





LEAP 2025 Science 3-8

Technical Report: 2023–2024

Prepared by DRC, Pearson, and WestEd





EXECUTIVE SUMMARY

The Louisiana Educational Assessment Program 2025 (LEAP 2025) is composed of tests that are carefully constructed to fairly assess the achievement of Louisiana students. This technical report provides information on the operational test administrations, scoring activities, analyses, and results of the spring 2024 administration of the LEAP 2025 Science tests that included both operational and field test items.

While this technical report and its associated materials have been produced in a way that can help educators understand the technical characteristics of the assessment used to measure student achievement, the information is primarily intended for use by those who evaluate tests, interpret scores, or use test results in making educational decisions. It is assumed that the reader has technical knowledge of test construction and measurement procedures, as stated the *Standards for Educational and Psychological Testing* (AERA et al., 2014).

The chapters of this technical report outline general information about the administration and scoring activities of the LEAP 2025 assessments, classical test theory (CTT) and item response theory (IRT) analysis results, 2024 test results, demographic characteristics of students, reliability and validity, and the interpretation of the scores on the tests. Additionally, because of conditions related to COVID-19, please use caution when making any inferences from the statistical results of the spring 2024 administration.

Table of Contents

EXECUTIVE SUMMARY	2
1. Introduction	7
Summary of the 2023–2024 Activities	7
2. Assessment Frameworks	9
3. Overview of the Test Development Process	11
Item Development Plan	11
Proposal and Review of Topics and Sources	30
Performance Expectation Bundling	30
Phenomena Selection and Outline Development	31
Matching Phenomena to Item Sets and Tasks and Foci to Standalone Items	32
Outline and Stimuli Development	34
Item Writing and Review Process	35
Data Review Process and Results	42
4. Construction of Test Forms with Embedded Field Test	46
Test Design	46
Initial Construction	52
Operational Form	52
Field Test Versions	55
Revision and Review	56

Psychometric Approval of Operational Forms	56
LDOE Review	57
Test Forms and Accessible Versions	57
Online and Paper Forms	57
Accommodated Print Versions	57
Form Versions for Students with Visual Impairments	58
5. Test Administration	59
Training of School Systems	59
Ancillary Materials	59
Return Material Forms and Guidelines	68
Security Checklists	69
Interpretive Guides	69
Time	69
Online Forms Administration, Grades 3–8	69
Paper-Based Forms Administration, Grade 3	70
Accessibility and Accommodations	70
Testing Windows	71
Test Security Procedures	71
Data Forensic Analyses	72
Alerts for Disturbing Content	73

6.	Scoring Activities	74
(Constructed-Response and Extended-Response Scoring	75
7.	Data Analysis	89
	Classical Item Statistics	89
	Differential Item Functioning	90
	Measurement Models	95
	Calibration and Linking	95
	Operational Item Parameters	103
	Item Fit	103
	Dimensionality and Local Item Independence	105
	Scaling	106
	Test Characteristic Curve	107
	Test Information Curve, Score Distribution, and IRT Difficulty Distribution	109
	Field Test Data Review	119
8.	Test Results and Score Reports	120
	Demographic Characteristics of Students	120
	Test Results	121
	Effect Size	128
	Score Reports	128
	Achievement Level Policy Definitions	131
	Timing Analysis	133
9.	Reliability	134
	Internal Consistency Reliability Estimation	134

	Classical Standard Error of Measurement	135
	Conditional Standard Error of Measurement and Cut Scores	136
	Student Classification Accuracy and Consistency	139
10.	Validity	141
	Evidence for Construct-Related Validity	142
	Internal Structure of Reporting Categories	142
	Content-Related Evidence	142
	Dimensionality and Principal Component Analysis	143
	Item Development and Field-Test Analysis	143
	Mode Effect Study	145
Ref	ferences	. 147
Apı	pendix A: Training Agendas	. 150
Apı	pendix B: Test Summary	. 169
Apı	pendix C: Item Analysis Summary Report	. 185
Apı	pendix D: Dimensionality	. 241
Apı	pendix E: Scale Distribution and Statistical Report	. 247
Apı	pendix F: Reliability and Classification Accuracy	260
Apı	pendix G: Accommodated Print and Braille Creation	.269
Apı	pendix H: On-Going Quality Control	. 272

1. Introduction

The Louisiana Department of Education (LDOE) has a long and distinguished history in the development and administration of assessments that support its state accountability system and are aligned to the Louisiana Student Standards. Per state law, the LDOE is to administer statewide science assessments in grades 3–8 and high school Biology annually. Fulfilling the directive of the Louisiana State Board of Elementary and Secondary Education (BESE), the LDOE must deliver high-quality, Louisiana-specific standards-based assessments. The LDOE and the BESE are committed to the development of rigorous assessments as one component of their comprehensive plan—Louisiana Believes—designed to ensure that every Louisiana student is on track to be successful in postsecondary education and the workforce.

The purpose of this technical report is to describe the process for the operational administration of the statewide summative science assessments for grades 3–8 as part of the Louisiana Educational Assessment Program 2025 (LEAP 2025). This report outlines the testing procedures, including forms construction, administration, statistical analyses, scoring and analyses, and reporting of scores.

Summary of the 2023–2024 Activities

WestEd and Pearson, in partnership with the LDOE and Data Recognition Corporation (DRC), the administration vendor, developed a timeline to capture the major activities necessary to produce the spring 2024 Science grades 3–8 operational forms with embedded field test items (EFT). All tests were delivered in a computer-based format, with a paper-based option for grade 3. An accommodated paper-based format was available for students in grades 4–8 who are not physically able to test on a computer. Table 1.1 summarizes the key activities along with the months during which the activities were completed.

Table 1.1

Key Activities from August 2022 to August 2024

Date	Activity
August– December 2022	 Started item development planning for the spring 2024 test The LDOE approved the item development plans, proposed bundles, and standalone item topics WestEd updated the content development specifications, style guides, and training materials WestEd developed outlines for the stimulus review committees and began standalone item development The Technical Advisory Committee (TAC) meeting convened
January– February 2023	 The LDOE convened stimulus review committees The LDOE provided feedback and approval to begin set/task development
March–May 2023	 WestEd led in item writing and development LDOE staff reviewed proposed item sets, tasks, and standalones
June 2023	 WestEd and the LDOE convened Item Content/Bias Review Committees onsite in Baton Rouge The LDOE and WestEd staff held reconciliation meetings
July–October 2023	 Content was finalized and the LDOE approved Online content delivered to administration vendor Frameworks were finalized and the LDOE approved Conducted data review Operational and field test forms were selected and the LDOE approved The LDOE, WestEd, and DRC met for planning meeting
November– December 2023	 November TAC convened Accommodated print/braille forms and alt text constructed, the LDOE approved, and delivered to administration vendor The LDOE and WestEd staff reviewed proposed spring 2024 EFT selections in administration platform
February 2024	 The TAC convened The LDOE, WestEd, and DRC met for planning meeting
April-May 2024	Spring 2024 test was administered, including EFT
August 2024	Data Review Held

2. Assessment Frameworks

The assessment framework addresses:

- the test designs,
- test blueprints,
- range of standards to be covered,
- reporting categories,
- percentages of assessment items and score points by reporting category,
- projected testing times,
- the numbers of forms to be administered, and
- select psychometric analysis activities.

Measuring student proficiency of the full depth and breadth of the Louisiana Student Standards for Science (LSSS) requires assessments built from a range of item types. The choice of a specific item type is a function of efficient and effective measurement of the target content. Multiple-choice (MC) and multiple-select (MS) item types provide students an opportunity to select the correct answer or answers from a set of choices. MS items can elicit a greater depth of understanding than traditional MC items by requiring the selection of more than one correct response, efficiently scored by an automated scoring engine. Constructed-response (CR) and extended-response (ER) items allow students to develop an explanation, describe a model, design a solution, and/or otherwise apply and communicate scientific understanding as required by the Science and Engineering Practices (SEPs) and Crosscutting Concepts (CCCs). These types of student-produced responses are handscored by teams of trained readers. Technology-enhanced (TE) items allow students to apply and communicate scientific knowledge and understanding as required by the SEPs and CCCs in ways that may not be addressed by MC or MS item types, but in a manner more cost-effective and less time-consuming than CR and ER item types with automated engine scoring. TE items may ask students to develop models or to sort processes by dragging components into a valid order, construct viable explanations by selecting words or phrases from several drop-down menus, or complete other tasks. The complexity of the TE items reduces the probability of randomly guessing the correct answer. Two-part items involve the application of understanding different but related knowledge to a concept or supporting assertions with evidence.

For two-part items, students may construct an explanation and support the explanation with evidence or make a claim and evaluate evidence to support the claim. Another application of two-part items is to develop a model in part A and to evaluate the model in part B. A range of item types and applications allows greater test-taker engagement and provides a more authentic assessment experience.

The test design includes item sets, a task, and standalone items. A stimulus that describes a scientific phenomenon anchors each item set or task. A focus that details some aspects of a phenomenon provides the common anchor for standalone items. Item sets are composed of four items associated with a common stimulus. The item sets may include 1point selected-response items (single-select and/or MS formats), 1- and 2-point TE items, and 2-point two-part items (two-part independent [TPI] and/or two-part dependent [TPD] formats) tied to a common stimulus. For grades 5–8, item sets may include 1- or 2-point TE items. Three item sets include a two-point CR item. The assessment also includes one task. The task consists of five items tied to a common stimulus and includes 1-point selected-response items (both single-select and MS formats), 2-point two-part items (TPI and/or TPD formats), and a 9-point ER item for grades 5–8. The standalone items provide flexibility to meet the test blueprint and afford greater coverage of the standards while still requiring students to make connections among the three dimensions of the LSSS. All points associated with the task contribute to a student's overall score, but the ER item is not a component of the current blueprint and therefore not included in the proportional representation of content assessed by other parts of the test.

Because the assessment at grade 3 was administered primarily via paper, the item types were limited to selected-response (i.e., MC and MS), two-part (i.e., TPI and/or or TPD), and CR items. Assessments for grades 4–8 were administered primarily online, so TE items were viable at these grades. However, paper and pencil versions of the assessments for grades 4–8 were made available as accommodated forms for students who were unable to test online. For those forms, TE items were adapted for paper presentation to address the same content.

The Assessment Frameworks were reviewed by the LDOE content and psychometric staff to ensure that the test designs, blueprints, and form designs met the necessary content, reporting, and psychometric requirements.

3. Overview of the Test Development Process

Item Development Plan

Table 3.1a presents the acronyms used in item and test development.

Table 3.1a

Acronyms Used in Item and Test Development

Acronvm	Meaning
ARG	Engaging in Argument from Evidence
CCC	Crosscutting Concepts
C/E	Cause and Effect
DATA	Analyzing and Interpreting Data
DCI	Disciplinary Core Ideas
E/M	Energy and Matter
E/S	Constructing Explanations and Designing Solutions
INFO	Obtaining, Evaluating, and Communicating Information
INV	Planning and Carrying Out Investigations
LEAP	Louisiana Educational Assessment Program
LS	Life Science
LSSS	Louisiana Student Standards for Science
MCT	Using Mathematics and Computational Thinking
MOD	Developing and Using Models
PAT	Patterns
PE	Performance Expectation
Q/P	Asking Questions and Defining Problems

Table 3.1a

Acronyms Used in Item and Test Development (continued)

Acronym	Meaning	
S/C	Stability and Change	
SEP	Science and Engineering Practices	
S/F	Structure and Function	
SPQ	Scale, Proportion, and Quantity	
SYS	Systems and System Models	

The test blueprints that guided item development projections for grade 3 are presented in Tables 3.1b-3.1g.

Table 3.1b

Test Blueprint for LEAP 2025 Grade 3: DCI Domain Coverage

Grade 3: DCI Domain Coverage				
# of PEs in LSSS Relative % in LSSS % by Points of All It				
ESS	3	20%	15%–25%	
LS	8	53%	48%–58%	
PS	4	27%	22%–32%	
Total	15	100%		

Table 3.1c

Test Blueprint for LEAP 2025 Grade 3: Minimal PE Coverage

Grade 3: Minimal PE Coverage						
SEP CCC Min Items						
03-ESS2-1	SEP 4 – DATA	CCC 1 – PAT	1			
03-ESS2-2	SEP 8 – INFO	CCC 1 – PAT	1			
03-ESS3-1	SEP 7 – ARG	CCC 2 – C/E	1			
03-LS1-1	SEP 2 – MOD	CCC 1 – PAT	1			
03-LS2-1	SEP 7 – ARG	CCC 4 – SYS	1			
03-LS3-1	SEP 4 – DATA	CCC 1 – PAT	1			
03-LS3-2	SEP 6 – E/S	CCC 2 – C/E	1			
03-LS4-1	SEP 4 – DATA	CCC 3 – SPQ	1			
03-LS4-2	SEP 6 – E/S	CCC 2 – C/E	1			
03-LS4-3	SEP 7 – ARG	CCC 2 – C/E	1			
03-LS4-4	SEP 7 – ARG	CCC 4 – SYS	1			
03-PS2-1	SEP 3 – INV	CCC 2 – C/E	1			
03-PS2-2	SEP 3 – INV	CCC 1 – PAT	1			
03-PS2-3	SEP 1 – Q/P	CCC 2 – C/E	1			
03-PS2-4	SEP 1 – Q/P	CCC 1 – PAT	1			

Table 3.1d

Test Blueprint for LEAP 2025 Grade 3: CCC Coverage

Grade 3: CCC Coverage				
CCC Overall # in PEs in LSSS Relative % in LSSS		% by Points of CCC Items		
CCC 1 – PAT	6	40%	35%-45%	
CCC 2 – C/E	6	40%	35%-45%	
CCC 3 – SPQ	1	7%	5%-15%	
CCC 4 – SYS	2	13%	8%–18%	
CCC 5 – E/M	0	0%	0%	
CCC 6 –S/F	0	0%	0%	
CCC 7 – S/C	0	0%	0%	
Total	15	100%		

Table 3.1e
Test Blueprint for LEAP 2025 Grade 3: SEP Coverage

Grade 3: SEP Coverage				
SEP Overall	# in PEs in LSSS Relative % in LSSS		% by Points of SEP Items	
SEP 1 – Q/P	2	13%	8%-18%	
SEP 2 – MOD	1	7%	5%-15%	
SEP 3 – INV	2	13%	8%-20%	
SEP 4 – DATA	3	20%	15%-25%	
SEP 5 – MCT	0	0%	0%	
SEP 6 – E/S	2	13%	8%-18%	
SEP 7 – ARG	4	27%	22%-32%	
SEP 8 – INFO	1	7%	5%-15%	
Total	15	100%		

Table 3.1f
Test Blueprint for LEAP 2025 Grade 3: SEP Reporting Category Coverage

Grade 3: SEP reporting category Coverage						
Reporting Category	# PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	Min Points		
Reporting Category 1 (SEPs 1 & 3)	4	29%	24%-34%	7		
Reporting Category 2 (SEPs 4, 5, 7)	7	50%	45%-55%	7		
Reporting Category 3 (SEPs 2 & 6)	3	21%	16%-26%	7		
Total	14	100%				

Note: SEP 8 (Obtaining, evaluating, and communicating information) is assumed to be embedded within each reporting category (1–3), so SEP 8 is not being repeated across the reporting categories.

Table 3.1g
Test Blueprint for LEAP 2025 Grade 3: SEP Compared to CCC Ratio

Grade 3: SEP Compared to CCC Ratio				
	Relative Weight in LSSS	Minimum %		
SEPs	50%	30%		
CCCs	50%	30%		

The test blueprints that guided item development projections for grade 4 are presented in Tables 3.1h-3.1m.

Table 3.1h

Test Blueprint for LEAP 2025 Grade 4: DCI Domain Coverage

Grade 4: DCI Domain Coverage				
Domain	# of PEs in LSSS	Relative % in LSSS	% by Points of All Items	
ESS	6	43%	38%-48%	
LS	2	14%	9%–19%	
PS	6	43%	38%-48%	
Total	14	100%		

Table 3.1i
Test Blueprint for LEAP 2025 Grade 4: Minimal PE Coverage

Grade 4: Minimal PE Coverage Every PE will be included at least one time in a test				
PE	SEP	CCC	Min Items	
04-ESS1-1	SEP 6 – E/S	CCC 1 – PAT	1	
04-ESS2-1	SEP 3 – INV	CCC 2 – C/E	1	
04-ESS2-2	SEP 4 – DATA	CCC 1 – PAT	1	
04-ESS2-3	SEP 1 – Q/P	CCC 2 – C/E	1	
04-ESS3-1	SEP 8 – INFO	CCC 2 – C/E	1	
04-ESS3-2	SEP 6 – E/S	CCC 2 – C/E	1	
04-LS1-1	SEP 7 – ARG	CCC 4 – SYS	1	
04-LS1-2	SEP 6 – E/S	CCC 2 – C/E	1	
04-PS3-1	SEP 6 – E/S	CCC 5 – E/M	1	
04-PS3-2	SEP 3 – INV	CCC 5 – E/M	1	
04-PS3-3	SEP 1 – Q/P	CCC 5 – E/M	1	
04-PS3-4	SEP 6 – E/S	CCC 5 – E/M	1	
04-PS4-1	SEP 2 – MOD	CCC 1 - PAT	1	
04-PS4-2	SEP 2 – MOD	CCC 2 - C/E	1	

Table 3.1j
Test Blueprint for LEAP 2025 Grade 4: CCC Coverage

Grade 4: CCC Coverage				
CCC Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of CCC Items	
CCC 1 – PAT	3	21%	16%-26%	
CCC 2 – C/E	6	43%	38%-48%	
CCC 3 – SPQ	0	0%	0%	
CCC 4 – SYS	1	7%	5%-15%	
CCC 5 – E/M	4	29%	24%-34%	
CCC 6 – S/F	0	0%	0%	
CCC 7 – S/C	0	0%	0%	
Total	14	100%		

Table 3.1k *Test Blueprint for LEAP 2025 Grade 4: SEP Coverage*

Grade 4: SEP Coverage				
SEP Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	
SEP 1 – Q/P	2	14%	9%-19%	
SEP 2 – MOD	2	14%	9%–19%	
SEP 3 – INV	2	14%	9%–19%	
SEP 4 – DATA	1	7%	5%-15%	
SEP 5 – MCT	0	0%	0%	
SEP 6 – E/S	5	36%	31%-41%	
SEP 7 – ARG	1	7%	5%-15%	
SEP 8 – INFO	1	7%	5%-15%	
Total	14	100%		

Table 3.11
Test Blueprint for LEAP 2025 Grade 4: SEP Reporting Category Coverage

Grade 4: SEP Reporting Category Coverage					
SEP Reporting Category	# PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	Min Points	
Reporting Category 1 (SEPs 1 & 3)	4	31%	26%-36%	7	
Reporting Category 2 (SEPs 4, 5, 7)	2	15%	10%-20%	7	
Reporting Category 3 (SEPs 2 & 6)	7	54%	49%-59%	7	
Total	13	100%			

Note: SEP 8 (Obtaining, evaluating, and communicating information) is assumed to be embedded within each reporting category (1–3), so SEP 8 is not being repeated across the reporting category.

Table 3.1m

Test Blueprint for LEAP 2025 Grade 4: SEP Compared to CCC Ratio

Grade 4: SEP Compared to CCC Ratio				
	Relative Weight in LSSS	Minimum %		
SEPs	50%	30%		
CCCs	50%	30%		

The test blueprints that guided item development projections for grade 5 are presented in Tables 3.1n-3.1s.

Table 3.1n
Test Blueprint for LEAP 2025 Grade 5: DCI Domain Coverage

Grade 5: DCI Domain Coverage				
Domain	# of PEs in LSSS	Relative % in LSSS	% by Points of All Items	
ESS	5	38%	33%-43%	
LS	2	15%	9%-20%	
PS	6	46%	41%-51%	
Total	13	100%		

Table 3.10
Test Blueprint for LEAP 2025 Grade 5: Minimal PE Coverage

Grade 5: Minimal PE Coverage Every PE will be included at least one time in a test				
PE	SEP	CCC	Min Items	
05-ESS1-1	SEP 7 – ARG	CCC 3 – SPQ	1	
05-ESS1-2	SEP 4 – DATA	CCC 1 – PAT	1	
05-ESS2-1	SEP 2 – MOD	CCC 4 – SYS	1	
05-ESS2-2	SEP 5 – MCT	CCC 3 – SPQ	1	
05-ESS3-1	SEP 6 – E/S	CCC 4 – SYS	1	
05-LS1-1	SEP 1 – Q/P	CCC 5 – E/M	1	
05-LS2-1	SEP 2 – MOD	CCC 4 – SYS	1	
05-PS1-1	SEP 2 – MOD	CCC 3 – SPQ	1	
05-PS1-2	SEP 5 – MCT	CCC 5 – E/M	1	
05-PS1-3	SEP 3 – INV	CCC 3 – SPQ	1	
05-PS1-4	SEP 3 – INV	CCC 2 – C/E	1	
05-PS2-1	SEP 7 – ARG	CCC 2 – C/E	1	
05-PS3-1	SEP 2 – MOD	CCC 5 – E/M	1	

Table 3.1p

Test Blueprint for LEAP 2025 Grade 5: CCC Coverage

Grade 5: CCC Coverage				
CCC Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of CCC Items	
CCC 1 – PAT	1	8%	5%-15%	
CCC 2 – C/E	2	15%	9%-22%	
CCC 3 – SPQ	4	31%	22%-36%	
CCC 4 – SYS	3	23%	18%-28%	
CCC 5 – E/M	3	23%	18%–28%	
CCC 6 – S/F	0	0%	0%	
CCC 7 – S/C	0	0%	0%	
Total	13	100%		

Table 3.1q
Test Blueprint for LEAP 2025 Grade 5: SEP Coverage

Grade 5: SEP Coverage				
SEP Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	
SEP 1 – Q/P	1	8%	3%-13%	
SEP 2 – MOD	4	31%	26%-36%	
SEP 3 – INV	2	15%	10%–20%	
SEP 4 – DATA	1	8%	3%-13%	
SEP 5 – MCT	2	15%	10%-20%	
SEP 6 –E/S	1	8%	3%-15%	
SEP 7 – ARG	2	15%	10%–20%	
SEP 8 – INFO	0	0%	0%	
Total	13	100%		

Table 3.1r
Test Blueprint for LEAP 2025 Grade 5: SEP Reporting Category Coverage

Grade 5: SEP Reporting Category Coverage					
	# of PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	Min Points	
Reporting Category 1 (SEPs 1 & 3)	3	23%	18%-28%	7	
Reporting Category 2 (SEPs 4, 5, 7)	5	38%	32%-43%	7	
Reporting Category 3 (SEPs 2 & 6)	5	38%	33%-43%	7	
Total	13	100%			

Note: SEP 8 (Obtaining, evaluating, and communicating information) is assumed to be embedded within each reporting category (1–3), so SEP 8 is not being repeated across the reporting categories.

Table 3.1s

Test Blueprint for LEAP 2025 Grade 5: SEP Compared to CCC Ratio

Grade 5: SEP Compared to CCC Ratio			
	Relative Weight in LSSS	Minimum %	
SEPs	50%	30%	
CCCs	50%	30%	

The test blueprints that guided item development projections for grade 6 are presented in Tables 3.1t-3.1y.

Table 3.1t

Test Blueprint for LEAP 2025 Grade 6: DCI Domain Coverage

Grade 6: DCI Domain Coverage				
Domain	# of PEs in LSSS	Relative % in LSSS	% by Points of All Items	
ESS	4	21%	15–26%	
LS	5	26%	21%-31%	
PS	10	53%	48%-58%	
Total	19	100%		

Table 3.1u
Test Blueprint for LEAP 2025 Grade 6: Minimal PE Coverage

Grade 6: Minimal PE Coverage Every PE will be included at least one time in a test				
PE	SEP	CCC	Min Items	
06-MS-ESS1-1	SEP 2 – MOD	CCC 1 – PAT	1	
06-MS-ESS1-2	SEP 2 – MOD	CCC 4 – SYS	1	
06-MS-ESS1-3	SEP 4 – DATA	CCC 3 – SPQ	1	
06-MS-ESS3-4	SEP 7 – ARG	CCC 2 – C/E	1	
06-MS-LS1-1	SEP 3 – INV	CCC 3 – SPQ	1	
06-MS-LS1-2	SEP 2 – MOD	CCC 6 – S/F	1	
06-MS-LS2-1	SEP 4 – DATA	CCC 2 – C/E	1	
06-MS-LS2-2	SEP 6 – E/S	CCC 1 – PAT	1	
06-MS-LS2-3	SEP 2 – MOD	CCC 5 – E/M	1	
06-MS-PS1-1	SEP 2 – MOD	CCC 3 – SPQ	1	
06-MS-PS2-1	SEP 6 – E/S	CCC 4 – SYS	1	
06-MS-PS2-2	SEP 3 – INV	CCC 7 – S/C	1	
06-MS-PS2-3	SEP 1 – Q/P	CCC 2 – C/E	1	
06-MS-PS2-4	SEP 7 – ARG	CCC 4 – SYS	1	
06-MS-PS2-5	SEP 3 – INV	CCC 2 – C/E	1	
06-MS-PS4-2	SEP 2 – MOD	CCC 6 – S/F	1	
06-MS-PS3-1	SEP 4 – DATA	CCC 3 – SPQ	1	
06-MS-PS3-2	SEP 2 – MOD	CCC 4 – SYS	1	
06-MS-PS4-1	SEP 5 – MCT	CCC 1 – PAT	1	

Table 3.1v *Test Blueprint for LEAP 2025 Grade 6: CCC Coverage*

Grade 6: CCC Coverage				
CCC Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of CCC Items	
CCC 1 – PAT	3	16%	11%–21%	
CCC 2 – C/E	4	21%	16%–26%	
CCC 3 – SPQ	4	21%	16%–26%	
CCC 4 – SYS	4	21%	16%–26%	
CCC 5 – E/M	1	5%	5–10%	
CCC 6 – S/F	2	11%	6–16%	
CCC 7 – S/C	1	5%	5–10%	
Total	19	100%		

Table 3.1w
Test Blueprint for LEAP 2025 Grade 6: SEP Coverage

Grade 6: SEP Coverage				
SEP Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	
SEP 1 – Q/P	1	5%	5%-10%	
SEP 2 – MOD	7	37%	32%-42%	
SEP 3 – INV	3	16%	11%–21%	
SEP 4 – DATA	3	16%	11%–21%	
SEP 5 – MCT	1	5%	5%-10%	
SEP 6 – E/S	2	11%	5%-16%	
SEP 7 – ARG	2	11%	5%-16%	
SEP 8 – INFO	0	0%	0%	
Total	19	100%		

Table 3.1x
Test Blueprint for LEAP 2025 Grade 6: SEP Reporting Category Coverage

Grade 6: SEP Reporting Category Coverage						
SEP Reporting Category	# PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	Min Points		
Reporting Category 1 (SEPs 1 & 3)	4	21%	16%-26%	7		
Reporting Category 2 (SEPs 4, 5, 7)	6	32%	27%-37%	7		
Reporting Category 3 (SEPs 2 & 6)	9	47%	42%-52%	7		
Total	19	100%				

Note: SEP 8 (Obtaining, evaluating, and communicating information) is assumed to be embedded within each reporting category (1–3), so SEP 8 is not being repeated across the reporting categories.

Table 3.1y
Test Blueprint for LEAP 2025 Grade 6: SEP Compared to CCC Ratio

Grade 6: SEP Compared to CCC Ratio				
Relative Weight in LSSS Minimum %				
SEPs	50%	30%		
CCCs	50%	30%		

The test blueprints that guided item development projections for grade 7 are presented in Tables3.1z-3.1ee.

Table 3.1z
Test Blueprint for LEAP 2025 Grade 7: DCI Domain Coverage

Grade 7: DCI Domain Coverage				
Domain	# of PEs in LSSS	Relative % in LSSS	% by Points of All Items	
ESS	4	25%	20%-35%	
LS	8	50%	45%-55%	
PS	4	25%	20%-35%	
Total	16	100%		

Table 3.1aa
Test Blueprint for LEAP 2025 Grade 7: Minimal PE Coverage

Grade 7: Minimal PE Coverage Every PE will be included at least one time in a test				
PE	SEP	CCC	Min Items	
07-MS-ESS2-4	SEP 2 – MOD	CCC 5 – E/M	1	
07-MS-ESS2-5	SEP 3 – INV	CCC 2 – C/E	1	
07-MS-ESS2-6	SEP 2 – MOD	CCC 4 – SYS	1	
07-MS-ESS3-5	SEP 1 – Q/P	CCC 7 – S/C	1	
07-MS-LS1-3	SEP 7 – ARG	CCC 4 – SYS	1	
07-MS-LS1-6	SEP 6 – E/S	CCC 5 – E/M	1	
07-MS-LS1-7	SEP 2 – MOD	CCC 5 – E/M	1	
07-MS-LS2-4	SEP 7 – ARG	CCC 7 – S/C	1	
07-MS-LS2-5	SEP 6 – E/S	CCC 7 – S/C	1	
07-MS-LS3-2	SEP 2 – MOD	CCC 2 – C/E	1	
07-MS-LS4-4	SEP 6 – E/S	CCC 2 – C/E	1	
07-MS-LS4-5	SEP 8 – INFO	CCC 2 – C/E	1	
07-MS-PS1-2	SEP 4 – DATA	CCC 1 – PAT	1	
07-MS-PS1-4	SEP 2 – MOD	CCC 2 – C/E	1	
07-MS-PS1-5	SEP 2 – MOD	CCC 5 – E/M	1	
07-MS-PS3-4	SEP 3 – INV	CCC 3 – SPQ	1	

Table 3.1bb

Test Blueprint for LEAP 2025 Grade 7: CCC Coverage

Grade 7: CCC Coverage				
CCC Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of CCC Items	
CCC 1 – PAT	1	6%	1%–11%	
CCC 2 – C/E	5	31%	20%–36%	
CCC 3 – SPQ	1	6%	1%-11%	
CCC 4 – SYS	2	13%	8%-18%	
CCC 5 – E/M	4	25%	20%-32%	
CCC 6 – S/F	0	0%	0%	
CCC 7 – S/C	3	19%	14%-24%	
Total	16	100%		

Table 3.1cc
Test Blueprint for LEAP 2025 Grade 7: SEP Coverage

Grade 7: SEP Coverage				
SEP Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	
SEP 1 – Q/P	1	6%	5%-15%	
SEP 2 – MOD	6	38%	33%-43%	
SEP 3 – INV	2	13%	8%-18%	
SEP 4 – DATA	1	6%	5%-15%	
SEP 5 – MCT	0	0%	0%	
SEP 6 – E/S	3	19%	14%-24%	
SEP 7 – ARG	2	13%	8%-18%	
SEP 8 – INFO	1	6%	5%-15%	
Total	16	100%		

Table 3.1dd

Test Blueprint for LEAP 2025 Grade 7: SEP Reporting Category Coverage

Grade 7: SEP Reporting Category Coverage					
SEP Reporting Category	# PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	Min Points	
Reporting Category 1 (SEPs 1 & 3)	3	20%	15%-25%	7	
Reporting Category 2 (SEPs 4, 5, 7)	3	20%	15%-25%	7	
Reporting Category 3 (SEPs 2 & 6)	9	60%	55%-65%	7	
Total	15	100%			

Note: SEP 8 (Obtaining, evaluating, and communicating information) is assumed to be embedded within each reporting category (1–3), so SEP 8 is not being repeated across the reporting categories.

Table 3.1ee
Test Blueprint for LEAP 2025 Grade 7: SEP Compared to CCC Ratio

Grade 7: SEP Compared to CCC Ratio					
Relative Weight in LSSS Minimum %					
SEPs	50%	30%			
CCCs	50%	30%			

The test blueprints that guided item development projections for grade 8 are presented in Tables 3.1ff-3.1kk.

Table 3.1ff
Test Blueprint for LEAP 2025 Grade 8: DCI Domain Coverage

Grade 8: DCI Domain Coverage					
Domain	# of PEs in LSSS	Relative % in LSSS	% by Points of All Items		
ESS	7	37%	32%-42%		
LS	7	37%	32%-42%		
PS	5	26%	21%-31%		
Total	19	100%			

Table 3.1gg *Test Blueprint for LEAP 2025 Grade 8: Minimal PE Coverage*

Grade 8: Minimal PE Coverage Everv PE will be included at least one time in a test						
PE	SEP	CCC	Min Items			
08-MS-ESS1-4	SEP 6 – E/S	CCC 3 – SPO	1			
08-MS-ESS2-1	SEP 2 – MOD	CCC 7 – S/C	1			
08-MS-ESS2-2	SEP 6 – E/S	CCC 3 – SPQ	1			
08-MS-ESS2-3	SEP 4 – DATA	CCC 1 – PAT	1			
08-MS-ESS3-1	SEP 6 – E/S	CCC 2 – C/E	1			
08-MS-ESS3-2	SEP 4 – DATA	CCC 1 – PAT	1			
08-MS-ESS3-3	SEP 6 – E/S	CCC 2 – C/E	1			
08-MS-LS1-4	SEP 7 – ARG	CCC 2 – C/E	1			
08-MS-LS1-5	SEP 6 – E/S	CCC 2 – C/E	1			
08-MS-LS3-1	SEP 2 – MOD	CCC 6 – S/F	1			
08-MS-LS4-1	SEP 4 – DATA	CCC 1 – PAT	1			
08-MS-LS4-2	SEP 6 – E/S	CCC 1 – PAT	1			
08-MS-LS4-3	SEP 4 – DATA	CCC 1 – PAT	1			
08-MS-LS4-6	SEP 5 – MCT	CCC 2 – C/E	1			
08-MS-PS1-1	SEP 2 – MOD	CCC 3 – SPO	1			
08-MS-PS1-3	SEP 8 – INFO	CCC 6 – S/F	1			
08-MS-PS1-6	SEP 6 – E/S	CCC 5 – E/M	1			
08-MS-PS3-3	SEP 6 – E/S	CCC 5 – E/M	1			
08-MS-PS3-5	SEP 7 – ARG	CCC 5 – E/M	1			

Table 3.1hh
Test Blueprint for LEAP 2025 Grade 8: CCC Coverage

	Grade 8: CCC Coverage					
CCC Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of CCC Items			
CCC 1 – PAT	5	26%	21%–31%			
CCC 2 – C/E	5	26%	21%–31%			
CCC 3 – SPQ	3	16%	11%–21%			
CCC 4 – SYS	0	0%	0%			
CCC 5 – E/M	3	16%	11%–21%			
CCC 6 – S/F	2	11%	5%-16%			
CCC 7 – S/C	1	5%	1%–11%			
Total	19	100%				

Table 3.1ii
Test Blueprint for LEAP 2025 Grade 8: SEP Coverage

	Grade 8: SEP Coverage					
SEP Overall	# in PEs in LSSS	Relative % in LSSS	% by Points of SEP Items			
SEP 1 – Q/P	0	0%	0%			
SEP 2 – MOD	3	16%	11%–21%			
SEP 3 – INV	0	0%	0%			
SEP 4 – DATA	4	21%	16%-26%			
SEP 5 – MCT	1	5%	2%-15%			
SEP 6 – E/S	8	42%	37%-42%			
SEP 7 – ARG	2	11%	5%-16%			
SEP 8 – INFO	1	5%	5%-15%			
Total	19	100%				

Table 3.1jj
Test Blueprint for LEAP 2025 Grade 8: SEP Reporting Category Coverage

Grade 8: SEP Reporting Category Coverage						
SEP Reporting Category	# PEs in LSSS	Relative % in LSSS	% by Points of SEP Items	Min Points		
Investigate (SEPs 4, 6, 8)	6	31.5%	27%-37%	7		
Evaluate (SEPs 4, 5, 7)	6	31.5%	27%-37%	7		
Reason Scientifically (SEPs 2 & 6)	7	37%	32%-42%	7		
Total	19	100%				

Table 3.1kk

Test Blueprint for LEAP 2025 Grade 8: SEP Compared to CCC Ratio

Grade 8: SEP Compared to CCC Ratio					
Relative Weight in LSSS Minimum %					
SEPs	50%	30%			
CCCs	50%	30%			

The assessment item development plans were created in conjunction with LDOE content staff. The development plans allowed for item attrition throughout the item development process, including reviews by LDOE assessment staff and by a content and bias review committee consisting of Louisiana educators. In addition, the number of items to be field tested also allowed for item loss due to deviations from psychometric criteria for item statistics based on student performance.

The development plans and the content distribution determined the focus of the item sets, tasks, and standalone items to be developed. Tables 3.2 show the item development plans for the number of items developed by WestEd by reporting category for grades 3–8. There were no new items developed for grade 4.

Table 3.2

Number of New Items Developed for the Spring 2024 Field Test for Item Sets, Tasks, and Standalone Items

Grade	Development Type	Total Number of Sets or Tasks	1-pt SRs	1-pt TEs	2-pt TEs	TPD/ TPI	ER	CR	Total Number of Items (non-ER/CR)
	Item Sets	4	22	0	0	14	0	4	36
3	Tasks	0	0	0	0	0	0	0	0
	Standalone Items	n/a	3	0	0	5	0	0	8
	Item Sets	3	7	10	8	8	0	3	33
5	Tasks	2	4	6	7	3	4	0	20
	Standalone Items	n/a	4	3	2	7	0	0	16
	Item Sets	3	15	3	10	5	0	3	33
6	Tasks	3	10	5	12	3	6	0	30
	Standalone Items	n/a	8	3	5	0	0	0	16
	Item Sets	3	14	3	16	0	0	3	33
7	Tasks	1	4	1	3	2	2	0	10
	Standalone Items	n/a	7	0	6	3	0	0	16
	Item Sets	3	20	0	10	3	0	3	33
8	Tasks	2	7	5	5	3	4	0	20
	Standalone Items	n/a	5	4	3	4	0	0	16

The development plans also may include item sets and tasks that were revised and refield-tested. Table 3.3 shows the number of items that were revised for refield-testing in 2024. There were no revised and refield-tested items in Grades 5 or 6.

Table 3.3
Number of Items Revised and Refield-tested for the 2024 field test.

Grade	Development Type	Total Number of Sets or Tasks	1-pt SRs	1-pt TEs	2-pt TEs	TPD/ TPI	ER	CR	Total Number of Items (non-ER/CR)
3	Item Sets	1	1	0	0	3	0	0	4
4	Item Sets	1	0	0	0	2	0	1	2
7	Item Sets	1	2	2	0	0	0	0	4
/	Tasks	1	2	0	0	0	2	0	2
0	Item Sets	0	0	0	0	0	0	0	0
8	Tasks	2	2	1	0	0	4	0	3

Proposal and Review of Topics and Sources

Performance Expectation Bundling

In the previous item development cycle, WestEd used the 2017 LSSS to recommend how performance expectations could be bundled in a task or item set to ensure that the breadth of all dimensions of constituent PEs is assessed in a meaningful way. Key to this bundling was the need to ensure that paired PEs and phenomena achieved a "natural fit." Therefore, not all PEs were bundled, some PEs appeared in more than one bundle, and some PEs were bundled across content domains. In previous development, the LDOE and WestEd determined that some item sets and tasks would allow a "mix and match" approach in which the science and engineering practice (SEP) for one of the PEs in a bundle could be used to develop items aligned to the disciplinary core idea (DCI) and crosscutting concept (CCC) of the other PE in the bundle. This approach was discontinued beginning with the current cycle because it generated some items with a SEP alignment outside the reporting category for the PE the item aligned to and therefore did not fit the reporting category. Within each task or item set, each item was given a primary assignment to one PE (DCI, SEP, and/or CCC) in the bundle, and to two or three of the dimensions comprising the three-dimensional structure of the performance expectation. However, the items in each item set or task worked together to assess the multidimensional nature of the performance expectations bundle.

Each year additional PE bundles may be proposed to the LDOE. Table 3.4 shows the bundles approved by the LDOE by grade, as well as the number of approved bundles that then were targeted for development in the 2023-2024 development cycle.

Table 3.4

PE Bundling by Grade

Grade	Total Number of PE Bundles Approved	Number of Bundles Targeted for Development
3	21	4
4	19	0
5	22	5
6	20	6
7	23	4
8	22	5

Phenomena Selection and Outline Development

Phenomena describe observable events in nature and include relevant data, images, and text that provide students with the information they need to engage in the scientific practices described in the LSSS. The stimuli for the LEAP 2025 grades 3–8 assessments are anchored on scientific phenomena described by text, images, tables, graphs, models, and graphic organizers created by WestEd's Design Team.

Phenomena and bundles were chosen to represent the breadth of assessable science content. As part of the item development plan, all PEs were aligned to at least one standalone item or to an item in an item set.

After studying the LSSS, the content lead generated lists of bundled and associated phenomena for item sets.

When identifying a phenomenon, the content lead considered:

- the emphasis of each performance expectation, as described in the clarification statements for each performance expectation;
- whether a proposed phenomenon was rich enough to support the required number of items, including overage;

- whether the phenomenon fit with the "PE bundles" developed earlier to provide meaningful, three-dimensional assessment of performance expectations; and
- whether the phenomenon was well suited for an item set (rather than a task).

Phenomena were chosen to represent the breadth of content described by the LSSS. The process of determining phenomena and associated bundles was iterative and included the identification of phenomena that could be assessed with a particular bundle, as well as understanding the need to assess PEs that had not been assessed in the previous field test.

Matching Phenomena to Item Sets and Tasks and Foci to Standalone Items

Item sets and tasks were targeted for development for the 2023–2024 development cycle based on an analysis of the test bank for each grade. The development of item sets influenced the selection of phenomena. Like the tasks, the item sets are phenomenabased, but unlike the tasks, they are made up of independent items that do not necessarily build upon each other. Also, unlike the tasks, the items in the item sets do not scaffold to help discriminate student performance levels, do not require a specific order, and do not contain a three-dimensional extended-response (ER) item. Although an item set does not need to contain a constructed-response (CR) item, WestEd developed CRs for all item sets. Table 3.5 shows the total number of ERs and CRs developed per grade.

Table 3.5

Constructed-Response Item Development by Grade

Grade	Number of ERs Developed	Number of CRs Developed
3	0	4
4	0	0
5	4	3
6	6	3
7	2	3
8	4	3

For the item sets and tasks, WestEd offered a document containing descriptions of phenomena associated with bundles to the LDOE to review prior to item development. Table 3.6 shows the number of phenomena submitted to the LDOE for item sets and tasks at grades 3–8.

Table 3.6

Phenomena Submitted by Grade

Grade	Number of Phenomena Submitted for Item Sets	Number of Phenomena Submitted for Tasks
3	8	0
4	0	0
5	6	4
6	6	6
7	6	2
8	6	4

For the item sets and tasks, the LDOE identified four phenomena at grades 3 and 7, five phenomena at grades 5, and 8, and six phenomena at grade 6 to be developed into stimuli. Upon approval of the phenomena, WestEd submitted item outlines containing stimuli and item descriptions to the LDOE. Once the item outlines were approved, item development for the item sets began.

In contrast to item sets and tasks, standalone items reflected independent content and are supported by a focus. A focus differs from a phenomenon in that it explores only

certain key aspects of an event and is typically supported by less data. As stated previously, the standalone items were included within the blueprints to provide greater coverage of the standards assessed and to provide flexibility in meeting the blueprints and test characteristic curve targets across test administrations. The WestEd content lead developed the foci for standalone items, based on standards that lacked coverage across the item sets and tasks. Consequently, these items were developed last. For standalone items, WestEd submitted the items and corresponding foci simultaneously; there was no separate focus approval phase for these items.

Outline and Stimuli Development

WestEd used both experienced internal and external science assessment editors to develop the phenomena-based stimuli for item sets. Before the editors began the process, the WestEd content lead trained them on the process of conducting an effective internet search for science articles on the LDOE's objectives, as well as training in universal design and bias and sensitivity issues. For an outline of the training, see Appendix A for the LEAP 2025 Grades 3–8 Training Agenda (2019–2024).

To support the outline development process, writers were given the LSSS. They were also provided specific item set templates that described the PE bundle to be written to, as well as the point value, item types, dimensional alignment of each of the items in the set, and whether the dimensions of the bundled PEs could be mixed or matched. The outline contained space for writers to enter the primary sources they used in researching their phenomenon and writing their stimulus, space for the writers to include a draft of the stimulus and its supporting data, as well as space to describe each item and its metadata. Writers submitted their item outlines to the editors, who finalized the item set outlines before they were submitted to the content lead and manager for senior review. After this review, the outlines were submitted to the LDOE.

Evaluating the Reading Level of Stimuli. WestEd performed Lexile and ATOS analyses on each stimulus to obtain quantitative measures of the readability of the texts. The Lexile Analyzer, developed by MetaMetrics, analyzes the semantic and syntactic features of a text and assigns it a Lexile measure. MetaMetrics also provides grade-level ranges corresponding to Lexile ranges. It should be noted that the grade-level ranges include overlap across grade levels. The ATOS text analysis tool, developed by Renaissance

Learning, considers the most important predictors of text complexity, including average sentence length and average word length, and uses a graded vocabulary list of more than 100,000 words to analyze word difficulty level. It reports on a grade-level scale. In addition to the Lexile and ATOS measures, the LSSS were used as an additional measure of grade-level appropriateness. WestEd and the LDOE also drew on the professional experience of educators, during Content and Bias Committee review, to verify that sources would be accessible to students, and made changes based on their feedback. Most of the stimuli developed for the assessments were found to be below or at grade level; however, some of the science vocabulary was evaluated as above grade level. In those cases, additional support such as parenthetical definitions (glossing) was included for necessary science content words that were above grade level and for words or phrases that were thought to be sources of potential confusion for students. The appropriateness of the stimuli for both content and readability was an explicit part of the content review process with Louisiana teachers.

Item Writing and Review Process

WestEd employed a cadre of item writers for the grades 3–8 assessments. All writers' resumes were approved by the LDOE before engaging in any item development activities. As the first step in the item writing process, the WestEd content lead provided a webinar training to all writers in February 2022. For an outline of the information covered, see Appendix A for the LEAP 2025 Grades 3–8 Item Outline Development Training Agenda. In the training, writers were provided context for the assessment, including LDOE expectations, the LSSS, and a review of best practices for item development. The item writers were provided the approved item topics and drafts of the stimuli, as well as item outlines that provided explanations of the phenomena underlying the item sets. Item writers were also provided with alignment to the Science and Engineering Practices, Crosscutting Concepts, and Disciplinary Core Ideas of the LSSS, and guidance on how each item set should be developed. The use of item set overviews allowed WestEd to provide direction for the items developed during the development cycle. For standalone development, item writers were provided with assignments that indicated the number of items to write to each performance expectation, as well as the specific dimensions to align to for each item.

The item writing assignments for each set also specified the set type, the item types (e.g., SR, MS, TE, TPI, TPD, CR, ER), and the number of items to be written, as well as potential item stems to be used for each item. Significant attention was devoted to understanding how to write TE items as well as scoring guides for CR items. Although all the writers were science writers with experience in writing three-dimensional items, WestEd also gave instructions in basic assessment item writing principles. Writers were instructed to make certain that the vocabulary and context of the items were grade-level appropriate, to ensure that the distracters were incorrect but plausible, and to avoid cueing and outliers in the items. Writers were also provided training in universal design and bias/sensitivity. A variety of items were presented and reviewed using universal design and bias/sensitivity lenses. This training also included an overview of these topics, (see Appendix A for the LEAP 2025 Grades 3–8 Item Writer Training Agenda). WestEd provided training and feedback to the writers throughout the development cycle, as the LDOE and WestEd gained a clearer understanding of how the stimuli, items, and sets worked together.

WestEd provided additional training to a subset of editors outlining the specific responsibilities for those who served as editors for the grades 3–8 assessments. For an outline of the information covered, see Appendix A for the LEAP 2025 Grades 3–8 Editor Training Agenda. Items went through two rounds of content editing that examined characteristics of items including alignment to the dimensions of the performance expectations of the LSSS, content accuracy, cognitive complexity, and quality of distractors. Items then went through one round of proofreading, which focused on grammar, usage, and consistent style of graphics, and a final round of review before being submitted to the LDOE for their first round of review.

Item Development Platform. Items were developed in Assessment Banking and Building solutions for Interoperable assessment (ABBI), Pearson's proprietary item development platform. In addition to the items and stimuli, the platform captured item metadata and allowed viewers to preview items using Pearson's format viewer (TestNav 8). In this view, items appeared together with all of the associated stimuli in the set. The ability to examine the items and stimuli as a set was critical in the item review and in the evaluation of the sets' content and cognitive demands on students.

Style Guidelines. Style guidelines continue to be based on documentation established with the LEAP 2025 Biology and Science assessments. This documentation was amended and updated as the development cycle progressed. When questions of style arose that

were unanswered by existing documentation, WestEd consulted the LDOE, and approved changes were added to the project style guide.

LDOE Content Review. As writing and editing for batches of item sets and standalone items were completed, these batches were sent to the LDOE for review by the LDOE Science Assessment Coordinators; Assessment Content Supervisor for Math, Science, and Small Populations; and Science Program Coordinators. Feedback from the LDOE review was implemented before the content and bias review meetings.

Content and Bias Review. After the completion of item development, WestEd coordinated content and bias review meetings. The meetings were led by facilitators from the LDOE, WestEd, and Pearson. Participants included current classroom teachers, retired teachers, content specialists, and school administrators. For the content and bias review meeting, participants completed nondisclosure agreements as part of the activities. The recruitment process, conducted by LDOE staff, also included participants from regions across the state. Participants represent the population of Louisiana students served—including special education, English Learners, students with disabilities—as well as the diverse geographic and demographic composition of the state. Table 3.7 provides the demographic characteristics of the review committee.

Table 3.7
Representation of Educators Participating in 2023–2024 Content and Bias Reviews

Grade Level	3	5	6	7	8
Classroom Teacher	8	9	7	6	8
Instructional Lead/Supervisor	1	0	1	2	0
School Administrator	0	0	0	0	0
Special Education Teacher	0	1	1	1	1
Visually or Hearing-Impaired Teacher	0	0	0	0	0
Other Staff	1	0	0	1	1
Black or African American	2	2	3	4	3
White	7	8	7	5	7
Male	0	0	1	1	2
Female	10	10	9	9	8
Total Participants	10	10	10	10	10

Before the committee members began the item review process, they received an orientation from the LDOE about the LEAP 2025 science assessments, and the WestEd content lead provided training on the criteria for evaluating items for content and bias considerations and the use of ABBI for item review. The committee members individually reviewed PE, SEP, DCI, and CCC alignment for each item and recorded the degree of alignment for each dimension and overall alignment on a worksheet on a scale of 0 (not aligned) to 3 (well aligned), referring to LSSS <u>Appendix A</u> (Learning Progressions). An item was considered to have a high degree of alignment if it aligned to the bullet listed in the PE. An item was considered to have a lower degree of alignment if it aligned to another bullet listed in the learning progression for that SEP or CCC. Committee members also recorded whether the science for each item was accurate and whether each item was free of bias. Areas of concern considered included opportunity and access, portrayal of groups represented, and protecting privacy and avoiding offensive content.

After the review of each item, each member voted in ABBI on whether to accept, accept with edits, or reject each item, recording comments for any item where they noted issues with science accuracy or bias. (If participants skipped an item or chose not to record a decision for a given item, the system registered the response as "No Vote" for that individual review. "No Vote" was recorded as the consensus rating when an initial group decision on an item was not reached, and the committee failed to return to that item and register a final vote to accept, revise, or reject the item.) Participants used Pearson laptops to access ABBI and only had access to ABBI during meeting times. Participants were locked out of ABBI when the meeting was not in progress. WestEd monitored participants to be sure that they did not use their cell phones at the table. WestEd also collected all materials at the end of each day, including notepads provided to the participants to write notes on as they reviewed the items.

Following the individual reviewers' votes, the group came together to view and discuss each stimulus and item as it was projected on-screen, with the goal of achieving consensus. The WestEd and Pearson facilitators compiled detailed notes about committee decisions for implementation after the review.

Results of Content Review. The results of the reviewers' individual judgments were captured in ABBI. Tables 3.8a-e provides these results, based on the participants' individual votes on each item following their initial review.

Table 3.8a *Grade 3 Vote Totals Based on Individual Votes Following Initial Review*

Item Type	N Items	Accept	Accept with Edits	No Vote	Reject	Total
CR	4	40	0	0	0	40
ER	0	0	0	0	0	0
MC	25	238	11	1	0	250
MS	0	0	0	0	0	0
TE	0	0	0	0	0	0
TPD	12	113	6	1	0	120
TPI	7	65	5	0	0	70
All Grade 3	48	456	22	2	0	480

Table 3.8b Grade 5 Vote Totals Based on Individual Votes Following Initial Review

Item Type	N Items	Accept	Accept with Edits	No Vote	Reject	Total
CR	3	26	4	0	0	30
ER	4	38	2	0	0	40
MC	14	135	5	0	0	140
MS	1	10	0	0	0	10
TE	36	331	27	0	2	360
TPD	16	151	8	1	0	160
TPI	2	20	0	0	0	20
All Grade 5	76	711	46	1	2	760

Table 3.8c Grade 6 Vote Totals Based on Individual Votes Following Initial Review

Item Type	N Items	Accept	Accept with Edits	No Vote	Reject	Total
CR	3	28	2	0	0	30
ER	6	48	12	0	0	60
MC	31	281	26	0	3	310
MS	1	10	0	0	0	10
TE	39	336	50	2	2	390
TPD	1	7	3	0	0	10
TPI	7	56	13	0	1	70
All Grade 6	88	766	106	2	6	880

Table 3.8d Grade 7 Vote Totals Based on Individual Votes Following Initial Review

Item Type	N Items	Accept	Accept with Edits	No Vote	Reject	Total
CR	3	25	5	0	0	30
ER	2	19	1	0	0	20
MC	22	182	36	2	0	220
MS	3	25	4	1	0	30
TE	29	243	44	3	0	290
TPD	0	0	0	0	0	0
TPI	5	48	2	0	0	50
All Grade 7	64	542	92	6	0	640

Table 3.8e

Grade 8 Vote Totals Based on Individual Votes Following Initial Review

Item Type	N Items	Accept	Accept with Edits	No Vote	Reject	Total
CR	3	27	2	1	0	30
ER	4	21	19	0	0	40
MC	30	255	43	2	0	300
MS	2	17	3	0	0	20
TE	27	203	63	3	1	270
TPD	3	22	8	0	0	30
TPI	7	53	17	0	0	70
All Grade 8	76	598	155	6	1	760

At the end of the meeting, consensus votes for each grade were compiled. There were no rejected items or item sets in any grade. All other items reviewed at each grade were either accepted as is or accepted with edits.

Post-Review Finalization. After the content and bias review, the WestEd staff implemented the committee's feedback and then met virtually with LDOE staff for reconciliation. WestEd provided records of all implemented changes to the LDOE prior to the virtual reconciliation meetings. During the reconciliation meeting, content leads from the LDOE and WestEd reviewed items to ensure that the items reflected the content, clarity, and style appropriate for inclusion in the field test. Following the reconciliation meetings, which focused on the finalization of item content, the LDOE and WestEd content leads worked together to finalize the scoring guides for CR and ER items through a separate series of communications. Once all content considerations were resolved, all items and stimuli went through a final formal fact-check by content editors and two additional rounds of proofreading. Any changes resulting from these reviews were submitted to the LDOE for approval.

Data Review Process and Results

During data review of the spring 2024 FT items, content experts and psychometric support staff reviewed field-tested items with accompanying data to make judgments about the appropriateness of items for use on future operational test forms. Statistically flagged items were not rejected on the sole basis of statistics; only items with identifiable flaws based on content were rejected.

The data review meetings began with a refresher presentation to data review. The presentation included a review of item statistics (difficulty, discrimination, DIF, score distributions), appropriate interpretations and inferences, what would be considered reasonable values, and how the values might differ across item types.

Facilitators from Pearson and WestEd led the data review. Statistical information was evaluated for each item to determine whether the item functioned as intended. Each item's suitability for future operational tests was then evaluated in the context of the field-test statistics. Judgments to accept, accept with edits (or "revise/refield-test"), or reject were then recorded for each item. If the decision was to edit or to reject an item, additional information was captured to document the reason for the decision. Table 3.9 summarizes the disposition of field-tested items from data review.

Table 3.9 FT Item Dispositions by Item Type, 2024 Data Review

			Number of Items				
Grade	ltem Type	Accept	Edits Accepted	Reject	Total	% of Total	
	CR	2	0	0	2	4	
	MC	25	1	0	26	55	
	MS	0	0	0	0	0	
3	TE	0	0	0	0	0	
	TPI	6	1	0	7	15	
	TPD	9	3	0	12	26	
	Total	42	5	0	47	100	
	CR	1	0	0	1	20	
4	MC	3	0	0	3	60	
	MS	0	0	0	0	0	

			N	umber of Iten	ns	
Grade	ltem Type	Accept	Edits Accepted	Reject	Total	% of Total
	TE	0	0	0	0	0
	TPI	1	0	0	1	20
	TPD	0	0	0	0	0
	Total	5	0	0	5	100
	CR	1	0	0	1	1.5
	ER	0	0	0	0	0
	MC	11	0	2	13	20
5	MS	1	0	0	1	1.5
5	TE	31	1	2	34	52
	TPI	3	0	0	3	5
	TPD	12	0	1	13	20
	Total	59	1	5	65	100
	CR	2	0	0	2	3
	ER	0	0	0	0	0
	MC	14	5	8	27	40
	MS	0	1	0	1	2
6	TE	14	9	8	31	46
	TPI	5	0	0	5	7
	TPD	0	1	0	1	2
	Total	35	16	16	67	100
	CR	2	1	0	3	4
	ER	0	0	0	0	0
	MC	14	11	0	25	35
7	MS	3	0	0	3	4
7	TE	21	9	1	31	44
	TPI	6	1	0	7	10
	TPD	1	1	0	2	3
	Total	47	23	1	71	100
	CR	3	0	0	3	4
0	ER	0	0	0	0	0
8	MC	22	2	1	25	36
	MS	3	0	0	3	4

			Number of Items				
Grade	ltem Type	Accept	Edits Accepted	Reject	Total	% of Total	
	TE	18	4	4	26	38	
	TPI	3	2	1	6	9	
	TPD	3	2	1	6	9	
	Total	52	10	7	69	100	

Table 3.9

FT Item Dispositions by Item Type, 2023 Data Review (continued)

			N	umber of Iten	ns	
Grade	ltem Type	Accept	Edits Accepted	Reject	Total	% of Total
	CR	2	2	0	4	7
	ER	0	0	0	0	0
	MC	10	3	2	15	25
6	MS	1	3	0	4	7
0	TE	15	13	0	28	46
	TPI	3	0	0	3	5
	TPD	3	2	2	7	11
	Total	34	23	4	61	100
	CR	0	0	0	0	0
	ER	0	0	0	0	0
	MC	0	0	0	0	0
7	MS	0	0	0	0	0
/	TE	0	0	0	0	0
	TPI	0	0	0	0	0
	TPD	0	0	0	0	0
	Total	0	0	0	0	0
	CR	2	0	0	2	3
	ER	0	0	0	0	0
	MC	12	6	4	22	37
8	MS	2	0	0	2	3
δ	TE	16	12	0	28	47
	TPI	1	1	0	2	3
	TPD	2	2	0	4	7
	Total	35	21	14	60	100

Following the data review meeting, LDOE content specialists considered the item level data review outcomes to determine which sets and tasks could be used operationally or rejected unless revised/re-field tested. The reconciliation decisions were the final decisions. It should be noted that the training presentation agenda for data review is included in <u>Appendix A: Training Agendas</u>.

4. Construction of Test Forms with Embedded Field Test

Test Design

To assess the integrated nature of the content, practices, and crosscutting concepts of the LSSS, the LEAP 2025 grades 3–8 science assessments involved set-based designs. The tests included item sets and, for grades 5–8, a task on each form, each anchored by a common stimulus or stimuli. Additionally, standalone items were included to support meeting the specific targets of the test blueprints. Table 4.1a shows the Test Design for Science Grade 3.

Table 4.1a *Test Design for Science Grade 3*

Test Session	Numbers of Items
Session 1: One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
OP Standalone Items	4 OP Standalone SR Items 1 OP Standalone TPD/TPI Items
One FT Item Set	2 FT Item Set SR Items 1–2 FT Item Set TPD/TPI Item 0–1 FT Item Set CR Items
FT Standalone Items	0–2 FT Standalone SR Items 0–2 FT Standalone TPD/TPI Items
Session 2: One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
OP Standalone Items	6 OP Standalone SR Items 1 OP Standalone TPD/TPI Items
Total Items Field Tested Across Forms for Grade 3	3 FT Standalone SR Items 4 FT Standalone TPD/TPI Items 20 FT Item Set SR Items 15 FT Item Set TPD/TPI Items 2 Item Set CR Items

Note: Students do not complete more than one CR per item set. There were a total of 3 operational CR items per form.

Table 4.1b shows the Test Design for Science Grade 4.

Table 4.1b

Test Design for Science Grade 4

Test Session	Numbers of Items
Session 1: One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
OP Standalone Items	2 OP Standalone SR Items 1 OP Standalone TPD/TPI Items
FT Standalone Item	0–1 FT Standalone SR Items 0–1 FT Standalone TPD/TPI Items
Session 2: One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One FT Item Set	2 FT Task SR Items 2 FT Task TPD/TPI Items 0–1 FT Item Set CR Items
Total Items Field Tested Across Forms for Grade 4	1 FT Standalone SR items 0 FT Standalone TPD/TPI Items 2 FT Item Set SR Items 1 FT Item Set TPD/TPI Item 1 Item Set CR Item

Note: Students did not complete more than one CR per item set. There were a total of 3 operational CRs per form. Item sets field tested included one item set developed in 2018.

Table 4.1c shows the Test Design for Science Grades 5–8.

Table 4.1c

Test Design for Science Grades 5–8

Test Session	Numbers of Items
Session 1: One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
OP Standalone Items	2 OP Standalone SR Items 1 OP Standalone TPD/TPI Items
Session 2: One OP Task	2 OP Task SR Items 2 OP Task TPD/TPI Items 1 OP Task ER Item
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
One OP Item Set	2 OP Item Set SR Items 1–2 OP Item Set TPD/TPI Items 0–1 OP Item Set CR Items
OP Standalone Items	1 OP Standalone SR Item 2 OP Standalone TPD/TPI Items
Session 3: One FT Item Set or Task	2 FT Item Set SR Items 2 FT Item Set TPD/TPI Items 0–1 FT Item Set CR Items OR 2 FT Item Set SR Items 2 FT Item Set TPD/TPI Items 1 FT Item Set ER Item
FT Standalone Items	0–2 FT Standalone SR Items 0–2 FT Standalone TPD/TPI Items

Test Session	Numbers of Items
Total Items Field Tested Across Forms for Grade 5	4 FT Standalone SR Items 5 FT Standalone TE Items 7 FT Standalone TPD/TPI Items 6 FT Item Set SR Items 16 FT Item Set TE Items 6 FT Item Set TPD/TPI Items 3 FT Item Set CR Items 4 FT Task SR Items 13 FT Task TE Items 3 FT Task TE Items 4 FT Task ER Items
Total Items Field Tested Across Forms for Grade 6	8 FT Standalone SR Items 8 FT Standalone TE Items 0 FT Standalone TPD/TPI Items 11 FT Item Set SR Items 7 FT Item Set TE Items 3 FT Item Set TPD/TPI Items 3 FT Item Set CR Items 9 FT Task SR Items 16 FT Task TE Items 3 FT Task TPD/TPI Items
Total Items Field Tested Across Forms for Grade 7	8 FT Standalone SR Items 6 FT Standalone TE Items 5 FT Standalone TPD/TPI Items 13 FT Item Set SR Items 18 FT Item Set TE Items 2 FT Item Set TPD/TPI Items 3 FT Item Set CR Items 7 FT Task SR Items 7 FT Task TE Items 2 FT Task TE Items 4 FT Task ER Items

Test Session	Numbers of Items				
Total Items Field Tested Across Forms for Grade 8	5 FT Standalone SR Items 7 FT Standalone TE Items 4 FT Standalone TPD/TPI Items 12 FT Item Set SR Items 7 FT Item Set TE Items 2 FT Item Set TPD/TPI Items 3 FT Item Set CR Items 11 FT Task SR Items 12 FT Task TE Items 6 FT Task TPD/TPI Items 8 FT Task ER Items				

Note: Students do not complete more than one CR per item set. There were a total of three operational CR items per form.

Initial Construction

The purpose of the spring 2024 forms construction activities was to create operational forms using the spring 2018, spring 2019, spring 2022, and spring 2023 field test items that were approved for operational use and to embed field test items in the spring 2024 forms for potential use in future operational assessments. This section describes the process used to create operational and field test forms.

Operational Form

Data review-approved items, field tested in spring 2018, 2019, 2022, or 2023 were available for use on the spring 2024 operational assessments.

For each of grades 3–8, WestEd completed item selection for one operational (OP) form for the spring 2023 administration. WestEd worked with the LDOE content staff to select items for the forms following the data review meeting in September and submitted these forms to Pearson psychometricians for consideration before formal submission to the LDOE for approval.

For grades 3 and 4, a combination of item sets and standalone items were chosen that would ensure that the relative distribution of score points by reporting category would meet the blueprints for the operational assessment while avoiding similar content and topics across the balance of items and item types. For grades 5–8, the WestEd content lead selected the task first and followed with a combination of item sets and standalone items that would ensure that the relative distribution of score points by reporting category would meet the blueprints for the operational assessment while avoiding similar content and topics across the balance of items and item types. Tables 4.2a–f provide the operational test composition for grades 3–8 for spring 2024.

Table 4.2a *LEAP 2025 Grade 3: Operational Test Composition*

ltem Sets/ltem Types	Total Sets	Total Items per Set	Total Points per Set	SR	CR, Two- Part	Total Items	Total Points
4-Item Set	6	4	6	12	12	24	36
Standalone items	1	12	14	10	2	12	14
Totals	_	_	-	22	14	36	50

Table 4.2b
LEAP 2025 Grade 4: Operational Test Composition

Item Sets/Item Types	Total Sets	Total Items per Set	Total Points per Set	SR	CR, Two- Part	Total Items	Total Points
4-Item Set	7	4	6	14	16	28	42
Standalone items	1	8	10	16	2	8	10
Totals	-	-	-	20	18	36	52

Table 4.2c

LEAP 2025 Grade 5: Operational Test Composition

ltem Sets/Item Types	Total Sets	Total Items per Set	Total Points per Set	SR, 1- pt TE	CR, 2-pt TE, Two-Part	ER	Total Items	Total Points
4-Item Set	5	4	6	10	10		20	30
Standalone items	1	12	16				12	16
Task	1	5	15	2	2	1	5	15
Totals	-	-	-	12	12	1	37	61

Table 4.2d LEAP 2025 Grade 6: Operational Test Composition

ltem Sets/ltem Types	Total Sets	Total Items per Set	Total Points per Set	SR, 1-pt TE	CR, 2-pt TE, Two-Part	ER	Total Items	Total Points
4-Item Set	5	4	6	10	10		20	30
Standalone items	1	12	16				12	16
Task	1	5	15	2	2	1	5	15
Totals	_	_	_	12	12	1	37	61

Table 4.2e LEAP 2025 Grade 7: Operational Test Composition

ltem Sets/ltem Types	Total Sets	Total Items per Set		SR, 1- pt TE	CR, 2-pt TE, Two-Part	ER	Total Items	Total Points
4-Item Set	5	4	6	10	10		20	30
Standalone items	1	12	16				12	16
Task	1	5	15	2	2	1	5	15
Totals	_	_	_	12	12	1	37	61

Table 4.2f
LEAP 2025 Grade 8: Operational Test Composition

ltem Sets/Item Types	Total Sets	Total Items per Set	Total Points per Set	SR, 1- pt TE	CR, 2-pt TE, Two-Part	ER	Total Items	Total Points
4-Item Set	5	4	6	10	10		20	30
Standalone items	1	12	16				12	16
Task	1	5	15	2	2	1	5	15
Totals	_	1	1	12	12	1	37	61

Field Test Versions

The number of field test versions administered in spring 2024 varied by grade. These data are shown in Table 4.4.

Table 4.4

Spring 2023 Field Test Versions Administered by Grade

Grade	Number of Versions
3	10
4	1
5	14
6	14
7	14
8	14

In some cases, the number of field test slots exceeded the number of items available for field testing. As a result, some items were repeated among field test versions. One or two versions of each item set were field tested as needed.

For grade 3, one field test item set and one field test standalone item were embedded within session 1 of the operational form. For grade 4, one field test standalone item was embedded in session 1 and a field test item set was embedded in session 2. For grades 5-8, one item set and five standalone items were embedded in session 3.

In addition to content balance, the WestEd content lead was careful to avoid cueing and clanging between items. Cueing occurs when content in one item provides clues to the answer of another item. Clanging refers to overlap or similarity of content. Because content was purposefully distributed across the forms, cueing and clanging were intended to have been avoided; however, developers also conducted a separate review of the forms to check for inadvertent cueing or clanging.

Following the final item placement by the WestEd content lead, test maps containing each item's unique identification number (UIN) were created. The test maps captured details about each proposed form, including test session, item sequence, unique item number, and associated item metadata. Item descriptions were also included for each item, to aid in the review of the selection and placement of individual items.

Revision and Review

Psychometric Approval of Operational Forms

Prior to submitting the forms to the LDOE staff for review, Pearson psychometricians and WestEd content specialists participated in an iterative process of reviewing and revising the forms. The psychometric review consisted of comparisons of the expected representation and the actual representation of reporting categories, science and engineering practices, disciplinary core ideas, crosscutting concepts, performance expectations, and item types on the operations forms including SR, CR, TPI, and TPD at grades 3 and 4; and SR, CR, TE, TPI, TPD, and ER at grades 5–8.

The answer keys for MC items were also examined to determine whether any forms had significantly non-uniform distributions of correct responses (A, B, C, and D). Spreadsheets were used to generate frequency tables of reporting categories, science and engineering practices, disciplinary core ideas, crosscutting concepts, performance expectations, item types, and MC answer keys for each form and across forms. Deviations from the blueprint were identified and addressed. Test characteristic curves (TCC) based on item response theoretic models were applied to data, and conditional standard errors of measurement were computed for each iteration during the test construction process to evaluate how well a proposed test form matched psychometric targets. Psychometric approval from Pearson was provided for all forms prior to submission to the LDOE for their review. Criteria to flag items based on scoring point can be found in Table 4.4.

Table 4.4
Summary of Flagging Criteria to Select/Flag Items: Classical Analysis and IRT

	P-value		P-B	DIF		IRT	
Point	Low Bound	Upper Bound	Lower Bound	Exclude	a	b	С
1	0.25	0.90	0.20	(0.30-3.50	-3.00-3.00	< 0.35
2 and higher	0.25	0.90	0.20		0.30-3.50	-3.00-3.00	N/A

Note: Detailed information can be found from the 2021–2023 Framework and Test Construction Document. It should be noted that these values are psychometric recommendations. Actual item decision occurs by content staff based on these recommendation criteria.

LDOE Review

Following the psychometric reviews, the test maps and constructed sets were delivered to the LDOE for approval. Forms were reviewed by both LDOE content and psychometric staff. Based on the LDOE review, sets or standalone items were replaced and the sequence of answer choices (for field test items) and the sequence of items within sets were revised as requested. Following these changes, the overall balance of answer choices and key runs was re-evaluated and final adjustments were made to achieve the appropriate balance.

Finalized test maps were used to create PDF versions of paper forms, which were reviewed by WestEd's proofreaders before the items were transferred from ABBI to DRC.

Test Forms and Accessible Versions

Online and Paper Forms

The LEAP 2025 science assessments for grades 3–8 were administered as computer-based tests (CBT) with a paper-based option for grade 3 (selected at the school system level) and an accommodated print form only for a student who required a paper-based accommodation for grades 4–8.

Accommodated Print Versions

For grades 4–8, the accommodated print form was selected based on the field test version that contained the fewest and least complex technology-enhanced items. This version was identified as Version 1. The technology-enhanced items in this version were converted to a paper and pencil format that allowed students to record their responses, or have their responses transcribed into the test booklet. In addition, alternate text was written for all stimuli and items containing graphics. Detailed information can be found in <u>Appendix G</u>, Accommodated Print and Braille Creation.

Form Versions for Students with Visual Impairments

Braille and large-print test form versions were constructed for each grade to enable students with visual impairments to participate in the LEAP 2025 assessments. Version 1 of the grade 3 paper-based test form served as the basis for braille and large-print development. Braille forms for grades 4–8 were based on the accommodated print forms for operational items in Version 1. There were no large-print versions of the grades 4–8 accommodated print forms. Instead, students who needed a large-print version in grades 4–8 used larger-sized monitors and/or the magnification features of the online testing system. All online test content had been developed to scale in relation to the available area on larger monitors while maintaining the correct aspect ratio. Specific recommendations on how to transcribe items into braille were provided by the braille publisher to produce the braille version of the LEAP 2025 assessments and the test administrator's notes that accompanied the braille forms. The goal was to maximize the number of items that could be transcribed into braille.

For students who were administered a large-print or braille version, examiners were instructed to transcribe students' responses from the large-print or braille version into a consumable test booklet for grade 3, and the online testing system (INSIGHT) for grades 4–8, exactly as the students responded. Detailed information can be found in <u>Appendix G</u>, Accommodated Print and Braille Creation.

5. Test Administration

This chapter describes processes and activities implemented and information disseminated to help ensure standardized test administration procedures and, thus, uniform test administration conditions for students. According to the American Educational Research Association (AERA), American Psychological Association (APA), and National Council on Measurement in Education (NCME) (2014) *Standards for Educational and Psychological Testing* (hereafter the *Standards*), "The usefulness and interpretability of test scores require that a test be administered and scored according to the developer's instructions" (111). This chapter examines how test administration procedures implemented for the Louisiana Education Assessment Program 2025 (LEAP 2025) strengthen and support the intended score interpretations and reduce constructivelevant variance that could threaten the validity of score interpretations.

Training of School Systems

To ensure that the LEAP 2025 assessments are administered and scored in accordance with the department's policies, the LDOE takes a primary role in communicating with and training school system personnel. The LDOE provides train-the-trainer opportunities for the district test coordinators, who in turn convey test administration training to schools within their school systems. The LDOE conducts quality-assurance visits during testing to ensure adherence to the standardized administration of the tests.

The district test coordinators are responsible for the schools within their systems. They disseminate information to each school, offer assistance with test administration, and serve as liaisons between the LDOE and their school systems. The LDOE also provides assistance with and interpretation of assessment data and test results.

Ancillary Materials

Ancillary materials for LEAP 2025 test administration contribute to the body of evidence of the validity of score interpretation. This section examines how the test materials address the *Standards* related to test administration procedures.

For the spring test administration, Data Recognition Corporation (DRC) produces two administration manuals: the *LEAP 2025 Grade 3 Paper-Based Test Administration Manual* (TAM) and the *LEAP 2025 Grades 3–8 Computer-Based Test Administration Manual* (TAM). The TAMs provide detailed instructions for administering the LEAP assessments. The manuals include information on test security, test administrator responsibilities, test preparation, administration of tests (computer-based or paper-based), and post-test procedures.

Table of Contents for LEAP 2025 Paper-Based Test Administration Manual (TAM)

- Notes and Reminders
- Test Administrator Pre-Administration Oath of Security and Confidentiality
 Statement
- Test Administrator Post-Administration Oath of Security and Confidentiality Statement
- Overview
- Test Security
 - Secure Test Materials
 - Testing Irregularities and Security Breaches
 - Testing Environment
 - Violations of Test Security
 - Answer Change Analysis
 - Voiding Student Tests
- Test Administrator Responsibilities
- Test Administration Checklists
 - Before Testing
 - During Testing
 - After Testing (Daily)
 - After Testing (Last Day)
- Test Administrators' Frequently Asked Questions
- Test Materials
 - Receipt of Test Materials
- Testing Guidelines
 - Testing Eligibility
 - Test Schedule
 - Extended Time for Testing
- Testing Times

- Makeup Testing
- Testing Conditions
- Special Populations and Accommodations
 - o IDEA Special Education Students
 - Students with One or More Disabilities According to Section 504
 - o Gifted and Talented Special Education Students
 - o Test Accommodations for Special Education and Section 504 Students
 - Special Considerations for Deaf and Hard-of-Hearing Students
 - English Learners (ELs)
- Hand-Coded Consumable Test Booklets
- Students Absent from Testing
- Consumable Test Booklet Coding
 - Coding the Demographic Section
- Sample Grade 3 English Language Arts Consumable Test Booklet
- General Instructions for LEAP 2025
 - Student Marking/Erasing on Consumable Test Booklet
 - o Reading Directions to Students
 - Special Instructions
- Directions for Administering LEAP 2025 Tests
- Post-Test Procedures
 - o Test Administrator Oath of Security and Confidentiality Statement
 - Used and Unused Consumable Test Booklets (Defined)
 - Transferring Student Responses
 - o Returning Test Materials to the School Test Coordinator
- Index

Table of Contents for LEAP 2025 Computer-Based Test Administration Manual (TAM)

- Notes and Reminders
- Test Administrator Pre-Administration Oath of Security and Confidentiality Statement
- Test Administrator Post-Administration Oath of Security and Confidentiality Statement
- Overview
- Test Security
 - Secure Test Materials
 - Testing Irregularities and Security Breaches
 - Testing Environment
 - Violations of Test Security
 - Voiding Student Tests
- Test Administrator Responsibilities
 - Software Tools and Features for Test Administrators
- Test Administration Checklists
 - Before Testing
 - During Testing
 - After Testing (Daily)
 - After Testing (Last Day)
- Test Administrators' Frequently Asked Questions
- Test Materials
 - Receipt of Test Materials
- Testing Guidelines
 - Testing Eligibility
 - Testing Schedule
 - Extended Time for Testing
- Testing Times for Grades 3-8
 - Makeup Testing
 - Testing Conditions
- Online Tools Training
- Student Tutorials
 - Student Tutorials
- Special Populations and Accommodations
 - o IDEA Special Education Students

- Students with One or More Disabilities According to Section 504
- Gifted and Talented Special Education Students
- Test Accommodations for Special Education and Section 504 Students
- Special Considerations for Deaf and Hard-of-Hearing Students
- English Learners (ELs)
- General Instructions
 - Reading Directions to Students
- LEAP 2025: Grades 3–8 English Language Arts (All Sessions)
- LEAP 2025: Grades 3–8 Mathematics (All Sessions)
- LEAP 2025: Grades 3–8 Science (Sessions 1–2)
- LEAP 2025: Grades 5–8 Science Session 3 Select Schools Only
- LEAP 2025: Grades 3–8 Social Studies (All Sessions)
- Post-Test Procedures
 - Test Administrator Post-Administration Oath of Security and Confidentiality
 Statement
 - Returning Test Materials to the School Test Coordinator
- Index

DRC also produces test coordinator manuals for paper- and computer-based test administrations. The TCMs provide detailed instructions for district and school test coordinators' responsibilities for distributing, collecting, and returning test materials to DRC for scoring.

Table of Contents for LEAP 2025 Paper-Based Testing Test Coordinators Manual (TCM)

- Key Dates
- Resources Available in DRC INSIGHT Portal.
- Alerts
- Pre-Administration Oath of Security and Confidentiality Statement
- Post-Administration Oath of Security and Confidentiality Statement
- General Information
 - Test Security
 - Key Definitions
 - Violations of Test Security
 - Answer Change Analysis
 - Voiding Student Tests

- Testing Guidelines
 - Testing Eligibility
 - Testing Conditions
 - Test Schedule
 - Extended Time for Testing
 - Extended Breaks
 - Makeup Testing
 - Test Administration Resources
- Testing Times for Grade 3
- District Test Coordinator
 - Conduct Training Session
 - Receive Test Materials
 - Large-Print and Braille Test Materials and Communication Assistance Scripts
 (CAS)
 - Accommodated Materials
 - Verify and Distribute Test Materials to School Test Coordinators
 - Request Additional Test Materials and Bar-Code Labels
 - Collect Materials from Schools After Testing
 - Used and Unused Consumable Test Booklets (Defined)
 - Unscorable Documents and Unscorable Document Labels
- Directions for Returning Test Materials to DRC in May
 - Pickup 1: ELA and Mathematics Scorable Test Materials
 - o Pickup 2: Science and Social Studies Scorable Test Materials
 - Pickup 3: Nonscorable Test Materials
 - Final Checklist for Returning Test Materials to DRC
- School Test Coordinator
 - Receive and Verify Test Materials
 - Conduct Test Administration and Security Training Session
 - Supervise Application of Bar-Code Labels and Coding of Consumable Test Booklets
 - o Soiled, Damaged, and Other Unscorable Consumable Test Booklets
 - o Verify and Distribute Materials to Test Administrators
 - Supervise Test Administration
 - Collect Test Materials
 - Used and Unused Consumable Test Booklets (Defined)

- Coding Responsibilities of Principals—Before Testing
- o Coding Responsibilities of Principals—Before or After Testing
- Coding Responsibilities of Principals—After Testing
- Directions for Returning Test Materials to District Test Coordinator
 - o Pickup 1: ELA and Mathematics Scorable Test Materials
 - o Pickup 2: Science and Social Studies Scorable Test Materials
 - Pickup 3: Nonscorable Test Materials
 - Final Checklist for Returning Test Materials to DTC
- Void Notification
- Index

Table of Contents for LEAP 2025 Computer-Based Testing Test Coordinators Manual (TCM)

- Key Dates
- Resources Available in DRC INSIGHT Portal
- Alerts
- Pre-Administration Oath of Security and Confidentiality Statement
- Post-Administration Oath of Security and Confidentiality Statement
- General Information
 - DRC INSIGHT Portal and INSIGHT
- Test Security
 - Key Definitions
 - Violations of Test Security
- Testing Guidelines
 - Testing Eligibility
 - Testing Conditions
 - Testing Schedule
 - Extended Time for Testing
 - Extended Breaks
 - Accommodations
 - Makeup Testing
 - o Test Administration Resources
- Testing Times for Grades 3 through 8
- Roles and Responsibilities
 - District Test Coordinator
 - School Test Coordinator

- Technology Coordinator
- Managing Test Tickets
 - Student Transfers
 - Locked Test Tickets
 - Technical Issues
 - Invalidating Test Tickets
- Resources for Online Testing
 - Test Administration Manuals
 - DRC INSIGHT Portal User Guide
 - LEAP 2025 Accommodations and Accessibility Features User Guide
 - INSIGHT Technology User Guide
 - Online Tools Training (OTT)
 - Student Tutorials
- Void Notification

The LDOE assessment staff review, provide feedback, and give final approval for these manuals. The manuals are inclusive of grades 3–8 English Language Arts (ELA), Mathematics, Social Studies, and Science.

The *Standards* contain multiple references relevant to test administration. Information in the TAMs addresses these in the following manner.

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

The TAMs provide instructions for activities that happen before, during, and after testing with sufficient detail and clarity to support reliable test administrations by qualified test administrators. To ensure uniform administration conditions throughout the state, instructions in the TAMs describe the following: general rules of paper and online testing; assessment duration, timing, and sequencing information; and the materials required for testing.

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

To ensure the usefulness and interpretability of test scores and to minimize sources of construct-irrelevant variance, it was essential that the LEAP 2025 tests were administered according to the prescribed TAMs. It should be noted that adhering to the test schedule is also a critical component. The TCMs included instructions for scheduling the test within the state testing window. The TAMs and TCMs also contained the schedule for timing each test session.

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

Department staff release annual test security reports that describe a wide range of improper activities that may occur during testing, including the following: copying and reviewing test questions with students; cueing students during testing, verbally or with written materials on the classroom walls; cueing students nonverbally, such as by tapping or nodding the head; allowing students to correct or complete answers after tests have been submitted; splitting sessions into two parts; ignoring the standardized directions in the online assessment; paraphrasing parts of the test to students; changing or completing (or allowing other school personnel to change or complete) student answers; allowing accommodations that are not written in the Individualized Education Program (IEP), Individual Accommodation Plan/504 Plan (IAP), or English Learner Plan (EL plan); allowing accommodations for students who do not have an IEP, IAP, or EL plan; or defining terms on the test.

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

The TAMs outline the steps that teachers should take to prepare the classroom testing environment for administering the LEAP 2025 test. These include the following:

- Determine the layout of the classroom environment.
- Plan seating arrangements. Allow enough space between students to prevent the sharing of answers.
- Eliminate distractions such as bells or telephones.

- Use a Do Not Disturb sign on the door of the testing room.
- Make sure classroom maps, charts, and any other materials that relate to the content and processes of the test are covered or removed or are out of the students' view.

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

The TAMs present instructions for post-test activities to ensure that online tests are submitted and printed test materials are handled properly to maintain the integrity of student information and test scores. Detailed instructions guide test examiners in submitting all online test records. For students who were administered a large-print or braille version of the LEAP 2025 assessment, examiners are instructed to transcribe students' responses from the large-print or braille test book into the online testing system (INSIGHT) exactly as they responded in the large-print or braille test book.

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

Throughout the manuals, test coordinators and examiners are reminded of test security requirements and procedures to maintain test security. Specific actions that are direct violations of test security are noted. Detailed information about test security procedures is presented under "Test Security" in the manuals.

Return Material Forms and Guidelines

The paper-based TCM instructs test coordinators regarding procedures for organizing and packing materials and returning them to DRC for secure inventory purposes. The LDOE assessment staff have opportunities to review, provide feedback, and give final approval of the guidelines. The purpose of the instructions is to ensure that secure test materials are properly accounted for and organized appropriately for the return shipment.

Security Checklists

As soon as printed test materials are received by a school system, the district test coordinator ensures that the first and last security barcodes on the tests match the packing list they received. The district test coordinator then packages the tests to be sent to schools. Upon returning test books to DRC, school and district test coordinators are required to complete and submit an accountability form that details the number of test books or printed test forms returned. This form also requires that systems/schools document nonstandard situations, including lost, damaged, destroyed, extra, or missing test books.

Interpretive Guides

Essential to making valid interpretations of test scores is an understanding of what the test scores mean and how to interpret score reports. The Interpretive Guide is written for Louisiana teachers and administrators who receive the LEAP 2025 score reports. https://www.louisianabelieves.com/resources/search/assessment

Time

Each session of each content area test is timed to provide sufficient time for students to attempt all items. Only students with extended time accommodation were permitted to exceed the established time limits of any given session. The manuals provide examiners with timing guidelines for the assessments.

Online Forms Administration, Grades 3–8

The online forms are administered via DRC's INSIGHT online assessment system. School system and school personnel set up test sessions via DRC's INSIGHT portal and print test tickets. Students enter their ticket information to access the test in INSIGHT. In addition, students have access to the Online Tools Training (OTT) before the testing window, which allows them to practice using tools and features within INSIGHT. Tutorials with online

video clips that demonstrate features of the system are also available to students before testing.

Paper-Based Forms Administration, Grade 3

Schools with testers at grade 3 had the option to participate in either paper-based or computer-based testing for the spring assessment. DRC prints and ships paper materials to the sites that choose paper-based testing. These materials are returned to DRC after testing for processing and scoring with the online tests.

Accessibility and Accommodations

Accessibility features and accommodations include Access for All, Accessibility Features, and Accommodations.

- Access for All features are available to all students taking an assessment.
- Accessibility Features are available to students when deemed appropriate by a team of educators.
- Accommodations must appear in a student's IEP/IAP/EL plan.

Accommodations may be used with students who qualify under the Individuals with Disabilities Education Act (IDEA) and have an IEP or Section 504 of the Americans with Disabilities Act and have an IAP, or who are identified as English Learners (ELs) and have an EL plan.

Accommodations must be specified in the qualifying student's individual plan and must be consistent with accommodations used during daily classroom instruction and testing. The use of any accommodation must be indicated on the student information sheet at the time of test administration. AERA, APA, and NCME Standard 6.2 states:

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

In compliance with this standard, the TAM contains the list of Universal Tools, Designated Supports, and Accommodations permissible for the LEAP assessments. The following accommodations were provided by DRC for this administration:

- Braille
- Text-to-Speech
- Directions in Native Language

The following additional access and accommodation features were also available:

- Answers Recorded
- Extended Time
- Transferred Answers
- Individual/Small Group Administration
- Tests Read Aloud
- English/Native Language Word-to-Word Dictionary
- Directions Read Aloud/Clarified in Native Language
- Text-to-Speech for online testers
- Human Read Aloud
- Directions in Native Language

For more details about these accommodations, please refer to the <u>LEAP 2025 Accessibility</u> and <u>Accommodations Manual</u>.

Testing Windows

The computer-based testing window was available from April 15 through May 17, 2024. Paper-based testing occurred from April 17 through April 22, 2024.

Test Security Procedures

Maintaining the security of all test materials is crucial to preventing the possibility of random or systematic errors, such as unauthorized exposure of test items that would affect the valid interpretation of test scores. Several test security measures are

implemented for the LEAP 2025 assessments. Test security procedures are discussed throughout the TCMs and TAMs.

Test coordinators and administrators are instructed to keep all test materials in locked storage, except during actual test administration, and access to secure materials must be restricted to authorized individuals only (e.g., test administrators and the school test coordinator). During the testing sessions, test administrators are directly responsible for the security of the LEAP 2025 assessment and must account for all test materials and supervise the test administrators at all times.

Data Forensic Analyses

Due to the importance of the LEAP 2025 assessment, it is prudent to ensure that the results from the assessments are based on effective instruction and true student achievement. To help ensure that scores are related to actual learning and that results are valid, data forensic analyses take place to assist in separating meaningful gains from spurious gains. It is important to note that although the results of the analyses may be used to identify potential problems within a school, the identification of a problem is not an accusation of misconduct.

Multiple methods are incorporated into the forensic analysis. The following methods are applied:

- Response Change Analysis
- Score Fluctuation Analysis
- Web Monitoring
- Plagiarism Detection
- Alerts for Disturbing Content

Response Change Analysis. Students make changes to answer choices when taking the LEAP 2025 assessments, and this behavior is expected. Unfortunately, changes to student answers are sometimes influenced by school personnel who want to improve performance. Therefore, the response change analysis is conducted to identify school-and test administrator-level response change patterns that are statistically improbable when compared to the expected pattern at the state level.

Score Fluctuation Analysis. It is anticipated that performance on the LEAP 2025 assessments will improve over time for legitimate reasons such as changes in the curriculum and improvement in instruction. However, large and unexpected score changes may be a sign of testing impropriety. The LDOE applied an approach where the state's level of change in performance from one year to the next is compared to schools' and test administrators' change in performance during the same time frame. Schools and test administrators are identified when the level of change is statistically unexpected.

Web Monitoring. The content of the LEAP 2025 assessments should not appear outside the boundaries of the forms administered. To protect Louisiana test content, the internet is monitored for postings that contain, or appear to contain, potentially exposed and/or copied test content. When test content is verified, steps are taken to quickly remove the infringing content.

Plagiarism Detection. The LDOE monitors for two different plagiarism situations: copying from student to student and copying from an outside source, such as Wikipedia or another internet source. Instances of plagiarism are identified by human scorers and artificial intelligence. Alerts are set to identify responses that may indicate the possibility of teacher interference or plagiarism. Alerted responses are given additional review so that the appropriate response can be taken.

Alerts for Disturbing Content

Scorers for the LEAP 2025 assessments also have the ability to apply an alert flag to student responses that may indicate disturbing content (e.g., possible physical or emotional abuse, suicidal ideation, threats of harm to themselves or others). All alerted responses are automatically routed to the scoring director, who reviews and forwards appropriate responses to senior project staff for review. If it is concluded that a response warrants an alert, project management will contact the LDOE to take the necessary action. At no point during this process do scorers or staff have access to demographic information for any students participating in the assessment.

6. Scoring Activities

Directory of Test Specification (DOTS) Process. DRC creates a DOTS file, based on the approved test selection. The DOTS is a document containing information about each item on a test form, such as item identifier, item sequence, answer key, score points, subtype, session, alignment, and prior use of item. WestEd reviews and confirms the contents of the DOTS file as part of test review rounds. The DOTS file is then provided to the LDOE for review and final approval. Once approved the information contained in the DOTS is used in scoring the test and in reporting.

Selected-Response (SR) Item Keycheck. SR items for Social Studies include multiple-choice (MC) and multiple-select (MS) questions. Pearson calculates MC and MS item statistics and flags items if item statistics fall outside expected ranges. For example, items are flagged if few students select the correct response (*p*-value less than 0.25), if the item does not discriminate well between students of lower and higher ability (point-biserial correlation less than 0.20), or if many students (more than 40%) select a certain incorrect response. Lists of flagged MC and MS items, with the reasons for flagging, are provided to the LDOE and WestEd content staff for key verification. The staff reviews the list of flagged MC and MS items to confirm that the answer keys are accurate. The scoring of MC and MS items is also evaluated at data review.

Scoring of Technology-Enhanced (TE) Items. All TE items are processed through DRC's autoscoring engine and scored according to the assigned scoring rules established during content creation by WestEd in conjunction with the LDOE. DRC ensures that all rubrics and scoring rules are verified for accuracy before scoring any TE items. DRC has an established adjudication process for TE items to verify that correct answers are identified. DRC's TE scoring process includes the following procedures:

- A scoring rubric is created for each TE item. The rubric describes the one and only correct answer for dichotomously scored items (i.e., items scored as either right or wrong). If partial credit is possible, the rubric describes in detail the type of response that could receive credit for each score point.
- The information from each scoring rubric is entered into the scoring system within the item banking system so that the truth resides in one place along with the item image and other metadata. This scoring information designates specific information that varies by item type. For example, for a drag-and-drop

item, the information includes which objects are to be placed in each drop region to receive credit.

- The information is then verified by another autoscoring expert.
- After testing starts, reports are generated that show every response, how many students gave that response, and the score the scoring system provided for that response.
- The scoring is then checked against the scoring rubric using two levels of verification.
- If any discrepancies are found, the scoring information is modified and verified again. The scoring process is then rerun. This checking and modification process continues until no other issues are found.
- As a final check, a final report is generated that shows all student responses, their frequencies, and their received scores.

In the case of braille and accommodated print test forms, student responses to TE items are transcribed into the online system by a test administrator.

Adjudication. TE items and other eligible items identified in the test map are automatically scored as tests are processed. TE items are scored according to scoring rules in the DOTS, which includes scoring information for all item types.

The adjudication process focuses on detecting possible errors in scoring TE and MS items. DRC provides a report listing the frequency distributions of TE item responses and MS items. Members of the LDOE and WestEd content staff examine the TE and MS response distributions and the auto-frequency reports to evaluate whether the items are scored appropriately. When scoring issues are identified, WestEd content staff and the LDOE recommend changes to the scoring algorithm. Any changes to the scoring algorithm are based on the LDOE's decisions. DRC, in turn, applies the approved scoring changes to any affected items.

Constructed-Response and Extended-Response Scoring

Constructed-response items are scored by human raters trained by DRC. Extended-response items are scored by Project Essay Grade (PEG), an Artificial Intelligence (AI) scoring engine. Ten percent of the responses are scored twice to monitor and maintain

inter-rater reliability. Scoring supervisors also conduct read-behinds and review all nonscores and alerts. Handscoring processing rules are detailed in the LEAP 2025 Spring 2024 Handscoring/AI Documentation document.

Selection of Scoring Evaluators. Standard 4.20 states the following:

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

The following sections explain how scorers are selected and trained for the LEAP 2025 handscoring process and describe how the scorers are monitored throughout the handscoring process.

Recruitment and Interview Process. DRC strives to develop a highly qualified, experienced core of evaluators to appropriately maintain the integrity of all projects. All readers hired by DRC to score 2023–2024 LEAP 2025 test responses had at least a four-year college degree.

DRC has a human resources director dedicated solely to recruiting and retaining the handscoring staff. Applications for reader positions are screened by the handscoring project manager, the human resources director, or recruiting staff to create a large pool of potential readers. In the screening process, preference is given to candidates with previous experience scoring large-scale assessments and with degrees emphasizing the appropriate content areas. At the personal interview, reader candidates are asked to demonstrate their proficiency in writing by responding to a DRC writing topic and their proficiency in mathematics by solving word problems with correct work shown. These steps result in a highly qualified and diverse workforce. DRC personnel files for readers and team leaders include evaluations for each project completed. DRC uses these evaluations to place individuals on projects that best fit their professional backgrounds, their college degrees, and their performances on similar projects at DRC. Once placed, all readers go through rigorous training and qualifying procedures specific to the project on which they are placed. Any scorer who does not complete this training and does not

demonstrate the ability to apply the scoring criteria by qualifying at the end of the process is not allowed to score live student responses.

Security. Whether training and scoring are conducted within a DRC facility or done remotely, security is essential to the handscoring process. When users log into DRC's secure, web-based scoring application, ScoreBoard, they are required to read and accept the security policy before they are allowed to access any project. For each project, scorers are also required to read and sign non-disclosure agreements, and during training emphasis is always given to what security means, the importance of maintaining security, and how this is accomplished.

Readers only have access to student responses they are qualified to score. Each scorer is assigned a unique username and password to access DRC's imaging system and must qualify before viewing any live student responses. DRC maintains full control of who may access the system and which item each scorer may score. No demographic data is available to scorers at any time.

Each DRC scoring center is a secure facility. Access to scoring centers is limited to badge-wearing staff and to visitors accompanied by authorized staff. All readers are made aware that no scoring materials may leave the scoring center. To prevent the unauthorized duplication of secure materials, cell phone/camera use within the scoring rooms is strictly forbidden. Readers only have access to student responses they are qualified to score.

In a remote environment, security reminders are given on a daily basis. Similar to the work that occurs within DRC scoring sites, in a remote environment, education about security expectations is the best way to maintain security of any project materials. DRC requires scorers working remotely to work in a private environment away from other people (including family members). Restrictions are in place that define the hours during the day scorers log into the system. If any type of security breach were to occur, immediate action would be taken to secure materials, and the employee would be terminated. DRC has the same policy within the scoring centers.

Handscoring Training Process. Standard 6.9 specifies:

Those responsible for test scoring should establish and document quality control processes and criteria. Adequate training should be provided. The quality of scoring

should be monitored and documented. Any systematic source of scoring errors should be documented and corrected (118).

Training Material Development. DRC scoring supervisors train scorers using the LDOE-approved training materials. These materials are developed by DRC and LDOE staff from a selection scored by Louisiana educators at rangefinding and include the following:

- Prompts and associated sources
- Rubrics
- Anchor sets
- Practice sets
- Qualifying sets

Training and Qualifying Procedures. Handscoring involves training and qualifying team leaders and evaluators, monitoring scoring accuracy and production, and ensuring security of both the test materials and the scoring facilities. The LDOE reviews training materials and oversees the training process.

Qualifying Standards. Scorers demonstrate their ability to apply the scoring criteria by qualifying (i.e., scoring with acceptable agreement with true scores on qualifying sets). After each qualifying set is scored, the DRC scoring director responsible for training leads the scorers in a discussion of the set.

Any scorer who does not qualify by the end of the qualifying process for an item is not allowed to score live student responses.

Monitoring the Scoring Process. Standard 6.8 states:

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

The following section explains the monitoring procedures that DRC uses to ensure that handscoring evaluators follow established scoring criteria while items are being scored. Detailed scoring rubrics, which specify the criteria for scoring, are available for all constructed- and extended-response items.

Reader Monitoring Procedures. Throughout the handscoring process, DRC project managers, scoring directors, and team leaders review the statistics that are generated daily. DRC uses one team leader for every 10 to 12 readers. If scoring concerns are apparent among individual scorers or if a scorer needs clarification on the scoring rules, team leaders address those issues on an individual basis. DRC supervisors typically monitor one out of five of the scorer's readings, making adjustments to that ratio as needed. If a supervisor disagrees with a reader's scores during monitoring, the supervisor provides retraining in the form of direct feedback to the reader, using rubric language and applicable training responses.

Validity Sets and Inter-Rater Reliability. In addition to the feedback that supervisors provide to readers during regular read-behinds and the continuous monitoring of interrater reliability and score point distributions, DRC also conducts validity scoring using the LDOE-approved validity responses identified by the DRC scoring supervisors during live scoring for newly operational items. Validity responses are inserted among the live student responses.

The validity responses are added to DRC's image handscoring system prior to the beginning of scoring. Validity reports compare readers' scores to predetermined scores and are used to help detect potential room drift as well as individual scorer drift. This data is used to make decisions regarding the retraining and/or release of scorers, as well as the rescoring of responses.

Approximately 10% of all live student responses are scored by a second reader to establish inter-rater reliability statistics for all constructed- and extended-response items. This procedure is called a "double-blind read" because the second reader does not know the first reader's score. DRC monitors inter-rater reliability based on the responses that are scored by two readers. If a scorer falls below the expected rate of agreement, the team leader or scoring director retrains the scorer. If a scorer fails to improve after retraining and feedback, DRC removes the scorer from the project. In this situation, DRC removes all scores assigned by the scorer in question. The responses are then reassigned and rescored.

To monitor inter-rater reliability, DRC produces scoring summary reports daily. DRC's scoring summary reports display exact, adjacent, and nonadjacent agreement rates for each reader. These rates are calculated based on responses that are scored by two readers, and their definitions are included below.

- Percentage Exact (%EX)—total number of responses by reader where scores are the same, divided by the number of responses that were scored twice
- Percentage Adjacent (%AD)—total number of responses by reader where scores are one point apart, divided by the number of responses that were scored twice
- Percentage Nonadjacent (%NA)—total number of responses by reader where scores are more than one point apart, divided by the number of responses that were scored twice

Each reader is required to maintain a level of exact agreement on validity responses and on inter-rater reliability. Additionally, readers are required to maintain a low rate of nonadjacent agreement.

Calibration Sets. DRC pulls calibration responses for items. DRC uses these sets to perform calibration across the entire scorer population for an item if trends are detected (e.g., low agreement between certain score points if a certain type of response is missing from initial training). These calibrations are designed to help refocus scorers on how to properly use the scoring guidelines. They are selected to help illustrate particular points and familiarize scorers with the types of responses commonly seen during operational scoring. After readers score a calibration set, the scoring director reviews it from the front of the room, using rubric language and scoring concepts exemplified by the anchor responses to explain the reasoning behind each response's score.

Reports and Reader Feedback. Reader performance and intervention information are recorded in reader feedback logs. These logs track information about actions taken with individual readers to ensure scoring consistency in regard to reliability, score point distribution, and validity performance. In addition to the reader feedback logs, DRC provides the LDOE with handscoring quality control reports for review throughout the scoring window.

Inter-Rater Reliability. DRC and LDOE have agreed to expectations around inter-rater reliability and validity agreements as shown in Table 6.1.

Table 6.1 *Inter-Rater Reliability for Operational Constructed-Response Items*

Agreement Rate Expectations for Validity and Inter-Rater Reliability – LEAP 2025						
Content Area/Course	Score Point Range	Perfect Agreement	Perfect Agreement + Adjacent			
Grades 3-8 Science	0-2 Rubric	80%	95%			
CR items						
Grades 5-8 Science	0-1 Rubric	90%	100%			
Composite	0-2 Rubric	80%	95%			
(multi-part) ER items	0-3 Rubric	70%	95%			
	0-4 Rubric	70%	95%			
	0-5 Rubric	70%	95%			
	0-6 Rubric	60%	93%			
	0-7 Rubric	60%	93%			
	0-8 Rubric	60%	90%			
Grades 5-8 Science	0-9 Rubric	60%	90%			
Comprehensive						
(single part) ER items						

A minimum of 10% of the responses for constructed- and extended-response items are scored independently by a second reader. This is the case regardless of whether the first reader is a human rater or AI. The statistics for inter-rater reliability are calculated for all items at all grades. To determine the reliability of scoring, the percentage of perfect agreement and adjacent agreement between the first and second scores is examined.

Tables 6.2–6.9 provide the inter-rater reliability and score point distributions by grade level for the constructed-response and extended-response items administered in the spring 2024 forms.

Table 6.2 *Inter-Rater Reliability for Operational Constructed-Response Items*

			Inter-Rate	r Reliability*	
Grade	ltem	2x	Exact Agreement (%)	Adjacent Agreement (%)	Nonadjacent (%)
	ltem 1	≥16,460	88	12	0
3**	ltem 2	≥15,980	88	11	0
	ltem 3	≥16,070	92	8	0
	ltem 1	≥11,570	94	6	0
4	ltem 2	≥16,870	93	7	0
	ltem 3	≥20,510	96	4	0
	ltem 1	≥11,800	97	3	0
5	ltem 2	≥12,340	91	8	1
	ltem 3	≥12,740	90	9	1
	ltem 1	≥12,990	89	10	1
6	ltem 2	≥11,950	90	9	0
	ltem 3	≥11,540	86	14	0
	ltem 1	≥11,800	92	4	4
7	ltem 2	≥13,390	91	9	0
	ltem 3	≥11,760	86	13	0
	ltem 1	≥13,390	89	9	2
8	ltem 2	≥13,280	94	6	1
	ltem 3	≥16,710	90	10	1

^{*} The percent may not add up to 100% due to rounding.

^{**} Grade 3 report combines both online and paper forms.

Table 6.3

Score Point Distributions for Operational Constructed-Response Items

			S	core Point I	Distribution	1*	
Grade	ltem	Total	"0" Rating (%)	"1" Rating (%)	"2" Rating (%)	Blank (%)	Nonscore Codes (%)**
	ltem 1	≥60,240	35	36	14	5	10
3***	ltem 2	≥59,950	43	33	9	5	10
	ltem 3	≥60,000	51	30	4	6	9
	ltem 1	≥54,460	36	47	13	0	2
4	ltem 2	≥56,890	57	20	9	0	14
	ltem 3	≥58,700	45	28	7	0	20
	ltem 1	≥54,160	74	6	15	0	4
5	ltem 2	≥54,280	57	17	21	0	5
	ltem 3	≥54,570	59	21	13	0	5
	ltem 1	≥54,100	78	11	3	0	7
6	ltem 2	≥53,550	73	19	2	0	5
	ltem 3	≥53,330	66	24	5	0	4
	ltem 1	≥53,390	71	7	17	0	4
7	ltem 2	≥54,120	61	20	14	0	5
	ltem 3	≥53,230	32	53	9	0	5
	ltem 1	≥54,610	57	23	13	0	7
8	ltem 2	≥54,470	64	17	11	0	8
	Item 3	≥55,940	51	21	14	0	13

^{*} The percent may not add up to 100% due to rounding.

^{**} Nonscore codes include Foreign language (F), Insufficient (I), Don't Understand (N), Refusal (R), Off Topic (T), and Unintelligible (U). Responses that cannot be assigned a score based on the rubric are assigned a nonscore code and count as zero points toward student scores.

^{***} Grade 3 report combines both online and paper forms.

Table 6.4

Inter-Rater Reliability for Operational Extended-Response Items

		Inter-Rater Reliability*							
Grade	2x	Part	Exact Agreement (%)	Adjacent Agreement (%)	Nonadjacent (%)				
		Part A	94	6	0				
5	≥12,690	Part B	93	6	1				
		Part C	93	6	1				
6	≥15,910	Part A	98	2	0				
O	215,910	Part B	91	6	3				
7	≥15,370	N/A	86	12	2				
8	\16 220	Part A	84	14	2				
0	≥16,220	Part B	79	17	4				

^{*} The percent may not add up to 100% due to rounding.

Table 6.5

Score Point Distributions for Operational Extended-Response Items

				Score Point Distribution*										
														Non
			"0"	"1"	"2"	"3"	"4"	"5"	"6"	"7"	"8"	"9"	Blank	score Codes
Grade	Total	Part	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)	(%)**
		Part A	35	41	7	10							0	5
5	≥54,400	Part B	64	9	14	8							0	5
		Part C	75	8	8	3							0	5
6	>EE 600	Part A	69	17	4								0	10
6	≥55,600	Part B	35	7	14	7	13	5	7	3			0	10
7	≥55,240	N/A	11	9	12	16	16	11	11	2	1	1	0	10
8	>EE 640	Part A	35	25	21	7							0	11
0	≥55,640	Part B	12	14	19	19	14	7	2				0	11

^{*} The percent may not add up to 100% due to rounding.

^{**} Nonscore codes include Foreign language (F), Insufficient (I), Don't Understand (N), Refusal (R), Off Topic (T), and Unintelligible (U). Responses that cannot be assigned a score based on the rubric are assigned a nonscore code and count as zero points toward student scores.

Table 6.6 *Inter-Rater Reliability for Field Test Constructed-Response Items*

			Inter-Rate	r Reliability*	
Grade	ltem	2x	Exact Agreement (%)	Adjacent Agreement (%)	Nonadjacent (%)
2	ltem 1	≥410	95	4	0
3	ltem 2	≥340	88	12	0
4	ltem 1	≥340	95	5	0
5	ltem 1	≥390	92	5	4
6	ltem 1	≥370	88	10	2
6	ltem 2	≥380	88	12	0
	ltem 1	≥340	80	15	5
7	ltem 2	≥380	92	5	4
	Item 3	≥380	82	11	7
	Item 1	≥400	85	15	1
8	Item 2	≥420	87	11	1
	ltem 3	≥340	86	11	3

^{*} The percent may not add up to 100% due to rounding.

Table 6.7

Score Point Distributions for Field Test Constructed-Response Items

			S	core Point I	Distribution	1 *	
Grade	ltem	Total	"0" Rating (%)	"1" Rating (%)	"2" Rating (%)	Blank (%)	Nonscore Codes (%)**
3	ltem 1	≥1,700	44	31	16	2	6
3	ltem 2	≥1,670	65	20	5	7	3
4	ltem 1	≥1,340	64	19	8	0	9
5	ltem 1	≥1,690	64	4	26	0	5
6	ltem 1	≥1,680	39	32	24	0	4
О	ltem 2	≥1,680	28	38	29	0	5
	ltem 1	≥1,670	73	13	12	0	3
7	ltem 2	≥1,680	80	7	8	0	5
	Item 3	≥1,680	64	16	15	0	4
	ltem 1	≥1,690	26	35	33	0	5
8	ltem 2	≥1,700	61	19	12	0	8
	ltem 3	≥1,660	38	32	27	0	2

^{*} The percent may not add up to 100% due to rounding.

^{**} Nonscore codes include Foreign language (F), Insufficient (I), Don't Understand (N), Refusal (R), Off Topic (T), and Unintelligible (U). Responses that cannot be assigned a score based on the rubric are assigned a nonscore code and count as zero points toward student scores.

Table 6.8 *Inter-Rater Reliability for Field Test Extended-Response Items*

			Inter-Rater	Reliability*	
Grade 8	2x	Part	Exact Agreement (%)	Adjacent Agreement (%)	Nonadjacent (%)
		Part A	92	8	0
Item 1	≥480	Part B	88	12	0
		Part C	87	11	2
		Part A	91	8	2
Item 2	≥390	Part B	93	5	2
		Part C	80	16	4
		Part A	83	13	4
Item 3	≥410	Part B	95	5	0
		Part C	90	10	0

^{*} The percent may not add up to 100% due to rounding.

Table 6.9

Score Point Distributions for Field Test Extended-Response Items

				Score Point Distribution*										
Grade 8	Total	Part	"0" (%)	"1" (%)	"2" (%)	"3" (%)	"4" (%)	"5" (%)	"6" (%)	"7" (%)	"8" (%)	"9" (%)	Blank (%)	Non score Codes (%)**
		Part A	49	30	12								0	9
Item 1	≥1,740	Part B	67	19	4	0							0	9
		Part C	34	40	12	4							0	9
		Part A	66	17	10	2							0	5
Item 2	≥1,690	Part B	64	24	6	1							0	5
		Part C	46	16	21	11							0	5
	Part A	48	22	14	6	2						0	7	
Item 3	≥1,700	Part B	73	7	8	4							0	7
		Part C	59	23	11								0	7

^{*} The percent may not add up to 100% due to rounding.

^{**} Nonscore codes include Foreign language (F), Insufficient (I), Don't Understand (N), Refusal (R), Off Topic (T), and Unintelligible (U). Responses that cannot be assigned a score based on the rubric are assigned a nonscore code and count as zero points toward student scores.

7. Data Analysis

Classical Item Statistics

This section describes the classical item analysis for data obtained from the operational LEAP 2025 Science tests. The classical analysis includes statistical analysis based on the following types of items: multiple-choice/multiple-select items, rule-based machine-scored items such as technology-enhanced items, and handscored items such as constructed-and extended-response items. For each operational item, the statistical analysis produces item difficulty (*p*-value) and item discrimination (point-biserial).

Tables and figures that provide the additional information on classical item statistics for the Spring 2024 test can be found in Appendix C: Item Analysis Summary Report. Tables C.1–C.4 show the summaries of classical item statistics. As a measure of item difficulty, p (or "the *p*-value") indicates the average proportion of total points earned on an item. For example, if p = 0.50 on an MC item, then half of the examinees earned a score of 1. If p = 0.50 on a CR item, then examinees earned half of the possible points on average (e.g., 1 out of 2 possible points). A measure of point-biserial correlation indicates a measure of item discrimination. Items with higher item-total correlations provide better information about how well items discriminate between lower- and higher-performing students. It should be also noted that a corrected point-biserial correlation indicates the correlation between an item score and the total test score, where the item score is not included in the total score. The results can be found in Tables C.2–C4. By the way, the statistical analysis results for operational and field test (FT) items are stored in Pearson's Assessment Banking and Building solutions for Interoperable assessment (ABBI) system.

Differential Item Functioning

Differential item functioning (DIF) analyses are intended to statistically signal potential item bias. DIF is defined as a difference between similar-ability groups' (e.g., males or females that attain the same total test score) probability of getting an item correct.

Because test scores can reflect many sources of variation, the test developers' task is to create assessments that measure the intended knowledge and skills without introducing construct-irrelevant variance. When tests measure something other than what they are intended to measure, test scores may reflect those extraneous elements in addition to what the test is purported to measure. If this occurs, these tests can be called biased (Angoff, 1993; Camilli & Shepard, 1994; Green, 1975; Zumbo, 1999). Different cultural and socioeconomic experiences are among some factors that can confound test scores intended to reflect the measured construct.

One DIF methodology applied to dichotomous items was the Mantel–Haenszel (*MH*) DIF statistic (Holland & Thayer, 1988; Mantel & Haenszel, 1959). The MH method is a frequently used method that offers efficient statistical power (Clauser & Mazor, 1998). The *MH* chi-square statistic is

$$MH_{\chi^2} = \frac{\left(\sum_k F_k - \sum_k E(F_k)\right)^2}{\sum_k Var(F_k)},$$

where F_k is the sum of scores for the focal group at the k_{th} level of the matching variable (Zwick, Donoghue, & Grima, 1993). Note that the MH statistic is sensitive to N such that larger sample sizes increase the value of the chi-square.

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

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

where $N_{\rm rlk}$ is the number of correct responses in the reference group at ability level k, $N_{\rm f0k}$ is the number of incorrect responses in the focal group at ability level k, $N_{\rm k}$ is the

total number of responses, N_{flk} is the number of correct responses in the focal group at ability level k, and N_{r0k} is the number of incorrect responses in the reference group at ability level k. The MH DIF statistic is based on a $2\times2\times M$ (2 groups \times 2 item scores \times M strata) frequency table, in which students in the reference (male or white) and focal (female or black) groups are matched on their total raw scores.

The ΔMH DIF is then computed as

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

Positive values of ΔMH DIF indicate items that favor the focal group (i.e., positive DIF items are differentially easier for the focal group); negative values of ΔMH DIF indicate items that favor the reference group (i.e., negative DIF items are differentially easier for the reference group). Ninety-five percent confidence intervals for ΔMH DIF are used to conduct statistical tests.

The MH chi-square statistic and the ΔMH DIF were used in combination to identify operational test items exhibiting strong, weak, or no DIF (Zieky, 1993). Table 7.1 defines the DIF categories for dichotomous items.

Table 7.1

DIF Categories for Dichotomous Items

DIF Category	Criteria
A (negligible)	ΔMH DIF is not significantly different (p <0.05) from 0.0 or is less than 1.0.
B (slight to moderate)	 ΔMH DIF is significantly different (p <0.05) from 0.0 but not from 1.0, andis at least 1.0; OR ΔMH DIF is significantly different (p <0.05) from 1.0 (p <0.05) but is less than 1.5. Positive values are classified as "B+" and negative values as "B"
C (moderate to large)	\mid $\Delta MH\ DIF\mid$ is significantly different (p <0.05) than 1.0 and is at least 1.5. Positive values are classified as "C+" and negative values as "C"

For polytomous items, the standardized mean difference (SMD) (Dorans & Schmitt, 1991; Zwick, Thayer, & Mazzeo, 1997) and the Mantel χ P2P statistic (Mantel, 1963) are used to identify items with DIF. SMD estimates the average difference in performance between the reference group and the focal group while controlling for student ability. To calculate

the SMD, let M represent the matching variable (total test score). For all M = m, identify the students with raw score m and calculate the expected item score for the reference group (E_{rm}) and the focal group (E_{fm}). DIF is defined as $D_m = E_{fm} - E_{rm}$, and SMD is a weighted average of D_m using the weights $w_m = N_{fm}$ (the number of students in the focal group with raw score m), which gives the greatest weight at score levels most frequently attained by students in the focal group.

$$SMD = \frac{\sum_{m} w_{m} (E_{fm} - E_{rm})}{\sum_{m} w_{m}} = \frac{\sum_{m} w_{m} D_{m}}{\sum_{m} w_{m}}$$

The *SMD* is converted to an effect-size metric by dividing it by the standard deviation of item scores for the total group. A negative *SMD* value indicates an item on which the focal group has a lower mean than the reference group, conditioned on the matching variable. On the other hand, a positive *SMD* value indicates an item on which the reference group has a lower mean than the focal group, conditioned on the matching variable.

The *MH DIF* statistic is based on a $2\times(T+1)\times M$ (2 groups \times T+1 item scores \times M strata) frequency table, where students in the reference and focal groups are matched on their total raw scores (T = maximum score for the item). The Mantel χ^2 statistic is defined by the following equation:

Mantel
$$\chi^2 = \frac{\left(\sum_m \sum_t N_{rtm} Y_t - \sum_m \frac{N_{r+m}}{N_{r+m}} \sum_t N_{+tm} Y_t\right)^2}{\sum_m Var(\sum_t N_{rtm} Y_t)}$$
.

The p-value associated with the Mantel χ^2 statistic and the *SMD* (on an effect-size metric) are used to determine DIF classifications. Table 7.2 defines the DIF categories for polytomous items.

Table 7.2

DIF Categories for Polytomous Items

DIF Category	Criteria
A (negligible)	Mantel $\chi^2 p$ -value > 0.05 or $ SMD/SD \le 0.17$
B (slight to moderate)	Mantel $\chi^2 p$ -value < 0.05 and 0.17< SMD/SD \leq 0.25
C (moderate to large)	Mantel χ^2 p-value < 0.05 and $ SMD/SD > 0.25$

Four DIF analyses were conducted for the operational test items only: Female/Male, African American/White, Hispanic/White, and Economically Disadvantaged/Not Economically Disadvantaged. That is, item score data were used to detect items on which female or male students performed unexpectedly well or unexpectedly poorly, given their performance on the full assessment. The same methods were used to detect items on which African American or White students, Hispanic or White students and Economically Disadvantaged or Not Economically Disadvantaged performed unexpectedly well or unexpectedly poorly, given their performance on the full assessment. The last two columns of Tables 7.3.1-7.3.4 provide the number of items flagged for DIF. Items flagged with A-DIF show negligible DIF, items flagged with B-DIF are said to exhibit slight to moderate DIF, and items with C-DIF are said to exhibit moderate to large DIF. None of the operational test items were flagged for C-DIF by any analyses. Note that DIF flags for dichotomous items are based on the *MH* statistics while DIF flags for polytomous items are based on the combination of Mantel χ^2 p-value and SMD statistics.

Table 7.3.1

Summary of Female/Male DIF Flags by Grade

Grade	Α	[B+],[B-]	[C+],[C-]
3	36	[0],[0]	[0],[0]
4	36	[0],[0]	[0],[0]
5	38	[1],[0]	[0],[0]
6	37	[1],[0]	[0],[0]
7	37	[0],[0]	[0],[0]
8	36	[0],[1]	[0],[0]

Table 7.3.2 Summary of African American/White DIF Flags by Grade

Grade	Α	[B+],[B-]	[C+],[C-]
3	36	[0],[0]	[0],[0]
4	36	[0],[0]	[0],[0]
5	39	[0],[0]	[0],[0]
6	37	[0],[1]	[0],[0]
7	36	[0],[1]	[0],[0]
8	36	[0],[1]	[0],[0]

Table 7.3.3

Summary of Hispanic/White DIF Flags by Grade

Grade	А	[B+],[B-]	[C+],[C-]
3	36	[0],[0]	[0],[0]
4	36	[0],[0]	[0],[0]
5	39	[0],[0]	[0],[0]
6	37	[0],[1]	[0],[0]
7	36	[1],[0]	[0],[0]
8	36	[0],[1]	[0],[0]

Table 7.3.4 Summary of Economically Disadvantaged/Not Economically Disadvantaged DIF Flags by Grade

Grade	А	[B+],[B-]	[C+],[C-]
3	36	[0],[0]	[0],[0]
4	36	[0],[0]	[0],[0]
5	39	[0],[0]	[0],[0]
6	38	[0],[0]	[0],[0]
7	37	[0],[0]	[0],[0]
8	37	[0],[0]	[0],[0]

Measurement Models

IRTPRO, a software application for item calibration and test scoring, was used to estimate IRT parameters from LEAP 2025 data. MC, MS, and some TE items (i.e., one-point) were scored dichotomously (0/1), so the three-parameter logistic model (3PL) was applied to those data:

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

In that model, $p_i(\theta_j)$ 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.

Since the Science tests also included polytomous items scored higher than 1 point, the generalized partial credit model (GPCM) (Muraki, 1992) was used to estimate the parameters of these items:

$$p_{im}(\theta_{j}) = \frac{\exp[\sum_{k=0}^{m} Da_{i}(\theta_{j} - b_{i} + d_{ik})]}{\sum_{v=0}^{M_{i}-1} \exp[Da_{i}(\theta_{j} - b_{i} + d_{iv})]'}$$

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 Mi is the number of score categories of item i with possible item scores as consecutive integers from 0 to Mi – 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, etc.).

Calibration and Linking

LEAP 2025 Science assessments are standards-based assessments that have been constructed to align to the LSSS, as defined by the LDOE and Louisiana educators. For each course, the content standards specify the subject matter students should know and the skills they should be able to perform. In addition, performance standards specify how much of the content standards students need to master in order to achieve proficiency. Constructing tests to content standards enables the tests to assess the same constructs from one year to the next.

Item Response Theory (IRT) models were used in the item calibration for the LEAP 2025 Science tests. All calibration activities were independently replicated by Pearson staff as an added quality-control check.

The most common and straightforward way to score a test is to simply use the sum of points a student earned on the test, namely, the raw score. Although the raw score is conceptually simple, it can be interpreted only in terms of a particular set of items. When new test forms are administered in subsequent administrations, other types of derived scores must be used to compensate for any differences in the difficulty of the items and to allow direct comparisons of student performance between administrations. Thus, the primary purpose of form equating is to establish score equivalency between two (or more) forms. Equivalency is established by first building the forms to be equated according to content specifications. Then the form scores are placed on the same scale (by equating), such that students performing on two scaled assessments at the same level of underlying achievement should receive the same scale score on both forms, although they may not receive the same number-correct score (or raw score). LDOE and Pearson strive to maintain equivalent samples or use near-census samples over the years, minimizing the potential differences caused by the different samples.

Tables 7.4.1-7.4.6 provide scale scores at selected percentiles that can be used to compare the distributional characteristics of the Spring 2024 test form to previous administrations. Although these scale scores are rounded values, there were differences in the scale score values for a given percentile across the forms. These variations could arise for several reasons: (1) differences in the proficiency (i.e., achievement) of the students in the samples or growth in student achievement across years; (2) unevenness in the respective distributions that combine with the number-correct-to-scale- score scoring method, leaving "gaps" in the scale; or (3) other sources of equating error. In general, however, the test characteristic function equating techniques will "level" the equated forms through the raw-to-scale- score adjustment.

Table 7.4.1

Comparisons of Scale Scores at Selected Percentiles: Grade 3 Operational Forms

Percentile	2019 Spring Form A	2021 Spring Form A	2022 Spring Form B	2023 Spring Form C	2024 Spring Form D
99	791	787	791	790	789
95	775	773	777	773	773
90	765	762	765	765	765
85	760	755	759	757	760
80	755	750	751	753	756
75	750	745	748	748	751
70	745	740	743	743	747
65	742	734	737	738	742
60	737	731	734	733	737
55	734	725	731	730	735
50	731	722	725	727	730
45	728	719	721	721	727
40	722	715	718	718	724
35	719	712	714	714	718
30	715	703	709	710	714
25	712	698	705	705	711
20	703	693	700	700	702
15	698	687	694	694	698
10	693	679	687	687	692
5	679	669	679	663	677
1	650	650	650	650	650

Table 7.4.2 *Comparisons of Scale Scores at Selected Percentiles: Grade 4 Operational Forms*

Percentile	2019 Spring Form A	2021 Spring Form A	2022 Spring Form B	2023 Spring Form C	2024 Spring Form D
99	798	798	803	809	803
95	782	779	782	789	786
90	774	770	771	776	773
85	766	762	764	770	768
80	764	756	759	765	760
75	758	751	754	757	755
70	753	748	749	754	750
65	751	742	747	749	745
60	748	739	741	744	740
55	743	734	739	739	736
50	740	731	733	737	733
45	737	725	730	732	728
40	734	721	727	729	722
35	728	718	723	723	719
30	725	712	720	720	713
25	722	707	716	717	709
20	716	703	711	714	706
15	708	695	701	706	702
10	704	690	695	701	692
5	690	678	687	690	686
1	668	651	664	672	660

Table 7.4.3

Comparisons of Scale Scores at Selected Percentiles: Grade 5 Operational Forms

Percentile	2019 Spring Form A	2021 Spring Form A	2022 Spring Form B	2023 Spring Form C	2024 Spring Form D
99	807	807	804	813	815
95	788	785	785	791	791
90	776	773	774	779	780
85	768	765	766	771	769
80	762	760	761	763	764
75	757	752	756	758	757
70	752	747	750	752	752
65	747	742	745	745	747
60	745	737	739	739	742
55	740	735	733	737	737
50	735	729	730	731	731
45	732	723	724	725	726
40	726	717	718	719	720
35	723	714	714	715	717
30	717	707	706	708	710
25	714	703	702	700	706
20	707	694	693	695	699
15	698	689	688	690	689
10	689	677	676	677	684
5	677	671	660	670	671
1	654	650	650	650	650

Table 7.4.4

Comparisons of Scale Scores at Selected Percentiles: Grade 6 Operational Forms

Percentile	2019 Spring Form A	2021 Spring Form A	2022 Spring Form B	2023 Spring Form C	2024 Spring Form D
99	797	794	800	793	804
95	779	776	778	773	780
90	769	767	766	763	768
85	763	758	758	756	759
80	758	753	753	749	754
75	753	749	747	745	747
70	749	744	741	740	742
65	744	739	736	735	737
60	742	734	730	730	732
55	736	731	727	725	727
50	734	725	721	722	722
45	728	722	717	716	720
40	725	719	714	713	715
35	722	716	706	709	712
30	719	709	702	706	706
25	712	704	697	698	703
20	709	700	692	693	696
15	704	695	687	687	692
10	695	683	680	681	688
5	683	676	665	664	677
1	657	650	650	650	654

Table 7.4.5

Comparisons of Scale Scores at Selected Percentiles: Grade 7 Operational Forms

Percentile	2019 Spring Form A	2021 Spring Form A	2022 Spring Form B	2023 Spring Form C	2024 Spring Form D
99	809	805	812	802	810
95	786	783	784	783	790
90	775	770	773	774	775
85	767	762	765	765	768
80	759	754	757	760	760
75	754	748	751	754	756
70	751	743	746	750	751
65	746	740	743	745	747
60	743	735	737	740	740
55	737	732	735	735	738
50	735	726	729	730	733
45	729	723	726	728	728
40	726	717	723	722	723
35	723	714	717	716	718
30	717	711	713	713	715
25	714	707	710	706	709
20	707	699	702	702	702
15	703	695	698	694	698
10	695	690	688	689	690
5	685	679	681	677	679
1	662	651	653	650	651

Table 7.4.6

Comparisons of Scale Scores at Selected Percentiles: Grade 8 Operational Forms

Percentile	2019 Spring Form A	2021 Spring Form A	2022 Spring Form B	2023 Spring Form C	2024 Spring Form D
99	803	799	802	802	801
95	784	778	781	785	783
90	773	768	773	774	771
85	766	761	765	767	766
80	761	756	758	759	759
75	756	750	754	755	754
70	752	745	749	750	749
65	747	743	744	745	744
60	743	738	740	740	739
55	741	733	735	737	734
50	736	729	730	732	729
45	731	726	728	727	726
40	729	721	723	724	720
35	723	718	717	717	718
30	721	712	711	714	711
25	715	708	708	710	708
20	708	701	701	706	701
15	705	697	697	698	697
10	697	687	687	693	688
5	682	675	682	680	677
1	658	650	658	662	653

Operational Item Parameters

Appendix C summarizes the distributions of item parameters and provides the graphical displays of the distributions of IRT parameter estimates for each grade. TPI, TPD, CR, and ER items have no *c* parameters because they are polytomous items and are therefore modeled using the GPCM. The number of item parameters associated with the ER items reflect item parameter estimates associated with particular "part scores" that comprise the total ER item. By the way, it should be noted that statistical results of FT items can be found at Pearson ABBI.

Item Fit

IRT scaling algorithms attempt to find item parameters (numerical characteristics) that create a match between observed patterns of item responses and theoretical response patterns defined by the selected IRT models. The Q_1 statistic (Yen, 1981) is used as an index for how well theoretical item curves match observed item responses. Q_1 is computed by first conducting an IRT item parameter estimation, then estimating students' achievement using the estimated item parameters, and, finally, using students' achievement scores in combination with estimated item parameters to compute expected performance on each item. Differences between expected item performance and observed item performance are then compared at 10 selected equal intervals across the range of student achievement. Q_1 is computed as a ratio involving expected and observed item performance. Q_1 is interpretable as a chi-square (χ^2) statistic, which is a statistical test that determines whether the data (observed item performance) fit the hypothesis (the expected item performance). Q_1 for each item type has varying degrees of freedom because the different item types have different numbers of IRT parameters. Therefore, Q_1 is not directly comparable across item types. An adjustment or linear transformation (translation to a Z-score, Z_{Q_i}) is made for different numbers of item parameters and sample size to create a more comparable statistic.

It should be noted that Yen's Q_1 statistic (Yen, 1981) was calculated to evaluate item fit for both operational and field test items by comparing observed and expected item performance. MAP (maximum *a posteriori*) estimates from IRTPRO were used as student ability estimates. For dichotomous items, Q_1 is computed as

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

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

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

where $P_i(\hat{\theta}_a)$ is the item characteristic function for item i and examinee a. The summation is taken over examinees 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 are transformed to ZQ_1 and are compared to a criterion $ZQ_{1,crit}$ to determine whether fit is acceptable. The conversion formulas are

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

and

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

where *df* is the degrees of freedom (the number of intervals minus the number of independent item parameters). Items are categorized as exhibiting either fit or misfit.

A summary of IRT item parameter statistics and item fit for operational items is displayed in <u>Appendix D: Dimensionality</u>.

Dimensionality and Local Item Independence

By fitting all items simultaneously to the same achievement scale, IRT is operating under the assumption that there is a single predominant construct that underlies the performance of all items. Under this assumption, item performance should be related to achievement and, additionally, any relationship of performance between pairs of items should be explained or accounted for by variance in students' levels of achievement. This is the "local item independence" assumption of unidimensional IRT and is associated with a test for unidimensionality called the Q_3 statistic (Yen, 1984).

Computation of the Q_3 statistic starts with expected student performance on each item, which is calculated using item parameters and estimated achievement scores. Then, for each student and each item, the difference between expected and observed item performance is calculated. The difference is the remainder in performance after accounting for underlying achievement. If performance on an item is driven by a predominant achievement construct, then the residual will be small (as tested by the Q_1 statistic), and the correlation between residuals of the item pairs will also be small. These correlations are analogous to partial correlations or the relationship between two variables (items) after accounting for the effects of a third variable (underlying achievement). The correlation among IRT residuals is the Q_3 statistic.

When calculating the level of local item dependence for two items (i and j), the Q_3 statistic is

$$Q_3 = r_{d_i d_j}.$$

The correlation between d_i and d_j values is the correlation of the residuals—that is, the difference between expected and observed scores for each item. For test taker k,

$$d_{ik} = u_{ik} - P_i(\theta_k),$$

where u_{ik} is the score of the kth test taker on item i and $P_i(\theta_k)$ represents the probability of test taker k responding correctly to item i.

With n items, there are n(n-1)/2 Q_3 statistics. If an assessment consists of 48 items, for example, there are 1,128 Q_3 values. The Q_3 values should all be small. Summaries of the distributions of Q_3 are provided in <u>Appendix D: Dimensionality</u>. Specifically, Q_3 data are

summarized by minimum, 5th percentile, median, 95th percentile, and maximum values for LEAP 2025 Science grades 3 through 8. To add perspective to the meaning of Q_3 distributions, the average zero-order correlation (simple intercorrelation) among item responses is also shown. If the achievement construct accounts for the relationships between items, Q_3 values should be much smaller than the zero-order correlations. The Q_3 summary tables in the dimensionality reports in Appendix D show for all grades and subjects that at least 90% (between the 5th and 95th percentiles) of the items are expectedly small. These data, coupled with the Q_1 data, indicate that the unidimensional IRT model provides a reasonable solution to capture the essence of student science achievement defined by the selected set of items for each grade level.

Scaling

Based on the panelist recommendations and LDOE approval, the scale is set using two cut scores, Basic and Mastery, with fixed scale score points of 725 and 750, respectively. The scale scores for Approaching Basic and Advanced vary by grade level. The highest obtainable scale score (HOSS) and lowest obtainable scale score (LOSS) for the scale determined by the LDOE are 650 and 850.

IRT ability estimates (θ s) are transformed to the reporting scale with a linear transformation equation of the form

$$SS = A\theta + B$$
,

where SS is scale score, θ is IRT ability, A is a slope coefficient, and B is an intercept. The slope can be calculated as

$$A = \frac{SS_{Mastery} - SS_{Basic}}{\theta_{Mastery} - \theta_{Basic}},$$

where $\theta_{Mastery}$ is the Mastery cut score on the theta scale, and θ_{Basic} is the Basic cut score on the theta scale. $SS_{Mastery}$ and SS_{Basic} are the Mastery and Basic scale score cuts, respectively. With A calculated, B are derived from the equation

$$SS_{Mastery} = A\theta_{Mastery} + B$$
,

which are rearranged as

$$B = SS_{Mastery} - A\theta_{Mastery} \text{ or } B = SS_{Mastery} - \frac{SS_{Mastery} - SS_{Basic}}{\theta_{Mastery} - \theta_{Basic}} \theta_{Mastery}.$$

Thus, the general equation for converting θ s to scale scores is

$$SS = \left(\frac{SS_{Mastery} - SS_{Basic}}{\theta_{Mastery} - \theta_{Basic}}\right)\theta + \left(SS_{Mastery} - \frac{SS_{Mastery} - SS_{Basic}}{\theta_{Mastery} - \theta_{Basic}}\theta_{Mastery}\right).$$

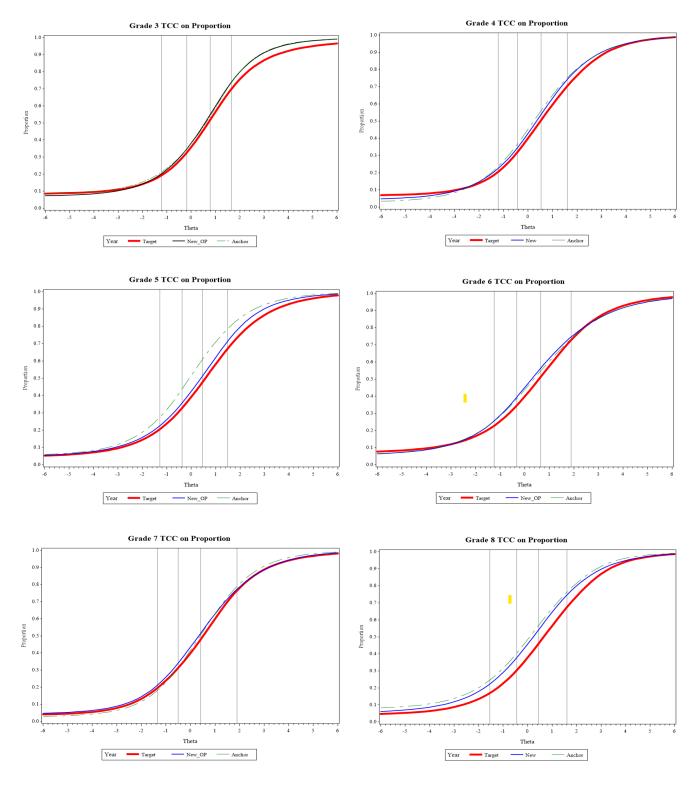
The scaling constants A and B are calculated, and the Advanced cut score and the Approaching Basic cut score on the θ scale are transformed to the reporting scale, rounded to the nearest integer. At this point, the score ranges associated with the five achievement levels are determined. The same scaling constants A and B are used to convert student ability estimates to the reporting scale until new achievement level standards are set. Descriptive Statistics and Frequency Distribution of LEAP 2025 Science Scale Scores can be found in <u>Appendix E: Scale Distribution and Statistical Report</u>.

Test Characteristic Curve

Additional evidence of comparability can be found by reviewing the test characteristic curves (TCCs) across administrations of the LEAP 2025 Science assessments, as can be seen in the following figure. As seen from Plot 7.1 below, the TCCs between two years were similar across ability ranges. By the way, it should be noted that while the base form for all grades was the 2019 operational form, grade 4 used the 2022 operational test form. In addition, although the vertical lines are in theta scale, they indicate performance cuts. Each theta cut corresponding to the scale score of a performance-level cut (e.g., 704, 725, 750, and 778 for grade 4).

Plot 7.1

Test Characteristic Curve



Test Information Curve, Score Distribution, and IRT Difficulty Distribution

In this section, student's Science test score distribution, IRT item difficulty (i.e., b-parameter) distribution, and item information curve are presented. Compared to the base year (i.e., 2023 Science test), the 2024 Science tests generally provide more test information around the lower and middle range of theta than other rages, as can be observed in Tables 7.5.1-7.5.6 and Plot 7.2; it should be noted that the primary goal of 2024 test construction was to improve measurement accuracy for lower performing students. Although the vertical lines are in theta scale, they indicate performance cuts. That is, each theta cut corresponds to the scale score of a performance-level cut.

Table 7.5.1

SPR 2024 Student's Score and IRT B-Parameter Distribution: Grade 3

Percent of Students' Theta	Theta Range	Number of Items of IRT-B
0.00	theta < -3.5	0
1.70	-3.5 ≤ theta < -3.0	0
0.00	-3.0 ≤ theta < -2.5	0
3.78	-2.5 ≤ theta < -2.0	0
3.02	-2.0 ≤ theta < -1.5	0
11.60	-1.5 ≤ theta < -1.0	0
12.32	-1.0 ≤ theta < -0.5	1
18.07	-0.5 ≤ theta < 0.0	5
15.76	0.0 ≤ theta < 0.5	6
13.57	0.5 ≤ theta < 1.0	14
12.36	1.0 ≤ theta < 1.5	5
4.94	1.5 ≤ theta < 2.0	5
2.30	2.0 ≤ theta < 2.5	0
0.43	2.5 ≤ theta < 3.0	0
0.14	3.0 ≤ theta < 3.5	0
0.02	3.5 ≤ theta	0
-3.09	Minimum	-0.57
4.86	Maximum	1.85
-0.02	Mean	0.75
1.13	SD	0.62
≥50,070	Total	36

Table 7.5.2 SPR 2024 Student's Score and IRT B-Parameter Distribution: Grade 4

Percent of Students' Theta	Theta Range	Number of Items of IRT-B
0.43	theta < -3.5	0
0.00	-3.5 ≤ theta < -3.0	0
0.66	-3.0 ≤ theta < -2.5	0
3.10	-2.5 ≤ theta < -2.0	0
9.73	-2.0 ≤ theta < -1.5	1
12.22	-1.5 ≤ theta < -1.0	1
14.70	-1.0 ≤ theta < -0.5	2
14.95	-0.5 ≤ theta < 0.0	5
13.39	0.0 ≤ theta < 0.5	12
11.56	0.5 ≤ theta < 1.0	5
9.77	1.0 ≤ theta < 1.5	7
5.59	1.5 ≤ theta < 2.0	2
2.44	2.0 ≤ theta < 2.5	1
0.86	2.5 ≤ theta < 3.0	0
0.47	3.0 ≤ theta < 3.5	0
0.13	3.5 ≤ theta	0
-3.51	Minimum	-1.54
4.35	Maximum	2.02
-0.13	Mean	0.41
1.19	SD	0.81
≥48,780	Total	36

Table 7.5.3

SPR 2024 Student's Score and IRT B-Parameter Distribution: Grade 5

Percent of Students' Theta	Theta Range	Number of Items of IRT-B
0.00	theta < -3.5	0
1.34	-3.5 ≤ theta < -3.0	0
1.13	-3.0 ≤ theta < -2.5	0
4.01	-2.5 ≤ theta < -2.0	1
8.78	-2.0 ≤ theta < -1.5	1
9.43	-1.5 ≤ theta < -1.0	1
15.56	-1.0 ≤ theta < -0.5	4
13.82	-0.5 ≤ theta < 0.0	10
14.97	0.0 ≤ theta < 0.5	5
12.70	0.5 ≤ theta < 1.0	8
9.55	1.0 ≤ theta < 1.5	8
5.30	1.5 ≤ theta < 2.0	0
2.31	2.0 ≤ theta < 2.5	1
0.59	2.5 ≤ theta < 3.0	0
0.33	3.0 ≤ theta < 3.5	0
0.18	3.5 ≤ theta	0
-3.17	Minimum	-2.12
3.98	Maximum	2.14
-0.17	Mean	0.23
1.23	SD	0.92
≥48,360	Total	39

Table 7.5.4

SPR 2024 Student's Score and IRT B-Parameter Distribution: Grade 6

Percent of Students' Theta	Theta Range	Number of Items of IRT-B
0.70	theta < -3.5	0
0.80	-3.5 ≤ theta < -3.0	0
1.14	-3.0 ≤ theta < -2.5	0
4.22	-2.5 ≤ theta < -2.0	0
10.17	-2.0 ≤ theta < -1.5	1
14.26	-1.5 ≤ theta < -1.0	0
15.97	-1.0 ≤ theta < -0.5	7
13.74	-0.5 ≤ theta < 0.0	6
12.12	0.0 ≤ theta < 0.5	10
10.19	0.5 ≤ theta < 1.0	4
8.36	1.0 ≤ theta < 1.5	3
4.52	1.5 ≤ theta < 2.0	2
2.13	2.0 ≤ theta < 2.5	3
1.07	2.5 ≤ theta < 3.0	1
0.40	3.0 ≤ theta < 3.5	1
0.20	3.5 ≤ theta	0
-3.57	Minimum	-1.65
5.01	Maximum	3.00
-0.31	Mean	0.42
1.27	SD	1.12
≥47,810	Total	38

Table 7.5.5

SPR 2024 Student's Score and IRT B-Parameter Distribution: Grade 7

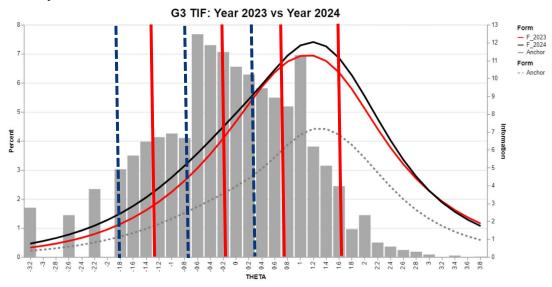
Percent of	Theta Range	Number of
0.66	theta < -3.5	0
0.72	-3.5 ≤ theta < -3.0	0
1.24	-3.0 ≤ theta < -2.5	0
3.74	-2.5 ≤ theta < -2.0	0
8.24	-2.0 ≤ theta < -1.5	0
12.18	-1.5 ≤ theta < -1.0	2
14.69	-1.0 ≤ theta < -0.5	1
16.05	-0.5 ≤ theta < 0.0	7
14.38	0.0 ≤ theta < 0.5	10
12.06	0.5 ≤ theta < 1.0	7
7.51	1.0 ≤ theta < 1.5	5
4.33	1.5 ≤ theta < 2.0	4
2.77	2.0 ≤ theta < 2.5	1
1.01	2.5 ≤ theta < 3.0	0
0.27	3.0 ≤ theta < 3.5	0
0.14	3.5 ≤ theta	0
-3.89	Minimum	-1.24
4.22	Maximum	2.33
-0.23	Mean	0.49
1.25	SD	0.82
≥47,950	Total	37

Table 7.5.6

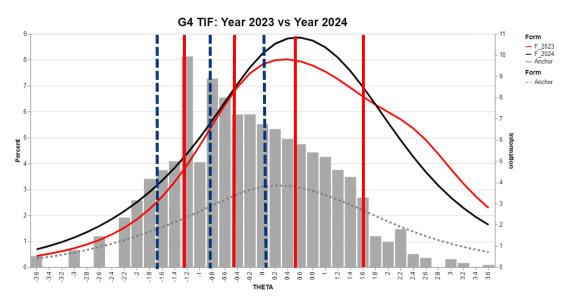
SPR 2024 Student's Score and IRT B-Parameter Distribution: Grade 8

Percent of	Theta Range	Number of
0.00	theta < -3.5	0
1.14	-3.5 ≤ theta < -3.0	0
0.90	-3.0 ≤ theta < -2.5	0
3.30	-2.5 ≤ theta < -2.0	0
8.42	-2.0 ≤ theta < -1.5	1
13.79	-1.5 ≤ theta < -1.0	2
16.68	-1.0 ≤ theta < -0.5	2
14.91	-0.5 ≤ theta < 0.0	3
15.38	0.0 ≤ theta < 0.5	13
10.45	0.5 ≤ theta < 1.0	7
8.47	1.0 ≤ theta < 1.5	6
4.90	1.5 ≤ theta < 2.0	2
1.15	2.0 ≤ theta < 2.5	0
0.34	2.5 ≤ theta < 3.0	1
0.14	3.0 ≤ theta < 3.5	0
0.03	3.5 ≤ theta	0
-3.49	Minimum	-1.63
4.02	Maximum	2.55
-0.26	Mean	0.42
1.14	SD	0.86
≥48,220	Total	37

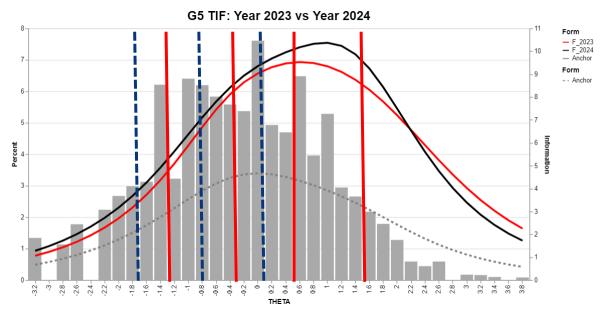
Plot 7.2 Test Information Curve



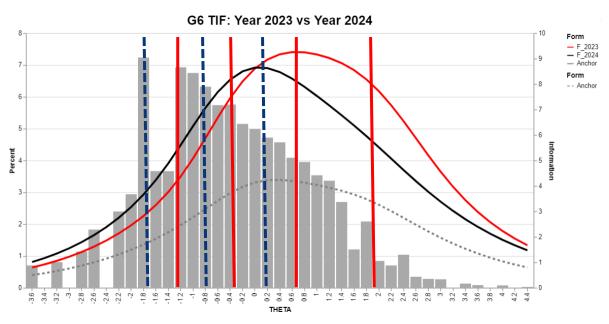
Note: The solid red straight lines representing performance level cuts are on theta scale. The theta cuts correspond to the following scale scores: 698, 725, 750, and 773; The straight dotted lines indicate growth measurement cuts and were created after the non-proficient achievement levels (i.e., Unsatisfactory, Approaching Basic, and Basic) were cut in half.



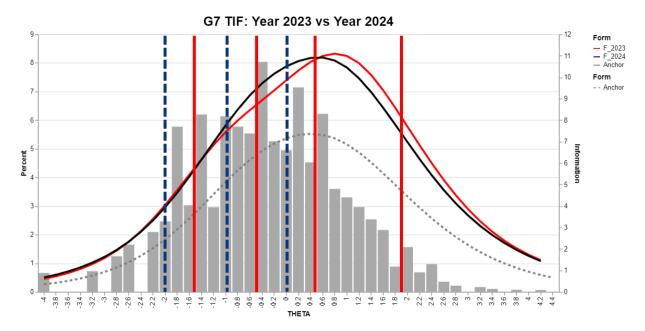
Note: The solid red straight lines representing performance level cuts are on theta scale The theta cuts correspond to the following scale scores: 704, 725, 750, and 778; The straight dotted lines indicate growth measurement cuts and were created after the non-proficient achievement levels (i.e., Unsatisfactory, Approaching Basic, and Basic) were cut in half.



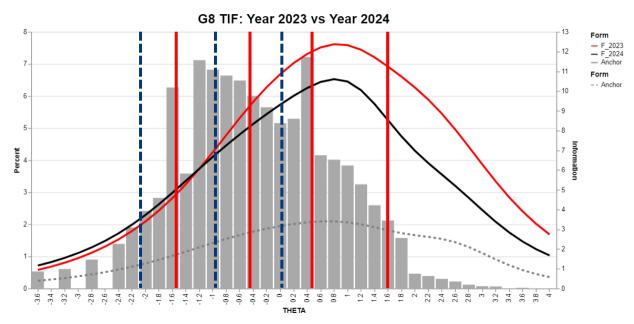
Note: The solid red straight lines representing performance level cuts are on theta scale. The theta cuts correspond to the following scale scores: 698, 725, 750, and 781; The straight dotted lines indicate growth measurement cuts and were created after the non-proficient achievement levels (i.e., Unsatisfactory, Approaching Basic, and Basic) were cut in half.



Note: The solid red straight lines representing performance level cuts are on theta scale. The theta cuts correspond to the following scale scores: 701, 725, 750, and 782; The straight dotted lines indicate growth measurement cuts and were created after the non-proficient achievement levels (i.e., Unsatisfactory, Approaching Basic, and Basic) were cut in half.



Note: The solid red straight lines representing performance level cuts are on theta scale. The theta cuts correspond to the following scale scores: 702, 725, 750, and 790; The straight dotted lines indicate growth measurement cuts and were created after the non-proficient achievement levels (i.e., Unsatisfactory, Approaching Basic, and Basic) were cut in half.



Note: The solid red straight lines representing performance level cuts are on theta scale. The theta cuts correspond to the following scale scores: 694, 725, 750, and 782; The straight dotted lines indicate growth measurement cuts and were created after the non-proficient achievement levels (i.e., Unsatisfactory, Approaching Basic, and Basic) were cut in half.

Field Test Data Review

The process used to complete the field test item equating is an anchored item equating process. In this process the item parameters from the 2024 operational items were fixed as constant (i.e., to calculate Stocking-Lord equating constant) and the item parameters for the field test items were freely calibrated, placing the item parameters for the field test items on the same scale as the operational items.

The data review meeting began with a refresher presentation to data review. The presentation included a review of item statistics (difficulty, discrimination, DIF, score distributions) based on CTT and IRT, appropriate interpretations and inferences, what would be considered reasonable values, and how the values might differ across item types. The result of such reviews is to determine if items are eligible to be placed in the item bank for future test construction or if items need to be updated and field tested again.

While the 2024 FT data review aimed to improve SEM accuracy for lower-performing students, we accepted FT items with slightly lower IRT-a parameters, even if they were relatively easy (e.g., p-values of 0.7 or 0.8), provided there were no content flaws. It should be noted that all the results of spring 2024 data review are saved in Pearson ABBI. It should be noted that the training presentation agenda for data evaluation is included in Appendix A: Training Agendas.

8. Test Results and Score Reports

This section provides the Spring LEAP 2025 Science test results including the scale score and performance levels. Presenting the results by performance level helps translate the numerical scale scores into descriptive categories reflecting student achievement levels (i.e., Level 1: Unsatisfactory, Level 2: Approaching Basic, Level 3: Basic, Level 4: Mastery, and Level 5: Advanced). Tables 8.1–8.6 present evidence of the score reliability and validity for the LEAP 2025 Science 3–8 tests.

Demographic Characteristics of Students

The operational Science tests were administered to all eligible students in the appropriate grade level during Spring 2024. Grade 3 results combine both online and paper forms. Spring 2024 operational score results were based on the following student characteristics:

- Gender
 - o Female
 - Male
- Race
 - African American
 - American Indian or Alaska Native
 - o Asian, Hispanic/Latino
 - Native Hawaiian or Other Pacific Islander
 - Two or More Races
 - White
- Education Classification
- Economic Status
- English Learner (EL)
- Migrant Status
- Homeless Status
- Military Affiliation
- Foster Care Status

Test Results

For the Spring 2024 Science tests, the lowest obtainable scale score (LOSS) on the tests is 650 and the highest obtainable scale score (HOSS) is 850. Scale score means and standard deviations as well as the percentages of students in each performance level are reported for the state and disaggregated into various demographic groups. In addition to the descriptive statistics presented in the following tables, scale score frequency distributions are presented in Appendix E: Scale Distribution and Statistical Report.

Table 8.1

LEAP 2025 State Test Results for Spring 2024: Grade 3

Calamand	Carlo anno ann dala	S	Scale Score				% at Performance Level***					
Category*	Subgroup**	N	Mean	SD	1	2	3	4	5			
Total	Total		729.25	29.31	12	31	31	19	6			
Candar	Female	≥24,800	729.49	28.73	11	31	32	19	6			
Gender	Male	≥25,270	729.01	29.87	13	32	30	19	7			
	African American	≥20,400	719.79	27.41	17	40	29	12	≤5			
	AI/AN	≥260	729.15	26.46	11	31	34	20	≤5			
	Asian	≥760	744.04	31.13	8	17	29	28	19			
Ethnicity	Hispanic/Latino	≥6,000	722.26	29.27	17	36	29	14	≤5			
	NHPI	≥30	741.85	34.4	12	12	32	26	18			
	Two or More	≥2,000	734.92	28.08	9	27	34	23	8			
	White	≥20,580	739.56	27.39	6	23	34	27	11			
Economically	No	≥12,380	745.56	26.66	≤5	17	32	32	15			
Disadvantaged	Yes	≥37,450	724	28.08	15	36	31	15	≤5			
English Learner	No	≥46,880	730.64	29.04	11	30	32	20	7			
English Learner	Yes	≥3,210	708.87	25.39	27	47	21	≤5	≤5			
Education Classification	Regular	≥43,700	731.11	28.97	11	30	32	20	7			
Education Classification	Special	≥6,380	716.52	28.45	21	43	24	10	≤5			
Section 504	No	≥46,210	729.75	29.44	12	31	31	20	7			
Section 504	Yes	≥3,870	723.2	27.07	14	40	30	13	≤5			
Migrapt	No	≥49,990	729.28	29.3	12	31	31	19	6			
Migrant	Yes	≥90	710.92	28.1	29	42	18	8	≤5			
Homeless Status	No	≥48,890	729.57	29.25	12	31	31	19	6			
Homeless status	Yes	≥1,190	716.07	28.84	22	40	26	11	≤5			
Military Affiliation	No	≥49,120	729.02	29.29	12	32	31	19	6			
Military Affiliation	Yes	≥960	740.63	27.94	7	20	33	29	11			
Fostor Caro Status	No	≥49,890	729.28	29.31	12	31	31	19	6			
Foster Care Status	Yes	≥190	720.61	28.97	19	37	25	17	≤5			

^{*} Six students had invalid gender status. 28 students had missing ethnicity status. 247 students lacked economic status information; ** Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander; *** Level 1 = Unsatisfactory. Level 2 = Approaching Basic. Level 3 = Basic. Level 4 = Mastery. Level 5 = Advanced. The overall performance level may not add up to 100% due to rounding.

Table 8.2

LEAP 2025 State Test Results for Spring 2024: Grade 4

Catagografi	Cubananath	Sc	ale Score		% at	Perfo	rmano	e Leve	e ***
Category*	Subgroup**	N	Mean	SD	1	2	3	4	5
Total		≥48,790	732.83	30.98	18	23	28	23	8
Gender	Female	≥24,040	731.49	29.48	18	24	30	22	6
Gender	Male	≥24,740	734.13	32.33	18	22	27	24	9
	African American	≥20,120	721.24	27.04	26	30	27	14	≤5
	AI/AN	≥250	736.98	28.52	14	17	35	28	7
	Asian	≥760	752.88	31.95	7	13	21	37	21
Ethnicity	Hispanic/Latino	≥5,540	724.92	29.74	25	25	28	17	≤5
	NHPI	≥50	743.22	29.33	12	18	22	34	14
	Two or More	≥1,940	737.29	29.72	14	19	32	26	9
	White	≥20,090	745.36	29.86	9	15	29	33	14
Economically	No	≥12,600	751.18	29.56	6	12	26	37	18
Disadvantaged	Yes	≥35,900	726.56	28.83	22	27	29	18	≤5
English Loarner	No	≥45,990	734.27	30.81	17	22	29	24	8
English Learner	Yes	≥2,800	709.27	23.49	42	33	20	≤5	≤5
Education Classification	Regular	≥42,410	735.39	30.47	15	22	29	25	8
Education Classification	Special	≥6,370	715.86	28.96	37	29	21	10	≤5
Costion FO4	No	≥44,320	733.52	31.17	18	22	28	24	8
Section 504	Yes	≥4,470	726	28.16	21	29	28	17	≤5
Migrant	No	≥48,720	732.84	30.99	18	23	28	23	8
Migrant	Yes	≥60	724.4	28.52	23	32	25	15	≤5
Llamalaga Ctatus	No	≥47,740	733.2	30.95	18	23	28	23	8
Homeless Status	Yes	≥1,040	716.13	27.88	35	29	24	11	≤5
Military Affiliation	No	≥47,850	732.52	30.94	18	23	28	23	8
Military Affiliation	Yes	≥940	748.91	28.71	7	12	28	38	16
Fostor Cara Status	No	≥48,600	732.86	30.99	18	23	28	23	8
Foster Care Status	Yes	≥190	724.88	28.7	23	27	31	14	≤5

^{* 23} students had missing ethnicity status. 284 students lacked economic status information.

^{**} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

^{***} Level 1 = Unsatisfactory. Level 2 = Approaching Basic. Level 3 = Basic. Level 4 = Mastery. Level 5 = Advanced. The overall performance level may not add up to 100% due to rounding.

Table 8.3

LEAP 2025 State Test Results for Spring 2024: Grade 5

Catamanut	Colomonath	Sc	Scale Score				% at Performance Level***					
Category*	Subgroup**	N	Mean	SD	1	2	3	4	5			
Total		≥48,360	731.09	36.77	18	25	23	25	9			
Condor	Female	≥23,670	730.75	35.33	17	26	25	24	8			
Gender	Male	≥24,690	731.43	38.1	19	24	22	25	10			
	African American	≥19,880	717.73	32.87	26	33	23	16	≤5			
	AI/AN	≥270	736.73	34.65	11	25	27	27	10			
	Asian	≥850	758.07	37.95	6	13	18	34	29			
Ethnicity	Hispanic/Latino	≥5,380	721.61	36.91	26	27	22	19	≤5			
	NHPI	≥30	745.9	37.38	10	18	21	36	15			
	Two or More	≥1,850	737.52	35.2	12	23	26	28	11			
	White	≥20,050	745.05	34.6	9	18	25	34	14			
Economically	No	≥12,650	752.74	33.88	6	14	23	39	20			
Disadvantaged	Yes	≥35,470	723.56	34.57	22	29	24	20	≤5			
Facilials Lagrange	No	≥45,860	732.88	36.36	17	24	24	26	9			
English Learner	Yes	≥2,500	698.47	27.96	48	34	14	≤5	≤5			
Education Classification	Regular	≥42,470	734.62	35.62	15	24	25	27	10			
Education Classification	Special	≥5,890	705.68	34.92	44	29	14	11	≤5			
Continue FOA	No	≥43,510	732.14	36.91	18	24	24	26	9			
Section 504	Yes	≥4,850	721.77	34.14	24	31	23	17	≤5			
Minusont	No	≥48,300	731.11	36.77	18	25	23	25	9			
Migrant	Yes	≥60	718.54	34.56	31	25	25	17	≤5			
Lla manda na Chatana	No	≥47,310	731.49	36.75	18	25	24	25	9			
Homeless Status	Yes	≥1,050	713.44	33.46	32	30	22	13	≤5			
Military Affiliation	No	≥47,450	730.79	36.77	18	25	23	24	9			
Military Affiliation	Yes	≥900	746.91	33.18	7	18	25	36	15			
Footon Como Status	No	≥48,180	731.14	36.78	18	25	23	25	9			
Foster Care Status	Yes	≥170	719.56	33.3	26	30	23	18	≤5			

^{* 18} students had missing ethnicity status. 238 students lacked economic status information.

^{**} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

^{***} Level 1 = Unsatisfactory. Level 2 = Approaching Basic. Level 3 = Basic. Level 4 = Mastery. Level 5 = Advanced. The overall performance level may not add up to 100% due to rounding.

Table 8.4

LEAP 2025 State Test Results for Spring 2024: Grade 6

Cotomount	Cook and cooking	Sc	ale Scor	е	% at	Perfo	rmano	e Leve	el***
Category*	Subgroup**	N	Mean	SD	1	2	3	4	5
Total	_	≥47,820	725.46	32.58	24	26	27	18	≤5
Gender	Female	≥23,330	724.52	30.88	24	27	28	17	≤5
Gender	Male	≥24,490	726.35	34.09	25	24	26	19	6
	African American	≥20,020	713.57	28.01	35	31	24	9	≤5
	AI/AN	≥260	728.21	29.74	18	26	34	19	≤5
	Asian	≥810	752.54	35.12	8	16	20	37	20
Ethnicity	Hispanic/Latino	≥5,490	718.11	32.61	33	26	24	14	≤5
	NHPI	≥30	730.15	31.05	15	26	31	26	≤5
	Two or More	≥1,670	732.08	31.24	16	26	31	21	6
	White	≥19,500	738.01	31.55	13	20	32	27	8
Economically	No	≥12,740	744.49	31.65	9	16	32	32	11
Disadvantaged	Yes	≥34,810	718.65	30.04	30	29	26	13	≤5
English Lagrage	No	≥45,480	726.91	32.31	23	26	28	19	≤5
English Learner	Yes	≥2,340	697.34	23.86	59	28	10	≤5	≤5
Education Classification	Regular	≥42,590	728.03	32.18	21	25	29	19	≤5
Education Classification	Special	≥5,230	704.51	27.92	50	28	15	6	≤5
Section 504	No	≥42,830	726.38	32.67	23	25	28	19	≤5
Section 504	Yes	≥4,980	717.59	30.66	32	30	23	11	≤5
Minuset	No	≥47,740	725.48	32.57	24	26	27	18	≤5
Migrant	Yes	≥80	716.19	33.19	33	28	23	11	≤5
Lla ma ala an Chahun	No	≥46,790	725.83	32.58	24	26	27	18	≤5
Homeless Status	Yes	≥1,030	708.59	27.43	44	30	17	8	≤5
Military Affiliation	No	≥46,850	725.12	32.5	25	26	27	18	≤5
Military Affiliation	Yes	≥960	741.9	31.87	10	19	32	29	10
Fostor Cara Status	No	≥47,660	725.51	32.57	24	26	27	18	≤5
Foster Care Status	Yes	≥160	710.92	31.23	42	30	18	8	≤5

^{* 14} students had missing ethnicity status. 266 students lacked economic status information.

^{**} AI/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

^{***} Level 1 = Unsatisfactory. Level 2 = Approaching Basic. Level 3 = Basic. Level 4 = Mastery. Level 5 = Advanced. The overall performance level may not add up to 100% due to rounding.

Table 8.5

LEAP 2025 State Test Results for Spring 2024: Grade 7

Category*	Subgroup**	S	cale Scor	е	% at Performance Level***					
		N	Mean	SD	1	2	3	4	5	
Total		≥47,960	732.37	33.91	18	24	28	25	≤5	
Gender	Female	≥23,510	732.8	32.4	16	25	30	25	≤5	
Gender	Male	≥24,450	731.95	35.29	19	23	27	26	6	
	African American	≥20,220	719.8	29.94	26	31	27	14	≤5	
	AI/AN	≥250	735.24	28.78	12	24	32	30	≤5	
	Asian	≥750	758.97	37.18	7	11	20	42	21	
Ethnicity	Hispanic/Latino	≥5,220	725.13	33.89	25	25	27	21	≤5	
	NHPI	≥30	733.94	32.74	11	28	31	25	6	
	Two or More	≥1,670	738.22	32.55	13	22	29	32	6	
	White	≥19,780	745.57	32.14	8	17	29	37	9	
Economically	No	≥13,220	752.13	32.04	6	13	27	41	12	
Disadvantaged	Yes	≥34,470	724.96	31.44	22	28	28	19	≤5	
For ellely Language	No	≥45,880	733.83	33.52	16	23	29	26	≤5	
English Learner	Yes	≥2,080	700.21	25.41	52	32	13	≤5	≤5	
Education Classification	Regular	≥42,820	735.24	33.14	15	23	29	27	6	
Education Classification	Special	≥5,130	708.46	30.62	42	31	17	9	≤5	
Caption 504	No	≥42,790	733.43	33.98	17	23	28	26	≤5	
Section 504	Yes	≥5,170	723.56	32	24	30	26	18	≤5	
Minnest	No	≥47,890	732.38	33.91	18	24	28	25	≤5	
Migrant	Yes	≥70	720.24	28.98	34	18	30	17	≤5	
Llamada a Chahua	No	≥46,960	732.72	33.92	17	24	28	26	≤5	
Homeless Status	Yes	≥1,000	715.66	28.85	30	33	25	12	≤5	
NATIONAL ACCIDENT	No	≥47,040	732.01	33.84	18	24	28	25	≤5	
Military Affiliation	Yes	≥920	750.42	32.23	7	13	27	41	11	
Factor Core Status	No	≥47,810	732.42	33.9	18	24	28	25	≤5	
Foster Care Status	Yes	≥140	713.51	30.97	37	24	28	9	≤5	

^{* 13} students had missing ethnicity status. 267 students lacked economic status information.

^{**} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

^{***} Level 1 = Unsatisfactory. Level 2 = Approaching Basic. Level 3 = Basic. Level 4 = Mastery. Level 5 = Advanced. The overall performance level may not add up to 100% due to rounding.

Table 8.6

LEAP 2025 State Test Results for Spring 2024: Grade 8

Category*	Subgroup**	S	cale Scor	е	%		erforn vel**		2
		N	Mean	SD	1	2	3	4	5
Total		≥48,220	730.2	31.79	14	30	28	23	≤5
Gender	Female	≥23,480	730.53	30.58	12	31	29	22	≤5
dender	Male	≥24,740	729.89	32.89	15	30	27	23	6
	African American	≥20,120	718.42	28.1	20	40	26	12	≤5
	AI/AN	≥270	732.83	29.59	11	30	32	21	6
	Asian	≥800	752.82	34.62	6	15	23	35	21
Ethnicity	Hispanic/Latino	≥5,500	720.86	32.57	23	32	25	17	≤5
	NHPI	≥30	740.44	31.2	≤5	28	31	31	8
	Two or More	≥1,660	736.55	30.24	8	28	29	29	6
	White	≥19,810	743.27	29.38	6	21	31	33	9
Economically	No	≥13,190	748.59	29.02	≤5	17	30	38	12
Disadvantaged	Yes	≥34,740	723.41	29.93	17	36	27	17	≤5
Footish Loomon	No	≥45,900	731.87	31.21	12	30	29	24	6
English Learner	Yes	≥2,320	697.21	24.21	48	38	12	≤5	≤5
Education Classification	Regular	≥43,380	732.68	31.26	12	29	29	24	6
Education Classification	Special	≥4,840	708.02	27.64	33	43	16	7	≤5
Costion FO4	No	≥43,000	731.25	31.89	13	30	28	23	6
Section 504	Yes	≥5,220	721.53	29.53	19	38	26	15	≤5
Migrant	No	≥48,150	730.22	31.78	14	30	28	23	≤5
Migrant	Yes	≥70	716.38	32.13	29	29	26	16	≤5
Homeless Status	No	≥47,260	730.56	31.73	13	30	28	23	≤5
Homeless status	Yes	≥960	712.4	29.36	29	39	22	10	≤5
Military Affiliation	No	≥47,330	729.91	31.75	14	31	28	22	≤5
Military Affiliation	Yes	≥890	745.8	29.59	≤5	20	28	36	11
Foster Care Status	No	≥48,060	730.25	31.78	14	30	28	23	≤5
roster Care Status	Yes	≥160	715.01	28.82	27	36	25	12	≤5

^{* 11} students had missing ethnicity status. 285 students lacked economic status information.

^{**} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

^{***} Level 1 = Unsatisfactory. Level 2 = Approaching Basic. Level 3 = Basic. Level 4 = Mastery. Level 5 = Advanced. The overall performance level may not add up to 100% due to rounding.

Effect Size

One way to evaluate the magnitude of the standardized mean difference (SMD) is to calculate the ES. Cohen's *d* was used to calculate the ES and is given by the following formula:

$$d = \frac{\overline{x_a} - \frac{1}{x_b}}{\sqrt{\frac{(n_a - 1)s_a^2 + (n_b - 1)s_b^2}{(n_a + n_b) - 2}}},$$

where is x_a is the mean score of group A, x_b is the mean score of group B, s_a^2 is the variance of group A, s_b^2 is the variance of group B, n_a is the number of students in group A, and n_b is the number of students in group B.

Cohen's d, then, expresses the difference in group means in terms of the standard deviation. Cohen (1988) offered guidelines for interpreting the meaning of the d statistic: d = 0.20 is a small ES, d = 0.50 is a medium ES, and d = 0.80 is a large ES. Based on Cohen's (1988) guidelines, certain trends are observable in Tables B.6.1–B.6.6. Although no big difference in Science tests was seen between females and males, mean raw scores and ESs show that Asian and White students tend to outperform other ethnicity groups. There were clear performance differences among regular education, gifted/talented education, and special education students in Education Classification and Non-English Learner and English Learner in EL status. Performance differences were also observed from Economically Disadvantaged status, Homeless status, Foster Care status, and Military Affiliation status.

Score Reports

Score reports are the primary means of communicating test scores to appropriate school system personnel (e.g., testing coordinators or superintendents), teachers, and parents. Interpretations of test scores from each administration are disseminated in two ways: the individual score report and the LEAP Interpretive Guide. The LDOE and DRC strive to create documents that will be accessible to parents, teachers, and all other stakeholders. The Individual Student-Level Report (ISR) is the primary means for sharing student test results with parents. As such, it is a standalone document from which parents can glean information that is relevant to understanding their children's test scores. For more

information about the test, parents are provided the <u>Parent Guide to the LEAP 2025</u> Student Reports. In the 2021–2022 administration year, student reports for each school were posted by subject, then downloaded and printed from eDIRECT by the school systems and schools. eDIRECT is DRC's secure online system that provides schools and districts access to student tests and reports.

School Roster Report. A School Roster Report, which provides summary information about student performance on the LEAP 2025 Grades 3–8 Science tests, is available to school systems and schools through eDIRECT. Total test scores and achievement level indicators are shown for the test of interest. Category and subcategory performance ratings are also reported for students. At the school level, the percentage of students at each achievement level and rating by category and subcategory are summarized. More details can be found in the <u>LEAP 2025 Grades 3-8 Interpretive Guide (iGUIDE) Spring 2022</u>.

Individual Student-Level Report. The ISR is another type of report available through the eDIRECT system. ISRs may be downloaded and printed by schools to be sent home to parents. At the top of the page, overall student performance is reported by scale score and achievement level. In the middle of the page, category and subcategory performance indicators are reported. When a student does not receive a scale score, their achievement level will be left blank. ISRs for students whose scores were invalidated will display a blank scale score for a given course.

LEAP 2025 Grades 3-8 Interpretive Guide (iGUIDE) Spring 2022. The LEAP 2025 Grades 3-8 Interpretive Guide (iGUIDE) Spring 2022 was written to help Louisiana school system and school administrators, teachers, parents, and the general public understand the LEAP Science Grades 3-8 tests. The LEAP 2025 Grades 3-8 Interpretive Guide (iGUIDE) Spring 2022 was developed collaboratively by DRC and LDOE staff. LDOE staff had opportunities to review the guide, provide feedback, and give final approval. The elements of the table of contents are provided below:

- Introduction to the Interpretive Guide
 - Overview
 - Purpose of the Interpretive Guide
 - Test Design
 - Scoring
 - Item Types and Scoring
 - o Interpreting Scores and Achievement Levels
 - Scale Score
 - Achievement Level Definitions
 - Student Rating by Reporting Category and Subcategory
- Student-Level Reports
 - o Sample Student Report: Explanation of Results and Terms
 - o Sample Student Report A
 - o Sample Student Report B
 - o Sample Student Report C
 - o Sample Student Report D
- School Roster Report
 - o Sample School Roster Report: Explanation of Results and Terms
 - o Sample Science School Roster Report

Achievement Level Policy Definitions

Achievement level policy definitions for the LEAP 2025 Science tests are shown in Table 8.7. The titles and descriptions of the achievement levels were defined to be part of a cohesive assessment system, and the achievement levels indicate a student's ability to demonstrate proficiency on the LSSS defined for a specific course. The standard-setting section of the LEAP 2025 Science 2018-2019 technical report contains comprehensive information.

Table 8.7

Achievement Level Policy Definitions for LEAP 2025

Achievement Level	Achievement Level Policy Definition
Advanced	Students performing at this level have exceeded college and career readiness expectations and are well prepared for the next level of studies in this content area.
Mastery	Students performing at this level have met college and career readiness expectations and are prepared for the next level of studies in this content area.
Basic	Students performing at this level have nearly met college and career expectations and may need additional support to be fully prepared for the next level of studies in this content area.
Approaching Basic Students performing at this level have partially met college and care readiness expectations and will need much support to be prepared for next level of studies in this content area.	
Unsatisfactory	Students performing at this level have not yet met the college and career readiness expectations and will need extensive support to be prepared for the next level of studies in this content area.

It should be noted that the overall purpose of reporting test results is to communicate information on student performance to stakeholders. These results are presented in the context of score reports that aid the user in understanding the meaning of the test scores. The reports and ancillary information address multiple best practices of the testing industry. Table 8.8 shows the cut of each performance level, and the CSEM for each performance level can be found at Table 9.1 in Chapter 9, Reliability. The standard-setting section of the LEAP 2025 Science 2018–2019 technical report contains comprehensive information.

Table 8.8 *Performance Level Cuts: Science G3-8*

Grade	Approaching Basic	Basic	Mastery	Advanced
3	698	725	750	773
4	704	725	750	778
5	698	725	750	781
6	701	725	750	782
7	702	725	750	790
8	694	725	750	782

Timing Analysis

Since how much time test takers took for each item is available for the CBT administrations, average timing on test by session was calculated for the Spring 2024 Science G3-8 tests. It should be noted that only students with an extended time accommodation were permitted to exceed the established time limits of any given session. The following table summarizes the number of students included in this analysis, the number of items (including both operational and field test items) the students were administered, the average amount of minutes spent across all items by session, and the standard deviation. Since there are extreme test times on both ends (some are very short, and some are very long), the median is included as it is less influenced by these extremes. In this circumstance, it is a more useful description of expected values than the mean.

Table 8.9 SPR 2024 Timing Analysis by Session (Time in Minutes): Science G3-8

Grade	Test Session	Number of Test Items	Total Number of Students	Test Time Mean	Test Time SD	Test Time Median
3	1	18	≥24,000	52.19	17.95	52.27
	2	23	≥24,120	47.25	19.16	46.93
4	1	21	≥46,840	54.55	17.80	53.40
	2	20	≥46,890	54.14	19.04	53.41
5	1	18	≥46,360	51.68	17.12	49.90
	2	19	≥46,660	53.45	17.70	52.55
	3	9	≥34,100	25.57	10.55	24.00
6	1	18	≥45,580	55.29	19.04	53.44
	2	19	≥46,290	47.80	16.84	45.78
	3	9	≥33,650	26.13	10.02	24.98
7	1	18	≥45,300	56.35	81.69	54.26
	2	19	≥46,570	42.93	15.41	41.15
	3	9	≥33,860	24.30	9.92	22.74
8	1	18	≥45,830	51.41	25.90	49.35
	2	19	≥46,720	50.95	17.94	49.63
	3	9	≥33,830	27.23	10.95	25.87

9. Reliability

Internal Consistency Reliability Estimation

Internal consistency methods use data from 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 that require multiple test administrations. One of the most frequently used internal consistency reliability estimates is coefficient alpha (Cronbach, 1951). Coefficient alpha is based on the assumption that inter-item covariances constitute 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 coefficient alpha is

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

where N is the number of items on the test, $s_{\gamma_i}^2$ is the sample variance of the i_{th} item or component, and s_X^2 is the observed score variance for the test. Coefficient alpha is appropriate for use when the items on the test are reasonably homogeneous. The homogeneity of LEAP 2025 Science tests is evidenced through a dimensionality analysis. Dimensionality analyses results are discussed in "<u>Chapter 7. Data Analysis</u>."

The reliability and classification accuracy reports in <u>Appendix F: Reliability and Classification Accuracy</u> provide coefficient alpha and IRT model-based or "marginal reliability" (Thissen, Chen, & Bock, 2003) for the total test.

While coefficient alpha values were between 0.861 and 0.897, the marginal alpha values were between 0.85 and 0.91 for the Science tests. Marginal reliability is described as "an average reliability over levels of θ or theta" (Thissen, 1990). Marginal reliability may be reproduced by squaring and subtracting from 1 each of the 31 "posterior standard deviations" (SEMs) in the IRTPRO output file. Since the variance of the population is 1, each of these values represents the reliability at each of the 31 θ s. Marginal reliability is the

average of these computations weighted by the normal probabilities for each of the 31 quadrature intervals. The formula for marginal reliability is

$$\overline{\rho} = \frac{s_{\theta}^2 - E(SEM_{\theta}^2)}{s_{\theta}^2} ,$$

where s_{θ}^2 is the variance of a given θ (is 1 for standardized θ) and $E(SEM_{\theta}^2)$ is the average error variance or the mean of the squared posterior standard deviations by weighting population density. Marginal reliability can be interpreted in the same way as traditional internal consistency reliability estimates such as coefficient alpha.

Additional reliabilities were calculated on various demographics using the population of students. (Please refer to Table F.1.) Included with coefficient alpha in the tables are the number of students responding to the test, the mean score obtained by this group of students, and the standard deviation of the scores obtained for this group.

Coefficient alpha estimates are computed for the entire test and each subscale by reporting category. Subscore reliability will generally be lower than total score reliability because reliability is influenced by the number of items as well as their covariation. In some cases, the number of items associated with a subscore is small (10 or fewer). Subscore results must be interpreted carefully when these measures reflect the limited number of items associated with the score.

Classical Standard Error of Measurement

The classical standard error of measurement (SEM) represents the amount of variance in a score that results from random factors other than what the assessment is intended to measure. Because underlying traits such as academic achievement cannot be measured with perfect precision, the SEM is used to quantify the margin of uncertainty in test scores. For example, factors such as chance error and differential testing conditions can cause a student's observed score (the score achieved on a test) to fluctuate above or below his or her true score (the student's expected score). The SEM is calculated using both the standard deviation and the reliability of test scores, as follows:

$$SEM = \sigma_x \sqrt{(1 - P'_{xx})}$$

where P'_{xx} is the reliability estimate and σ_x is the standard deviation of raw scores on the test. A standard error provides some sense of the uncertainty or error in the estimate of the true score using the observed score. For example, suppose a student achieves a raw score of 50 on a test with an SEM of 3. Placing a one-SEM band around this student's score would result in a raw score range of 47 to 53. If the student took the test 100 times and 100 similar raw score ranges were computed, about 68 of those score ranges would include the student's true score.

It is important to note that the SEM provides an estimate of the average test score error for all students regardless of their individual proficiency levels. It is generally accepted that the SEM varies across the range of student proficiencies (Peterson, Kolen, & Hoover, 1989). For this reason, it is useful to report test-level SEM, and SEMs for 2024 Science between 3.31 and 3.93, as seen from Table B.4. In addition, SEMs by student group can be found in Appendix F.

Conditional Standard Error of Measurement and Cut Scores

It is important to note that the SEM index provides only an estimate of the average test score error for all students regardless of their individual levels of proficiency. By comparison, conditional standard error of measurement (CSEM) provides a reliability estimate at each score point on a test. Like the SEM, the CSEM reflects the amount of variance in a score resulting from random factors other than what the assessment is designed to measure, but it provides an estimate conditional on proficiency. The CSEM is usually smallest, and thus scores are most reliable, near the middle of the score distribution. Typically, achievement tests included relatively large numbers of moderately difficult items. Because these items are usually well matched to a students' ability, they provide the most reliable estimates of ability. It is desirable, for an achievement test where students are classified into pass/fail categories, that the CSEM be lowest at the cut score for passing. The CSEMs at the four cut scores of each grade that define the performance levels are presented in Table 9.1. The standard-setting section of the LEAP 2025 Science 2018-2019 technical report contains comprehensive information.

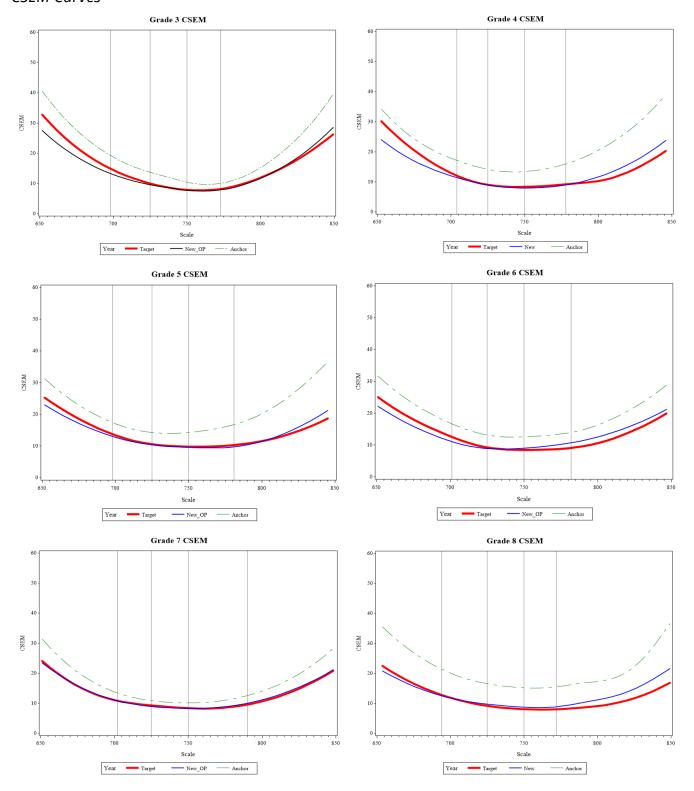
Table 9.1

Conditional Standard Errors of Cut Scores: Spring 2024 LEAP

	Approaching Basic		proaching Basic Basic		Mastery		Advanced	
Grade	Cut Score	CSEM	Cut Score	CSEM	Cut Score	CSEM	Cut Score	CSEM
3	698	13	725	10	750	8	773	8
4	704	12	725	9	750	8	778	9
5	698	13	725	10	750	10	781	10
6	701	11	725	9	750	9	782	11
7	702	11	725	9	750	8	790	10
8	694	13	725	10	750	9	782	10

IRT methods are used for estimating CSEM and are presented in the following graph. With fixed-form assessments, the estimates of measurement error tend to be higher at the low and high ends of the scale-score range (i.e., theta-scale range), where few items measure the ability levels. Generally, there are few students with extreme scores, and these score levels cannot be estimated as accurately as levels toward the middle of the ability range. The middle of the ability range, where cut scores are located, shows lower measurement error than the low and high ends of the ability ranges. Plot 9.1 below demonstrates that irrespective of grades, the tests are designed to minimize measurement error in the middle of the scale-score range, where the majority of students are located. In addition, we can also observe less SEMs for the 2024 assessment in the lower and middle range of the scale scores compared to the 2023 assessment; the primary goal of building the 2024 test forms was to increase less measurement errors for the lower-performing students.

Plot 9.1 *CSEM Curves*



Student Classification Accuracy and Consistency

Students are classified into one of five performance levels based on their scale scores. It is important to know the reliability of student scores in any examination; assessing the reliability of the classification decisions based on these scores is of even greater importance. Classification decision reliability is estimated by the probabilities of correct and consistent classification of students. Procedures were used from Livingston and Lewis (1995) and Lee, Hanson, and Brennan (2000) to derive accuracy and consistency classification measures.

Accuracy of Classification. According to Livingston and Lewis (1995, p. 180), the classification accuracy is "the extent to which the actual classifications of the test takers . . . agree with those that would be made on the basis of their true scores, if their true scores could somehow be known." Accuracy estimates are calculated from crosstabulations between "classifications based on an observable variable (scores on a test) and classifications based on an unobservable variable (the test takers' true scores)." True score is also referred to as a hypothetical mean of scores from all possible forms of the test if they could be somehow obtained (Young & Yoon, 1998).

Consistency of Classification. Classification consistency is "the agreement between classifications based on two non-overlapping, equally difficult forms of the test" (Livingston & Lewis, 1995, p. 180). Consistency is estimated using actual response data from a test and the test's reliability to statistically model two parallel forms of the test and compare the classifications on those alternate forms.

Accuracy and Consistency Indices. Three types of accuracy and consistency indices were generated: *overall, conditional-on-level*, and *cut point*, provided in <u>Appendix F:</u> <u>Reliability and Classification Accuracy.</u> The *overall accuracy* of performance-level classifications is computed as a sum of the proportions on the diagonal of the joint distribution of true score and observed score levels. It is a proportion (or percentage) of correct classification across all the levels. While the overall accuracy indices were between 0.675 and 0.716, the overall consistency indices were 0.566 and 0.613 for the 2024 Science tests.

Another way to express overall consistency is to use Cohen's Kappa (κ) coefficient (Cohen, 1960). The overall coefficient Kappa when applying all cutoff scores together is

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

where P is the probability of consistent classification, and P_c is the probability of consistent classification by chance (Lee, Hanson, & Brennan, 2000). P is the sum of the diagonal elements, and P_c is the sum of the squared row totals. The PChance indices were between 0.217 and 0.246 for the 2024 Science tests.

Kappa is a measure of "how much agreement exists beyond chance alone" (Fleiss, 1973), which means that it provides the proportion of consistent classifications between two forms after removing the proportion of consistent classifications expected by chance alone. The Kappa indices were between 0.425 and 0.495 for the 2024 Science tests.

Consistency conditional-on-level is computed as the ratio between the proportion of correct classifications at the selected level (diagonal entry) and the proportion of all the students classified into that level (marginal entry).

Accuracy conditional-on-level is analogously computed. The only difference is that in the consistency table, both row and column marginal sums are the same, whereas in the accuracy table, the sum that is based on true status is used as a total for computing accuracy conditional on level.

Perhaps the most important indices for accountability systems are those for the accuracy and consistency of classification decisions made at specific cut points. To evaluate decisions at specific cut points, the joint distribution of all the performance levels is collapsed into a dichotomized distribution around that specific cut point.

10. Validity

"Validity is defined as ... the degree to which evidence and theory support the interpretations of test scores entailed by proposed users of tests" (AERA/APA/NCME, 2014). The purpose of test score validation is not to validate the test itself but to validate interpretations of the test scores for particular purposes or uses. Test score validation is not a quantifiable property but an ongoing process, beginning at initial conceptualization and continuing throughout the entire assessment process.

The 2022–2023 LEAP 2025 Science tests were designed and developed to provide fair and accurate scores that support appropriate, meaningful information for educational decisions. The knowledge, expertise, and professional judgment offered by Louisiana educators ultimately ensure that the content of the LEAP 2025 Science tests is an adequate and representative sample of appropriate content, and that the content is a legitimate basis upon which to derive valid conclusions about student achievement.

Chapters 2, 3, and 4 provide a general discussion of test book creation and the editing process, describing the selection of operational test items, the content distribution of embedded field test items, and the process to obtain approvals from the LDOE. The test design process and participation by Louisiana educators throughout the process—from item development, content review, and bias review to test selection—reinforce confidence in the content and design of LEAP 2025 to derive valid inferences about Louisiana student performance. The data review process and results are also discussed. Chapter 5 of the technical report describes the process, procedures, and policies that guide the administration of the LEAP 2025 assessments, including accommodations, test security, and detailed written procedures provided to test administrators and school personnel. Chapter 6 describes scoring processes and activities for the LEAP 2025 Science tests.

Chapter 7 describes classical data analysis and item response theoretic calibration, scaling, and equating methods, as well as processes and procedures to clean data to ensure replicable, iterative calibrations and scaling of the 2024 Science tests to derive scale scores from students' raw scores. Some references to introductory and advanced discussions of IRT are provided. Chapter 7 also describes an analysis of DIF. Complete tables of gender and ethnicity DIF results for all 2024 Science operational items are presented in Appendix C. Chapter 8 of the technical report summarizes the test results,

score distributions, score reports, and achievement level information. Chapter 9 addresses Cronbach's alpha and marginal alpha as measures of internal consistency and describes analysis procedures for classification consistency and classification accuracy. In addition, test validity is addressed in this chapter.

Evidence for Construct-Related Validity

Evidence for construct-related validity—the meaning of test scores and the inferences they support—is the central concept underlying the LEAP 2025 validation process. Validity evidence, from the design of the test to item development and scoring, is created throughout the entire assessment process. Therefore, evidence of validity is described throughout the LEAP 2025 technical report.

Internal Structure of Reporting Categories

The 2024 Science tests contain three reporting categories: *Investigate, Evaluate, and Reason Scientifically.* Table D.1 shows correlations among the reporting categories, and the moderate correlations were observed among the reporting categories; since we used distinct items for each reporting category, a moderate correlation was anticipated.

Content-Related Evidence

Content validity is frequently defined in terms of the sampling adequacy of test items. That is, content validity is the extent to which the items in a test adequately represent the domain of items or the construct of interest (Suen, 1990). Consequently, content validity provides judgmental evidence in support of the domain relevance and representativeness of the content in the test (Messick, 1989). It should be noted that the 2024 Science operational test forms were built exclusively using an ABBI bank program which contained both content and statistical information about both operational and field-tested items.

Dimensionality and Principal Component Analysis

Appendix D: Dimensionality provides information about principal component analysis of the Science tests. Measurement implies order and magnitude along a single dimension (Andrich, 2004). Consequently, in the case of scholastic achievement, a one-dimensional scale is required to reflect this idea of measurement (Andrich, 1988, 1989). However, unidimensionality cannot be strictly met in a real testing situation because students' cognitive, personality, and test-taking factors usually have a unique influence on their test performance to some level (Andrich, 2004; Hambleton, Swaminathan, & Rogers, 1991).

Consequently, what is required for unidimensionality to be met is an investigation of the presence of a dominant factor that influences test performance. This dominant factor is considered as the ability measured by the test (Andrich, 1988; Hambleton et al., 1991; Ryan, 1983).

To check the unidimensionality of the Spring 2024 assessment, the relative sizes of the eigenvalues associated with a principal component analysis of the item set were examined using the Statistical Analysis System (SAS) program. The first and second principal component eigenvalues were compared without rotation. Table D.2 and Plot D.1 summarize the results of the first and second principal component eigenvalues of the assessments. A general rule of thumb in exploratory factor analysis suggests that a set of items may represent as many factors as there are eigenvalues greater than 1 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 with IRT (Loehlin, 1987; Orlando, 2004). As seen from the tables and figures, the first component is substantially larger than the second eigenvalue for the 2024 Science tests.

Item Development and Field-Test Analysis

Test development for LEAP Science tests is ongoing and continuous. Content specialists, teachers from across Louisiana, WestEd/Pearson, and LDOE were greatly involved in developing and reviewing test items. Committees such as content review and bias review reviewed all of the items, which were finally stored in the item bank. Specifically, an internal review by LDOE and WestEd/Pearson staff for alignment and quality required a

great deal of time and energy. More specific information on item (test) development and review can be obtained in Chapter 3, Overview of the Test Development Process.

Various field test forms were used to administer the test items. Once these items were scored, the LDOE and WestEd/Pearson conducted additional item analysis and content review. Any field test items that exhibited statistical results that suggested potential problems were carefully reviewed by both LDOE and WestEd/Pearson content specialists. A determination was then made as to whether an item should be accepted, rejected, and revised/refield-tested. Information on statistical analyses for field test items can be obtained in Chapter 6, Data Analysis.

In summary, additional evidence consistent with the validity, reliability, and consistency of the LEAP 2025 Science assessment has been documented in the LEAP Grades 3–8 Science framework, test development plans, and the 2019 Science standard-setting technical report. Table 10.1 summarizes the sources of validity evidence and indicates where the evidence can be found in the technical report.

Mode Effect Study

It is important to evaluate fairness in test administration in addition to evaluating fairness by examining performance among subgroups. Since two modes (i.e., paper-based tests and computer-based tests) were administered for grade 3, the following techniques (i.e., mode effect analysis and equating) were applied to operational test data to investigate the item mode effect. The mode effect analysis has been conducted, and the results indicate no items exhibiting C category DIF, suggesting no mode effect between online and paper tests; *all items* exhibited *A category DIF*.

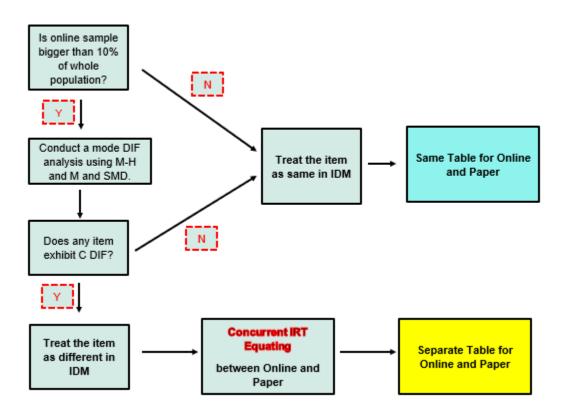


Figure 10.1 General overview of equating, including a mode effect analysis

Table 10.1 Evidence of Validity and the Corresponding Technical Report Chapter

Source of Validity	Related Information	Related Chapter/Source
		Chapter 3
	Item Development Process	LEAP 2025 Grades 3–8 Science
		Assessment Frameworks
	Test Blueprint and Item	Chapters 2 & 3
Evidence-Based on Test	Alignment to Curriculum and	Appendix A
Content	Standards	LEAP 2025 Grades 3–8 Science Assessment Frameworks
	Item Bias, Sensitivity, and Content Appropriateness	Chapter 3
	Accommodations	Chapter 4
	Field Test Analysis	Chapters 3, 7, & 9
Files Books Books	Data Review	LEAP 2025 Grades 3–8 Science
Evidence Based on Response Processes	Data Review	Assessment Frameworks
Processes	Classical Item Analysis	Chapter 7
	IRT Analysis	Chapter 7
	Differential Item Functioning	Chapter 7
Evidence Based on Internal	Reliability and Standard Errors of Measurement	Chapter 9
Structure	Correlation among Reporting Categories	Chapter 9
	Dimensionality Analysis	Chapter 9
Evidence Based on the	Scale Score and Performance Level Information	Chapter 8
Consequences of Testing	Test Interpretive Guide	Chapter 8

References

- American Educational Research Association (AERA), American Psychological Association (APA), & National Council on Measurement in Education (NCME). (2014). Standards for educational and psychological testing. AERA.
- Andrich, A. (1988). *Rasch models for measurement*. Sage Publications.
- Andrich, A. (1989). Distinctions between assumptions and requirements in measurement in the social sciences. In J. A. Keats, R. Taft, R. A. Heath, & H. H. Lovibond (Eds.). *Mathematical and theoretical systems*. ElsevierScience Publisher B.V.
- Andrich, A. (2004). *Modern measurement and analysis in social science*. Murdoch University, Perth. Western Australia.
- Angoff, W. (1993). Perspectives on differential item functioning methodology. In P. W. Holland & H. Warner (Eds.). *Differential item functioning* (pp. 3–24). Lawrence Erlbaum Associates.
- Camilli, G., & Shepard, L. A. (1994). *Methods for identifying biased test items*. Sage Publications.
- Campbell, D. T., & Fiske, D. W. (1959). Convergent and discriminant validation by the multitrait-multimethod matrix. *Psychological Bulletin*, *56*, 81–105.
- Clauser, B. E., & Mazor, K. M. (1998). Using statistical procedures to identify differentially functioning test items. *Educational Measurement: Issues and Practice*, 17, 31–44.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1), 37–47.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, *16*(3), 297–334.

- Dorans, N. J., & Schmitt, A. P. (1991). *Constructed response and differential item functioning: A pragmatic approach* (Research Report No. RR-91-47). Educational Testing Service.
- Fleiss, J. L. (1973). *Statistical methods for rates and proportions*. Wiley.
- Green, D. R. (1975). *Procedures for assessing bias in achievement tests*. Paper presented at the National Institute of Education (NIE) conference on Test Bias, Annapolis, MD
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Sage Publications.
- Holland, P. W., & Thayer, D. T. (1988). Differential item performance and the Mantel–Haenszel procedure. In H. Wainer & H. I. Braun (Eds.), *Test validity* (pp. 129–145). Lawrence Erlbaum Associates.
- Lee, W. C., Hanson, B. A., & Brennan, R. L. (2000). *Procedures for computing classification consistency and accuracy indices with multiple categories*. ACT Research Report Series, 2000(10). ACT.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, *32*(2), 179–197.
- Loehlin, J. C. (1987). Latent variable models. Lawrence Erlbaum Associates.
- Mantel, N. (1963). Chi-square tests with one degree of freedom: Extensions of the Mantel–Haenszel procedure. *Journal of the American Statistical Association, 58*, 690–700.
- Mantel, N., & Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute, 22*(4), 719–748.
- Messick, S. (1989). Meaning and values in test validation: The science and ethics of assessment. *Educational Researcher*, *18*, 5–11.

- Orlando, M. (2004). *Critical issues to address when applying item response theory (IRT) models*. Paper presented at the Drug Information Association, Bethesda, MD.
- Ryan, J. P. (1983). Introduction to latent trait analysis and item response theory. In W. E. Hathaway (Ed.), *Testing in the schools: New directions for testing and measurement* (p. 19). Jossey-Bass.
- Suen, H. K. (1990). *Principles of test theories*. Lawrence Erlbaum Associates.
- Young, M. J., & Yoon, B. (1998). *Estimating the consistency and accuracy of classifications in a standards-referenced assessment*. (CSE Technical Report 475). Los Angeles. Center for the Study of Evaluation, National Center for Research on Evaluation, Standards, and Student Testing. Graduate School of Education & Information Studies, University of California, Los Angeles.
- Zieky, M. (1993). DIF statistics in test development. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 337–348). Lawrence Erlbaum Associates.
- Zumbo, B. D. (1999). A handbook on the theory and methods of differential item functioning (DIF): Logistic regression modeling as a unitary framework for binary and Likert-type (ordinal) item scores. Directorate of Human Resources Research and Evaluation, Department of National Defense.
- Zwick, R., Donoghue, J. R., & Grima, A. (1993). Assessment of differential item functioning for performance tasks. *Journal of Educational Measurement*, *26*, 44–66.
- Zwick, R., Thayer, D. T., & Mazzeo, J. (1997). Descriptive and inferential procedures for assessing differential item functioning in polytomous items. *Applied Measurement inEducation*, *10*(4), 321–344.

Appendix A: Training Agendas

LEAP 2025 Grades 3–8 Item Outline Development Training Agenda

Item Development Cycle for 2019-2024 LEAP 2025 Assessment in Science

- I. Item Development Process
 - a. Overview
 - b. Steps in process
- II. Louisiana Student Standards for Science (LSSS)
 - a. New science standards were approved in early March 2017.
 - The LSSS represent the knowledge and skills needed for students to successfully transition to postsecondary education and the workplace.
 The standards call for students to:
 - 1. Apply content knowledge to real-world phenomena and to design solutions;
 - 2. Demonstrate the practices of scientists and engineers;
 - 3. Connect scientific learning to all disciplines of science; and
 - 4. Express ideas grounded in scientific evidence.
 - b. The Louisiana Student Standards are not the NGSS!
- III. Anatomy of the LSSS
 - a. Descriptor
 - b. Grade level
 - c. Standard
 - d. Domain
 - e. Topic number
 - f. Performance Expectation
 - i. Science and Engineering Practices
 - ii. Disciplinary Core Ideas
 - iii. Crosscutting Concepts
- IV. Outlines

- a. What outlines are
 - i. Definition and purpose
 - ii. Components
- b. What outlines are not
 - i. Characteristics
 - ii. Non-examples
- c. Outline assignments
 - i. Tasks

Components

- a. Stimulus
 - i. Purpose of graphics, data tables, and graphs
 - ii. Reading level
- b. Item types (G3, 4 vs. 5–EOC/Bio)
- c. Bundling of PEs
- ii. Item sets

Components

- a. Stimulus
- b. Item types (G3, 4 vs. 5-EOC/Bio)
- c. Bundling of PEs
- iii. Standalones
 - a. Purpose
 - b. Use of graphics, data tables, and graphs
 - c. Item types
 - d. Single PEs
- iv. Template
- V. Considerations
 - a. Tasks
 - i. Needed number of items and ERs
 - ii. Dimensionality
 - iii. Number of items seen by students vs. number of items developed
 - iv. Use of PEs
 - v. Use of scaffolding within the task
 - b. Item sets
 - i. Needed number of items and ERs
 - ii. Dimensionality

- iii. Interchangeability
- iv. Use of PEs (mix and match)
- v. Number of items seen by students vs. number of items developed
- c. Phenomena list (topics to avoid)
- d. Bias and sensitivity
 - i. Definitions
 - 1. Bias
 - 2. Sensitivity
 - 3. Stereotyping
 - 4. Fairness
 - ii. Rationale for removing bias and sensitivity
 - 1. Portrayal of groups within Louisiana's diverse population
 - 2. Protection of privacy and avoidance of offensive content
 - iii. Potential sources of bias
 - 1. Ethnicity
 - 2. Culture
 - 3. Religion
 - 4. Disability
 - 5. Gender/age stereotypes
 - 6. Geography
 - 7. Socioeconomic status
 - 8. Controversial issues or contexts
 - 9. English language proficiency
 - iv. Strategies to avoid bias
 - Include non-DCI-related information needed to understand stimulus/make stimulus accessible to students regardless of background.
 - 2. Use familiar language and contexts to avoid accessibility bias.
 - 3. Avoid issues and themes that demean, offend, or inaccurately portray any religion, ethnicity, culture, gender, social group, or disability.
 - 4. Avoid topics that will offend the privacy of values and beliefs of students, parents, or the public.

LEAP 2025 Grades 3–8 Item Writer Training Agenda

Item Development Cycle for 2019–2024 LEAP 2025 Assessment in Science

- I. Project Overview:
 - a. Purpose of LEAP project in science
 - b. Characteristics of assessment
 - i. Grade specific, ending the current practice of grade span assessments in grades 4 and 8;
 - Designed to be accessible for use by the widest possible range of students, including but not limited to students with disabilities and English Learners (ELs);
 - iii. Constructed to yield valid and reliable test results while reporting student performance to five achievement levels;
 - iv. Developed and/or reviewed with Louisiana educator and student involvement;
 - v. Non-computer-adaptive; and
 - vi. Administered online.
- II. Louisiana Student Standards for Science (LSSS)
 - a. New science standards were approved in early March 2017.
 - The LSSS represent the knowledge and skills needed for students to successfully transition to postsecondary education and the workplace.
 The standards call for students to:
 - Apply content knowledge to real-world phenomena and to design solutions;
 - 2. Demonstrate the practices of scientists and engineers;
 - 3. Connect scientific learning to all disciplines of science; and
 - 4. Express ideas grounded in scientific evidence.
 - b. The Louisiana Student Standards are not the NGSS!
- III. Anatomy of the LSSS
 - a. Descriptor
 - b. Grade level
 - c. Standard
 - d. Domain

- e. Topic number
- f. Performance Expectation
 - i. Science and Engineering Practices
 - ii. Disciplinary Core Ideas
 - iii. Crosscutting Concepts
- IV. More Acronyms
 - a. SEP key
 - i. 1. Q/P = Asking Questions and Defining Problems
 - ii. 2. MOD = Developing and Using Models
 - iii. 3. INV = Planning and Carrying Out Investigations
 - iv. 4. DATA = Analyzing and Interpreting Data
 - v. 5. MCT = Using Mathematics and Computational Thinking
 - vi. 6. E/S = Constructing Explanations and Designing Solutions
 - vii. 7. ARG = Engaging in Argument from Evidence
 - viii. 8. INFO = Obtaining, Evaluating, and Communicating Information
 - b. CCC key
 - i. PAT = Patterns
 - ii. C/E = Cause and Effect
 - iii. SPQ = Scale, Proportion, and Quantity
 - iv. SYS = Systems and System Models
 - v. E/M = Energy and Matter
 - vi. S/F = Structure and Function
 - vii. S/C = Stability and Change
 - c. "Acronyms Cheat Sheet"
- V. Multidimensional Standards à Multidimensional Assessment
 - a. Dimensions are never to be taught in isolation, and therefore are never tested in isolation.
 - b. The goal of a multidimensional assessment is to gather evidence that a student has proficiency in each of the three dimensions.
 - i. Every item must align to at least two of the three dimensions (with one exception for ERs—"mix and match").
 - ii. Assessment must reflect the different dimensional combinations.
 - 1. SEP and DCL
 - 2. DCI and CCC
 - 3. SEP and CCC (not content)

4. SEP, DCI, CCC

- VI. Aligning to Multiple Dimensions
 - a. SEP
 - i. Develop and model; Analyze data; Construct an explanation
 - b. DCI
 - c. CCC
 - i. Energy and Matter; Patterns; Scale, Proportion, and Quantity
- VII. Phenomena: Keystone of 3-D Assessments
 - a. Phenomena: Observable events that students can use the three dimensions to explain or make sense of
 - i. Links to phenomena websites are available in the "LEAP Phenomena and Context" document.
- VIII. Context: How Phenomena Are Presented
 - a. Contexts are the setting in which phenomena are presented (stimuli).
 - b. A single phenomenon can be presented in many different contexts.
 - c. Phenomena ≠ context; context ≠ phenomena
- IX. Contexts and Stimuli
 - a. Stimuli contain contexts in which phenomena are presented.
 - b. Contexts and stimuli should be unique and novel.
 - i. Non-textbook
 - ii. Think outside the box
 - c. Stimuli must be student friendly and grade appropriate.
 - i. Engaging to students
 - ii. Free of bias and sensitivity issues
 - 1. Definitions
 - a. Bias
 - b. Sensitivity
 - c. Stereotyping
 - d. Fairness
 - 2. Rationale for Removing Bias and Sensitivity
 - a. Portrayal of groups within Louisiana's diverse population
 - b. Protection of privacy and avoidance of offensive content
 - 3. Potential Sources of Bias
 - a. Ethnicity
 - b. Culture

- c. Religion
- d. Disability
- e. Gender/age stereotypes
- f. Geography
- g. Socioeconomic status
- h. Controversial issues or contexts
- i. English language proficiency

4. Strategies to Avoid Bias

- a. Include non-DCI related information needed to understand stimulus/make stimulus accessible to students regardless of background.
- b. Use familiar language and contexts to avoid accessibility bias.
- Avoid issues and themes that demean, offend, or inaccurately portray any religion, ethnicity, culture, gender, social group, or disability.
- d. Avoid topics that will offend the privacy of values and beliefs of students, parents, or the public.
- d. Phenomena, contexts, and stimuli need to be the right grain size.
- e. Goldilocks—provide only the information that is needed
- X. Phenomena and PE Bundles
 - a. PE bundle is usually 2 PEs, but 1-PE and 3-PE bundles are acceptable.
 - b. PE bundling is used in two of the three "item groupings" on LSSS assessment.
 - c. See "Phenomena and Context Overview" and "Contexts and Stimuli" documents for more information.
- XI. Assessment Design: Item Components
 - a. The LSSS assessment will consist of three distinct "components."
 - i. Tasks (PE bundles; phenomena)
 - ii. Item sets (PE bundles; phenomena)
 - iii. Standalone items (single PE only; foci)
- XII. Component: Task
 - a. Tasks (stimulus; four items + ER; dependency OK; phenomenon/PE bundle)
 - b. Tasks include a stimulus and a dependent set of four 1- or 2-point SRs and/or TE items, culminating with one 3-dimensional extended response.
 - c. Items in tasks may require a specific order.

- d. Information in one item may be used in another item (but NOT cue!).
- e. Items may be scaffolded to help discriminate student performance levels.
- f. All items help make sense of or explain a phenomenon.
- g. No CRs
- h. For ER: Can "mix and match" within dimensions from PE bundle as long as the ER aligns with one SEP, one DCI, and one CCC

XIII. Component: Item Set

- a. Item set (stimulus; four items total; CR possible; no inter-item dependency)
 - i. Item sets are composed of a stimulus and four 1- or 2-point SR, TE, and/or CR items.
 - ii. Some item sets will contain one 2-point CR.
 - iii. Item sets without a CR will contain one 2-point TE item (likely an evidence-based selected-response) [EBSR].
 - iv. Items are independent of one another, but all items must depend on the common stimulus.
 - v. Like tasks, the item set makes sense of or explains a phenomenon using a PE bundle. No ERs are included in item sets.

XIV. Component: Standalone Items

- a. Standalone items (single PE; no parts)
 - i. Standalone items will have a "focus" rather than a phenomenon upon which a stimulus is built. This is because a phenomenon is too large to explain or make sense of with one item.
 - ii. Item types include 1- and 2-point formats: no CRs or ERs.
- XV. Item Types: Selected-Response (SR) Formats
 - a. Multiple choice (MC) (1 point)
 - i. Four answer options with one and only one correct answer
 - b. Multiple select (MS) (1 point)
 - i. Five or six answer options with two or three correct answers
- XVI. Item Types: Open-Response Formats
 - a. Constructed response (CR) (2 points)
 - i. Students enter text into a response space
 - ii. Can be two parts
 - iii. Aligns to PE bundle
 - iv. 2-D or 3-D
 - v. Used in item sets ONLY (not all)

- b. Extended response (ER) (grades 3, 4: 6 points; grades 5–EOC: 9 points)
 - i. Students enter text into a response space
 - ii. Can be up to three parts
 - iii. 3-D: Aligns to one SEP, one DCI, and one CCC (mix and match from PE bundle)
 - iv. Can include additional stimulus
 - v. Can reference or depend on previous item in task
 - vi. Used in tasks ONLY

XVII. Item Types:

- a. Technology-enhanced items (TEIs)
 - i. TEIs are worth 1 or 2 points.
 - ii. Used in tasks, item sets, and standalone items
 - iii. TEI types (NO TEIs in grades 3 and 4!)
 - 1. Graphic Gap Match
 - Graphic Gap Match Response Interactions allow graphic gaps and graphic choices. This item type can also be used to create regular gap matches by creating the background in art.

2. Order Interaction

 An Order Interaction Response Interaction consists of choices that may be placed in order or sequence and is a drag-and-drop interaction type. Typically, this interaction type will have three or more choices. The test taker drags the options to the desired order.

3. Hot Spot

 A Hot Spot Response Interaction includes an art image or graphic. The initial state of this item type has no choices selected. This interaction type has a specific set of choices or hot spots that are defined within areas of the art image. One or more choices may be selected in this interaction.

4. Hot Text

 Hot Text Response Interactions include only text. The initial state of this item type has no choices selected. This interaction type has a specific set of hot text selections that are defined within areas of the text. One or more choices may be selected in this interaction.

5. Fill in the Blank (FIB)

- A Text Entry (FIB) Response Interaction includes a freeform field where the test taker enters text, without the ability to use the return or enter key. This interaction will not support multi-line responses.
- b. Evidence-based selected-response (EBSR): Combination of two questions; second question asks students to identify evidence used from the text to support their response to the first question.

XVIII. Development Process Overview

XIX. Universal Design

- a. Ensures that a fair test is developed that provides an accurate measure of what all assessed students know and can do without compromising reliability or validity
 - i. Use consistent naming and graphics conventions;
 - ii. Ensure reading level suitable for the grade level being tested;
 - iii. Replace low-frequency words with simple, common words;
 - iv. Avoid irregularly spelled words, words with ambiguous or multiple meanings, technical terms unless defined and integral to meaning, and concepts with multiple names, symbols, or representations;
 - v. Ensure clarity of noun-pronoun relationships (eliminate pronouns wherever possible);
 - vi. Simplify keys and legends;
 - vii. Use grade-appropriate content; and
 - viii. Avoid differential familiarity for any group, based on language, socioeconomic status, regional/geographic area, or prior knowledge or experience unrelated to the subject matter being tested (bias/sensitivity).
- b. See "Universal Design" for more information.

XX. Item Difficulty

- a. Item difficulty allows students to be placed along a learning progression and assigned to one of the FIVE proficiency levels (to be set at a future date).
 - i. Want a range of difficulty items among each item grouping
 - ii. Cognitive complexity is not difficulty.

b. See "Item Difficulty Overview" for more information.

XXI. Cognitive Complexity*

- a. Need for a range of items of varied cognitive complexity
- b. Existing models of cognitive complexity (e.g., DOK)
- c. Development of a model to address three-dimensional items of LEAP assessment*
- d. (*As the TAGS-M model was in development during the early portion of the 2018–2019 development cycle, item writers used their understanding of cognitive complexity to develop two- and three-dimensional items aligned to the PEs of the LSSS, targeting a broad range of cognitive complexities. These items were then coded by WestEd staff after the TAGS-M model was complete.)

XXII. Sourcing

- a. Sources are required for specific information, such as species, planets, stars, elements, or designs of existing solutions.
 - i. Sources are not needed for commonly known facts.
 - 1. Formula for photosynthesis
 - 2. The definition of speed
 - ii. If in doubt, source!
 - iii. Use reputable sources
 - iv. See "Sources" for more information.

XXIII. Graphics

- a. Graphics are used to convey ideas, data, and/or concepts in a simplified visual form.
 - i. Graphics are essential components of science and include:
 - 1. Tables, diagrams, models, graphs, images
 - ii. All graphics must be introduced appropriately with an introductory statement. Some graphics require only a brief introduction; some require a bit more, e.g.:
 - 1. The students' results are shown in the table below.
 - 2. Students made a scale drawing of their prototype. The scale drawing is shown below.
 - iii. Be aware that some graphics may be changed during production to control for colorblindness.

- iv. See "General Guidelines for Graphics" document for more information.
- v. Style guide
- XXIV. Development Process Overview
- XXV. Information Security
 - a. Do NOT email!
 - b. We will send/receive items and assignments using a secure system.
 - c. General questions about processes OK

LEAP 2025 Grades 3–8 Editor Training Agenda

Item Development Cycle for the 2018–2024 LEAP 2025 Science Assessment

- I. Item Set/Task/Standalone Item Overview
 - a. Criteria for review
- II. Item Development Process
 - a. One round of items slated for development in 2018–2019
 - b. All batches will go through four rounds of LDOE review at different stages of development before committee:
 - i. Outline review (item descriptions; graphic roughs)
 - ii. Item development
 - 1. R1 (fully fleshed-out items; functional TE items; graphics; sources)
 - 2. R2 (implementation of LDOE feedback; rewrites possible; revisions expected)
 - 3. R3 (final look before committee review—no editing, all comments are for committee review)
 - c. Committee review
- III. Process Overview for Intake/E1
- IV. Intake/E1 Rules for Returning Item Sets/Tasks/Standalone Item Submissions to Writers
- V. Feedback to Writers
- VI. Process Overview for Intake/E2
- VII. Intake/E1 Rules for Returning Item Sets/Tasks/Standalone Item Submissions to E1 Writer
- VIII. Use of the Style Guides
 - a. Social Studies/Science Content Style Guide
 - b. TEI Guide
 - c. Graphics Style Guide

LEAP 2025 Biology and Grades 3-8 Content and Bias Item Review Committee Training Agenda

Item Development Cycle for the 2022-2024 LEAP Science Assessment

- I. Welcome from LDOE
- II. Introductions
- III. Non-Disclosure Agreement
 - a. Test security and student confidentiality are of utmost importance to WestEd and the Louisiana Department of Education.
 - b. As a participant in the Science Content/Bias Item Review Meetings, you will have access to materials that must be regarded as secure.
 - c. All materials must be treated as confidential. You are not to disclose the content of these materials or copy or reproduce any of the materials, directly or indirectly.
 - d. By signing and submitting the form, you confirmed that you agree to adhere to these guidelines.
- IV. LEAP Test Development Process
- V. Purpose of Content and Bias Item Review
 - a. To ensure high-quality science tests that:
 - i. Reflect instructionally relevant content
 - ii. Provide valid information to students, parents, teachers, administrators, policymakers, and the public
 - iii. Are fair and appropriate for all students
- VI. What to Consider
 - a. Louisiana Student Standards for Science
 - b. Performance Expectation and the Phenomenon
 - c. Science Shifts
 - d. Components
 - i. Tasks
 - a) Based on a common stimulus
 - b) Items follow a prescribed order; items build on one another
 - c) For field testing, different versions of items included culminating with an extended-response (ER) item

ii. Item Sets

- a) Based on a common stimulus
- b) Items are not in a prescribed order
- c) 4 items on operational test; may have a constructedresponse (CR) item
- d) For field testing, extra items included (12 items developed to get 4)

iii. Standalone Items

VII. Item Types

VIII. Content alignment

- a. Alignment is the key element of content review.
 - i. Is the item providing an appropriate measure of the PE and its related dimensions?
 - ii. Item content alignment is the degree to which an item measures the intended PE and its related dimensions.
 - iii. Put another way: An item is determined to be aligned if the item allows the student to provide evidence of his or her understanding of the specified PE and its related dimensions.
- b. Additional considerations include:
 - i. Scoring/key accuracy
 - ii. Scientific accuracy
- IX. Principles of LSSS for Science Alignment
 - a. Items must be aligned to at least two of the three dimensions.
 - b. Multiple aspects of the item and the item's alignment need to be considered.
 - c. Relative degrees of alignment need to be evaluated.
 - d. Holistic (not analytic) judgments are used to determine acceptable alignment.
- X. Bias and Sensitivity Review
 - a. Items and stimuli should be free of bias and sensitivity concerns.
 - b. This helps to provide students with a fair opportunity to demonstrate their knowledge or skills, regardless of their backgrounds.
 - c. Bias is the presence of some language or content that prevents some members of a group from showing us their knowledge or skills in a particular content area.

- i. Result: Two individuals of the same ability but from different groups perform differently.
- d. What is sensitivity?
- e. Any reference in a stimulus or item that might cause a student to have an emotional reaction and prevent the student from showing us their knowledge and skills for a particular content area.
 - i. Result: Two individuals of the same ability but from different groups perform differently.
- f. If there are bias or sensitivity concerns for an item, the reviewer should be able to point to one of these areas as an area of concern.
 - i. Opportunity and Access
 - a) Problems:
 - i.) Not all Louisiana students have had the opportunity to visit different regions of the world, the US, or Louisiana.
 - ii.) Some students have stronger science skills than English skills.
 - b) Possible solutions:
 - i.) Include non-DCI information that makes a stimulus accessible to students from all backgrounds.
 - ii.) Avoid regional language or words with different meanings in different groups.
 - iii.) Avoid idioms and figurative language.
 - ii. Portrayal of Groups Represented
 - a) Problem:
 - i.) A group is stereotyped (portrayed consistently in a particular way, which may be offensive to members of that group).
 - b) Possible solution:
 - i.) Avoid issues and themes that demean, offend, or inaccurately portray a group, culture, ethnicity, disability.
 - iii. Protecting Privacy and Avoiding Offensive Content
 - a) Problem:

- i.) Some issues and contexts are controversial to particular groups.
- b) Possible solution:
 - i.) Avoid topics that will offend the privacy, values, and/or beliefs of students, parents, and the public.
- XI. Cognitive Complexity and Difficulty
 - a. Cognitive complexity ≠ difficulty
 - b. Cognitive complexity refers to the type and level of thinking and reasoning required of students to answer a test question.
 - c. Difficulty refers to the amount of time and/or effort needed to answer a test question (easy or hard) and can be measured in percentage answering question correctly.
 - d. Task Analysis Guide in Science (Tekkumru-Kisa, Stein & Schunn, 2014)—focused on instruction
 - e. Modified TAGS model is a tool for coding 2- and 3-dimensional items
 - f. Cognitive Complexity in TAGS model

XII. Content Review Decisions

- a. Yes ("Accept")
 - i. Item is acceptable as is
 - ii. Aligned
 - iii. Scientifically accurate
 - iv. Scoring information correct
 - v. Free of bias concerns
- b. No ("Accept with Edits" or "Reject")
 - i. Due to content concerns
 - ii. Metadata alignment with explanation
 - iii. Science accuracy concern with explanation
 - iv. Due to bias concerns
 - v. With explanation
- c. Reject when:
 - i. Complete alignment mismatch
 - ii. Unfixable context flaws
- d. Revise when:
 - i. Fixes can be made
 - ii. Item Alignment Information

XIII. Reviewing Items

- a. Review items in ABBI online
- b. Your facilitator will walk you through a few items to help you learn how to use this tool.
- c. Use the Review Tool for alignment decisions
- d. Vote in ABBI
- e. You will select from:
 - i. Accept
 - ii. Accept with Edits
 - iii. Reject
- f. "Accept with Edits" or "Reject" require comments/justification

XIV. Logistics

- a. Breaks will be announced by the facilitator
- b. ABBI access will be locked during non-meeting times
- c. Room will be locked over lunch
- d. At the conclusion of the meeting, you will receive email communications about:
 - i. Stipend
 - ii. Substitute Reimbursement Form
 - iii. Evaluation survey

LEAP 2025 Grades 3–8 Data Review Training Agenda

- I. What is a Data Review?
 - a. Statistical Definition: Classical Test Theory
 - 1. P-value
 - 2. Point-Biserial
 - 3. Option/Distribution Analysis
 - 4. Differential Item Function (DIF)
 - 5. Flagging Value

Statistics	Flagging Value
P-value	≤ 0.25 or > 0.95
Omit Percentage	> 4%
Point-biserial Correlation	< 0.20
Distractor Percentage	> 400/
(MC only)	- > 40%
Distractor Point-biserial Correlation (MC only)	> 0.00
DIF	B, C

- b. Statistical Definition: Item Response Theory (IRT)
 - 1. IRT Discrimination (a-parameter)
 - 2. IRT Difficulty (b-parameter)
 - 3. IRT Guessing (c-parameter)
 - 4. Q1 (Zq1)
 - 5. Item Fit Plot
 - 6. Flagging Value

Flagging Value for IRT Item Parameters									
a (Discrimination)	b (Difficulty)	c (Guessing)							
< 0.30	Lowe than -3.0 or Higher than 3.0	> 0.35							

- II. Judgement Task in ABBI
 - a. Accept
 - b. Accept with Edits
 - c. Reject

Appendix B: Test Summary

Science G3-8

Contents

Table B.1 Percentage of Points by Reporting Category (includes Task Items): Spring 2024 Operational SC G3-8

Table B.2 Standard Coverage: Spring 2024 Operational SC G3-8

Table B.3 Item Type Summary: Spring 2024 Operational SC G3-8

Table B.4 Raw Score Summary: Spring 2024 Operational SC G3-8

Table B.5 Raw Score Summary by Reporting Category: Spring 2024 Operational SC G3-8

Table B.6 Scale Score and Raw Score Summary: Spring 2024 Operational SC G3-8

Table B.1

Percentage of Points by Reporting Category (includes Task Items): Spring 2024 Operational SC G3-8

Reporting Category	G3	G4	G5	G6	G7	G8
N/A	8.0%	11.5%	-	-	3.3%	-
1 Investigate	28.0%	30.8%	19.7%	19.7%	14.8%	31.7%
2 Evaluate	46.0%	13.5%	47.5%	24.6%	31.1%	25.0%
3 Reason Scientifically	18.0%	44.2%	32.8%	55.7%	50.8%	43.3%

^{*} N/A indicates no reporting category.

Table B.2 Standard Coverage: Spring 2024 Operational SC G3-8

Grade 3

			No.	of Itei	ms		
Reporting Categ	ories	TPI	TPD	MS	МС	CR	% of Test
		N	N	N	N	N	
N/A	3-ESS2-2			1	1	1	8.33
	Sub-Total			1	1	1	8.33
1 Investigate	3-PS2-1		1		2		8.33
	3-PS2-2		1		1		5.56
	3-PS2-3				1		2.78
	3-PS2-4		1		2	1	11.11
	Sub-Total		3		6	1	27.78
2 Evaluate	3-ESS2-1		1		1		5.56
	3-ESS3-1		1		3		11.11
	3-LS2-1	1			1		5.56
	3-LS3-1	1			2		8.33
	3-LS4-1	1			2		8.33
	3-LS4-3				1	1	5.56
	3-LS4-4				1		2.78
	Sub-Total	3	2		11	1	47.22
3 Reason Scientifically	3-LS1-1	1			1		5.56
	3-LS3-2		2		1		8.33
	3-LS4-2				1		2.78
	Sub-Total	1	2		3		16.67
Total	-	4	7	1	21	3	100.00

^{*} N/A indicates no reporting category.

Grade 4

			N	lo. of	ltems			
Reporting Catego	ries	TPI	TPD	TEI	MS	МС	CR	% of Test
		N	N	N	N	N	N	
N/A	4-ESS2-1		1			1		5.56
	4-ESS3-1	1		1				5.56
	Sub-Total	1	1	1		1		11.11
1 Investigate	4-ESS2-1			1		1		5.56
	4-ESS2-3					2	1	8.33
	4-PS3-2		2		1	2		13.89
	4-PS3-3	1			1			5.56
	Sub-Total	1	2	1	2	5	1	33.33
2 Evaluate	4-ESS2-2	1						2.78
	4-LS1-1	2		1				8.33
	Sub-Total	3		1				11.11
3 Reason Scientifically	4-ESS1-1		1			1		5.56
	4-ESS3-2	1		1			1	8.33
	4-LS1-2					2		5.56
	4-PS3-1		1			1		5.56
	4-PS3-4					1		2.78
	4-PS4-1		1			1		5.56
	4-PS4-2	1				2	1	11.11
	Sub-Total	2	3	1		8	2	44.44
Total		7	6	4	2	14	3	100.00

^{*} N/A indicates no reporting category.

Grade 5

				No. o	of Ite	ms			
Reporting Category	ories	TPI	TPD	TEI	MS	МС	ER	CR	% of Test
		N	N	N	N	N	N	N	
1 Investigate	5-LS1-1			1				1	5.41
	5-PS1-3					1		1	5.41
	5-PS1-4			1		3			10.81
	Sub-Total			2		4		2	21.62
2 Evaluate	5-ESS1-1			1	1				5.41
	5-ESS1-2		1	1					5.41
	5-ESS2-2	1		1				1	8.11
	5-PS1-2		1	2		1	1		13.51
	5-PS2-1		1			1			5.41
	Sub-Total	1	3	5	1	2	1	1	37.84
3 Reason Scientifically	5-ESS2-1			4		1			13.51
	5-ESS3-1			1		1			5.41
	5-LS2-1	1				1			5.41
	5-PS1-1	1		3					10.81
	5-PS3-1			1		1			5.41
	Sub-Total	2		9		4			40.54
Total		3	3	16	1	10	1	3	100.00

 $[\]ensuremath{^{\star}}$ N/A indicates no reporting category.

Grade 6

				No. o	of Ite	ms			
Reporting Cate	egories	TPI	TPD	TEI	MS	МС	ER	CR	% of Test
		N	N	N	N	N	N	N	
1 Investigate	6-MS-LS1-1				1			1	5.41
	6-MS-PS2-2		1		1				5.41
	6-MS-PS2-3					2			5.41
	6-MS-PS2-5			1				1	5.41
	Sub-Total		1	1	2	2		2	21.62
2 Evaluate	6-MS-ESS1-3					1			2.70
	6-MS-ESS3-4					2		1	8.11
	6-MS-LS2-1		1	1					5.41
	6-MS-PS2-4					1			2.70
	6-MS-PS3-1			1					2.70
	6-MS-PS4-1			1		1			5.41
	Sub-Total		1	3		5		1	27.03
3 Reason Scientifically	6-MS-ESS1-1		1	2					8.11
	6-MS-ESS1-2					2			5.41
	6-MS-LS1-2			1		1			5.41
	6-MS-LS2-2		1		2		1		10.81
	6-MS-LS2-3	1				1			5.41
	6-MS-PS1-1			1					2.70
	6-MS-PS2-1					1			2.70
	6-MS-PS3-2	1		1					5.41
	6-MS-PS4-2		1			1			5.41
	Sub-Total	2	3	5	2	6	1		51.35
Total		2	5	9	4	13	1	3	100.00

^{*} N/A indicates no reporting category.

Grade 7

				No. o	of Ite	ms			
Reporting Cate	egories	TPI	TPD	TEI	MS	МС	ER	CR	% of Test
		N	N	N	N	N	N	N	
N/A	7-MS-LS4-5				1	1			5.41
	Sub-Total				1	1			5.41
1 Investigate	7-MS-ESS2-5				1				2.70
	7-MS-ESS3-5	1							2.70
	7-MS-PS3-4		2	1		1			10.81
	Sub-Total	1	2	1	1	1			16.22
2 Evaluate	7-MS-LS1-3			1		1	1		8.11
	7-MS-LS2-4			1		1		1	8.11
	7-MS-PS1-2			1	1				5.41
	Sub-Total			3	1	2	1	1	21.62
3 Reason Scientifically	7-MS-ESS2-4			3	1			1	13.51
	7-MS-ESS2-6	1							2.70
	7-MS-LS1-6				1				2.70
	7-MS-LS1-7			1		1			5.41
	7-MS-LS2-5			4		1			13.51
	7-MS-LS3-2	1	1			1			8.11
	7-MS-LS4-4				1			1	5.41
	7-MS-PS1-4					1			2.70
	7-MS-PS1-5			1					2.70
	Sub-Total	2	1	9	3	4		2	56.76
Total		3	3	13	6	8	1	3	100.00

^{*} N/A indicates no reporting category.

Grade 8

				No. o	of Ite	ms			
Reporting Cat	egories	TPI	TPD	TEI	MS	МС	ER	CR	% of Test
		N	N	N	N	N	N	N	
1 Investigate	8-MS-ESS3-2	1							2.70
	8-MS-ESS3-3			1		1			5.41
	8-MS-LS1-5			1				1	5.41
	8-MS-PS1-3			2		1			8.11
	8-MS-PS1-6				1				2.70
	8-MS-PS3-3	1						1	5.41
	Sub-Total	2		4	1	2		2	29.73
2 Evaluate	8-MS-ESS2-3			2		1			8.11
	8-MS-LS1-4			1					2.70
	8-MS-LS4-1			1		1			5.41
	8-MS-LS4-3					1			2.70
	8-MS-LS4-6			1		1			5.41
	8-MS-PS3-5			1		2			8.11
	Sub-Total			6		6			32.43
3 Reason Scientifically	8-MS-ESS1-4					1		1	5.41
	8-MS-ESS2-1			1					2.70
	8-MS-ESS2-2					2			5.41
	8-MS-ESS3-1			1		1	1		8.11
	8-MS-LS3-1			2					5.41
	8-MS-LS4-2		1	1					5.41
	8-MS-PS1-1					2			5.41
	Sub-Total		1	5		6	1	1	37.84
Total		2	1	15	1	14	1	3	100.00

^{*} N/A indicates no reporting category.

Table B.3 *Item Type Summary: Spring 2024 Operational SC G3-8*

Grade	МС	MS	TEI	CR	ER*	TPD	TPI
3	21	1	0	3	0	7	4
4	14	2	4	3	0	6	7
5	10	1	16	3	1	3	3
6	13	4	9	3	1	5	2
7	8	6	13	3	1	3	3
8	13	1	15	3	1	1	2

^{*}Classical analyses are calculated and estimated separately for each dimension of the ER item, and the result summarizes both dimensions.

Table B.4
Raw Score Summary: Spring 2024 Operational SC G3-8

Grade	N	Mean	SD	Min	Max	Mean_Pval	Mean_Pbis	Reliability*	SEM
3	≥50,060	20.55	9.50	0	49	0.41	0.42	0.87	3.41
4	≥48,770	22.00	10.41	0	51	0.44	0.47	0.90	3.31
5	≥48,350	25.42	11.87	0	60	0.46	0.46	0.90	3.72
6	≥47,790	25.06	11.24	1	59	0.44	0.44	0.88	3.93
7	≥47,920	24.95	12.11	0	61	0.42	0.48	0.90	3.85
8	≥48,170	25.65	10.81	0	57	0.43	0.42	0.88	3.74

^{*} Reliability is Cronbach's alpha.

Table B.5
Raw Score Summary by Reporting Category: Spring 2024 Operational SC G3-8

	Reporting								
Grade	Category	Mean	SD	Min	Max	Mean_Pval	Mean_Pbis	Reliability	SEM
3	Investigate	5.58	3.05	0	14	0.40	0.41	0.64	1.82
	Evaluate	9.75	4.66	0	23	0.43	0.41	0.76	2.31
	Reason Scientifically	4.23	2.35	0	9	0.47	0.46	0.57	1.54
4	Investigate	6.32	3.61	0	16	0.41	0.45	0.72	1.89
	Evaluate	3.18	1.80	0	7	0.47	0.51	0.54	1.22
	Reason Scientifically	10.27	4.86	0	23	0.47	0.47	0.80	2.16
5	Investigate	5.63	2.50	0	12	0.52	0.43	0.62	1.55
	Evaluate	10.41	6.28	0	29	0.40	0.51	0.83	2.59
	Reason Scientifically	9.38	4.19	0	20	0.49	0.43	0.74	2.12
6	Investigate	3.54	2.33	0	12	0.32	0.43	0.63	1.42
	Evaluate	6.69	3.43	0	15	0.47	0.47	0.71	1.85
	Reason Scientifically	14.82	6.58	0	34	0.48	0.44	0.77	3.15
7	Investigate	2.73	2.01	0	9	0.31	0.41	0.53	1.37
	Evaluate	6.85	3.83	0	19	0.38	0.49	0.60	2.42
	Reason Scientifically	14.53	6.99	0	31	0.47	0.50	0.86	2.59
8	Investigate	7.40	3.71	0	19	0.40	0.43	0.72	1.95
	Evaluate	7.75	2.99	0	15	0.47	0.38	0.72	1.58
	Reason Scientifically	10.50	5.26	0	26	0.43	0.45	0.67	3.01

Table B.6.1
Scale Score and Raw Score Summary: Spring 2024 Operational Science Grade 3

Category	Subgroup*	N	Percent	Scale Score Mean	Scale Score SD	Raw Score Mean	Raw Score SD	Effect Size
7	≥50,070	100.00	729.27	29.28	20.55	9.50	-	
Caradan	Female	≥24,800	49.53	729.50	28.71	20.58	9.33	-0.01
Gender	Male	≥25,260	50.46	729.04	29.83	20.51	9.67	-
	African American	≥20,390	40.72	719.83	27.38	17.37	8.24	0.74
	AI/AN	≥260	0.53	729.15	26.46	20.31	8.74	0.39
	Asian	≥760	1.53	744.04	31.13	25.80	10.60	-0.19
Ethnicity	Hispanic/Latino	≥6,000	11.99	722.29	29.25	18.31	8.98	0.60
	NHPI	≥30	0.07	741.85	34.40	25.12	11.10	-0.12
	Two or More Races	≥2,000	4.00	734.92	28.08	22.39	9.46	0.17
	White	≥20,570	41.10	739.57	27.37	23.97	9.48	-
Economically	No	≥12,380	24.73	745.60	26.61	26.12	9.45	-0.82
Disadvantaged	Yes	≥37,440	74.78	724.02	28.06	18.75	8.77	-
English Loornor	No	≥46,860	93.59	730.67	29.01	20.99	9.51	-0.74
English Learner	Yes	≥3,200	6.41	708.89	25.37	14.07	6.62	-
Education Classification	Regular	≥43,690	87.26	731.13	28.95	21.14	9.50	-0.50
	Special	≥6,380	12.74	716.56	28.42	16.49	8.43	-
Section 504	No	≥46,190	92.26	729.78	29.41	20.73	9.55	-0.25
Section 504	Yes	≥3,870	7.74	723.24	27.03	18.34	8.56	-
Migrapt	No	≥49,970	99.80	729.31	29.28	20.56	9.50	-0.60
Migrant	Yes	≥90	0.20	710.92	28.10	14.85	7.87	-
Homeless Status	No	≥48,870	97.61	729.59	29.22	20.65	9.50	-0.45
Tiorneless status	Yes	≥1,190	2.39	716.07	28.84	16.43	8.31	-
Military Affiliation	No	≥49,100	98.07	729.05	29.27	20.47	9.48	0.42
	Yes	≥960	1.93	740.63	27.94	24.43	9.59	-
Foster Care	No	≥49,880	99.61	729.30	29.28	20.56	9.50	-0.29
Status	Yes	≥190	0.39	720.61	28.97	17.83	8.87	-

^{*} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

Table B.6.2
Scale Score and Raw Score Summary: Spring 2024 Operational Science Grade 4

Category	Subgroup*	N	Percent	Scale Score Mean	Scale Score SD	Raw Score Mean	Raw Score SD	Effect Size
Total		≥48,780	100.00	732.84	30.98	22.00	10.41	-
Gender	Female	≥24,030	49.27	731.50	29.48	21.50	9.95	0.09
Gender	Male	≥24,740	50.73	734.14	32.32	22.48	10.81	-
	African American	≥20,120	41.25	721.25	27.03	18.02	8.73	0.86
	AI/AN	≥250	0.52	736.98	28.52	23.34	9.79	0.28
	Asian	≥760	1.57	752.88	31.95	28.86	10.84	-0.25
Ethnicity	Hispanic/Latino	≥5,530	11.35	724.93	29.75	19.39	9.72	0.68
	NHPI	≥50	0.10	743.22	29.33	25.58	10.59	0.07
	Two or More Races	≥1,940	3.98	737.29	29.72	23.45	10.18	0.27
	White	≥20,080	41.18	745.37	29.85	26.28	10.33	-
Economically	No	≥12,600	25.83	751.19	29.56	28.31	10.21	-0.87
Disadvantaged	Yes	≥35,890	73.59	726.57	28.83	19.84	9.54	-
Foodbale Language	No	≥45,980	94.26	734.27	30.80	22.47	10.40	-0.80
English Learner	Yes	≥2,800	5.74	709.28	23.49	14.25	6.79	-
Education	Regular	≥42,410	86.95	735.39	30.47	22.83	10.35	-0.62
Classification	Special	≥6,360	13.05	715.87	28.97	16.47	9.04	-
Cartina FOA	No	≥44,310	90.83	733.53	31.17	22.24	10.48	-0.26
Section 504	Yes	≥4,470	9.17	726.00	28.16	19.56	9.36	-
Missant	No	≥48,710	99.87	732.85	30.98	22.00	10.41	-0.29
Migrant	Yes	≥60	0.13	724.40	28.52	18.98	9.77	-
Harrada a Chahaa	No	≥47,730	97.86	733.21	30.95	22.12	10.41	-0.54
Homeless Status	Yes	≥1,040	2.14	716.13	27.88	16.56	8.59	-
Military Affiliation	No	≥47,840	98.07	732.52	30.94	21.89	10.39	0.55
	Yes	≥940	1.93	748.91	28.71	27.58	9.95	-
Foster Care	No	≥48,590	99.61	732.87	30.99	22.01	10.41	-0.27
Status	Yes	≥190	0.39	724.88	28.70	19.21	9.41	-

^{*} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

Table B.6.3
Scale Score and Raw Score Summary: Spring 2024 Operational Science Grade 5

Category	Subgroup*	N	Percent	Scale Score Mean	Scale Score SD	Raw Score Mean	Raw Score SD	Effect Size
Т	otal	≥48,360	100.00	731.10	36.77	25.42	11.87	-
Caradan	Female	≥23,670	48.95	730.75	35.33	25.21	11.48	0.04
Gender	Male	≥24,680	51.05	731.44	38.09	25.63	12.24	-
	African American	≥19,880	41.11	717.75	32.86	21.01	10.08	0.82
	AI/AN	≥270	0.58	736.73	34.65	27.08	11.50	0.25
	Asian	≥850	1.76	758.07	37.95	34.47	12.68	-0.38
Ethnicity	Hispanic/Latino	≥5,380	11.13	721.61	36.91	22.49	11.41	0.64
	NHPI	≥30	0.08	745.90	37.38	30.64	12.43	-0.06
	Two or More Races	≥1,850	3.84	737.57	35.15	27.40	11.63	0.22
	White	≥20,050	41.46	745.06	34.60	29.98	11.69	-
Economically	No	≥12,640	26.16	752.75	33.87	32.61	11.56	-0.88
Disadvantaged	Yes	≥35,470	73.35	723.57	34.57	22.91	10.89	-
English Loomer	No	≥45,850	94.82	732.89	36.35	25.97	11.83	-0.91
English Learner	Yes	≥2,500	5.18	698.47	27.96	15.39	7.25	-
Education	Regular	≥42,470	87.83	734.63	35.61	26.48	11.71	-0.75
Classification	Special	≥5,880	12.17	705.69	34.92	17.79	10.11	-
Costion FO4	No	≥43,500	89.96	732.14	36.90	25.77	11.94	-0.30
Section 504	Yes	≥4,850	10.04	721.80	34.12	22.28	10.76	-
Migrant	No	≥48,290	99.87	731.12	36.77	25.43	11.88	-0.34
Migrant	Yes	≥60	0.13	718.54	34.56	21.45	10.51	-
Llomologa Status	No	≥47,300	97.82	731.50	36.74	25.55	11.89	-0.48
Homeless Status	Yes	≥1,050	2.18	713.50	33.42	19.83	9.87	-
Military	No	≥47,450	98.12	730.80	36.77	25.32	11.86	0.44
Affiliation	Yes	≥900	1.88	746.91	33.18	30.59	11.34	-
Foster Care	No	≥48,180	99.63	731.15	36.77	25.44	11.88	-0.32
Status	Yes	≥170	0.37	719.56	33.30	21.61	10.37	-

^{*} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

Table B.6.4
Scale Score and Raw Score Summary: Spring 2024 Operational Science Grade 6

Category	Subgroup*	N	Percent	Scale Score Mean	Scale Score SD	Raw Score Mean	Raw Score SD	Effect Size
Т	otal	≥47,810	100.00	725.47	32.57	25.06	11.24	-
Canadan	Female	≥23,320	48.78	724.53	30.87	24.67	10.74	0.07
Gender	Male	≥24,490	51.22	726.36	34.09	25.43	11.69	-
	African American	≥20,010	41.87	713.58	28.00	20.87	9.40	0.83
	AI/AN	≥260	0.55	728.21	29.74	25.96	10.57	0.31
	Asian	≥810	1.70	752.54	35.12	34.43	12.09	-0.45
Ethnicity	Hispanic/Latino	≥5,480	11.48	718.11	32.61	22.64	10.98	0.61
	NHPI	≥30	0.08	730.15	31.05	26.82	10.97	0.24
	Two or More Races	≥1,670	3.51	732.08	31.24	27.26	11.01	0.20
	White	≥19,500	40.78	738.01	31.54	29.44	11.14	-
Economically	No	≥12,740	26.65	744.48	31.64	31.74	11.10	-0.86
Disadvantaged	Yes	≥34,800	72.79	718.66	30.03	22.66	10.27	-
	No	≥45,470	95.10	726.91	32.30	25.54	11.21	-0.89
English Learner	Yes	≥2,340	4.90	697.35	23.84	15.67	7.00	-
Education	Regular	≥42,580	89.07	728.04	32.17	25.92	11.19	-0.72
Classification	Special	≥5,220	10.93	704.48	27.88	17.99	8.90	-
5 11 504	No	≥42,820	89.57	726.38	32.66	25.38	11.28	-0.28
Section 504	Yes	≥4,980	10.43	717.59	30.66	22.25	10.47	-
	No	≥47,730	99.83	725.48	32.56	25.06	11.24	-0.28
Migrant	Yes	≥80	0.17	716.19	33.19	21.93	11.09	-
	No	≥46,780	97.84	725.84	32.57	25.18	11.25	-0.53
Homeless Status	Yes	≥1,030	2.16	708.59	27.43	19.24	9.05	-
Military	No	≥46,840	97.98	725.13	32.49	24.94	11.21	0.52
Affiliation	Yes	≥960	2.02	741.90	31.87	30.78	11.22	-
Foster Care	No	≥47,650	99.66	725.52	32.56	25.07	11.24	-0.44
Status	Yes	≥160	0.34	711.09	31.25	20.15	10.17	-

^{*} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

Table B.6.5
Scale Score and Raw Score Summary: Spring 2024 Operational Science Grade 7

Category	Subgroup*	N	Percent	Scale Score Mean	Scale Score SD	Raw Score Mean	Raw Score SD	Effect Size
Т	otal	≥47,950	100.00	732.38	33.89	24.95	12.11	-
Condor	Female	≥23,500	49.01	732.82	32.39	25.02	11.72	-0.01
Gender	Male	≥24,440	50.99	731.96	35.28	24.89	12.48	-
	African American	≥20,220	42.17	719.82	29.93	20.34	10.27	0.85
	AI/AN	≥250	0.53	735.24	28.78	25.77	10.78	0.34
	Asian	≥750	1.57	758.98	37.21	34.70	13.09	-0.42
Ethnicity	Hispanic/Latino	≥5,210	10.88	725.12	33.88	22.50	11.72	0.61
	NHPI	≥30	0.08	733.94	32.74	25.56	11.79	0.35
	Two or More Races	≥1,670	3.48	738.28	32.49	27.06	11.93	0.23
	White	≥19,780	41.26	745.58	32.13	29.75	11.87	-
Economically	No	≥13,220	27.59	752.13	32.04	32.18	11.85	-0.88
Disadvantaged	Yes	≥34,450	71.86	724.97	31.43	22.24	11.02	-
English Loomer	No	≥45,860	95.65	733.84	33.50	25.45	12.05	-0.96
English Learner	Yes	≥2,080	4.35	700.21	25.42	14.00	7.30	-
Education	Regular	≥42,820	89.30	735.25	33.13	25.94	11.99	-0.78
Classification	Special	≥5,130	10.70	708.42	30.58	16.74	9.78	-
Costion FOA	No	≥42,770	89.21	733.44	33.96	25.34	12.16	-0.30
Section 504	Yes	≥5,170	10.79	723.58	31.98	21.73	11.20	-
Migraph	No	≥47,880	99.85	732.40	33.90	24.96	12.11	-0.37
Migrant	Yes	≥70	0.15	720.24	28.98	20.51	10.24	-
Llowedees Status	No	≥46,950	97.91	732.74	33.90	25.08	12.13	-0.51
Homeless Status	Yes	≥1,000	2.09	715.66	28.85	18.91	9.68	-
Military	No	≥47,030	98.08	732.03	33.83	24.82	12.08	0.56
Affiliation	Yes	≥920	1.92	750.42	32.23	31.62	11.96	-
Foster Care	No	≥47,800	99.69	732.44	33.89	24.97	12.11	-0.54
Status	Yes	≥140	0.31	713.51	30.97	18.38	10.17	-

^{*} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

Table B.6.6

Scale Score and Raw Score Summary: Spring 2024 Operational Science Grade 8

Category	Subgroup*	N	Percent	Scale Score Mean	Scale Score SD	Raw Score Mean	Raw Score SD	Effect Size
Т	otal	≥48,220	100.00	730.20	31.79	25.65	10.81	-
Caradan	Female	≥23,470	48.69	730.53	30.58	25.69	10.48	-0.01
Gender	Male	≥24,740	51.31	729.89	32.89	25.61	11.12	-
	African American	≥20,110	41.72	718.42	28.09	21.54	9.17	0.88
	AI/AN	≥270	0.56	732.83	29.59	26.42	10.21	0.36
	Asian	≥800	1.67	752.84	34.63	33.53	11.92	-0.32
Ethnicity	Hispanic/Latino	≥5,500	11.41	720.86	32.57	22.70	10.54	0.71
	NHPI	≥30	0.07	740.44	31.20	29.11	10.74	0.10
	Two or More Races	≥1,660	3.45	736.55	30.24	27.76	10.55	0.23
	White	≥19,810	41.09	743.27	29.38	30.14	10.43	-
Economically	No	≥13,190	27.37	748.59	29.02	32.03	10.36	-0.87
Disadvantaged	Yes	≥34,730	72.04	723.41	29.93	23.29	9.97	-
Enablish Lagrana	No	≥45,890	95.18	731.88	31.21	26.18	10.72	-1.04
English Learner	Yes	≥2,320	4.82	697.21	24.21	15.20	6.48	-
Education	Regular	≥43,380	89.97	732.68	31.25	26.47	10.74	-0.77
Classification	Special	≥4,830	10.03	708.00	27.63	18.32	8.47	-
Continu FOA	No	≥42,990	89.17	731.26	31.89	26.02	10.87	-0.32
Section 504	Yes	≥5,220	10.83	721.53	29.53	22.60	9.84	-
Migraph	No	≥48,140	99.84	730.23	31.78	25.66	10.81	-0.40
Migrant	Yes	≥70	0.16	716.38	32.13	21.30	10.27	-
Lla ma ala an Chahan	No	≥47,250	98.00	730.57	31.73	25.77	10.81	-0.55
Homeless Status	Yes	≥960	2.00	712.40	29.36	19.81	9.10	-
Military	No	≥47,320	98.14	729.91	31.75	25.55	10.79	0.51
Affiliation	Yes	≥890	1.86	745.80	29.59	31.06	10.58	-
Foster Care	No	≥48,050	99.65	730.26	31.78	25.67	10.81	-0.47
Status	Yes	≥160	0.35	715.01	28.82	20.55	9.16	-

^{*} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

Appendix C: Item Analysis Summary Report

Summary Statistics Reports

Contents

Table C.1 P-Value Summary by Grade: Spring 2024 Operational SC G3-8

Plot C.1 P-Value Summary by Grade: Spring 2024 Operational SC G3-8

Table C.2 Item-Total Correlation Summary by Grade: Spring 2024 Operational SC G3-8

Plot C.2 Item-Total Correlation Summary by Grade: Spring 2024 Operational SC G3-8

Table C.3 Corrected Point-Biserial Correlation Summary by Grade: Spring 2024 Operational SC G3–8

Plot C.3 Corrected Point-Biserial Correlation Summary by Grade: Spring 2024 Operational SC G3–8

Table C.4 Item-Total Correlation Summary by Reporting Category and Grade: Spring 2024 Operational SC G3–8

Table C.5.1 IRT-A Parameter Summary by Reporting Category: SC G3-8

Table C.5.2 IRT-B Parameter Summary by Reporting Category: SC G3-8

Table C.5.3 IRT Parameter Summary: Spring 2024 Operational SC G3-8

Plot C.5.1 IRT Parameter Summary: Spring 2024 Operational SC G3-8: A-Parameter

Plot C.5.2 IRT Parameter Summary: Spring 2024 Operational SC G3-8: B-Parameter

Plot C.5.3 IRT Parameter Summary: Spring 2024 Operational SC G3-8: C-Parameter

Table C.6 Statistically Flagged Items by Item Type: Spring 2024 Operational SC G3-8

Table C.1.1

P-Value Summary by Grade: Spring 2024 Operational SC G3-8

Grade	No. of Items	0 ≤ p < 0.2	0.2 ≤ p < 0.4	0.4 ≤ p < 0.6	0.6 ≤ p < 0.8	0.8 ≤ p ≤ 1.0
3	36	0	16	19	1	0
4	36	1	13	18	4	0
5	39	2	15	14	8	0
6	38	5	8	18	7	0
7	37	1	16	15	5	0
8	37	2	12	20	3	0

Plot C.1.1

P-Value Summary by Grade: Spring 2024 Operational SC G3-8

Box and Whisker Plot

P-Value: Science

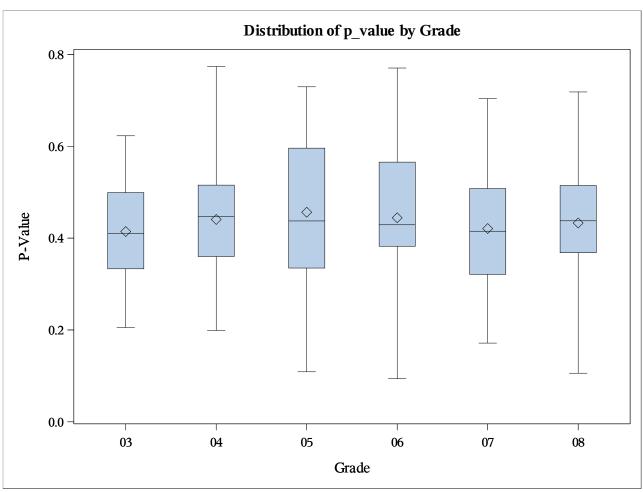


Table C.1.2

P-Value Summary by Item Type: Spring 2024 Operational SC G3-8

		No. of		25th		75th	
Grade	Type	Items	Minimum	Percentile	Median	Percentile	Maximum
	CR	3	0.205	0.205	0.276	0.346	0.346
	MC	21	0.272	0.342	0.446	0.495	0.553
3	MS	1	0.304	0.304	0.304	0.304	0.304
	TPD	7	0.311	0.311	0.357	0.484	0.550
	TPI	4	0.413	0.462	0.533	0.589	0.623
	CR	3	0.199	0.199	0.232	0.375	0.375
	MC	14	0.318	0.467	0.514	0.588	0.774
4	MS	2	0.285	0.285	0.354	0.424	0.424
4	TEI	4	0.228	0.245	0.436	0.615	0.621
	TPD	6	0.230	0.299	0.360	0.449	0.516
	TPI	7	0.374	0.413	0.445	0.489	0.499
	CR	3	0.187	0.187	0.244	0.299	0.299
	ER	3	0.109	0.109	0.201	0.296	0.296
5	MC	10	0.335	0.512	0.642	0.711	0.730
	MS	1	0.369	0.369	0.369	0.369	0.369
	TEI	16	0.246	0.375	0.524	0.589	0.693
	TPD	3	0.262	0.262	0.359	0.518	0.518
	TPI	3	0.401	0.401	0.437	0.544	0.544
	CR	3	0.094	0.094	0.116	0.176	0.176
	ER	2	0.127	0.127	0.208	0.289	0.289
	MC	13	0.281	0.497	0.588	0.645	0.770
6	MS	4	0.256	0.322	0.424	0.462	0.465
	TEI	9	0.362	0.410	0.422	0.524	0.728
	TPD	5	0.177	0.397	0.410	0.433	0.566
	TPI	2	0.382	0.382	0.486	0.590	0.590
	CR	3	0.206	0.206	0.240	0.363	0.363
	ER	1	0.333	0.333	0.333	0.333	0.333
	MC	8	0.311	0.399	0.497	0.580	0.613
7	MS	6	0.171	0.257	0.334	0.398	0.508
	TEI	13	0.241	0.401	0.439	0.513	0.704
	TPD	3	0.202	0.202	0.232	0.524	0.524
	TPI	3	0.321	0.321	0.458	0.681	0.681

Table C.1.2

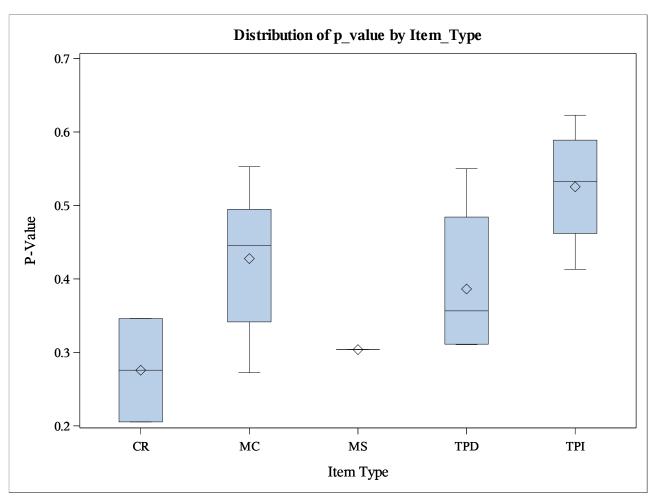
P-Value Summary by Item Type: Spring 2024 Operational SC G3-8 (continued)

		No. of		25th		75th	
Grade	Туре	Items	Minimum	Percentile	Median	Percentile	Maximum
	CR	3	0.197	0.197	0.251	0.255	0.255
	ER	2	0.307	0.307	0.339	0.372	0.372
	MC	13	0.269	0.371	0.473	0.516	0.561
8	MS	1	0.496	0.496	0.496	0.496	0.496
	TEI	15	0.106	0.429	0.458	0.579	0.718
	TPD	1	0.377	0.377	0.377	0.377	0.377
	TPI	2	0.413	0.413	0.414	0.416	0.416

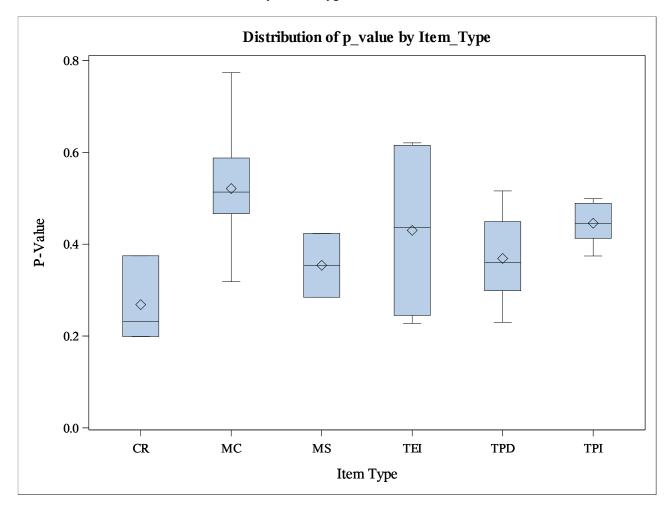
Plot C.1.2

P-Value Summary by Item Type: Spring 2024 Operational SC G3-8

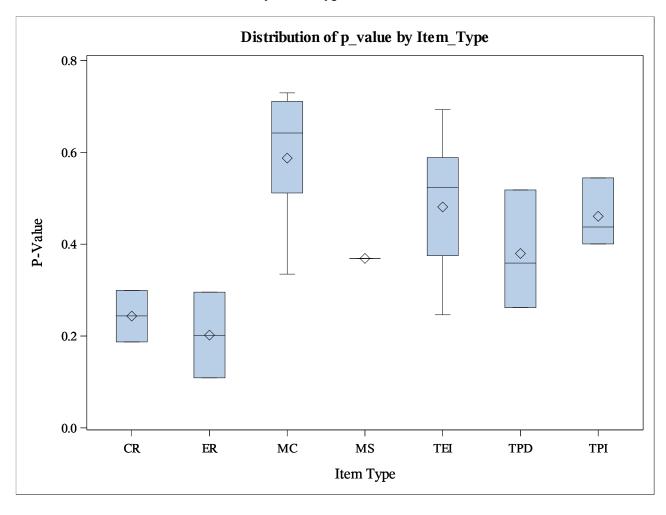
Box and Whisker Plot P-Value by Item Type: Science Grade 3



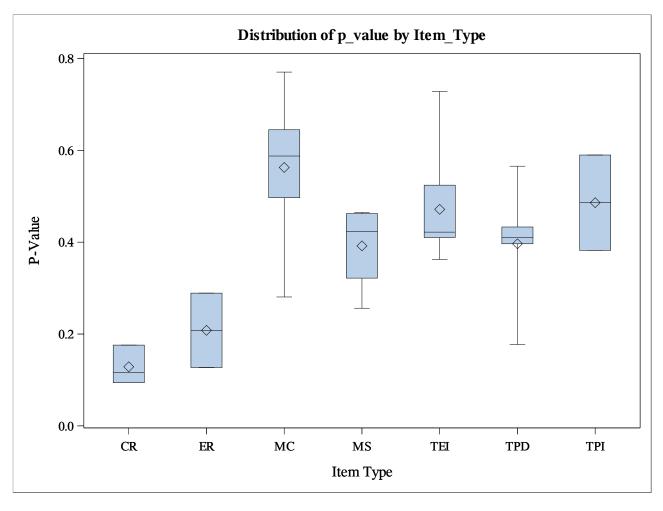
Box and Whisker Plot
P-Value by Item Type: Science Grade 4



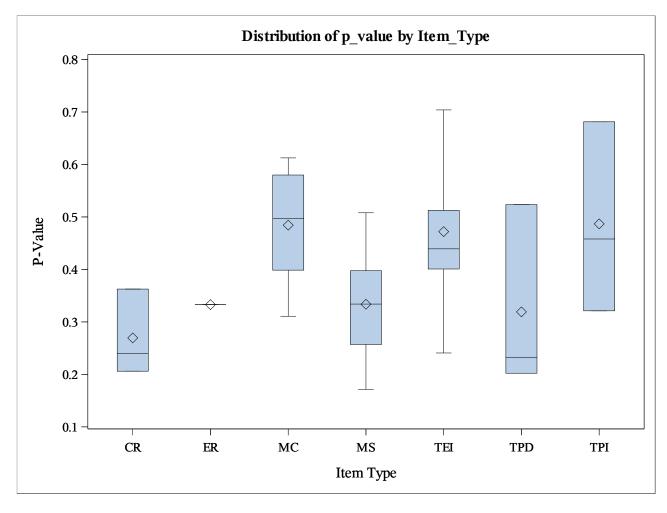
Box and Whisker Plot
P-Value by Item Type: Science Grade 5



Box and Whisker Plot
P-Value by Item Type: Science Grade 6



Box and Whisker Plot
P-Value by Item Type: Science Grade 7



Box and Whisker Plot
P-Value by Item Type: Science Grade 8

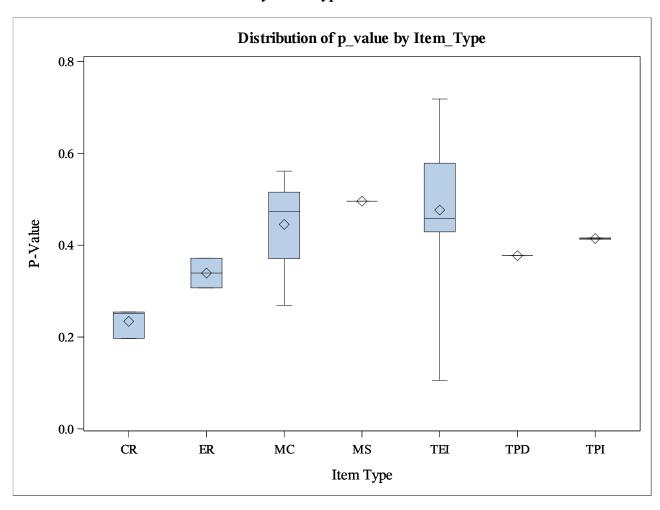


Table C.2.1 *Item-Total Correlation by Grade: Spring 2024 Operational SC G3–8*

	No. of						
Grade	Items	r < 0	$0.0 \le r < 0.2$	0.2 ≤ r < 0.3	0.3 ≤ r < 0.4	0.4 ≤ r < 0.5	r ≥ 0.5
3	36	0	0	6	9	13	8
4	36	0	0	2	5	14	15
5	39	0	0	0	10	17	12
6	38	0	0	2	10	16	10
7	37	0	0	2	6	11	18
8	38	1	1	6	8	13	9

Plot C.2.1 *Item-Total Correlation by Grade: Spring 2024 Operational SC G3–8*

Box and Whisker Plot Point-Biserial Correlation: Science

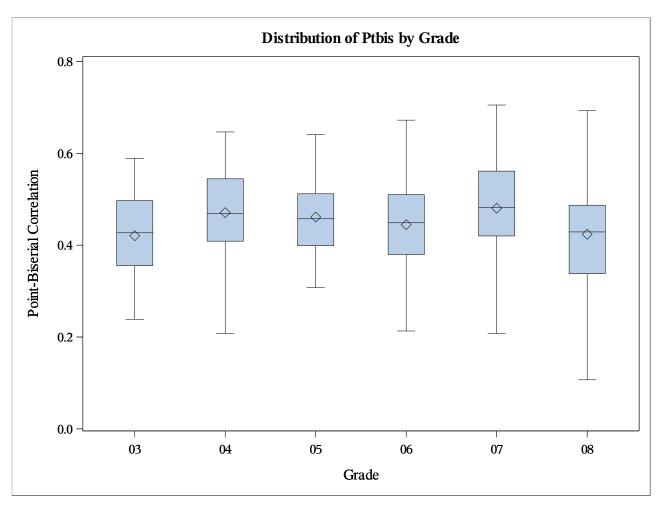


Table C.2.2 *Item-Total Correlation Summary by Item Type: Spring 2024 Operational SC G3–8 (continued)*

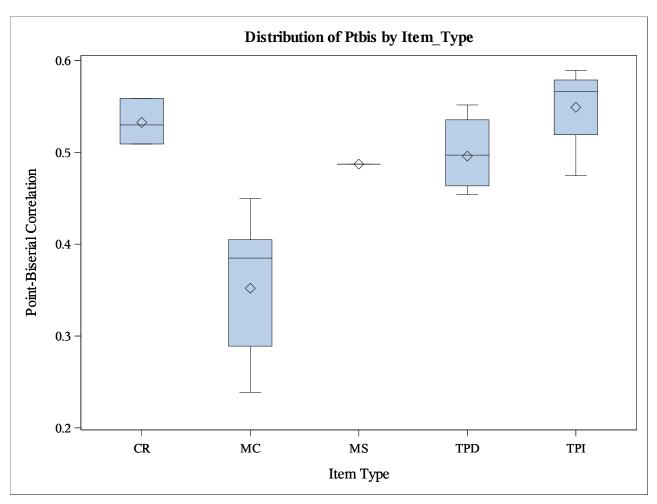
Grade	Туре	No. of Items	Minimum	25th Percentile	Median	75th Percentile	Maximum
	CR	3	0.509	0.509	0.530	0.559	0.559
	MC	21	0.239	0.289	0.385	0.405	0.450
3	MS	1	0.487	0.487	0.487	0.487	0.487
	TPD	7	0.454	0.463	0.497	0.535	0.552
	TPI	4	0.475	0.519	0.566	0.579	0.589
	CR	3	0.547	0.547	0.595	0.611	0.611
	MC	14	0.208	0.391	0.423	0.520	0.542
4	MS	2	0.411	0.411	0.450	0.490	0.490
4	TEI	4	0.299	0.323	0.403	0.469	0.479
	TPD	6	0.355	0.427	0.498	0.572	0.581
	TPI	7	0.405	0.462	0.542	0.596	0.647
	CR	3	0.506	0.506	0.512	0.600	0.600
	ER	3	0.502	0.502	0.599	0.611	0.611
	MC	10	0.324	0.372	0.391	0.428	0.478
5	MS	1	0.541	0.541	0.541	0.541	0.541
	TEI	16	0.308	0.401	0.455	0.488	0.641
	TPD	3	0.427	0.427	0.480	0.601	0.601
	TPI	3	0.411	0.411	0.455	0.560	0.560
	CR	3	0.366	0.366	0.385	0.512	0.512
	ER	2	0.403	0.403	0.538	0.672	0.672
	MC	13	0.254	0.389	0.422	0.452	0.512
6	MS	4	0.380	0.430	0.495	0.521	0.531
	TEI	9	0.213	0.376	0.458	0.477	0.567
	TPD	5	0.318	0.373	0.479	0.483	0.594
	TPI	2	0.518	0.518	0.571	0.625	0.625
	CR	3	0.517	0.517	0.561	0.585	0.585
	ER	1	0.705	0.705	0.705	0.705	0.705
	MC	8	0.208	0.340	0.369	0.413	0.550
7	MS	6	0.289	0.384	0.472	0.544	0.552
	TEI	13	0.322	0.454	0.482	0.585	0.629
	TPD	3	0.420	0.420	0.561	0.593	0.593
	TPI	3	0.449	0.449	0.531	0.561	0.561

Table C.2.2 *Item-Total Correlation Summary by Item Type: Spring 2024 Operational SC G3–8 (continued)*

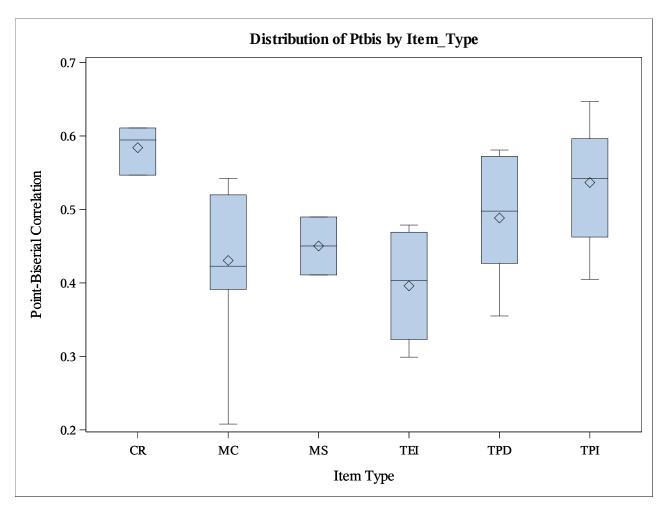
Grade	Туре	No. of Items	Minimum	25th Percentile	Median	75th Percentile	Maximum
	CR	3	0.404	0.404	0.444	0.548	0.548
	ER	2	0.666	0.666	0.679	0.693	0.693
	MC	13	0.246	0.310	0.368	0.422	0.487
8	MS	1	0.271	0.271	0.271	0.271	0.271
	TEI	15	0.107	0.317	0.452	0.519	0.646
	TPD	1	0.469	0.469	0.469	0.469	0.469
	TPI	2	0.552	0.552	0.590	0.628	0.628

Plot C.2.2 Item-Total Correlation Summary by Item Type: Spring 2024 Operational SC G3–8

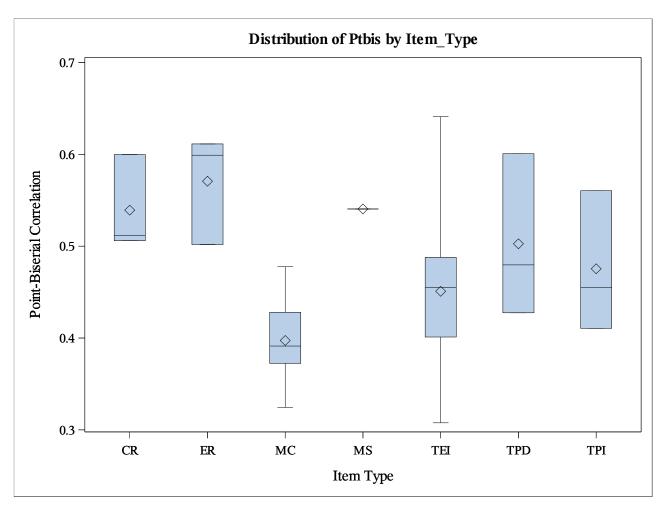
Box and Whisker Plot
Point-Biserial Correlation: Science Grade 3



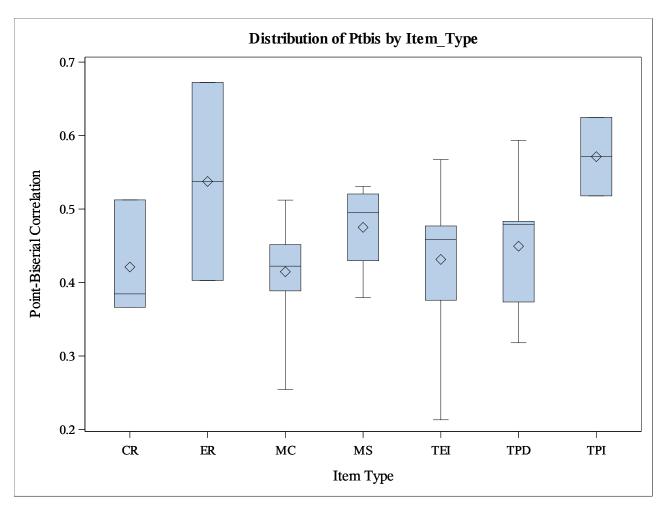
Box and Whisker Plot
Point-Biserial Correlation: Science Grade 4



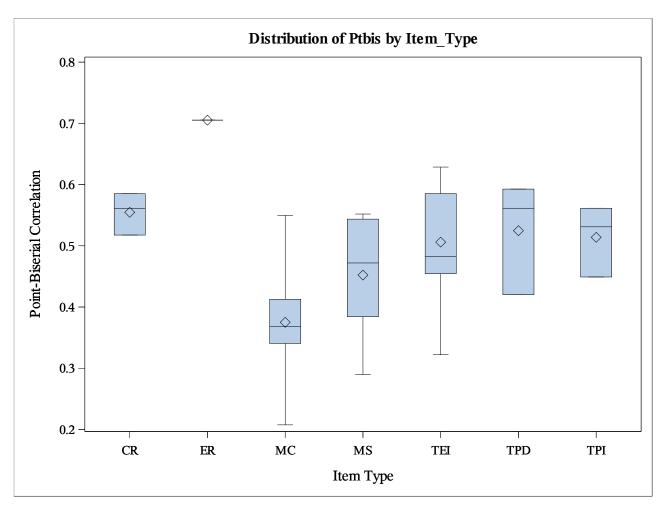
Box and Whisker Plot
Point-Biserial Correlation: Science Grade 5



Box and Whisker Plot
Point-Biserial Correlation: Science Grade 6



Box and Whisker Plot
Point-Biserial Correlation: Science Grade 7



Box and Whisker Plot
Point-Biserial Correlation: Science Grade 8

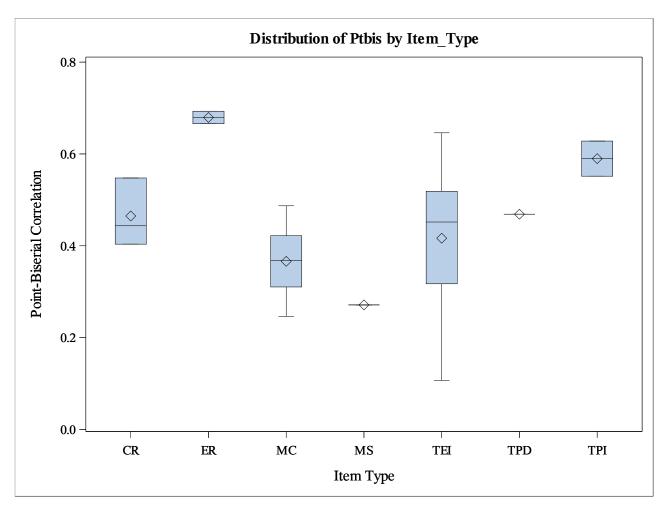


Table C.3.1

Corrected Point-Biserial Correlation* Summary by Grade: Spring 2024 Operational SC G3-8

	No. of						
Grade	Items	r < 0	0.0 ≤ r < 0.2	0.2 ≤ r < 0.3	0.3 ≤ r < 0.4	0.4 ≤ r < 0.5	r ≥ 0.5
3	36	0	3	6	14	9	4
4	36	0	1	3	11	12	9
5	39	0	0	3	14	14	8
6	38	0	1	3	13	17	4
7	37	0	1	3	8	12	13
8	38	1	2	9	11	9	6

^{*} Corrected point-biserial correlation, which was slightly more robust than point-biserial correlation, calculates the relationship between the item score and the total test score after removing the item score from the total test score.

Plot C.3.1 Corrected Point-Biserial Correlation Summary by Grade: Spring 2024 Operational SC G3–8



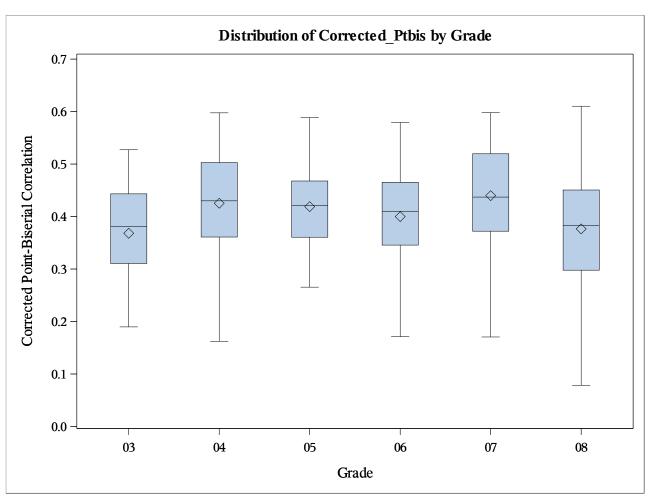


Table C.3.2

Corrected Point-Biserial Correlation* Summary by Item Type: Spring 2024 Operational SC G3-8

		No. of		25th		75th	
Grade	Туре	Items	Minimum	Percentile	Median	Percentile	Maximum
	CR	3	0.462	0.462	0.472	0.508	0.508
	MC	21	0.190	0.240	0.339	0.360	0.406
3	MS	1	0.449	0.449	0.449	0.449	0.449
	TPD	7	0.385	0.391	0.422	0.468	0.479
	TPI	4	0.413	0.460	0.508	0.518	0.527
	CR	3	0.499	0.499	0.552	0.570	0.570
	MC	14	0.162	0.350	0.384	0.483	0.507
4	MS	2	0.374	0.374	0.413	0.452	0.452
4	TEI	4	0.256	0.282	0.365	0.434	0.447
	TPD	6	0.290	0.361	0.434	0.511	0.533
	TPI	7	0.345	0.409	0.494	0.545	0.598
	CR	3	0.457	0.457	0.465	0.553	0.553
	ER	3	0.453	0.453	0.544	0.555	0.555
	MC	10	0.288	0.337	0.356	0.393	0.446
5	MS	1	0.511	0.511	0.511	0.511	0.511
	TEI	16	0.265	0.368	0.412	0.449	0.588
	TPD	3	0.365	0.365	0.424	0.544	0.544
	TPI	3	0.360	0.360	0.404	0.515	0.515
	CR	3	0.329	0.329	0.348	0.473	0.473
	ER	2	0.363	0.363	0.453	0.542	0.542
	MC	13	0.216	0.350	0.387	0.421	0.478
6	MS	4	0.346	0.395	0.460	0.487	0.499
	TEI	9	0.171	0.337	0.404	0.422	0.512
	TPD	5	0.273	0.324	0.413	0.422	0.539
	TPI	2	0.465	0.465	0.522	0.579	0.579
	CR	3	0.469	0.469	0.518	0.550	0.550
	ER	1	0.598	0.598	0.598	0.598	0.598
	MC	8	0.171	0.303	0.333	0.378	0.520
7	MS	6	0.256	0.351	0.444	0.515	0.522
	TEI	13	0.290	0.421	0.437	0.556	0.590
	TPD	3	0.372	0.372	0.514	0.539	0.539
	TPI	3	0.400	0.400	0.486	0.520	0.520

Table C.3.2

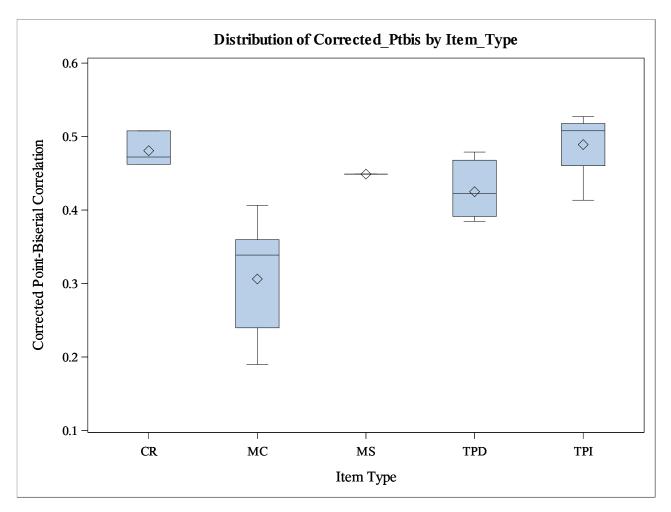
Corrected Point-Biserial Correlation* Summary by Item Type: Spring 2024 Operational SC G3-8 (continued)

		No. of		25th		75th	
Grade	Type	Items	Minimum	Percentile	Median	Percentile	Maximum
8	CR	3	0.349	0.349	0.388	0.497	0.497
	ER	2	0.599	0.599	0.605	0.610	0.610
	MC	13	0.202	0.268	0.328	0.383	0.451
	MS	1	0.228	0.228	0.228	0.228	0.228
	TEI	15	0.079	0.275	0.393	0.457	0.598
	TPD	1	0.412	0.412	0.412	0.412	0.412
	TPI	2	0.500	0.500	0.541	0.583	0.583

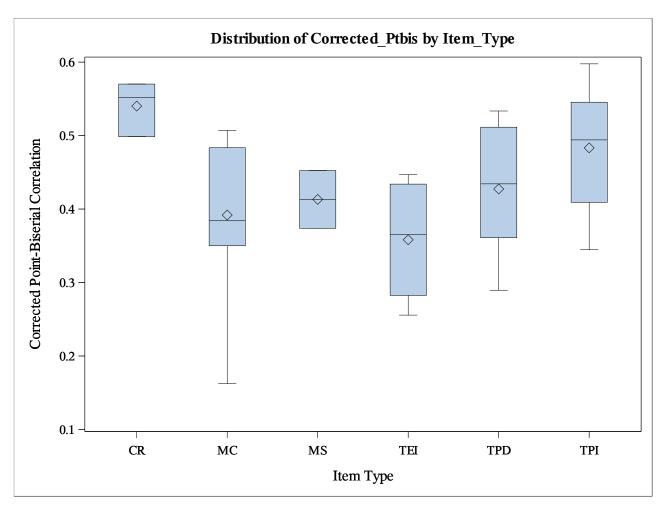
^{*} Corrected point-biserial correlation, which was slightly more robust than point-biserial correlation, calculates the relationship between the item score and the total test score after removing the itemscore from the total test score.

Plot C.3.2 Corrected Point-Biserial Correlation Summary by Item Type: Spring 2024 Operational SC G3–8

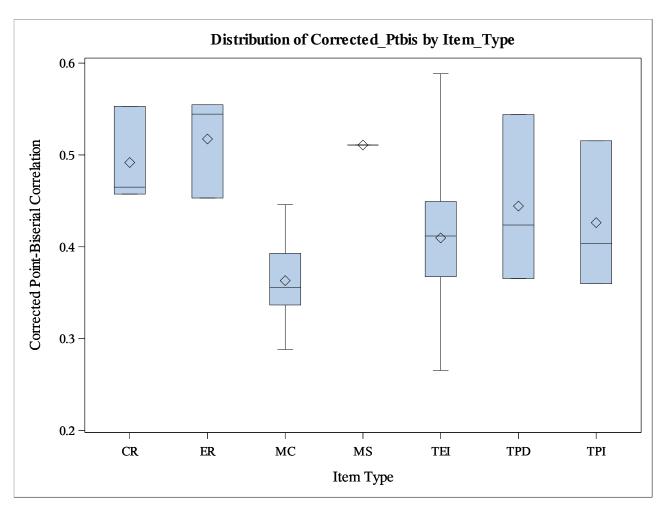
Box and Whisker Plot
Corrected Point-Biserial Correlation: Science Grade 3



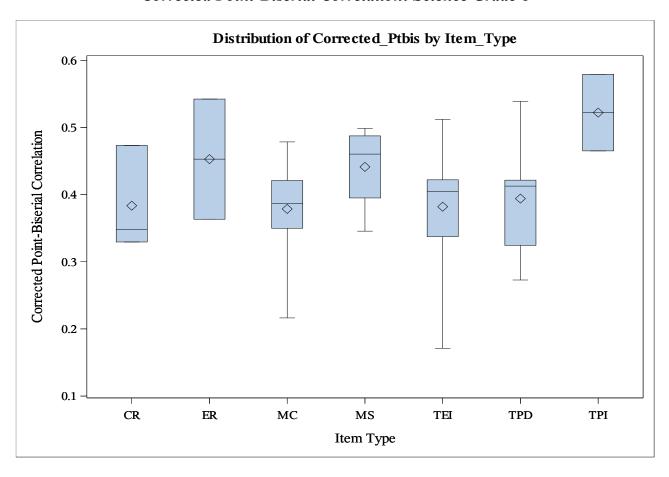
Box and Whisker Plot
Corrected Point-Biserial Correlation: Science Grade 4



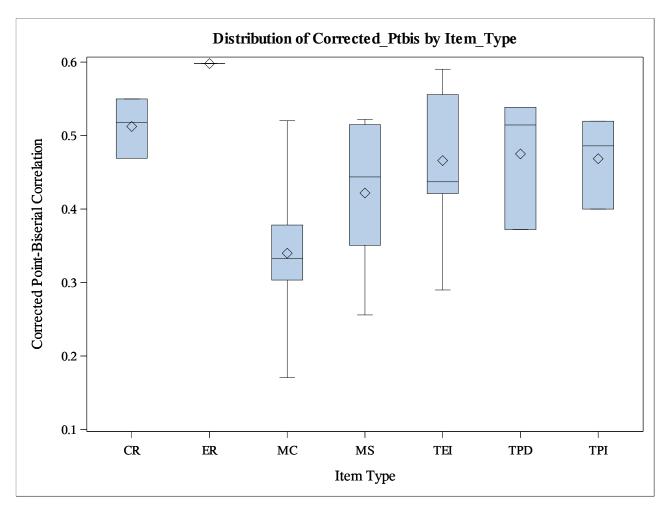
Box and Whisker Plot
Corrected Point-Biserial Correlation: Science Grade 5



Box and Whisker Plot Corrected Point-Biserial Correlation: Science Grade 6



Box and Whisker Plot
Corrected Point-Biserial Correlation: Science Grade 7



Box and Whisker Plot Corrected Point-Biserial Correlation: Science Grade 8

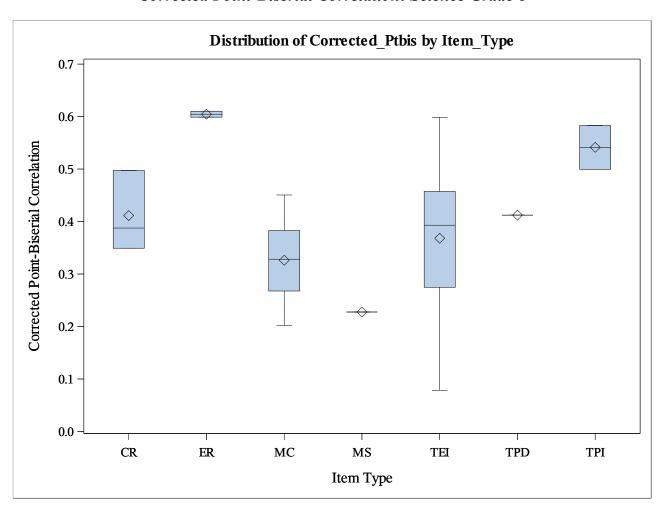


Table C.4.1 *Item-Total Correlation Summary by Reporting Category: Spring 2024 Operational SC G3–8*

	Reporting	No. of		25th		75th		
Grade	Category	Items	Minimum	Percentile	Median	Percentile	Maximum	
3	1 Investigate	10	0.241	0.332	0.431	0.472	0.559	
	2 Evaluate	17	0.240	0.327	0.405	0.497	0.589	
	3 Reason Scientifically	6	0.379	0.391	0.443	0.552	0.564	
4	1 Investigate	12	0.208	0.343	0.450	0.550	0.647	
	2 Evaluate	4	0.459	0.461	0.501	0.552	0.563	
	3 Reason Scientifically	16	0.347	0.407	0.446	0.556	0.596	
5	1 Investigate	8	0.308	0.356	0.391	0.509	0.584	
	2 Evaluate	16	0.324	0.449	0.491	0.599	0.641	
	3 Reason Scientifically	15	0.333	0.399	0.428	0.468	0.530	
6	1 Investigate	8	0.254	0.355	0.442	0.522	0.534	
	2 Evaluate	10	0.372	0.418	0.466	0.512	0.594	
	3 Reason Scientifically	20	0.213	0.384	0.445	0.479	0.672	
7	1 Investigate	6	0.208	0.322	0.435	0.517	0.561	
	2 Evaluate	8	0.335	0.424	0.467	0.557	0.705	
	3 Reason Scientifically	21	0.289	0.454	0.531	0.561	0.629	
8	1 Investigate	11	0.246	0.338	0.444	0.552	0.628	
	2 Evaluate	12	0.107	0.314	0.404	0.470	0.560	
	3 Reason Scientifically	14	0.224	0.341	0.445	0.548	0.693	

Table C.4.2.1

Item-Total Correlation Summary by Reporting Category and Item Type: Spring 2024 SC G3-4

			No. of		25th		75th	
Grade	Туре	Reporting Category	Items	Minimum	Percentile	Median	Percentile	Maximum
	CR	1 Investigate	1	0.559	0.559	0.559	0.559	0.559
	CK	2 Evaluate	1	0.530	0.530	0.530	0.530	0.530
		1 Investigate	6	0.241	0.306	0.358	0.412	0.450
	MC	2 Evaluate	11	0.240	0.272	0.388	0.405	0.441
3		3 Reason Scientifically	3	0.379	0.379	0.391	0.422	0.422
3		1 Investigate	3	0.454	0.454	0.472	0.535	0.535
	TPD	2 Evaluate	2	0.497	0.497	0.497	0.497	0.497
		3 Reason Scientifically	2	0.463	0.463	0.508	0.552	0.552
	TPI	2 Evaluate	3	0.475	0.475	0.568	0.589	0.589
		3 Reason Scientifically	1	0.564	0.564	0.564	0.564	0.564
	CR	1 Investigate	1	0.611	0.611	0.611	0.611	0.611
		3 Reason Scientifically	2	0.547	0.547	0.571	0.595	0.595
	МС	1 Investigate	5	0.208	0.330	0.408	0.520	0.528
		3 Reason Scientifically	8	0.377	0.400	0.423	0.469	0.542
	MS	1 Investigate	2	0.411	0.411	0.450	0.490	0.490
		1 Investigate	1	0.299	0.299	0.299	0.299	0.299
4	TEI	2 Evaluate	1	0.459	0.459	0.459	0.459	0.459
		3 Reason Scientifically	1	0.347	0.347	0.347	0.347	0.347
	TPD	1 Investigate	2	0.355	0.355	0.464	0.572	0.572
	טאו	3 Reason Scientifically	3	0.429	0.429	0.566	0.581	0.581
	TPI	1 Investigate	1	0.647	0.647	0.647	0.647	0.647
		2 Evaluate	3	0.462	0.462	0.541	0.563	0.563
		3 Reason Scientifically	2	0.405	0.405	0.501	0.596	0.596

Table C.4.2.2

Item-Total Correlation Summary by Reporting Category and Item Type: Spring 2024 SC G5-6

		relation summary by	No. of	<u> </u>	25th	,	75th	
Grade	Туре	Reporting Category	Items	Minimum		Median	Percentile	Maximum
	CD	1 Investigate	2	0.506	0.506	0.509	0.512	0.512
	CR	2 Evaluate	1	0.600	0.600	0.600	0.600	0.600
	ER	2 Evaluate	3	0.502	0.502	0.599	0.611	0.611
		1 Investigate	4	0.337	0.356	0.381	0.391	0.395
	MC	2 Evaluate	2	0.324	0.324	0.393	0.461	0.461
		3 Reason Scientifically	4	0.372	0.394	0.422	0.453	0.478
5	MS	2 Evaluate	1	0.541	0.541	0.541	0.541	0.541
		1 Investigate	2	0.308	0.308	0.446	0.584	0.584
	TEI	2 Evaluate	5	0.403	0.437	0.460	0.477	0.641
		3 Reason Scientifically	9	0.333	0.399	0.453	0.468	0.530
	TPD	2 Evaluate	3	0.427	0.427	0.480	0.601	0.601
	TPI	2 Evaluate	1	0.560	0.560	0.560	0.560	0.560
		3 Reason Scientifically	2	0.411	0.411	0.433	0.455	0.455
	CR	1 Investigate	2	0.366	0.366	0.439	0.512	0.512
	CIX	2 Evaluate	1	0.385	0.385	0.385	0.385	0.385
	ER	3 Reason Scientifically	2	0.403	0.403	0.538	0.672	0.672
		1 Investigate	2	0.254	0.254	0.299	0.343	0.343
	MC	2 Evaluate	5	0.372	0.418	0.422	0.456	0.512
		3 Reason Scientifically	6	0.389	0.422	0.445	0.452	0.461
	MS	1 Investigate	2	0.510	0.510	0.521	0.531	0.531
6	IVIO	3 Reason Scientifically	2	0.380	0.380	0.430	0.480	0.480
		1 Investigate	1	0.534	0.534	0.534	0.534	0.534
	TEI	2 Evaluate	3	0.476	0.476	0.477	0.567	0.567
		3 Reason Scientifically	5	0.213	0.339	0.376	0.442	0.458
	TPD	1 Investigate	1	0.373	0.373	0.373	0.373	0.373
		2 Evaluate	1	0.594	0.594	0.594	0.594	0.594
		3 Reason Scientifically	3	0.318	0.318	0.479	0.483	0.483
	TPI	3 Reason Scientifically	2	0.518	0.518	0.571	0.625	0.625

Table C.4.2.3

Item-Total Correlation Summary by Reporting Category and Item Type: Spring 2024 SC G7–8

Grade	Туре	Reporting Category	No. of Items	Minimum	25th Percentile	Median	75th Percentile	Maximum
	CD	2 Evaluate	1	0.585	0.585	0.585	0.585	0.585
	CR	3 Reason Scientifically	2	0.517	0.517	0.539	0.561	0.561
	ER	2 Evaluate	1	0.705	0.705	0.705	0.705	0.705
		1 Investigate	1	0.208	0.208	0.208	0.208	0.208
	MC	2 Evaluate	2	0.335	0.335	0.378	0.421	0.421
		3 Reason Scientifically	4	0.346	0.362	0.391	0.477	0.550
		1 Investigate	1	0.517	0.517	0.517	0.517	0.517
7	MS	2 Evaluate	1	0.427	0.427	0.427	0.427	0.427
/		3 Reason Scientifically	3	0.289	0.289	0.384	0.552	0.552
	TEI	1 Investigate	1	0.322	0.322	0.322	0.322	0.322
		2 Evaluate	3	0.451	0.451	0.482	0.528	0.528
		3 Reason Scientifically	9	0.454	0.457	0.534	0.601	0.629
	TPD	1 Investigate	2	0.420	0.420	0.491	0.561	0.561
		3 Reason Scientifically	1	0.593	0.593	0.593	0.593	0.593
	TPI	1 Investigate	1	0.449	0.449	0.449	0.449	0.449
		3 Reason Scientifically	2	0.531	0.531	0.546	0.561	0.561
	GD.	1 Investigate	2	0.404	0.404	0.424	0.444	0.444
	CR	3 Reason Scientifically	1	0.548	0.548	0.548	0.548	0.548
	ER	3 Reason Scientifically	2	0.666	0.666	0.679	0.693	0.693
		1 Investigate	2	0.246	0.246	0.292	0.338	0.338
	MC	2 Evaluate	6	0.310	0.398	0.420	0.456	0.487
0		3 Reason Scientifically	5	0.266	0.287	0.341	0.368	0.422
8	MS	1 Investigate	1	0.271	0.271	0.271	0.271	0.271
		1 Investigate	4	0.429	0.440	0.455	0.510	0.560
	TEI	2 Evaluate	6	0.107	0.275	0.346	0.485	0.560
		3 Reason Scientifically	5	0.224	0.365	0.476	0.519	0.646
	TPD	3 Reason Scientifically	1	0.469	0.469	0.469	0.469	0.469
	TPI	1 Investigate	2	0.552	0.552	0.590	0.628	0.628

Table C.5.1.1

IRT-A Parameter Summary by Reporting Category: SC G3

Grade	IRT-a Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	a < 0.0	0	0	0	0
	0.0 <= a < 0.2	0	0	0	0
	0.2 <= a < 0.4	1	3	1	5
	0.4 <= a < 0.6	4	6	1	11
	0.6 <= a < 0.8	3	4	2	10
	0.8 <= a < 1.0	0	2	2	5
	1.0 <= a < 1.2	1	0	0	2
	1.2 <= a < 1.4	1	0	0	1
3	1.4 <= a < 1.6	0	1	0	1
	1.6 <= a < 1.8	0	1	0	1
	1.8 <= a < 2.0	0	0	0	0
	2.0 <= a	0	0	0	0
	Minimum	0.33	0.26	0.31	0.26
	Maximum	1.36	1.64	0.91	1.64
	Mean	0.67	0.70	0.64	0.69
	SD	0.32	0.38	0.23	0.33
	Number of Items	10	17	6	36

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.1.2

IRT-A Parameter Summary by Reporting Category: SC G4

Grade	IRT-a Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	a < 0.0	0	0	0	0
	0.0 <= a < 0.2	0	0	0	0
	0.2 <= a < 0.4	2	0	2	5
	0.4 <= a < 0.6	3	3	6	13
	0.6 <= a < 0.8	5	1	6	12
	0.8 <= a < 1.0	1	0	1	2
	1.0 <= a < 1.2	1	0	1	4
	1.2 <= a < 1.4	0	0	0	0
4	1.4 <= a < 1.6	0	0	0	0
	1.6 <= a < 1.8	0	0	0	0
	1.8 <= a < 2.0	0	0	0	0
	2.0 <= a	0	0	0	0
	Minimum	0.23	0.41	0.27	0.23
	Maximum	1.01	0.63	1.00	1.16
	Mean	0.63	0.49	0.62	0.63
	SD	0.24	0.10	0.19	0.23
	Number of Items	12	4	16	36

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.3

IRT-A Parameter Summary by Reporting Category: SC G5

Grade	IRT-a Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	a < 0.0	0	0	0	0
	0.0 <= a < 0.2	0	0	0	0
	0.2 <= a < 0.4	1	2	5	8
	0.4 <= a < 0.6	6	8	4	18
	0.6 <= a < 0.8	1	3	4	8
	0.8 <= a < 1.0	0	1	2	3
	1.0 <= a < 1.2	0	1	0	1
	1.2 <= a < 1.4	0	1	0	1
5	1.4 <= a < 1.6	0	0	0	0
	1.6 <= a < 1.8	0	0	0	0
	1.8 <= a < 2.0	0	0	0	0
	2.0 <= a	0	0	0	0
	Minimum	0.30	0.26	0.30	0.26
	Maximum	0.63	1.31	0.99	1.31
	Mean	0.48	0.61	0.56	0.56
	SD	0.10	0.28	0.22	0.23
	Number of Items	8	16	15	39

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.1.4

IRT-A Parameter Summary by Reporting Category: SC G6

Grade	IRT-a Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	a < 0.0	0	0	0	0
	0.0 <= a < 0.2	0	0	1	1
	0.2 <= a < 0.4	1	2	4	7
	0.4 <= a < 0.6	3	5	7	15
	0.6 <= a < 0.8	3	0	4	7
	0.8 <= a < 1.0	1	3	2	6
	1.0 <= a < 1.2	0	0	1	1
	1.2 <= a < 1.4	0	0	1	1
6	1.4 <= a < 1.6	0	0	0	0
	1.6 <= a < 1.8	0	0	0	0
	1.8 <= a < 2.0	0	0	0	0
	2.0 <= a	0	0	0	0
	Minimum	0.30	0.31	0.19	0.19
	Maximum	0.97	0.98	1.22	1.22
	Mean	0.59	0.59	0.57	0.58
	SD	0.21	0.25	0.27	0.25
	Number of Items	8	10	20	38

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.1.5

IRT-A Parameter Summary by Reporting Category: SC G7

Grade	IRT-a Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	a < 0.0	0	0	0	0
	0.0 <= a < 0.2	0	0	0	0
	0.2 <= a < 0.4	3	2	1	6
	0.4 <= a < 0.6	1	3	8	12
	0.6 <= a < 0.8	2	3	4	9
	0.8 <= a < 1.0	0	0	6	7
	1.0 <= a < 1.2	0	0	1	2
	1.2 <= a < 1.4	0	0	1	1
7	1.4 <= a < 1.6	0	0	0	0
	1.6 <= a < 1.8	0	0	0	0
	1.8 <= a < 2.0	0	0	0	0
	2.0 <= a	0	0	0	0
	Minimum	0.22	0.25	0.36	0.22
	Maximum	0.72	0.78	1.25	1.25
	Mean	0.46	0.53	0.68	0.63
	SD	0.19	0.19	0.23	0.24
	Number of Items	6	8	21	37

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.1.6

IRT-A Parameter Summary by Reporting Category: SC G8

Grade	IRT-a Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	a < 0.0	0	0	0	0
	0.0 <= a < 0.2	0	0	0	0
	0.2 <= a < 0.4	5	2	4	11
	0.4 <= a < 0.6	4	2	7	13
	0.6 <= a < 0.8	1	4	1	6
	0.8 <= a < 1.0	0	1	1	2
	1.0 <= a < 1.2	1	1	1	3
	1.2 <= a < 1.4	0	1	0	1
8	1.4 <= a < 1.6	0	1	0	1
	1.6 <= a < 1.8	0	0	0	0
	1.8 <= a < 2.0	0	0	0	0
	2.0 <= a	0	0	0	0
	Minimum	0.29	0.37	0.31	0.29
	Maximum	1.11	1.51	1.13	1.51
	Mean	0.50	0.75	0.56	0.60
	SD	0.23	0.37	0.23	0.30
	Number of Items	11	12	14	37

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.2.1

IRT-B Parameter Summary by Reporting Category: SC G3

Grade	IRT-b Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	b < -3.5	0	0	0	0
	-3.5 <= b < -3.0	0	0	0	0
	-3.0 <= b < -2.5	0	0	0	0
	-2.5 <= b < -2.0	0	0	0	0
	-2.0 <= b < -1.5	0	0	0	0
	-1.5 <= b < -1.0	0	0	0	0
	-1.0 <= b < -0.5	0	0	1	1
	-0.5 <= b < 0.0	2	3	0	5
	0.0 <= b < 0.5	2	3	1	6
	0.5 <= b < 1.0	2	7	4	14
3	1.0 <= b < 1.5	3	2	0	5
	1.5 <= b < 2.0	1	2	0	5
	2.0 <= b < 2.5	0	0	0	0
	2.5 <= b < 3.0	0	0	0	0
	3.0 <= b < 3.5	0	0	0	0
	3.5 <= b	0	0	0	0
	Minimum	-0.24	-0.24	-0.57	-0.57
	Maximum	1.68	1.80	0.93	1.85
	Mean	0.78	0.71	0.47	0.75
	SD	0.63	0.57	0.59	0.62
	Number of Items	10	17	6	36

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.2.2

IRT-B Parameter Summary by Reporting Category: SC G4

Grade	IRT-b Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
0.000	b < -3.5	0	0	0	0
	-3.5 <= b < -3.0	0	0	0	0
	-3.0 <= b < -2.5	0	0	0	0
	-2.5 <= b < -2.0	0	0	0	0
	-2.0 <= b < -1.5	0	0	1	1
	-1.5 <= b < -1.0	0	0	1	1
	-1.0 <= b < -0.5	1	1	0	2
	-0.5 <= b < 0.0	1	1	3	5
	0.0 <= b < 0.5	4	1	5	12
	0.5 <= b < 1.0	1	1	3	5
4	1.0 <= b < 1.5	3	0	2	7
	1.5 <= b < 2.0	1	0	1	2
	2.0 <= b < 2.5	1	0	0	1
	2.5 <= b < 3.0	0	0	0	0
	3.0 <= b < 3.5	0	0	0	0
	3.5 <= b	0	0	0	0
	Minimum	-0.99	-0.63	-1.54	-1.54
	Maximum	2.02	0.82	1.59	2.02
	Mean	0.67	0.11	0.20	0.41
	SD	0.89	0.63	0.81	0.81
	Number of Items	12	4	16	36

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.2.3

IRT-B Parameter Summary by Reporting Category: SC G5

Grade	IRT-b Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	b < -3.5	0	0	0	0
	-3.5 <= b < -3.0	0	0	0	0
	-3.0 <= b < -2.5	0	0	0	0
	-2.5 <= b < -2.0	1	0	0	1
	-2.0 <= b < -1.5	1	0	0	1
	-1.5 <= b < -1.0	0	0	1	1
	-1.0 <= b < -0.5	2	2	0	4
	-0.5 <= b < 0.0	1	4	5	10
	0.0 <= b < 0.5	1	1	3	5
	0.5 <= b < 1.0	0	4	4	8
5	1.0 <= b < 1.5	2	4	2	8
	1.5 <= b < 2.0	0	0	0	0
	2.0 <= b < 2.5	0	1	0	1
	2.5 <= b < 3.0	0	0	0	0
	3.0 <= b < 3.5	0	0	0	0
	3.5 <= b	0	0	0	0
	Minimum	-2.12	-0.62	-1.35	-2.12
	Maximum	1.41	2.14	1.03	2.14
	Mean	-0.30	0.55	0.17	0.23
	SD	1.27	0.85	0.67	0.92
	Number of Items	8	16	15	39

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.2.4

IRT-B Parameter Summary by Reporting Category: SC G6

	l summary by Nep			Reason	Total Number
Grade	IRT-b Range	Investigate	Evaluate	Scientifically	of Items*
	b < -3.5	0	0	0	0
	-3.5 <= b < -3.0	0	0	0	0
	-3.0 <= b < -2.5	0	0	0	0
	-2.5 <= b < -2.0	0	0	0	0
	-2.0 <= b < -1.5	0	0	1	1
	-1.5 <= b < -1.0	0	0	0	0
	-1.0 <= b < -0.5	0	2	5	7
	-0.5 <= b < 0.0	2	3	1	6
	0.0 <= b < 0.5	1	2	7	10
	0.5 <= b < 1.0	1	2	1	4
6	1.0 <= b < 1.5	0	0	3	3
	1.5 <= b < 2.0	2	0	0	2
	2.0 <= b < 2.5	1	0	2	3
	2.5 <= b < 3.0	0	1	0	1
	3.0 <= b < 3.5	1	0	0	1
	3.5 <= b	0	0	0	0
	Minimum	-0.45	-0.99	-1.65	-1.65
	Maximum	3.00	2.58	2.48	3.00
	Mean	1.16	0.18	0.24	0.42
	SD	1.24	0.99	1.04	1.12
	Number of Items	8	10	20	38

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.2

IRT-B Parameter Summary by Reporting Category: SC G7

	ecc. Cammany by map			Reason	Total Number
Grade	IRT-b Range	Investigate	Evaluate	Scientifically	of Items*
	b < -3.5	0	0	0	0
	-3.5 <= b < -3.0	0	0	0	0
	-3.0 <= b < -2.5	0	0	0	0
	-2.5 <= b < -2.0	0	0	0	0
	-2.0 <= b < -1.5	0	0	0	0
	-1.5 <= b < -1.0	0	0	2	2
	-1.0 <= b < -0.5	0	0	1	1
	-0.5 <= b < 0.0	0	0	7	7
	0.0 <= b < 0.5	2	5	3	10
	0.5 <= b < 1.0	0	1	4	7
7	1.0 <= b < 1.5	1	1	3	5
	1.5 <= b < 2.0	2	1	1	4
	2.0 <= b < 2.5	1	0	0	1
	2.5 <= b < 3.0	0	0	0	0
	3.0 <= b < 3.5	0	0	0	0
	3.5 <= b	0	0	0	0
	Minimum	0.11	0.36	-1.24	-1.24
	Maximum	2.33	1.57	1.71	2.33
	Mean	1.23	0.68	0.19	0.49
	SD	0.89	0.44	0.82	0.82
	Number of Items	6	8	21	37

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

Table C.5.2.6 *IRT-B Parameter Summary by Reporting Category: SC G8*

Grade	IRT-b Range	Investigate	Evaluate	Reason Scientifically	Total Number of Items*
	b < -3.5	0	0	0	0
	-3.5 <= b < -3.0	0	0	0	0
	-3.0 <= b < -2.5	0	0	0	0
	-2.5 <= b < -2.0	0	0	0	0
	-2.0 <= b < -1.5	0	1	0	1
	-1.5 <= b < -1.0	0	2	0	2
	-1.0 <= b < -0.5	0	0	2	2
	-0.5 <= b < 0.0	2	1	0	3
	0.0 <= b < 0.5	5	4	4	13
	0.5 <= b < 1.0	0	1	6	7
8	1.0 <= b < 1.5	3	1	2	6
	1.5 <= b < 2.0	1	1	0	2
	2.0 <= b < 2.5	0	0	0	0
	2.5 <= b < 3.0	0	1	0	1
	3.0 <= b < 3.5	0	0	0	0
	3.5 <= b	0	0	0	0
	Minimum	-0.30	-1.63	-0.51	-1.63
	Maximum	1.70	2.55	1.44	2.55
	Mean	0.55	0.24	0.48	0.42
	SD	0.67	1.25	0.58	0.86
this The second	Number of Items	11	12	14	37

^{*}Note. The total number of items in each low includes those not assigned to any reporting category.

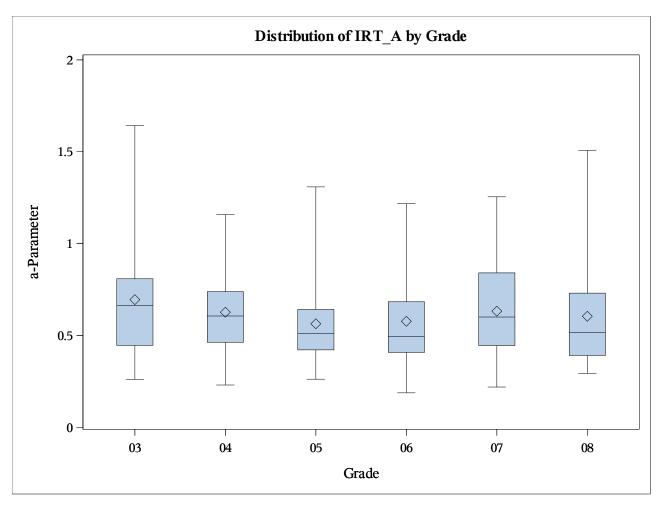
Table C.5.3

IRT Parameter Summary: Spring 2024 Operational SC G3-8

		No. of		25th		75th	
Grade	Parameter	Items	Minimum	Percentile	Median	Percentile	Maximum
	a	36	0.261	0.446	0.663	0.808	1.643
3	b	36	-0.568	0.324	0.821	1.108	1.848
	С	22	0.012	0.117	0.182	0.222	0.241
	а	36	0.23	0.463	0.606	0.738	1.158
4	b	36	-1.538	-0.023	0.268	1.087	2.015
	С	20	0.013	0.029	0.101	0.178	0.321
	а	39	0.262	0.421	0.511	0.642	1.309
5	b	39	-2.12	-0.392	0.278	0.865	2.136
	С	20	0.001	0.023	0.105	0.218	0.367
	а	38	0.187	0.408	0.494	0.683	1.217
6	b	38	-1.651	-0.453	0.151	1.019	3.005
	С	20	0	0.051	0.169	0.268	0.364
	а	37	0.219	0.445	0.6	0.84	1.255
7	b	37	-1.238	-0.002	0.443	1.053	2.333
	С	20	0.001	0.027	0.089	0.185	0.42
	а	37	0.293	0.391	0.516	0.73	1.508
8	b	37	-1.633	0.066	0.423	0.95	2.553
	С	19	0.016	0.06	0.135	0.241	0.327

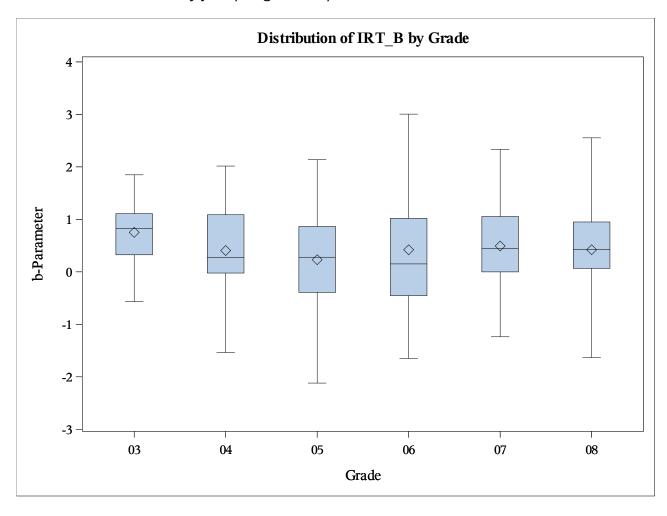
Plot C.5.1

IRT Item Parameter Summary for Spring 2024 Operational SC G3–8: A-Parameter



Plot C.5.2

IRT Item Parameter Summary for Spring 2024 Operational SC G3–8: B-Parameter



Plot C.5.3

IRT Item Parameter Summary for Spring 2024 Operational SC G3–8: C-Parameter

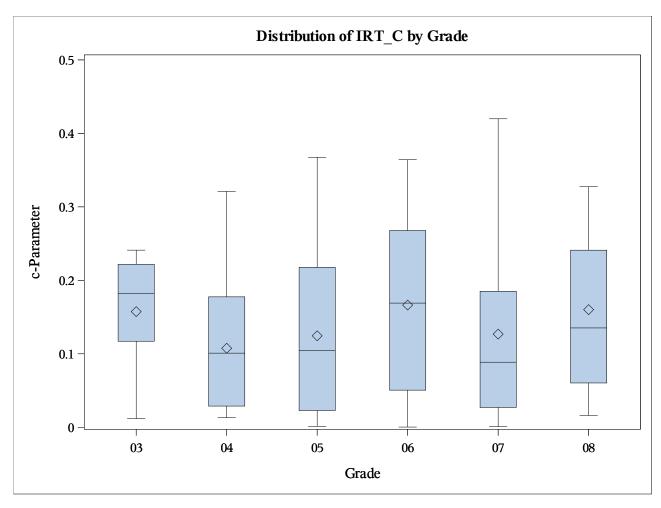


Table C.5.4

IRT Parameter Summary by Item Type: Spring 2024 Operational SC G3-8

			No. of		25th		75th	
Grade	Type	Parameter	Items	Minimum	Percentile	Median	Percentile	Maximum
	CR	а	3	0.564	0.564	0.663	0.67	0.67
		b	3	0.825	0.825	1.145	1.706	1.706
		a	21	0.261	0.663	0.744	0.973	1.643
	MC	b	21	-0.037	0.423	0.817	1.354	1.848
		С	21	0.012	0.133	0.185	0.222	0.241
3		a	1	0.803	0.803	0.803	0.803	0.803
5	MS	b	1	0.889	0.889	0.889	0.889	0.889
		С	1	0.013	0.013	0.013	0.013	0.013
	TPD	a	7	0.309	0.332	0.409	0.439	0.485
	IPD	b	7	-0.238	0.09	0.916	0.919	1.071
	TPI	a	4	0.452	0.508	0.567	0.61	0.649
		b	4	-0.568	-0.405	-0.134	0.27	0.565
	CR	a	3	0.571	0.571	0.711	0.756	0.756
	CIC	b	3	0.637	0.637	1.227	1.311	1.311
		a	14	0.403	0.612	0.747	0.945	1.093
	MC	b	14	-1.538	-0.05	0.104	0.676	1.978
		С	14	0.027	0.039	0.13	0.194	0.321
		a	2	0.558	0.558	0.613	0.668	0.668
	MS	b	2	0.308	0.308	0.732	1.156	1.156
4		С	2	0.019	0.019	0.023	0.027	0.027
		a	4	0.338	0.486	0.677	0.94	1.158
	TEI	b	4	-0.986	-0.81	0.298	1.412	1.593
		С	4	0.013	0.018	0.043	0.088	0.111
	TPD	a	6	0.23	0.271	0.346	0.514	0.599
	IFU	b	6	-0.209	0.109	0.746	1.234	2.015
	TPI	a	7	0.307	0.408	0.494	0.57	0.648
	161	b	7	-0.119	-0.078	0.159	0.396	0.819

Table C.5.4

IRT Parameter Summary by Item Type: Spring 2024 Operational SC G3-8 (continued)

			No. of		25th		75th	
Grade	Туре	Parameter	Items	Minimum	Percentile	Median	Percentile	Maximum
	CD	а	3	0.435	0.435	0.45	0.532	0.532
	CR	b	3	0.774	0.774	1.276	1.414	1.414
	ER	а	3	0.459	0.459	0.464	0.472	0.472
	EK	b	3	0.839	0.839	1.375	2.136	2.136
		a	10	0.447	0.543	0.613	0.882	1.309
	MC	b	10	-1.605	-0.617	-0.293	0.07	1.358
		С	10	0.027	0.1	0.186	0.26	0.295
		a	1	0.797	0.797	0.797	0.797	0.797
5	MS	b	1	0.464	0.464	0.464	0.464	0.464
		С	1	0.013	0.013	0.013	0.013	0.013
		a	16	0.298	0.397	0.512	0.66	1.001
	TEI	b	16	-2.12	-0.376	-0.008	0.761	1.388
		С	9	0.001	0.008	0.034	0.107	0.367
	TPD	a	3	0.262	0.262	0.336	0.433	0.433
	11.0	b	3	-0.226	-0.226	0.715	1.302	1.302
	TPI	a	3	0.309	0.309	0.342	0.52	0.52
	111	b	3	-0.407	-0.407	0.319	0.722	0.722
	CR	a	3	0.408	0.408	0.438	0.553	0.553
	CIN	b	3	1.74	1.74	2.576	3.005	3.005
	ER	a	2	0.187	0.187	0.303	0.418	0.418
	LIV	b	2	1.368	1.368	1.924	2.48	2.48
		a	13	0.448	0.632	0.868	0.972	1.217
	MC	b	13	-0.989	-0.655	-0.174	0.387	1.805
		С	13	0.035	0.15	0.204	0.289	0.364
		а	4	0.648	0.661	0.694	0.818	0.921
6	MS	b	4	-0.087	-0.041	0.033	0.72	1.379
		С	4	0*	0.007	0.034	0.088	0.122
		a	9	0.308	0.39	0.413	0.456	0.638
	TEI	b	9	-1.651	-0.453	0.157	0.634	2.069
		С	3	0.003	0.003	0.253	0.272	0.272
	TPD	a	5	0.289	0.292	0.304	0.309	0.499
	וויט	b	5	-0.618	0.093	0.338	1.019	2.266
	TDI	а	2	0.437	0.437	0.511	0.584	0.584
	TPI	b	2	-0.842	-0.842	-0.285	0.272	0.272

^{*}Actual c-parameter is 0.00029.

Table C.5.4

IRT Parameter Summary by Item Type: Spring 2024 Operational SC G3-8 (continued)

		Tilliary by ite	No. of		25th		75th	
Grade	Туре	Parameter	Items			Median		Maximum
		а	3	0.422	0.422	0.527	0.702	0.702
	CR	b	3	0.658	0.658	1.108	1.272	1.272
	רם	a	1	0.248	0.248	0.248	0.248	0.248
	ER	b	1	1.113	1.113	1.113	1.113	1.113
		а	8	0.219	0.487	0.683	0.937	1.255
	MC	b	8	-0.486	0.221	0.61	0.801	2.333
		С	8	0.001	0.052	0.185	0.301	0.42
		a	6	0.53	0.721	0.813	1.008	1.029
7	MS	b	6	-0.161	0.309	0.838	1.572	1.71
		С	6	0.02	0.021	0.05	0.075	0.167
		a	13	0.363	0.434	0.605	0.84	0.918
	TEI	b	13	-1.129	-0.397	0.326	0.417	1.761
		С	6	0.001	0.012	0.117	0.148	0.261
	TPD	a	3	0.329	0.329	0.445	0.481	0.481
		b	3	-0.326	-0.326	1.053	1.821	1.821
	TPI	a	3	0.361	0.361	0.509	0.544	0.544
	111	b	3	-1.238	-1.238	0.111	0.814	0.814
	CR	a	3	0.344	0.344	0.373	0.526	0.526
	CIX	b	3	0.95	0.95	1.275	1.705	1.705
	ER	a	2	0.357	0.357	0.474	0.592	0.592
	LIX	b	2	0.67	0.67	0.701	0.733	0.733
		a	13	0.405	0.493	0.789	1.046	1.508
	MC	b	13	-0.513	0.089	0.544	1.134	1.436
		С	13	0.016	0.131	0.194	0.241	0.327
		a	1	0.293	0.293	0.293	0.293	0.293
8	MS	b	1	-0.121	-0.121	-0.121	-0.121	-0.121
		С	1	0.020	0.020	0.020	0.020	0.020
		а	15	0.312	0.37	0.411	0.565	1.307
	TEI	b	15	-1.633	-0.506	0.095	0.452	2.553
		С	5	0.047	0.06	0.089	0.124	0.261
	TPD	a	1	0.391	0.391	0.391	0.391	0.391
	IFD	b	1	0.538	0.538	0.538	0.538	0.538
	TPI	a	2	0.516	0.516	0.612	0.708	0.708
	171	b	2	0.153	0.153	0.178	0.202	0.202

Table C.6
Statistically Flagged Operational Items: Spring 2024 Operational SC G3-8

Grade	Туре	No. of Items	N of Items Flagged for P-Value	N of Items Flagged for Point-Biserial Correlation	N of Items Flagged for DIF*	N of Items Flagged for Omitting
Graue	CR	3	1	0	0	0
	MC	21	0	0	0	0
3	MS	1	0	0	0	0
ی	TPD	7	0	0	0	0
	TPI	4	0	0	0	0
	CR	3	2	0	0	0
	MC	14	0	0	0	0
	MS	2	0	0	0	0
4	TEI	4	1	0	0	0
	TPD	6	1	0	0	0
	TPI	7	0	0	0	0
	CR	3	2	0	0	0
	ER*	1	1	0	1	0
	MC	10	0	0	0	0
5	MS	1	0	0	0	0
	TEI	16	1	0	0	0
	TPD	3	0	0	0	0
	TPI	3	0	0	0	0
	CR	3	3	0	0	0
	ER**	1	1	0	1	0
	MC	13	0	0	0	0
6	MS	4	0	0	1	0
	TEI	9	0	0	1	0
	TPD	5	1	0	0	0
	TPI	2	0	0	0	0

^{*} The number of flagged DIF items include both B and C DIF items.

^{**} Classical analyses were calculated and estimated separately for each dimension of the ER item, and the result summarize both dimensions.

Table C.6
Statistically Flagged Operational Items: Spring 2024 Operational SC G3-8 (continued)

Grade	Туре	No. of Items	N of Items Flagged for P-Value	N of Items Flagged for Point-Biserial Correlation	N of Items Flagged for DIF*	N of Items Flagged for Omitting
	CR	3	2	0	0	0
	ER**	1	0	0	0	0
	MC	8	0	0	0	0
7	MS	6	1	0	0	0
	TEI	13	1	0	1	0
	TPD	3	2	0	1	0
	TPI	3	0	0	0	0
	CR	3	1	0	0	0
	ER**	1	0	0	0	0
	MC	13	0	0	0	0
8	MS	1	0	0	0	0
	TEI	15	2	1	2	0
	TPD	1	0	0	0	0
	TPI	2	0	0	0	0

^{*} The number of flagged DIF items include both B and C DIF items.

^{**} Classical analyses were calculated and estimated separately for each dimension of the ER item, and the result summarize both dimensions.

Appendix D: Dimensionality

Dimensionality Reports: Science

Contents

Table D.1 Zq1 Statistics and Summary Data: Spring 2024 Operational SC G3-8

Table D.2 Q3 Statistics and Summary Data: Spring 2024 Operational SC G3-8

Table D.3 Reporting Category Intercorrelation Coefficients: Spring 2024 Operational SC G3-8

Table D.4 First and Second Eigenvalues: Spring 2024 Operational SC G3-8

Plot D.1 Principal Component Analysis: Spring 2024 Operational SC G3-8

Table D.1

Zq1 Statistics and Summary Data: Spring 2024 Operational SC G3-8

Grade	Туре	Minimum	25th Percentile	Median	75th Percentile	Maximum	No. of Items with Poor Fit
Grade	CR	92.82	92.82	128.81	380.88	380.88	1
	MC	10.74	21.82	27.43	40.25	194.67	1
3	MS	44.18	44.18	44.18	44.18	44.18	0
	TPD	141.93	234.23	290.10	315.00	539.84	7
	TPI	41.24	82.14	147.35	172.45	173.24	2
	CR	49.69	49.69	102.26	191.14	191.14	1
	MC	7.28	20.15	25.08	45.09	127.70	0
	MS	19.75	19.75	28.18	36.62	36.62	0
4	TEI	11.39	31.45	76.04	110.06	119.54	0
	TPD	42.35	70.17	104.13	334.20	343.84	2
	TPI	24.48	47.87	77.74	128.63	230.36	1
	CR	30.82	30.82	103.58	142.14	142.14	1
	ER	55.65	55.65	73.87	103.79	103.79	0
	MC	13.93	19.12	28.10	51.63	87.46	0
5	MS	34.17	34.17	34.17	34.17	34.17	0
	TEI	18.81	30.58	57.83	133.72	395.82	4
	TPD	14.77	14.77	57.27	508.37	508.37	1
	TPI	107.18	107.18	132.81	134.88	134.88	2
	CR	14.92	14.92	20.18	47.17	47.17	0
	ER	37.33	37.33	101.93	166.52	166.52	1
	MC	12.14	21.96	31.91	46.67	85.07	0
6	MS	24.09	30.31	48.33	96.69	133.24	1
	TEI	12.64	31.53	83.33	202.13	721.85	3
	TPD	61.51	113.88	267.89	292.78	583.74	3
	TPI	66.93	66.93	71.28	75.63	75.63	0
	CR	21.30	21.30	42.79	101.67	101.67	0
	ER	89.85	89.85	89.85	89.85	89.85	0
	MC	7.48	8.76	15.38	33.87	69.55	0
7	MS	11.08	29.37	30.15	30.46	43.67	0
	TEI	12.19	29.21	83.07	103.02	266.60	3
	TPD	44.73	44.73	74.08	259.92	259.92	1
	TPI	30.61	30.61	58.05	67.51	67.51	0

Table D.1

Zq1 Statistics and Summary Data: Spring 2024 Operational SC G3-8 (continued)

Grade	Туре	Minimum	25th Percentile	Median	75th Percentile	Maximum	No. of Items with Poor Fit
	CR	51.22	51.22	83.80	122.94	122.94	0
	ER	140.28	140.28	148.77	157.26	157.26	2
	MC	9.77	15.05	22.09	28.00	85.56	0
8	MS	66.48	66.48	66.48	66.48	66.48	0
	TEI	14.38	21.31	76.60	179.51	529.77	5
	TPD	215.83	215.83	215.83	215.83	215.83	1
	TPI	57.75	57.75	158.87	260.00	260.00	1

Table D.2

Q3 Statistics and Summary Data: Spring 2024 Operational SC G3-8 (Done)

	Average Zero		5th		95th	
Grade	Order Correlation	Minimum	Percentile	Median	Percentile	Maximum
3	0.155	-0.151	-0.079	-0.015	0.084	0.186
4	0.201	-0.131	-0.077	-0.015	0.096	0.189
5	0.194	-0.121	-0.071	-0.004	0.101	0.315
6	0.180	-0.190	-0.083	-0.013	0.102	0.232
7	0.213	-0.293	-0.099	-0.001	0.117	0.301
8	0.159	-0.182	-0.073	-0.014	0.114	0.398

Table D.3
Reporting Category Intercorrelation Coefficients: Spring 2024 Operational SC G3-8

Grade	Reporting Category	1 Investigate	2 Evaluate	3 Reason Scientifically
3	1 Investigate	1.00	-	-
	2 Evaluate	0.69	1.00	-
	3 Reason Scientifically	0.61	0.67	1.00
4	1 Investigate	1.00	-	-
	2 Evaluate	0.64	1.00	-
	3 Reason Scientifically	0.76	0.66	1.00
	1 Investigate	1.00	-	-
5	2 Evaluate	0.72	1.00	-
	3 Reason Scientifically	0.69	0.77	1.00
6	1 Investigate	1.00	-	-
	2 Evaluate	0.66	1.00	-
	3 Reason Scientifically	0.69	0.75	1.00
	1 Investigate	1.00	-	-
7	2 Evaluate	0.59	1.00	-
	3 Reason Scientifically	0.67	0.77	1.00
8	1 Investigate	1.00	-	-
	2 Evaluate	0.68	1.00	-
	3 Reason Scientifically	0.73	0.72	1.00

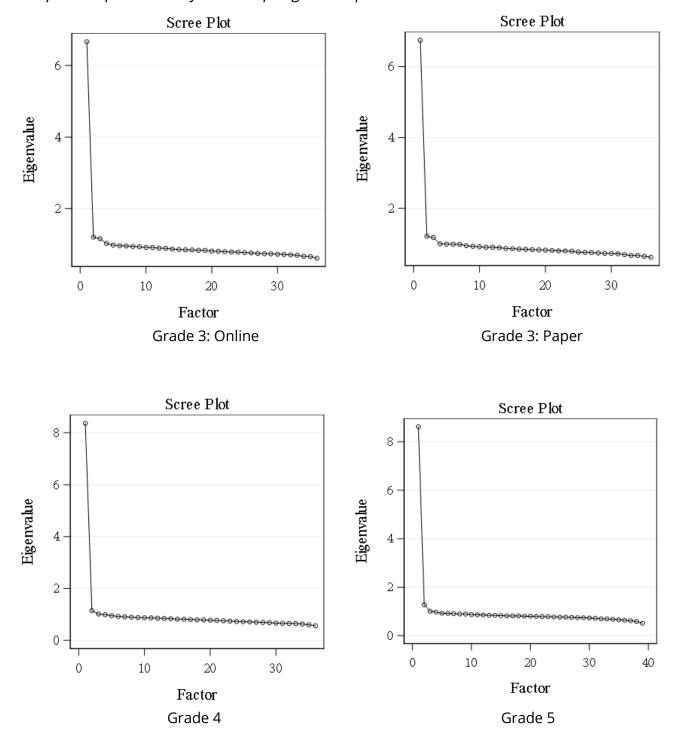
Table D.4

First and Second Eigenvalue: Spring 2024 Operational SC G3-8 (Done)

Grade	Mode	First Eigenvalue	Second Eigenvalue	Ratio
3	Online	6.664	1.202	5.544
3	Paper	6.748	1.208	5.615
4	Online	8.367	1.141	7.333
5	Online	8.621	1.276	6.756
6	Online	7.948	1.225	6.488
7	Online	9.040	1.180	7.661
8	Online	7.312	1.177	6.212

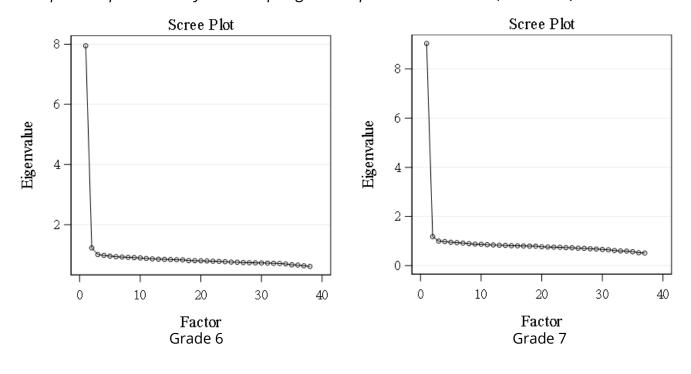
Plot D.1

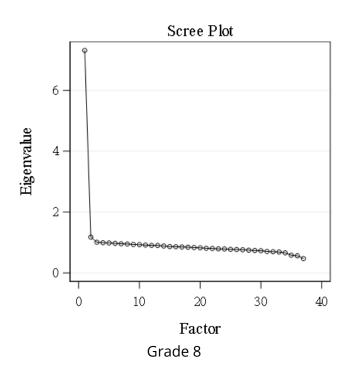
Principal Component Analysis Plot: Spring 2024 Operational SC G3-8



Plot D.1

Principal Component Analysis Plot: Spring 2024 Operational SC G3-8 (continued)



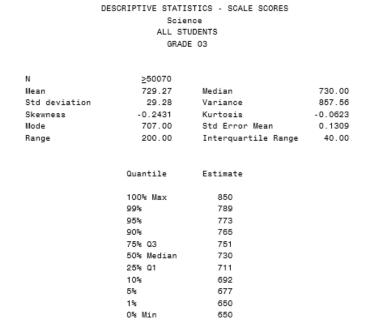


Appendix E: Scale Distribution and Statistical Report

Science

Contents
Table E.1.1 Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science Grade 3
Table E.1.2 Frequency Distribution of Scale Scores: Spring 2024 Operational Science Grade 3
Table E.2.1 Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science Grade 4
Table E.2.2 Frequency Distribution of Scale Scores: Spring 2024 Operational Science Grade 4
Table E.3.1 Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science Grade 5
Table E.3.2 Frequency Distribution of Scale Scores: Spring 2024 Operational Science Grade 5
Table E.4.1 Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science Grade 6
Table E.4.2 Frequency Distribution of Scale Scores: Spring 2024 Operational Science Grade 6
Table E.5.1 Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science Grade 7
Table E.5.2 Frequency Distribution of Scale Scores: Spring 2024 Operational Science Grade 7
Table E.6.1 Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science Grade 8
Table E.6.2 Frequency Distribution of Scale Scores: Spring 2024 Operational Science Grade 8

Table E.1.1
Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science: Grade 3



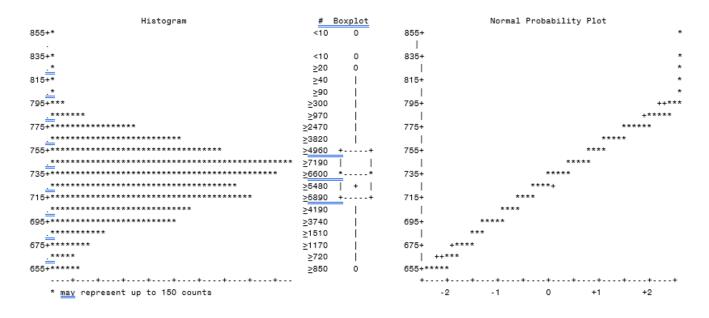
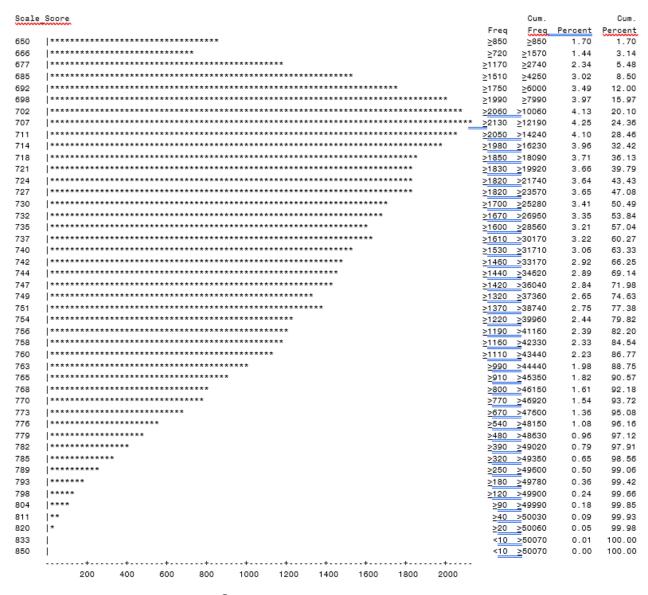


Table E.1.2
Frequency Distribution of Scale Scores: Spring 2024 Operational Science: Grade 3

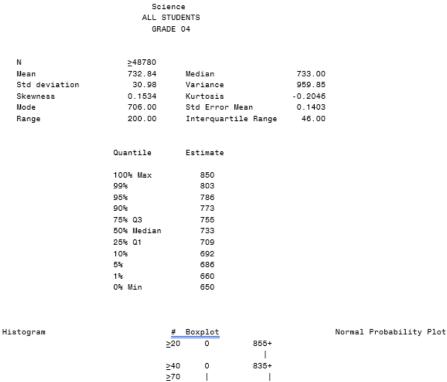
FREQUENCY DISTRIBUTION - SCALE SCORES
Science
ALL STUDENTS
GRADE 03



Frequency

Table E.2.1
Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science: Grade 4

DESCRIPTIVE STATISTICS - SCALE SCORES



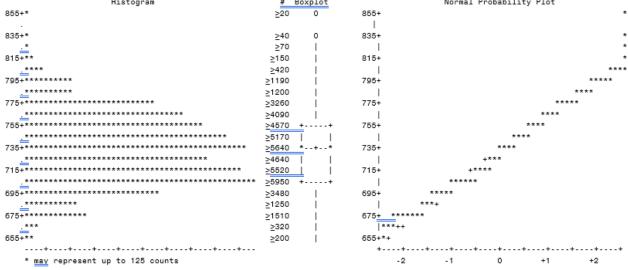
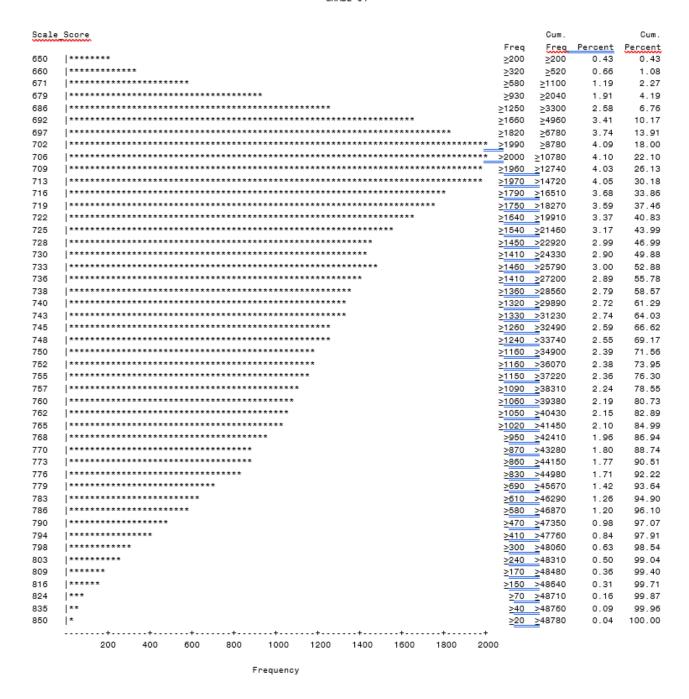


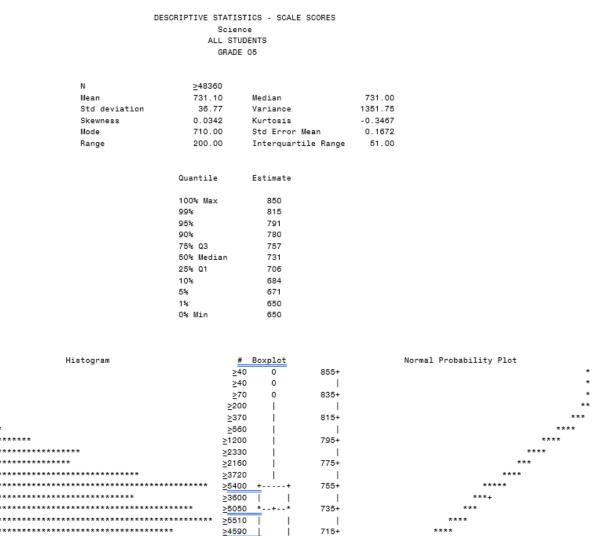
Table E.2.2 Frequency Distribution of Scale Scores: Spring 2024 Operational Science: Grade 4

FREQUENCY DISTRIBUTION - SCALE SCORES
Science
ALL STUDENTS
GRADE 04



251

Table E.3.1
Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science: Grade 5



695+

675+

655+*****

-2

- 1

0

≥3050 ≥3000

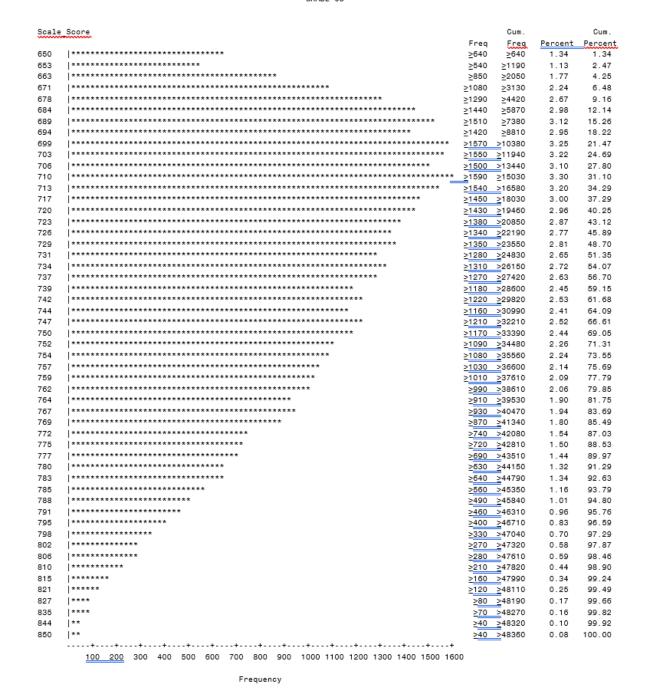
≥2950 ≥2370

≥850

855+*

* may represent up to 115 counts

Table E.3.2
Frequency Distribution of Scale Scores: Spring 2024 Operational Science: Grade 5



253

Table E.4.1
Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science: Grade 6

DESCRIPTIVE STATISTICS - SCALE SCORES
Science
ALL STUDENTS
GRADE 06

N	≥47810		
Mean	725.47	Median	722.00
Std deviation	32.57	Variance	1060.64
Skewness	0.2579	Kurtosis	-0.1429
Mode	692.00	Std Error Mean	0.1489
Range	200.00	Interquartile Range	44.00
	Quantile	Estimate	

QUANTILE	ESTIMAT
100% Max	850
99%	804
95%	780
90%	768
75% Q3	747
50% Median	722
25% Q1	703
10%	688
5%	677
1%	654
0% Min	650

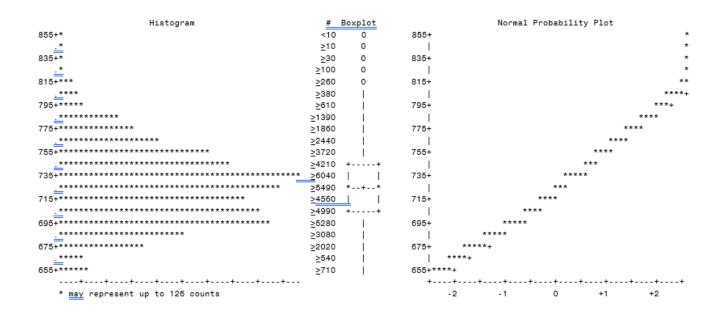


Table E.4.2
Frequency Distribution of Scale Scores: Spring 2024 Operational Science: Grade 6

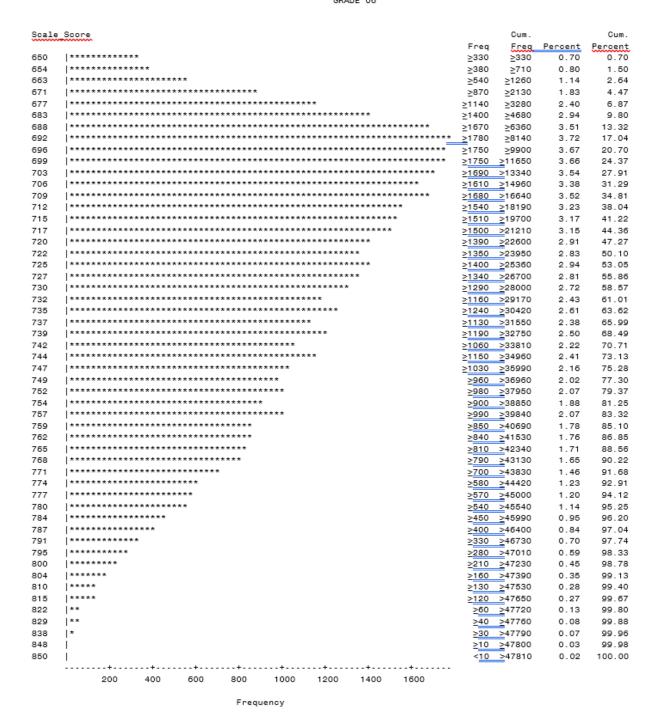
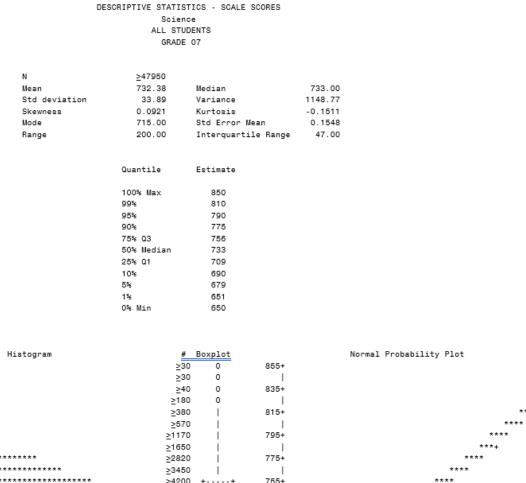


Table E.5.1
Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science: Grade 7



855+*

835+*

Table E.5.2
Frequency Distribution of Scale Scores: Spring 2024 Operational Science: Grade 7

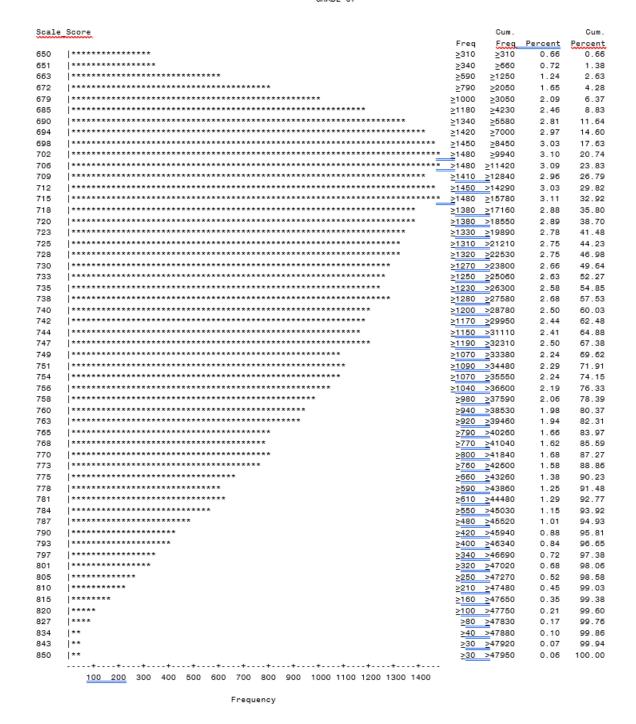


Table E.6.1
Scale Score Descriptive Statistics and Plots: Spring 2024 Operational Science: Grade 8

DESCRIPTIVE STATISTICS - SCALE SCORES
Science
ALL STUDENTS

GRADE 08 ≥48220 730.20 729.00 Std deviation 31.79 Variance 1010.33 -0.3326 Skewness 0.0377 Kurtosis Std Error Mean Mode 701.00 0.1447 Range 200.00 Interquartile Range 46.00

Quantile	Estimate
100% Max	850
99%	801
95%	783
90%	771
75% Q3	754
50% Median	729
25% Q1	708
10%	688
5%	677
1%	653
0% Min	650

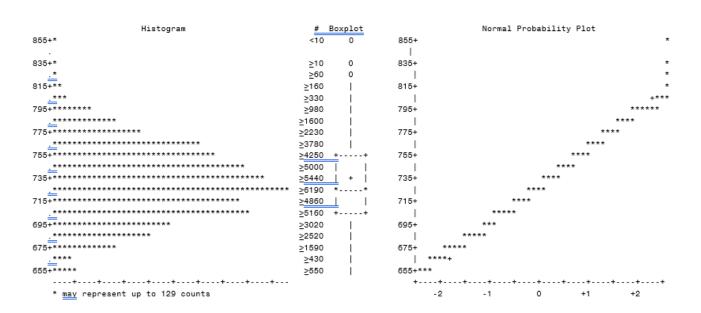


Table E.6.2
Frequency Distribution of Scale Scores: Spring 2024 Operational Science: Grade 8

Scale	Score		Cum.		Cum.
650	*******	Freq ≥250	Freg ≥250	Percent 0.53	Percent 0.53
653	*************	>290	≥550	0.61	1.14
663	***********	≥290 ≥430	≥980	0.01	2.04
671	******************	≥670	≥1650	1.39	3.44
677	*************************************	≥910	≥2570	1.91	5.34
683	************************************	>1160	>3740	2.42	7.76
688	*************************************	>1360	>5100	2.82	10.58
693		≥1530	≥6630	3.18	13.76
697		>1480	≥8120	3.09	16.85
701	 ***********************************	≥1720	≥9850	3.58	20.43
705			≥11570	3.57	24.00
708	*********************	>1710	>13280	3.55	27.55
711	*********************	>1660	>14950	3.46	31.01
715		>1620	>16570	3.36	34.37
718		>1570	>18140	3.26	37.63
720		>1620	>19770	3.38	41.01
723	********************	>1550	>21330	3.23	44.23
726	 ***********************************	>1570	>22900	3.26	47.49
729	****************	>1430	>24330	2.98	50.47
731		>1450	>25790	3.02	53.49
734	***********	>1350	>27140	2.81	56.30
736	**************	>1360	>28510	2.84	59.14
739	*************	>1260	>29780	2.63	61.77
741	************	>1210	>31000	2.53	64.30
744	**************	≥1290	≥32290	2.68	66.98
746	**********	>1250	>33550	2.61	69.59
749	*************	≥1230	<u>≥</u> 34780	2.55	72.14
751	**********	≥1140	≥35930	2.38	74.52
754	***********	>1090	>37030	2.28	76.80
756	********	≥1030	≥38060	2.14	78.94
759	********	>970	>39030	2.01	80.95
761	********	≥1020	≥40060	2.12	83.08
764	*************************	≥910	>40970	1.89	84.97
766	**********	≥980	<u>≥</u> 41950	2.04	87.01
769	*******************	≥860	>42820	1.79	88.80
771	************	≥810	<u>></u> 43630	1.69	90.49
774	****************	≥750	>44380	1.56	92.05
777	********************	≥670	≥45050	1.39	93.44
780	*************	≥580	<u>></u> 45630	1.20	94.65
783	**************	≥550	>46190	1.14	95.79
786	**************	≥470	≥46660	0.98	96.76
790	**********	≥420	>47080	0.88	97.64
793	**********	≥330	<u>></u> 47420	0.70	98.34
797	******	>220	>47640	0.46	98.80
801	******	≥180	≥47830	0.39	99.19
806	******	≥140	>47970	0.30	99.49
811	*****	≥100	<u>></u> 48080	0.22	99.71
816	***	≥50	>48140	0.12	99.83
822	**	≥30	<u>></u> 48170	0.07	99.91
829	**	≥30	≥48200	0.06	99.97
838	<u> </u> *	≥10	>48210	0.02	99.99
850		<10	≥48220	0.01	100.00
	++++++++		_		
	100 200 300 400 500 600 700 800 900 1000 1100 1200 1300 1400 1500 1600 170	0			

Frequency

Appendix F: Reliability and Classification Accuracy

Reliability and Classification Accuracy Reports Science

Contents

Tables F.1.1–F.1.2 Reliability and SEM for Overall and Subgroups: Spring 2024 Operational SC G3-8

Table F.2 Cronbach's Alpha and Marginal Reliability: Spring 2024 Operational SC G3-8

Table F.3.1–F.3.9 Classification Accuracy and Decision Consistency Matrices: Spring 2024 Operational SC G3-8

Table F.1.1
Reliability for Overall and Subgroups: Spring 2024 Operational Science

		Grade						
Category	Subgroup*	3	4	5	6	7	8	
All S	Students	0.871	0.899	0.902	0.878	0.899	0.880	
Condon	Female	0.865	0.889	0.894	0.865	0.892	0.872	
Gender	Male	0.876	0.907	0.910	0.890	0.905	0.886	
	African American	0.834	0.862	0.875	0.837	0.866	0.840	
	AI/AN	0.845	0.881	0.890	0.861	0.873	0.863	
	Asian	0.897	0.907	0.909	0.888	0.911	0.903	
Ethnicity	Hispanic/Latino	0.861	0.887	0.901	0.877	0.894	0.876	
	NHPI	0.907	0.899	0.904	0.866	0.890	0.878	
	Two or More	0.867	0.892	0.895	0.869	0.895	0.872	
	White	0.866	0.894	0.893	0.871	0.894	0.870	
Economically	No	0.866	0.892	0.888	0.869	0.894	0.869	
Disadvantaged	Yes	0.850	0.882	0.889	0.859	0.881	0.861	
Facilials Lagrages	No	0.870	0.898	0.901	0.877	0.898	0.877	
English Learner	Yes	0.772	0.791	0.804	0.760	0.781	0.736	
Education	Regular	0.870	0.896	0.898	0.876	0.896	0.877	
Classification	Special	0.849	0.878	0.891	0.841	0.866	0.830	
Costion FO4	No	0.872	0.900	0.903	0.879	0.899	0.880	
Section 504	Yes	0.846	0.878	0.889	0.869	0.885	0.860	
Minuset	No	0.871	0.899	0.902	0.878	0.899	0.880	
Migrant	Yes	0.829	0.888	0.883	0.876	0.870	0.869	
Llowedoss Ctatus	No	0.871	0.899	0.902	0.878	0.899	0.879	
Homeless Status	Yes	0.842	0.864	0.874	0.836	0.850	0.841	
Military	No	0.870	0.898	0.902	0.878	0.898	0.879	
Affiliation	Yes	0.871	0.884	0.885	0.873	0.896	0.874	
Foster Care	No	0.871	0.899	0.903	0.878	0.899	0.880	
Status	Yes	0.856	0.881	0.880	0.866	0.865	0.840	

^{*} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

Table F.1.2 SEM for Overall and Subgroups: Spring 2024 Operational Science

		Grade					
Category	Subgroup*	3	4	5	6	7	8
All S	Students	3.41	3.31	3.72	3.93	3.85	3.74
Candan	Female	3.43	3.32	3.74	3.95	3.85	3.75
Gender	Male	3.41	3.30	3.67	3.88	3.85	3.75
	African American	3.36	3.24	3.56	3.80	3.76	3.67
	AI/AN	3.44	3.38	3.81	3.94	3.84	3.78
	Asian	3.40	3.31	3.83	4.05	3.91	3.71
Ethnicity	Hispanic/Latino	3.35	3.27	3.59	3.85	3.82	3.71
	NHPI	3.39	3.37	3.85	4.02	3.91	3.75
	Two or More	3.45	3.35	3.77	3.98	3.87	3.77
	White	3.47	3.36	3.82	4.00	3.86	3.76
Economically	No	3.46	3.36	3.87	4.02	3.86	3.75
Disadvantaged	Yes	3.40	3.28	3.63	3.86	3.80	3.72
English Learner	No	3.43	3.32	3.72	3.93	3.85	3.76
English Learner	Yes	3.16	3.10	3.21	3.43	3.42	3.33
Education	Regular	3.43	3.34	3.74	3.94	3.87	3.77
Classification	Special	3.28	3.16	3.34	3.55	3.58	3.49
Costion FO4	No	3.42	3.31	3.72	3.92	3.86	3.77
Section 504	Yes	3.36	3.27	3.58	3.79	3.80	3.68
Migraph	No	3.41	3.31	3.72	3.93	3.85	3.74
Migrant	Yes	3.25	3.27	3.59	3.91	3.69	3.72
Llavasalasa Chahua	No	3.41	3.31	3.72	3.93	3.85	3.76
Homeless Status	Yes	3.30	3.17	3.50	3.66	3.75	3.63
Military	No	3.42	3.32	3.71	3.92	3.86	3.75
Affiliation	Yes	3.44	3.39	3.85	4.00	3.86	3.76
Foster Care	No	3.41	3.31	3.70	3.93	3.85	3.74
Status	Yes	3.37	3.25	3.59	3.72	3.74	3.66

^{*} Al/AN = American Indian or Alaska Native. NHPI = Native Hawaiian or Other Pacific Islander.

Table F.2 Cronbach's Alpha and Marginal Reliability: Spring 2024 Operational SC G3-8

Grade	Cronbach's Alpha	Marginal Reliability
3	0.871	0.87
4	0.899	0.90
5	0.902	0.90
6	0.878	0.89
7	0.899	0.91
8	0.880	0.89

Table F.3.1 Classification Accuracy Matrices: Spring 2024 Operational SC G3-8

		Unsatisfactory	Approaching				
Grade	Level	(1)	Basic (2)	Basic (3)	Mastery (4)	Advanced (5)	Total
	1	0.09	0.02	0.00	0.00	0.00	0.12
	2	0.03	0.23	0.05	0.00	0.00	0.30
3	3	0.00	0.06	0.21	0.05	0.00	0.32
3	4	0.00	0.00	0.05	0.13	0.05	0.24
_	5	0.00	0.00	0.00	0.01	0.01	0.02
	Total	0.12	0.31	0.31	0.19	0.06	1.00
	1	0.14	0.03	0.00	0.00	0.00	0.17
	2	0.04	0.14	0.04	0.00	0.00	0.22
4	3	0.00	0.05	0.20	0.05	0.00	0.31
4	4	0.00	0.00	0.04	0.17	0.02	0.23
	5	0.00	0.00	0.00	0.01	0.05	0.07
	Total	0.18	0.23	0.28	0.23	0.08	1.00
	1	0.15	0.03	0.00	0.00	0.00	0.18
	2	0.03	0.17	0.04	0.00	0.00	0.25
	3	0.00	0.05	0.15	0.05	0.00	0.24
5	4	0.00	0.00	0.04	0.18	0.03	0.26
	5	0.00	0.00	0.00	0.02	0.06	0.07
	Total	0.18	0.25	0.23	0.25	0.09	1.00
	1	0.20	0.04	0.00	0.00	0.00	0.24
	2	0.05	0.16	0.05	0.00	0.00	0.26
6	3	0.00	0.06	0.18	0.05	0.00	0.29
6	4	0.00	0.00	0.03	0.12	0.02	0.18
	5	0.00	0.00	0.00	0.01	0.03	0.04
	Total	0.24	0.26	0.27	0.18	0.05	1.00
	1	0.14	0.03	0.00	0.00	0.00	0.17
	2	0.03	0.16	0.05	0.00	0.00	0.24
_	3	0.00	0.05	0.19	0.05	0.00	0.29
7	4	0.00	0.00	0.04	0.20	0.02	0.26
	5	0.00	0.00	0.00	0.01	0.03	0.04
	Total	0.18	0.24	0.28	0.25	0.05	1.00
	1	0.10	0.03	0.00	0.00	0.00	0.13
	2	0.03	0.22	0.05	0.00	0.00	0.31
o	3	0.00	0.06	0.19	0.05	0.00	0.29
8	4	0.00	0.00	0.04	0.16	0.03	0.23
	5	0.00	0.00	0.00	0.01	0.03	0.04
	Total	0.14	0.30	0.28	0.23	0.05	1.00

Table F.3.2

Decision Consistency Matrices: Spring 2024 Operational SC G3-8

		Unsatisfactory	Approaching				
Grade	Level	(1)	Basic (2)	Basic (3)	Mastery (4)	Advanced (5)	Total
3	1	0.09	0.04	0.00	0.00	0.00	0.14
	2	0.03	0.19	0.07	0.00	0.00	0.30
	3	0.00	0.07	0.17	0.06	0.01	0.30
	4	0.00	0.01	0.07	0.10	0.04	0.21
	5	0.00	0.00	0.01	0.03	0.02	0.05
	Total	0.12	0.31	0.31	0.19	0.06	1.00
	1	0.14	0.05	0.01	0.00	0.00	0.19
	2	0.04	0.11	0.06	0.00	0.00	0.21
4	3	0.01	0.07	0.16	0.06	0.00	0.29
4	4	0.00	0.00	0.06	0.14	0.03	0.23
	5	0.00	0.00	0.00	0.03	0.05	0.08
	Total	0.18	0.23	0.28	0.23	0.08	1.00
	1	0.14	0.05	0.00	0.00	0.00	0.19
	2	0.04	0.14	0.06	0.01	0.00	0.24
5	3	0.00	0.06	0.11	0.06	0.00	0.23
Э	4	0.00	0.01	0.06	0.15	0.03	0.25
	5	0.00	0.00	0.00	0.03	0.05	0.09
	Total	0.18	0.25	0.23	0.25	0.09	1.00
	1	0.19	0.06	0.01	0.00	0.00	0.25
	2	0.05	0.13	0.07	0.00	0.00	0.25
6	3	0.01	0.07	0.15	0.05	0.00	0.27
0	4	0.00	0.00	0.05	0.10	0.02	0.18
	5	0.00	0.00	0.00	0.02	0.03	0.05
	Total	0.24	0.26	0.27	0.18	0.05	1.00
	1	0.13	0.05	0.00	0.00	0.00	0.19
	2	0.04	0.12	0.06	0.00	0.00	0.23
7	3	0.00	0.06	0.15	0.06	0.00	0.28
/	4	0.00	0.00	0.06	0.17	0.02	0.26
	5	0.00	0.00	0.00	0.02	0.03	0.05
	Total	0.18	0.24	0.28	0.25	0.05	1.00
	1	0.10	0.05	0.00	0.00	0.00	0.15
	2	0.04	0.18	0.07	0.01	0.00	0.29
8	3	0.00	0.07	0.15	0.06	0.00	0.28
٥	4	0.00	0.01	0.06	0.14	0.03	0.23
	5	0.00	0.00	0.00	0.02	0.03	0.05
	Total	0.14	0.30	0.28	0.23	0.05	1.00

Table F.3.3

Estimates of Accuracy and Consistency of Achievement Level Classification

Grade	Accuracy	Consistency	PChance	Карра
3	0.675	0.566	0.246	0.425
4	0.704	0.599	0.224	0.483
5	0.700	0.595	0.217	0.482
6	0.692	0.587	0.234	0.461
7	0.716	0.613	0.233	0.495
8	0.702	0.593	0.241	0.464

Table F.3.4

Accuracy of Classification at Each Achievement Level

Grade	Unsatisfactory (1)	Approaching Basic (2)	Basic (3)	Mastery (4)	Advanced (5)
3	0.792	0.751	0.647	0.561	0.625
4	0.816	0.635	0.653	0.726	0.797
5	0.835	0.682	0.601	0.703	0.752
6	0.829	0.615	0.636	0.697	0.775
7	0.825	0.655	0.658	0.755	0.773
8	0.801	0.724	0.632	0.700	0.734

Table F.3.5

Accuracy of Dichotomous Categorizations by Form (PAC Metric)

Grade	1 / 2+3+4+5	1+2 / 3+4+5	1+2+3 / 4+5	1+2+3+4 / 5
3	0.949	0.888	0.891	0.942
4	0.930	0.898	0.911	0.962
5	0.937	0.903	0.906	0.950
6	0.911	0.887	0.917	0.972
7	0.935	0.898	0.910	0.970
8	0.941	0.891	0.905	0.963

Table F.3.6

Consistency of Dichotomous Categorizations by Form (PAC Metric)

Grade	1 / 2+3+4+5	1+2 / 3+4+5	1+2+3 / 4+5	1+2+3+4 / 5
3	0.925	0.844	0.847	0.924
4	0.900	0.858	0.875	0.945
5	0.910	0.864	0.868	0.929
6	0.874	0.843	0.884	0.960
7	0.906	0.858	0.873	0.958
8	0.914	0.848	0.866	0.949

Table F.3.7
Kappa of Dichotomous Categorizations by Form (PAC Metric)

Grade	1 / 2+3+4+5	1+2 / 3+4+5	1+2+3 / 4+5	1+2+3+4 / 5
3	0.679	0.681	0.605	0.263
4	0.677	0.704	0.706	0.623
5	0.710	0.723	0.705	0.564
6	0.667	0.685	0.667	0.555
7	0.690	0.709	0.703	0.546
8	0.657	0.692	0.668	0.473

Table F.3.8

Accuracy of Dichotomous Categorizations: False Positive Rates (PAC Metric)

Grade	1 / 2+3+4+5	1+2 / 3+4+5	1+2+3 / 4+5	1+2+3+4 / 5
3	0.025	0.050	0.053	0.050
4	0.032	0.046	0.050	0.025
5	0.029	0.046	0.049	0.031
6	0.040	0.055	0.049	0.020
7	0.030	0.048	0.047	0.021
8	0.026	0.051	0.052	0.027

Table F.3.9

Accuracy of Dichotomous Categorizations: False Negative Rates (PAC Metric)

Grade	1 / 2+3+4+5	1+2 / 3+4+5	1+2+3 / 4+5	1+2+3+4 / 5
3	0.026	0.062	0.055	0.008
4	0.037	0.056	0.038	0.014
5	0.033	0.051	0.045	0.018
6	0.049	0.058	0.034	0.008
7	0.035	0.054	0.043	0.009
8	0.034	0.058	0.043	0.010

Appendix G: Accommodated Print and Braille Creation

Louisiana believes that all students requiring test accommodations should be presented with the same rigor as students taking tests without accommodations. To ensure this, Louisiana creates accommodated versions of the operational test form for each test administration, allowing all students to take the same items regardless of the need for an accommodated presentation. Careful consideration is given to all items that are used for Louisiana assessments for their ability to be faithfully represented in accommodated print (AP) and braille formats. Fairness for all populations, item integrity, and student-item interaction for technology-enhanced (TE) items are all factors when selecting the items that will appear on a Louisiana form. TE items are modified so that students who interact with an item on an AP or braille form will have a similar and equivalent experience to students who interact with that same item in the online environment. This maintains both the rigor and the content being assessed. Some examples of the modification process are provided below.

- Drag-and-drop items in the online environment require a student to place the answer options in an interactive table. For the AP and braille forms, the student is presented with a table with the same information as the interactive table (column or row headers, any completed cells, and blank spaces) and the answer options are listed below the table (similar to the online form in which the options are listed either below or to the right of the table). The directions are modified to ask the student to write the letter or number of the correct answer in its corresponding box. Students are also able to circle the text and draw arrows to indicate where it should be placed or add labels to the answer choices and write only the label in the box, as long as the intended response is clear to the test administrator who will transcribe the answers into the online system.
- Match interaction items in the online environment require a student to select a checkbox in one or more columns for each of multiple rows. In the AP and braille

forms, the student is provided with a table and asked to mark or select the correct answer in each row.

- Highlight-text items or item parts in the online environment require a student to
 click on the selected text, which highlights the selected word, phrase, or sentence.
 In the AP and braille forms, the text is presented in the same format and the
 student is asked to circle the answer. Where only certain words or phrases are
 selectable in the online system, those options are underlined in the AP and braille
 forms to indicate which words and/or phrases the student should select from.
- Drop-down menu items in the online environment have answer options in a drop-down menu format, oftentimes as part of a complete sentence. The AP and braille forms display the item with a blank line in place of the drop-down menu in the sentence, with all the answer options for the drop-down menu presented vertically below the sentence and lettered or numbered. The directions are then modified to ask the student to select the letter/number of the word/phrase that belongs in the blank.
- Short answer items in the online environment require a student to type the answer in a box. In the AP and braille forms, a box is provided for the student to write the response.
- Keypad input items in the online environment require a student to enter a numeric response including all rational and irrational numbers as well as expressions and equations. In the AP forms, a box is provided for the student to write the response.
 In the braille forms, students are asked to answer on the paper provided.
- Graphing items, including coordinate planes, number lines, line plots, and bar graphs, in the online environment require a student to complete a graph by

plotting points, adding Xs to create a line plot, or raising/lowering bars to create a bar graph or histogram. In the AP and braille forms, the student is provided with the same coordinate plane, number line, line plot, or bar graph as in the online item, including titles, axis labels, and keys, and is asked to complete the graph.

Displaying items similarly in accommodated print and braille forms and in the online environment (and allowing students to interact with the items in a similar manner) maintains item integrity by assessing a similar construct in a similar manner regardless of how a student encounters an item. This provides students who are unable to access the assessment online with an assessment at the same level of rigor as the online test.

AP forms are thoroughly reviewed by DRC and LDOE content experts alongside the online form, and braille forms are reviewed by an outside third-party braille expert against the AP form. Throughout the braille creation process, the braille vendor relies on the AP form and consults with the content experts at LDOE for additional clarification or modifications for specific items as needed. Students' responses to the accommodated print or braille test are captured in the same online test as used by the general population, either through use of a scribe or by themselves if able. This ensures a valid and reliable assessment for students who are unable to participate in the online assessment. Louisiana's sample sizes are too small for traditional studies of comparability for both AP and braille forms.

Appendix H: On-Going Quality Control

A system for monitoring, maintaining, and increasing the quality of its assessment system, including precise and technically sound criteria for the analyses of all of the assessments in its assessment system, is crucial and critical for keeping a high quality of assessments. Table H.1 outlines where information about monitoring, maintaining, and improving quality can be found within this report.

Table H.1

On-Going Quality Control

Related Information		Related Chapter/Source	
Test Materials	Item development quality procedures	Content alignment Cognitive complexity Bias, fairness, and sensitivity Technical design	Chapter 3
	Form development quality procedures	Test specifications Review of statistical quality of items	Chapter 4
Test Administration	Test administration training and procedures	Training and monitoring of test administrators Security Checklists Test Security Measurements	Chapter 5
	Monitoring test administrations	LDOE site audits Data Forensics Analysis Response-Change Analysis Web Monitoring Plagiarism Detection	Chapter 5

Table H.1

On-Going Quality Control (continued)

Related Information		Related Chapter/Source	
Scoring	Scorer recruitment, training and security procedures	Recruitment and interview process Security Training process, including material development and qualifying procedures	Chapter 6
	Monitoring scoring quality	Inter-rater reliability studies Validity Reader monitoring	Chapter 6
Psychometric Processes	Psychometric quality procedures	Specifications document for operational analysis	Pearson and the LDOE Internal documentation
	Monitoring psychometric quality	Key verification Calibration Scoring table generation Psychometric quality checks on the data	Chapter 7
	Cuts based on Performance-Level Setting	Quality-controlled procedures for performance-level setting Derivation of the cut scores	Chapter 8