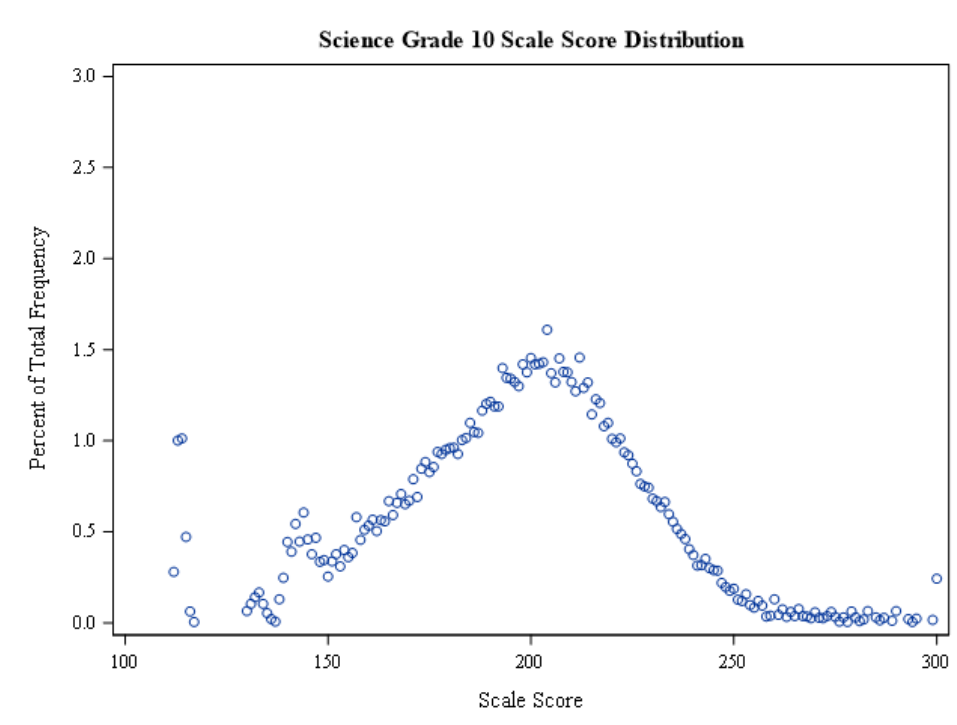


Figure I.8. Science Grade 10 Scale Score Distribution

Appendix J: Performance Level Distributions

Table J.1. English Grade 9 Performance Level Distribution

	Test Group	N	Below Proficient	Approaching Proficient	Proficient	Highly Proficient
All	Students Scored	46,511	14.43	42.91	38.78	3.87
Sex	Female	22,210	10.28	41.98	42.82	4.92
	Male	24,275	18.24	43.78	35.07	2.90
Ethnicity	Hispanic or Latino					
	Ethnicity	9,035	28.26	51.47	19.31	0.96
	Asian	816	12.25	37.87	42.52	7.35
	Native Hawaiian or Other Pacific Islander	699	26.90	56.65	15.88	0.57
	Black or African American	627	35.73	47.37	15.95	0.96
	American Indian or Alaska Native	463	31.97	54.21	13.61	0.22
	White	33,350	9.94	40.06	45.24	4.76
	Other	1,521	12.10	45.89	38.53	3.48
Limited English Proficiency	No	42,375	11.21	42.34	42.21	4.24
	Yes	4,136	47.49	48.77	3.68	0.07
Economic Disadvantage	No	33,773	10.27	40.45	44.50	4.78
	Yes	12,738	25.46	49.45	23.63	1.46
Special Education	No	41,883	10.95	42.57	42.21	4.26
	Yes	4,628	45.94	46.00	7.78	0.28

Table J.2. English Grade 10 Performance Level Distribution

	Test Group	N	Below Proficient	Approaching Proficient	Proficient	Highly Proficient
All	Students Scored	43,766	9.81	43.09	42.51	4.59
Sex	Female	20,668	6.21	41.33	47.06	5.40
	Male	23,056	13.05	44.70	38.39	3.86
Ethnicity	Hispanic or Latino					
	Ethnicity	8,441	18.15	57.64	23.26	0.96
	Asian	767	9.26	40.81	44.59	5.35
	Native Hawaiian or Other Pacific Islander	639	14.71	62.60	22.22	0.47
	Black or African American	557	21.36	54.94	22.44	1.26
	American Indian or Alaska Native	443	17.83	61.85	19.86	0.45
	White	31,577	7.24	38.35	48.64	5.77
	Other	1,342	8.57	43.89	43.52	4.02
	Limited English Proficiency	No	40,486	7.87	41.59	45.58
	Yes	3,280	33.78	61.59	4.57	0.06
Economic Disadvantage	No	32,770	7.81	39.26	47.38	5.55
	Yes	10,996	15.80	54.48	27.99	1.73
Special Education	No	39,754	7.54	41.60	45.84	5.02
	Yes	4,012	32.38	57.85	9.42	0.35

Table J.3. Reading Grade 9 Performance Level Distribution

	Test Group	N	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
All	Students Scored	46,680	12.62	43.35	32.43	11.61	
Sex	Female	22,256	9.61	43.16	34.30	12.94	
	Male	24,399	15.38	43.52	30.71	10.39	
Ethnicity	Hispanic or Latino						
	Ethnicity	9,176	23.25	53.29	19.46	4.00	
	Asian	816	9.19	37.38	36.76	16.67	
	Native Hawaiian or Other Pacific Islander	707	21.78	58.42	17.68	2.12	
	Black or African American	650	30.62	50.31	16.62	2.46	
	American Indian or Alaska Native	467	24.63	58.03	13.92	3.43	
	White	33,337	9.18	40.06	36.73	14.02	
	Other	1,527	10.15	44.01	33.14	12.70	
	Limited English Proficiency	No	42,448	10.01	42.20	35.05	12.75
		Yes	4,232	38.78	54.89	6.12	0.21
Economic Disadvantage	No	33,786	9.36	40.27	36.24	14.12	
	Yes	12,894	21.16	51.40	22.42	5.02	
Special Education	No	42,012	9.74	42.39	35.13	12.74	
	Yes	4,668	38.54	51.99	8.05	1.41	

Table J.4. Reading Grade 10 Performance Level Distribution

	Test Group	N	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
All	Students Scored	43,536	19.28	38.24	34.90	7.58	
Sex	Female	20,532	15.91	39.12	36.84	8.12	
	Male	22,966	22.32	37.47	33.13	7.08	
Ethnicity	Hispanic or Latino						
	Ethnicity	8,429	33.63	45.97	18.06	2.34	
	Asian	761	16.43	39.16	35.09	9.33	
	Native Hawaiian or Other Pacific Islander	630	36.67	44.60	18.10	0.63	
	Black or African American	561	39.57	39.57	18.00	2.85	
	American Indian or Alaska Native	439	32.12	49.20	17.77	0.91	
	White	31,373	14.68	35.85	40.22	9.25	
	Other	1,343	17.42	38.12	36.56	7.89	
	Limited English Proficiency	No	40,237	16.31	38.03	37.47	8.20
		Yes	3,299	55.50	40.89	3.55	0.06
Economic Disadvantage	No	32,595	15.67	36.44	38.90	9.00	
	Yes	10,941	30.03	43.63	22.97	3.36	
Special Education	No	39,524	16.14	38.06	37.52	8.27	
	Yes	4,012	50.20	40.03	9.02	0.75	

Table J.5. Mathematics Grade 9 Performance Level Distribution

	Test Group	N	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
All	Students Scored	45,163	23.21	44.01	26.71	6.07	
Sex	Female	21,400	21.29	47.68	26.71	4.32	
	Male	23,739	24.95	40.68	26.71	7.65	
Ethnicity	Hispanic or Latino						
	Ethnicity	8,736	43.62	44.22	10.87	1.28	
	Asian	794	16.37	41.31	28.84	13.48	
	Native Hawaiian or Other Pacific Islander	677	43.43	47.12	8.71	0.74	
	Black or African American	628	56.37	36.31	6.69	0.64	
	American Indian or Alaska Native	443	48.31	42.44	8.13	1.13	
	White	32,410	16.52	44.15	31.94	7.39	
	Other	1,475	22.03	43.39	26.71	7.86	
	Limited English Proficiency	No	41,098	19.27	44.96	29.11	6.67
		Yes	4,065	63.05	34.44	2.44	0.07
Economic Disadvantage	No	32,845	17.20	43.85	31.44	7.52	
	Yes	12,318	39.24	44.45	14.10	2.21	
Special Education	No	40,606	18.87	45.23	29.22	6.68	
	Yes	4,557	61.86	33.14	4.34	0.66	

Table J.6. Mathematics Grade 10 Performance Level Distribution

	Test Group	N	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
All	Students Scored	42,908	34.58	37.63	22.84	4.94	
Sex	Female	20,224	32.90	41.43	22.28	3.39	
	Male	22,652	36.06	34.24	23.36	6.33	
Ethnicity	Hispanic or Latino						
	Ethnicity	8,303	58.35	32.21	8.55	0.89	
	Asian	762	29.13	32.02	27.69	11.15	
	Native Hawaiian or Other Pacific Islander	632	59.81	34.02	5.85	0.32	
	Black or African American	544	65.81	25.74	6.99	1.47	
	American Indian or Alaska Native	427	59.25	33.72	7.03	0.00	
	White	30,924	26.94	39.54	27.43	6.09	
	Other	1,316	34.27	38.30	22.26	5.17	
	Limited English Proficiency	No	39,656	31.04	39.04	24.57	5.34
		Yes	3,252	77.74	20.48	1.72	0.06
Economic Disadvantage	No	32,188	28.66	38.90	26.43	6.01	
	Yes	10,720	52.36	33.83	12.08	1.73	
Special Education	No	38,941	30.18	39.60	24.82	5.40	
	Yes	3,967	77.74	18.30	3.45	0.50	

Table J.7. Science Grade 9 Performance Level Distribution

	Test Group	N	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
All	Students Scored	46,592	30.10	31.76	27.24	10.90	
Sex	Female	22,205	28.91	34.97	28.48	7.64	
	Male	24,362	31.20	28.83	26.11	13.86	
Ethnicity	Hispanic or Latino						
	Ethnicity	9,136	50.31	32.86	13.81	3.02	
	Asian	817	23.62	30.11	29.87	16.40	
	Native Hawaiian or Other Pacific Islander	695	54.39	33.09	11.51	1.01	
	Black or African American	652	63.19	27.76	8.13	0.92	
	American Indian or Alaska Native	470	53.19	33.83	10.64	2.34	
	White	33,298	23.32	31.48	31.84	13.36	
	Other	1,524	28.08	32.68	26.38	12.86	
	Limited English Proficiency	No	42,396	26.21	32.27	29.58	11.94
		Yes	4,196	69.35	26.64	3.62	0.38
Economic Disadvantage	No	33,716	24.45	31.46	30.99	13.10	
	Yes	12,876	44.89	32.55	17.44	5.13	
Special Education	No	41,936	26.31	32.35	29.47	11.87	
	Yes	4,656	64.20	26.44	7.22	2.15	

Table J.8. Science Grade 10 Performance Level Distribution

	Test Group	N	Below Proficient	Approaching Proficient	Proficient	Highly Proficient	
All	Students Scored	43,280	37.01	30.73	26.91	5.35	
Sex	Female	20,407	37.56	32.64	26.23	3.57	
	Male	22,833	36.50	29.05	27.50	6.95	
Ethnicity	Hispanic or Latino						
	Ethnicity	8,419	56.16	30.07	12.47	1.29	
	Asian	763	33.94	26.08	31.59	8.39	
	Native Hawaiian or Other Pacific Islander	631	58.00	32.17	9.03	0.79	
	Black or African American	564	61.70	24.65	12.59	1.06	
	American Indian or Alaska Native	435	55.63	32.64	11.26	0.46	
	White	31,149	30.77	31.06	31.55	6.61	
	Other	1,319	37.00	31.16	26.38	5.46	
	Limited English Proficiency	No	39,991	34.00	31.33	28.88	5.79
		Yes	3,289	73.61	23.53	2.86	0.00
Economic Disadvantage	No	32,378	32.22	30.99	30.42	6.37	
	Yes	10,902	51.21	29.99	16.48	2.32	
Special Education	No	39,341	33.90	31.28	29.00	5.81	
	Yes	3,939	67.99	25.26	5.97	0.79	

Appendix K: Principal Components Scree Plot

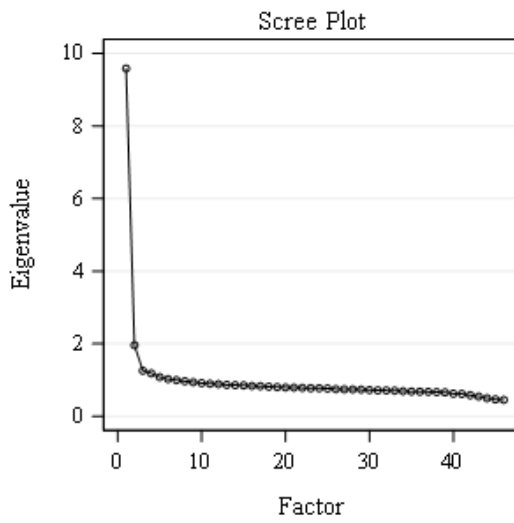


Figure K.1. English Grade 9 Principal Components Scree Plot

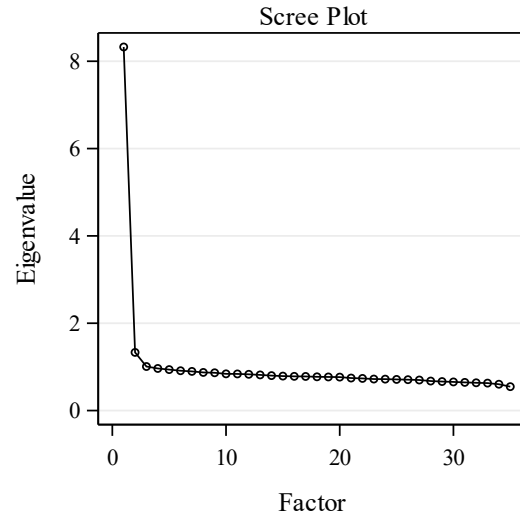


Figure K.3. Reading Grade 9 Principal Components Scree Plot

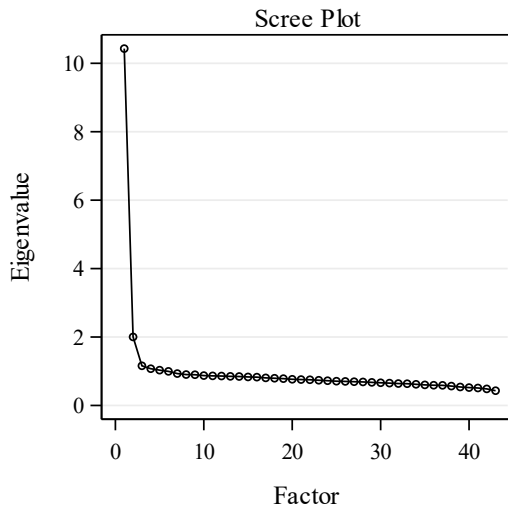


Figure K.2. English Grade 10 Principal Components Scree Plot

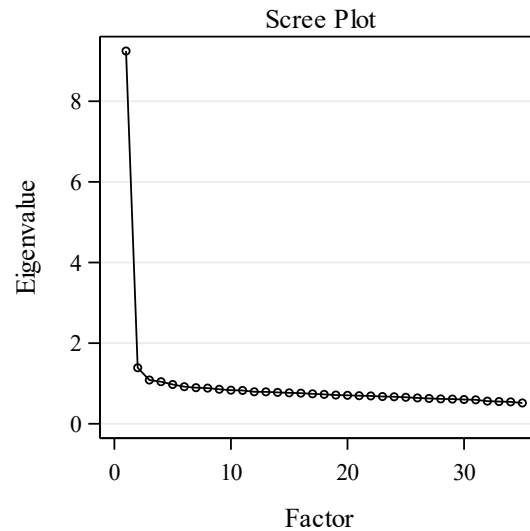


Figure K.4. Reading Grade 10 Principal Components Scree Plot

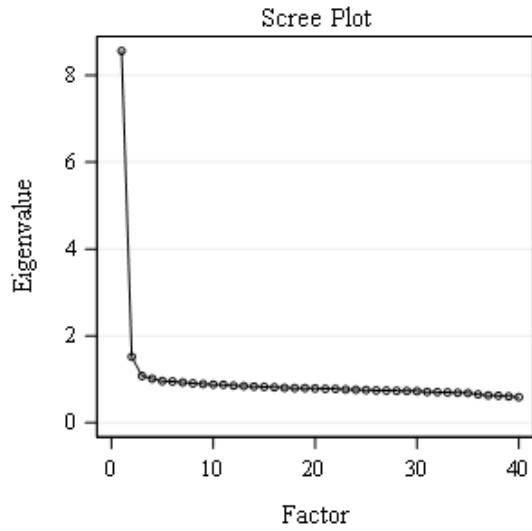


Figure K.5. Mathematics Grade 9 Principal Components Scree Plot

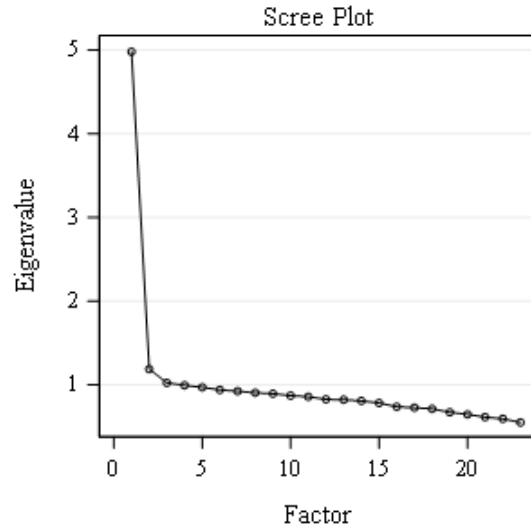


Figure K.7. Science Grade 9 Principal Components Scree Plot

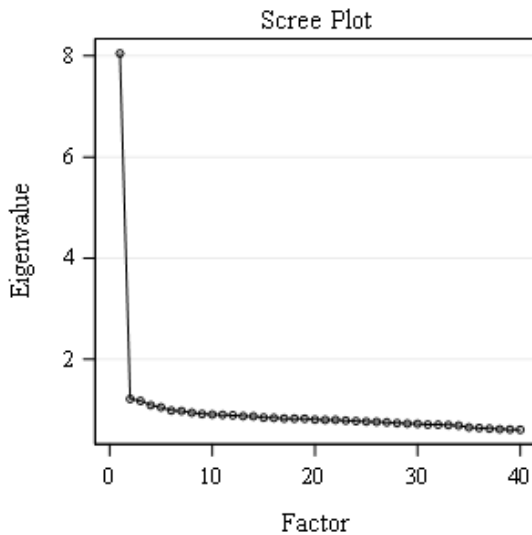


Figure K.6. Mathematics Grade 9 Principal Components Scree Plot

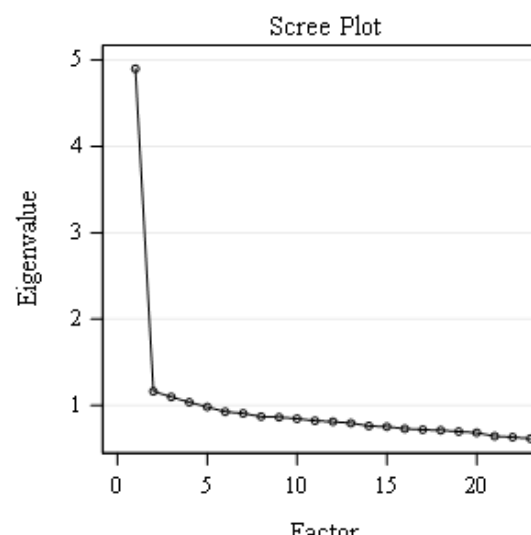


Figure K.8. Science Grade 10 Principal Components Scree Plot

Appendix L: Subscore Correlations

Table L.1. English Correlations of Total Score and Subscores

Grade	Subdomain	English Total	Conventions of Standard English	Knowledge of Language	Production of Writing
9	Total	1.00			
	Conventions of Standard English	0.95	1.00		
	Knowledge of Language	0.76	0.68	1.00	
	Production of Writing	0.88	0.79	0.64	1.00
10	Total	1.00			
	Conventions of Standard English	0.93	1.00		
	Knowledge of Language	0.85	0.79	1.00	
	Production of Writing	0.82	0.73	0.69	1.00

Table L.2. Reading Correlations of Total Score and Subscores

Grade	Subdomain	Reading Total	Key Ideas	Craft and Structure	Integration of Knowledge and Ideas
9	Total	1.00			
	Key Ideas	0.90	1.00		
	Craft and Structure	0.77	0.65	1.00	
	Integration of Knowledge and Ideas	0.92	0.79	0.68	1.00
10	Total	1.00			
	Key Ideas	0.89	1.00		
	Craft and Structure	0.71	0.62	1.00	
	Integration of Knowledge and Ideas	0.92	0.81	0.64	1.00

Table L.3. Mathematics Correlations of Total Score and Subscores

Grade	Subdomain	Math Total	Number and Quantity	Algebra	Functions	Geometry	Statistics and Probability
9	Total	1.00	—				
	Algebra	0.81	—	1.00			
	Functions	0.83	—	0.73	1.00		
	Geometry	0.83	—	0.69	0.70	1.00	
	Statistics and Probability	0.81	—	0.65	0.66	0.65	1.00
10	Total	1.00					
	Number and Quantity	0.77	1.00				
	Algebra	0.70	0.66	1.00			
	Functions	0.82	0.68	0.67	1.00		
	Geometry	0.70	0.58	0.55	0.58	1.00	
	Statistics and Probability	0.46	0.42	0.44	0.46	0.36	1.00

Table L.4. Science Correlations of Total Score and Subscores

Grade	Subdomain	Science Total	Gathering & Investigating	Developing Models	Using Mathematical Thinking	Construct Explanation
9	Total	1.00				
	Gathering & Investigating	0.74	1.00			
	Developing Models	0.76	0.52	1.00		
	Using Mathematical Thinking	0.75	0.50	0.51	1.00	
	Construct Explanation	0.74	0.50	0.52	0.59	1.00
10	Total	1.00				
	Gathering & Investigating	0.66	1.00			
	Developing Models	0.75	0.49	1.00		
	Using Mathematical Thinking	0.73	0.59	0.50	1.00	
	Construct Explanation	0.60	0.49	0.42	0.47	1.00

Appendix M: Item Drift

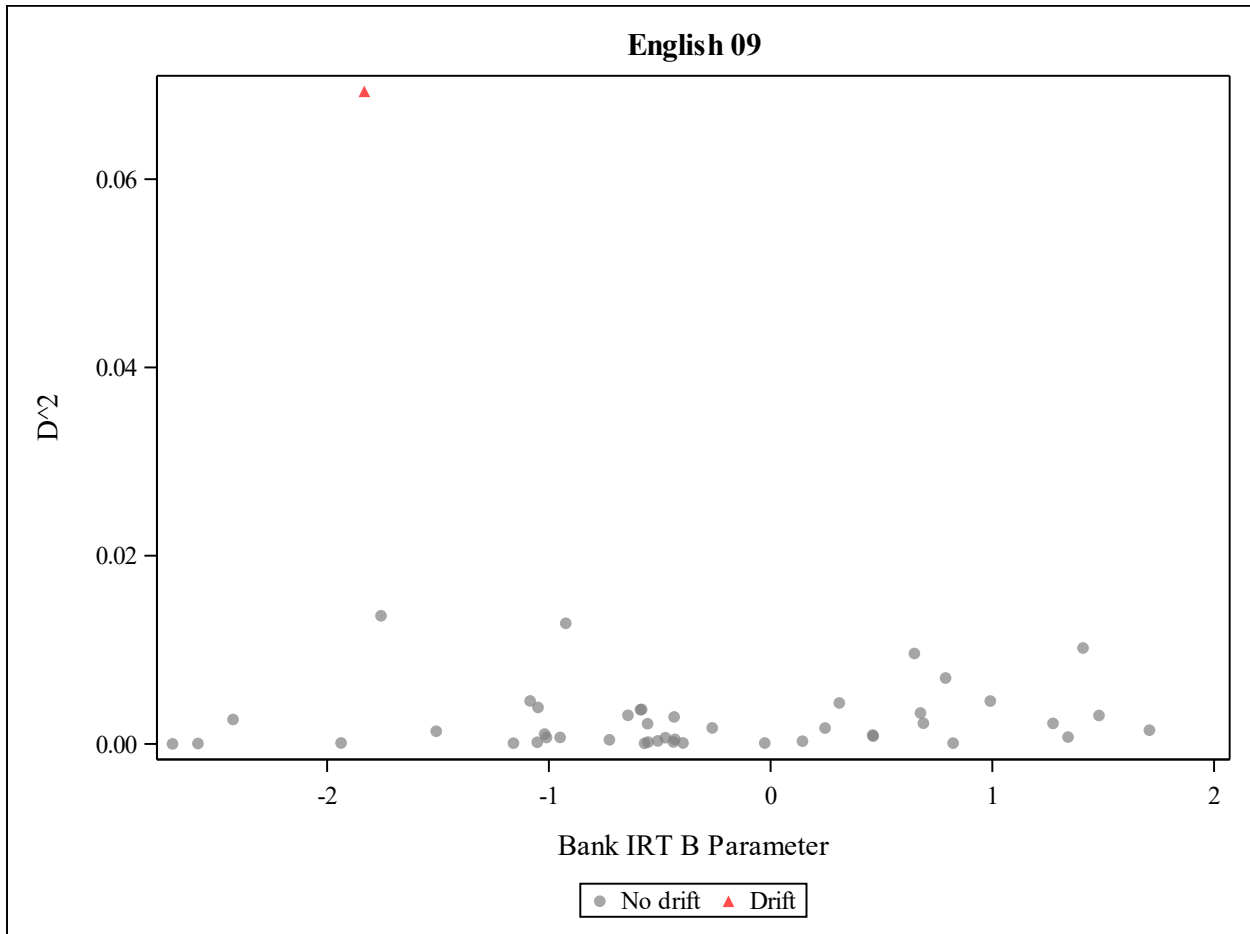


Figure M.1. English Grade 9 Item Drift

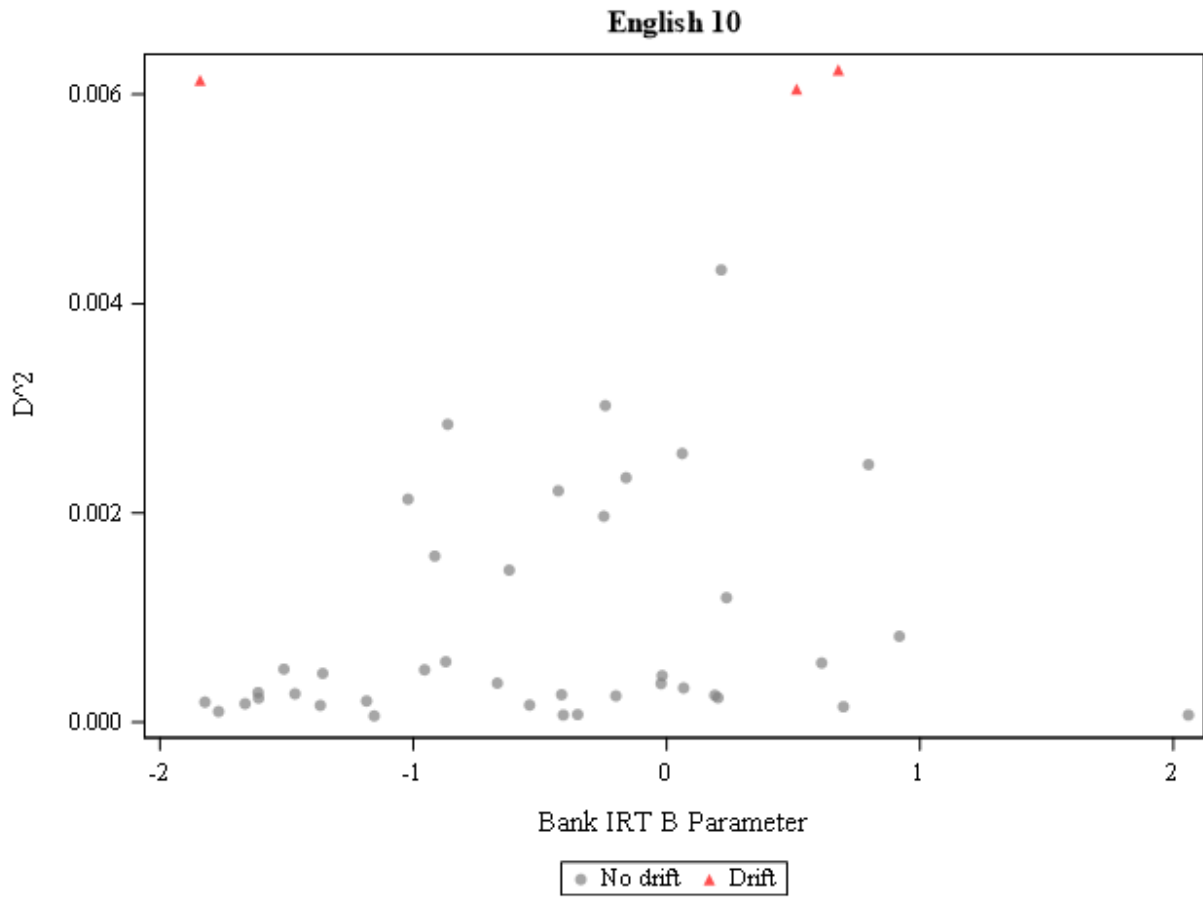


Figure M.2. English Grade 10 Item Drift

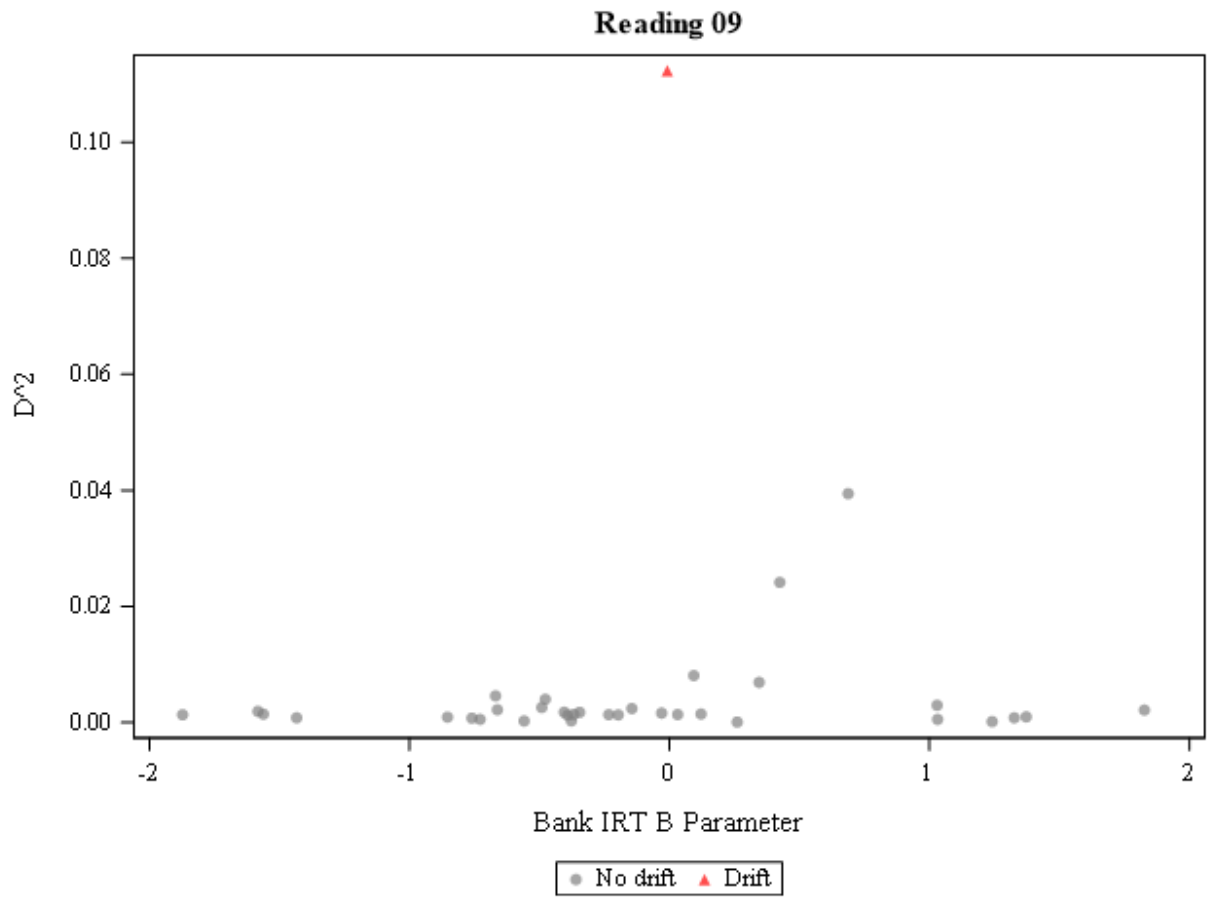


Figure M.3. Reading Grade 9 Item Drift

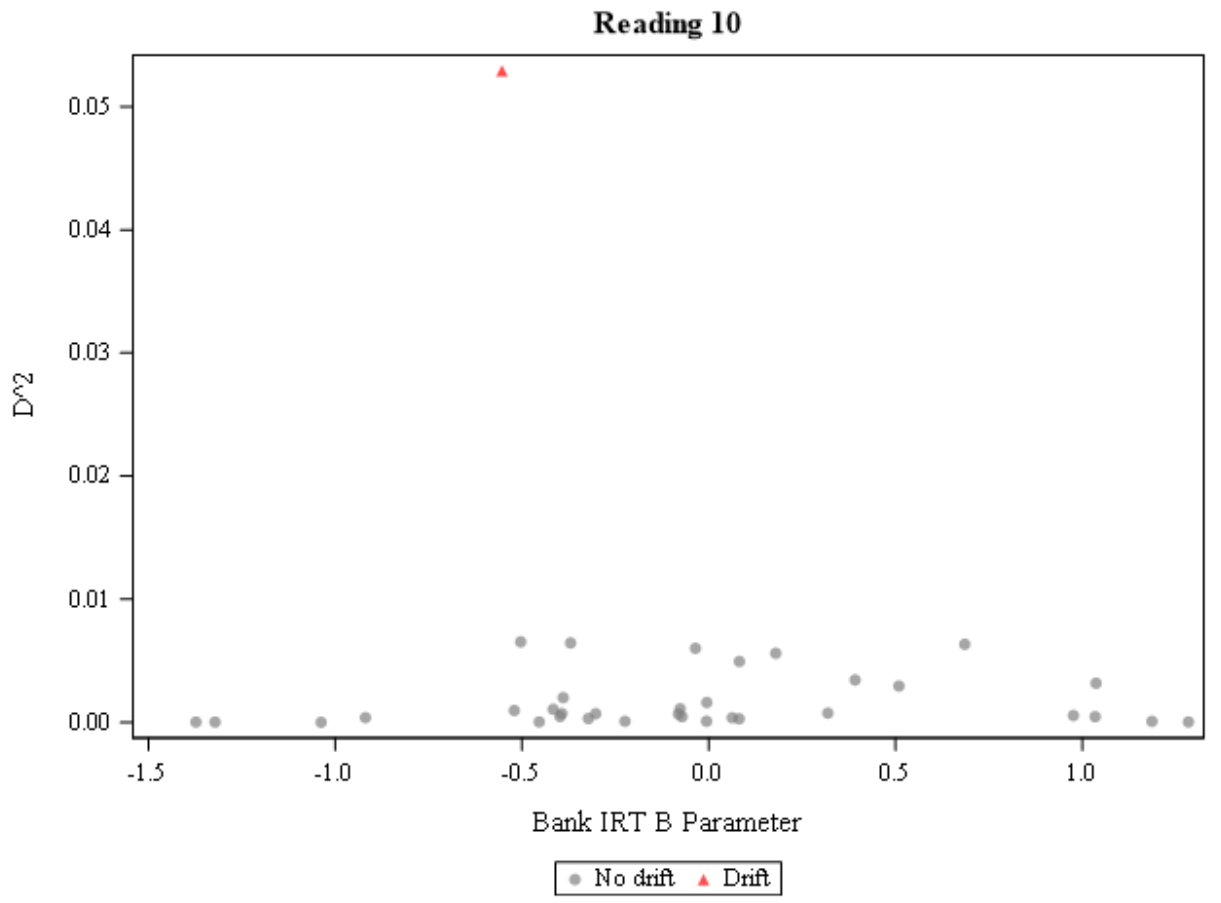


Figure M.4. Reading Grade 10 Item Drift

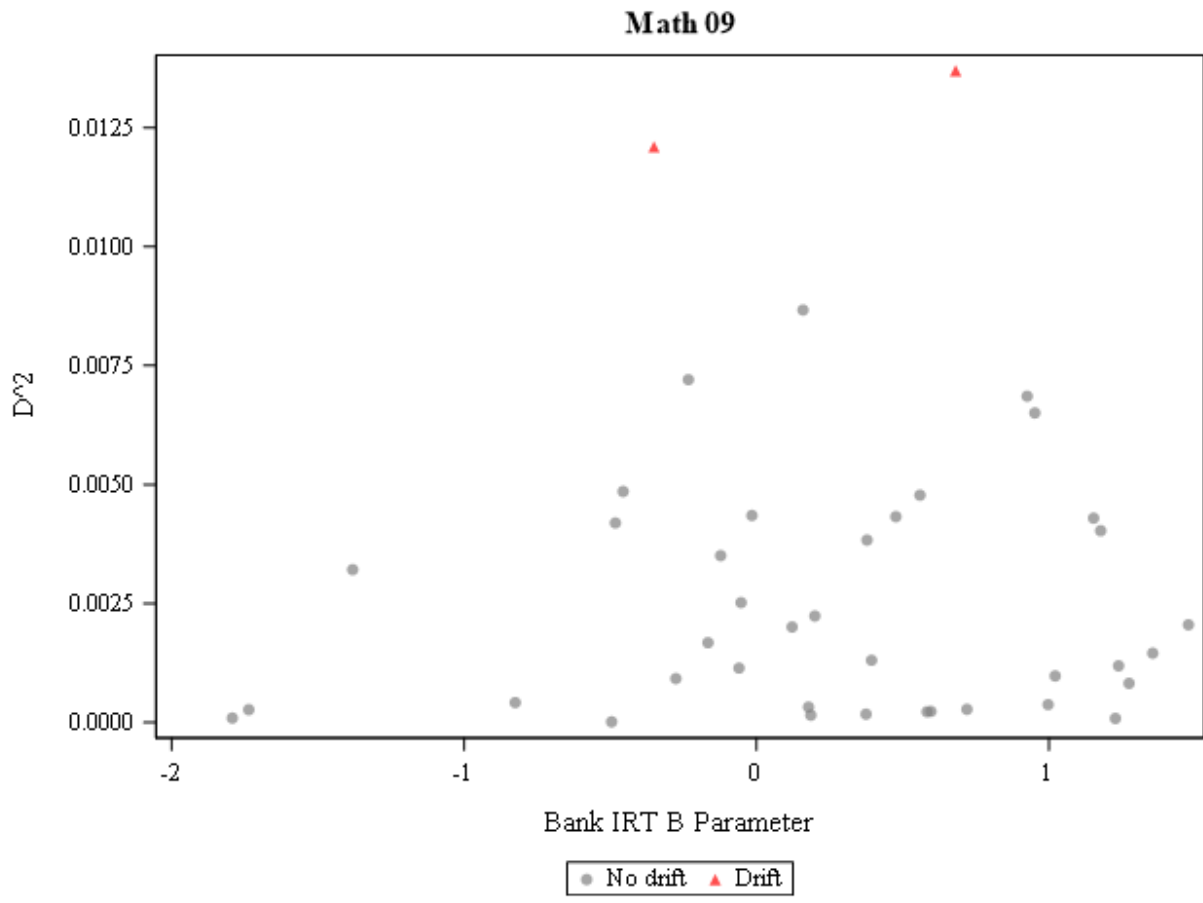


Figure M.5. Mathematics Grade 9 Item Drift

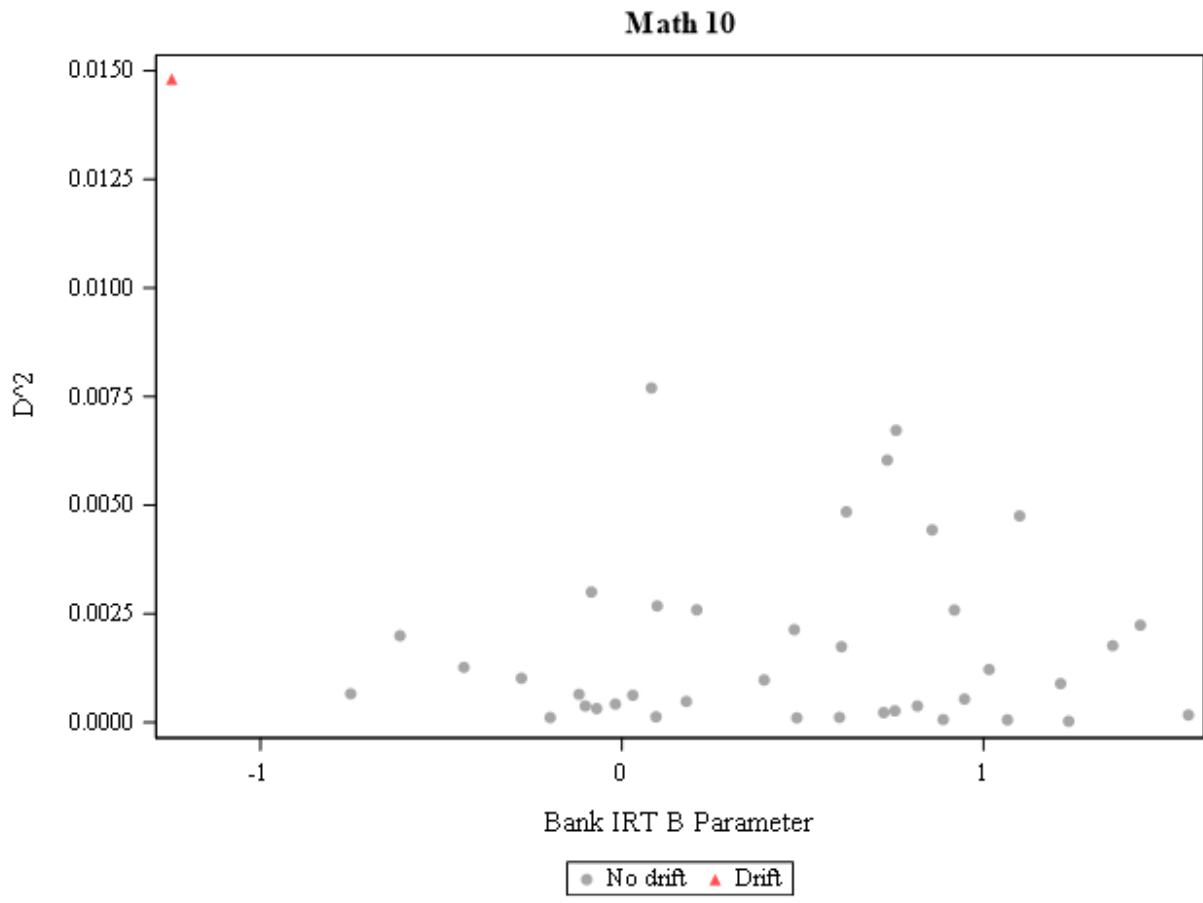


Figure M.6. Mathematics Grade 10 Item Drift

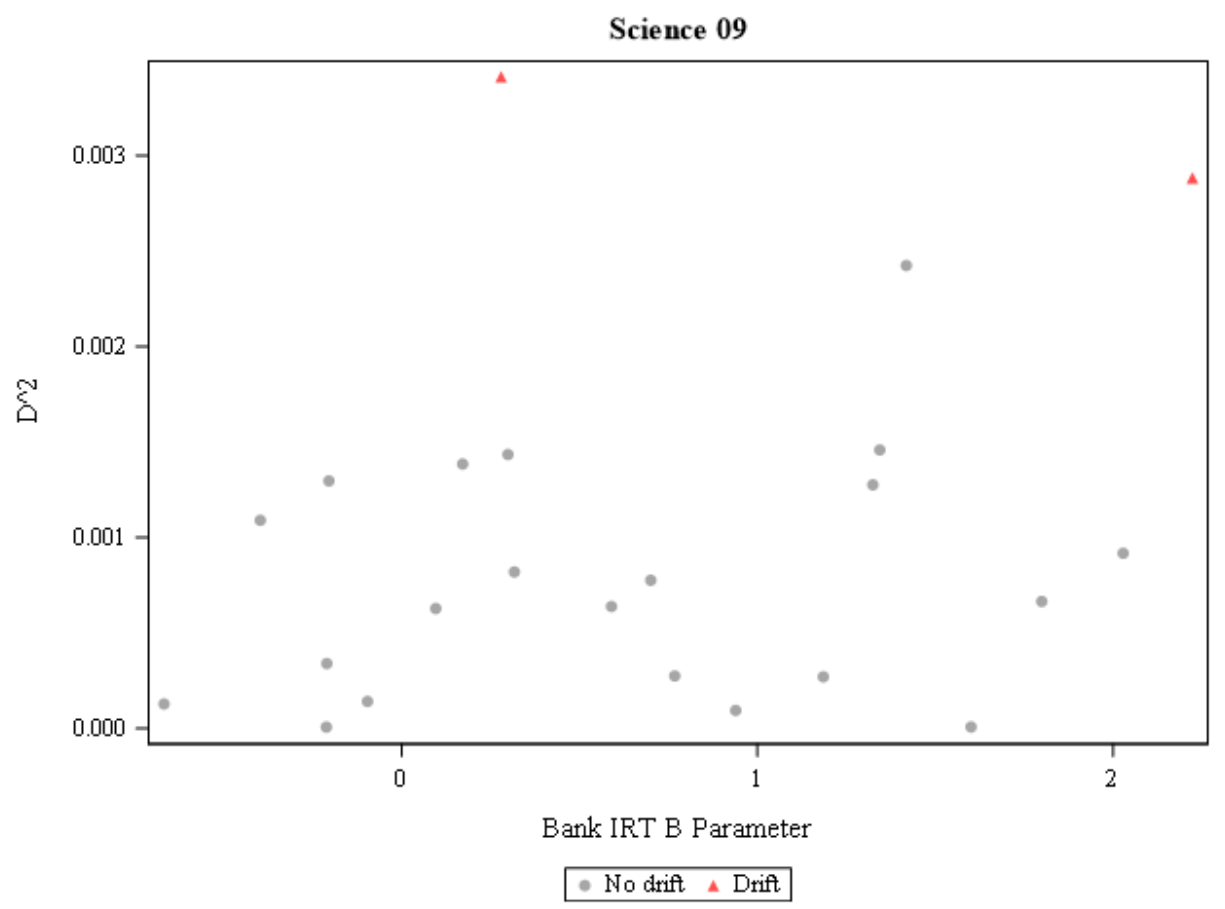


Figure M.7. Science Grade 9 Item Drift

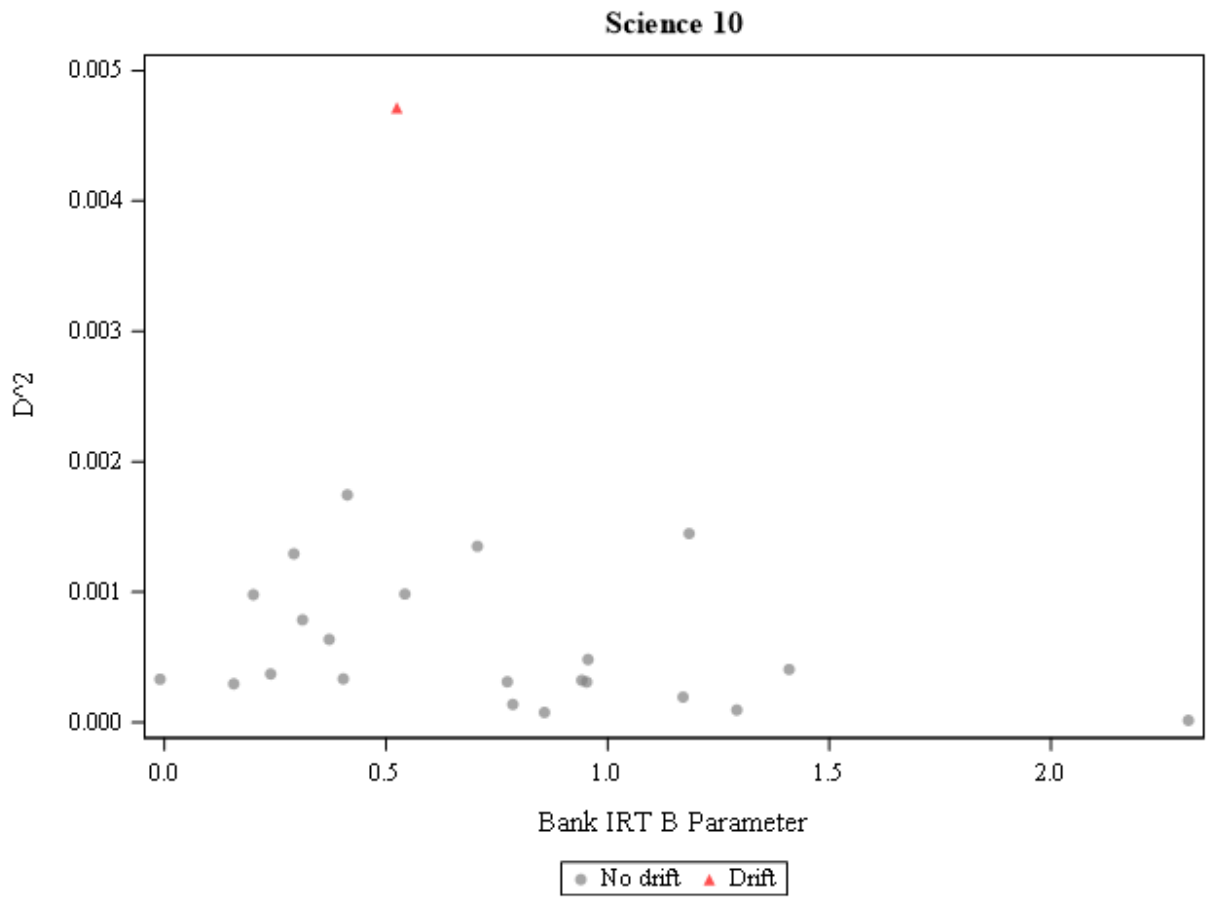


Figure M.8. Science Grade 10 Item Drift