

**Maine Through Year Assessment
Spring 2025 Technical Report**

nwea

Table of Contents

Section 1: Overview of the Maine Through Year Assessment	1
1.1. Structure of the Maine Through Year Assessment: A Balanced Assessment System ..	1
1.2. Intended Purposes and Uses of Test Results	2
1.3. Required Assessment and Policies for Including All Students	2
1.4. Meaningful Consultation.....	2
1.4.1. Schedule of Major Events.....	2
Section 2: Test Design and Content Development	4
2.1. Test Design & Development.....	4
2.2. Achievement Level Descriptors	5
2.3. Test Blueprints	6
2.3.1. Cognitive Complexity Blueprint Considerations	7
2.3.2. Reading Summative Blueprint Considerations.....	7
2.3.3. Mathematics Summative Blueprint Considerations.....	9
2.3.4. Fixed Forms	11
2.4. Item Types.....	11
2.5. Content Development	13
2.5.1. Item Development and Guidelines.....	13
2.5.2. Universal Design	14
2.5.3. Sensitivity and Fairness	15
2.6. Content and Bias Review Meeting	16
2.7. Field Testing and Data Review	17
Section 3: Administration and Security.....	20
3.1. Administration	20
3.2. Spring 2025 Administration	20
3.2.1. Student Population.....	20
3.2.2. Item Pool Characteristics.....	23
3.3. Constraint-Based Engine Adaptive Test Administration	26
3.3.1. Engine Evaluation	27
3.3.2. Blueprint Constraint Accuracy	28
3.3.3. Item Exposure Rates.....	30
3.3.5. Item Sequence	32
3.4. Paper Form Administration	32
3.4.1. Receiving and Taking Inventory of School Materials	32
3.4.2. Score Transcription	33
3.5. Assessment Security.....	33
3.5.1. Assessment Ethics and Appropriate Practice	34
3.5.2. Online Security.....	34
3.5.3. Student Assessment Security.....	34
3.5.4. Returning or Destroying Secure Materials	34
3.6. Systems for Protecting Data Integrity and Privacy.....	34
Section 4: Item Statistics, Calibration, and Scaling.....	36
4.1. Classical Item Statistics	36
4.1.1. Expected P Value.....	36

4.1.2. Item Discrimination (Item-Total Correlation)	38
4.2. IRT Calibration	41
4.3. IRT Model Assumptions	42
4.3.1. Local Independence	42
4.3.2. Model Fit	42
4.3.3. Unidimensionality	45
4.4. Scaling	45
Section 5: Technical Quality—Validity	55
5.1. Intended Purposes and Validity Evidence Framework	55
5.2. Purposes and Evidence	56
5.2.1. Test Purpose 1	56
5.2.2. Test Purpose 2	58
5.2.3. Test Purpose 3	58
5.2.4. Test Purpose 4	59
5.3. Interpretive Argument Claims	59
5.4. Summary of Validity Arguments	60
Section 6: Technical Quality—Other	62
6.1. Reliability	62
6.1.1. Marginal Reliability for Adaptive Tests	62
6.1.2. Classification Accuracy	63
6.1.3. Score Precision	69
6.2. Fairness and Accessibility	72
6.2.1. Logistic Regression (LR) DIF Method	73
6.2.2. DIF Results	74
6.3. Full Achievement Continuum	79
6.4. Scoring	80
6.4.1. Constructing the Maine Scale	80
6.4.2. Machine-Scored Items	80
6.4.3. Attemptedness Rule and Not-Tested Codes	80
6.5. Multiple Assessment Forms	81
6.6. Multiple Versions of an Assessment	81
6.7. Technical Analysis and Ongoing Maintenance	81
Section 7: Inclusion of All Students	82
7.1. Testing Population	82
7.2. Procedures for Including Students Who Utilize Accessibility Features	82
7.3. Procedures for Including Multilingual Learners	82
7.4. Accommodations	82
7.5. Monitoring Test Administration for Special Populations	84
7.5.1. Monitoring in Acacia	84
7.5.2. Maine DOE Administration Monitoring and Support	85
Section 8: Achievement Standards and Reporting	90
8.1. State Adoption of Achievement Standards	90
8.2. Achievement Standard Setting	90
8.3. Reporting	93

8.3.1. Achievement Level Descriptors	94
8.3.2. Setting the Cut Scores	94
8.3.3. Reports	94
8.3.4. Relations to Other Scores	97
References	98

List of Tables

Table 1.1. Schedule of Major Events for the Spring 2025 Administration	3
Table 2.1. Summary of Assessments by Content Area & Grade	4
Table 2.2. Maine’s Policy Achievement Level Descriptors.....	5
Table 2.3. Item Counts per DOK Level.....	7
Table 2.4. Instructional Areas for Maine Reading Summative Blueprints.....	7
Table 2.5. Approximate Summative Blueprint Percentages by Instructional Area: Reading, Grades 3–8 & HS.....	8
Table 2.6. Approximate Summative Blueprint Percentages by Text Type: Reading, Grades 3–8 & HS.....	8
Table 2.7. Approximate Reading Lexile Ranges, Grades 3–8 & HS.....	9
Table 2.8. Approximate Reading Word Count Ranges, Grades 3–8 & HS	9
Table 2.9. Instructional Areas for Maine Mathematics Summative Blueprints.....	9
Table 2.10. Approximate Summative Blueprint Percentages: Mathematics, Grades 3–5	11
Table 2.11. Approximate Summative Blueprint Percentages: Mathematics, Grades 6–8 & HS ..	11
Table 2.12. Online Item Types	11
Table 2.13. Item Type Percentages by Grade—Reading Summative Pools, Spring 2025.....	12
Table 2.14. Item Type Percentages by Grade—Mathematics Summative Pools, Spring 2025 ..	12
Table 2.15. October 2024 Content and Bias Review Results	17
Table 2.16. Data Review Flagging Criteria	17
Table 2.17. Data Review Results	19
Table 3.1. Demographic Information—Reading.....	21
Table 3.2. Demographic Information—Mathematics.....	21
Table 3.3. Ability Distribution—Summative Scale Scores	22
Table 3.4. Ability Distribution—Summative Theta.....	22
Table 3.5. Numbers of Items by Content and Instructional Areas.....	23
Table 3.6. ALD Distribution Across Instructional Areas—Reading.....	23
Table 3.7. ALD Distribution Across Instructional Areas—Mathematics.....	24
Table 3.8. Median Item Response Time by ALD	25
Table 3.9. Blueprint Constraint Accuracy by Reporting Category	28
Table 3.10. Operational Item Exposure Rates.....	30
Table 3.11. Field Test Item Exposure Rates	31
Table 3.12. Paper Form Summative Item Totals by Content and Grade	32
Table 4.1. Summary of <i>P</i> Values—Operational Items	37
Table 4.2. Summary of <i>P</i> Values—Field Test Items	38
Table 4.3. Summary of Item-Total Correlations—Operational Items.....	39
Table 4.4. Summary of Item-Total Correlations—Field Test Items	40
Table 4.5. Summary of IRT Item Statistics—Operational Items.....	41
Table 4.6. Summary of Mean-Square Infit and Outfit Statistics	43
Table 4.7. Correlations Among Instructional Area Scores—Reading.....	44
Table 4.8. Correlations Among Instructional Area Scores—Mathematics.....	44
Table 4.9. Maine Grade-Level Scale Properties	46
Table 4.10. Summative Scale Score Frequency Table—Reading	46
Table 4.11. Summative Scale Score Frequency Table—Mathematics	50
Table 4.12. Field Test Calibration Results.....	54

Table 5.1. Intended Test Purposes and Sources of Validity Evidence.....	56
Table 5.2. Interpretive Argument Claims—Evidence to Support the Essential Validity Elements	60
Table 6.1. Reliability Statistics—Reading	63
Table 6.2. Reliability Statistics—Mathematics	63
Table 6.3. Classification Accuracy by Achievement Level—Reading	65
Table 6.4. Classification Accuracy by Achievement Level—Mathematics.....	66
Table 6.5. Classification Accuracy by Achievement Level and Cut.....	67
Table 6.6. Classification Consistency by Achievement Level and Cut.....	68
Table 6.7. CSEMs at the Cut Scores.....	69
Table 6.8. CSEMs by Summative Score Decile.....	70
Table 6.9. Summary of CSEMs by Instructional Area—Reading	70
Table 6.10. Summary of CSEMs by Instructional Area—Mathematics	71
Table 6.11. LR DIF Categories.....	74
Table 6.12. DIF Analysis Results—Operational Items.....	74
Table 6.13. DIF Analysis Results—Field Test Items.....	77
Table 6.14. Available Not-Tested Codes	81
Table 7.1. Numbers of Students Who Were Assigned TTS	84
Table 8.1. Final Approved Cut Scores—Reading	92
Table 8.2. Impact Data Associated with Cut Scores—Reading	93
Table 8.3. Final Approved Cut Scores—Mathematics	93
Table 8.4. Impact Data Associated with Cut Scores—Mathematics	93
Table 8.5. Report Levels	94
Table 8.6. Correlations Among Spring 2025 Summative Scores and Fall 2024 Interim Scores	97

List of Figures

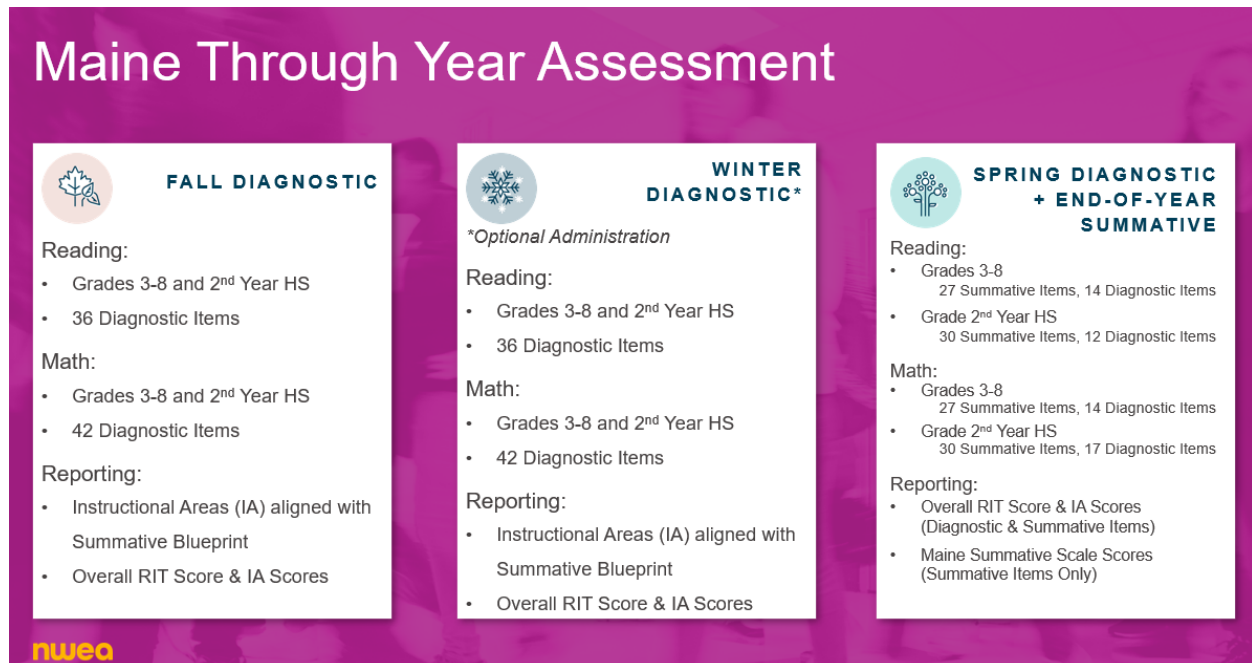
Figure 1.1. Structure of the Maine Through Year Assessment	1
Figure 2.1. Maine Blueprint Percentages—Mathematics, Grades 3–8 & 10 (HS).....	10
Figure 2.2. High School Reading Development Workflow	13
Figure 3.1. Adaptive Engine Overview	26
Figure 3.2. Student-Specific Plan Approach.....	27
Figure 7.1. Monitoring Testing Status in Acacia	85
Figure 7.2. 2025 Maine Educational Assessment Observation Form	86
Figure 8.1. Maine Through Year Assessment Embedded Standard Setting Iterative Processes	91
Figure 8.2. Individual Student Report.....	96

Section 1: Overview of the Maine Through Year Assessment

1.1. Structure of the Maine Through Year Assessment: A Balanced Assessment System

The Maine Through Year Assessment (MTYA) in reading and mathematics is a balanced assessment system that combines interim assessments administered multiple times throughout the academic year with an end-of-year summative assessment. The interim diagnostic assessments provide a measure of both student growth and achievement throughout the academic year. The spring administration is a combination of the end-of-year state summative assessment and the diagnostic assessment, with the summative assessment making up the majority of the administration. The state summative assessment reports student achievement according to grade-level state standards, specifically the Common Core State Standards (CCSS), and fulfills federal requirements for state assessments under the Every Student Succeeds Act (ESSA, 2015). In the spring, diagnostic items are used to help determine instructional area scores, or subscores.

Figure 1.1. Structure of the Maine Through Year Assessment



The MTYA assesses all publicly funded Maine students in grades 3 through 8 and the second year of high school (HS). The MTYA is an online adaptive test. For students with a need documented in an Individualized Education Plan (IEP) or 504 Plan, the test also offers three accommodated paper forms: paper/pencil standard print forms, large print forms, and braille forms.

This technical report documents the processes and procedures implemented to support the end-of-year summative portion of the MTYA. For the purposes of this technical report, only the spring administration of the assessment will be discussed. This technical report shows how the processes, methods applied, and results relate to the issues of validity and reliability and to the *Standards for Educational and Psychological Testing* (AERA et al., 2014). The complete technical report will be made available to the public by the Maine Department of Education at https://www.maine.gov/doe/Testing_Accountability/MECAS/NWEA no later than April 13, 2026.

The NWEA [MAP® Growth™ Technical Report](#) provides more information regarding elements of the diagnostic portions of the assessment, including item development and the computer adaptive test engine.

1.2. Intended Purposes and Uses of Test Results

The end-of-year summative portion of the MTYA has four primary purposes:

1. To report individual student achievement relative to the state-adopted content standards in reading and mathematics
2. To provide information to the public about school performance through the state's Every Student Succeeds Act (ESSA, 2015) reporting system, the ESSA Data Dashboard (<https://www.maine.gov/doe/dashboard>).
3. To support school identification within the state's ESSA compliant system of school identification and support
4. To provide a source of information for ongoing local program evaluation

The summative portion of the MTYA is designed to measure Maine's accountability standards, the Common Core State Standards (CCSS), in reading and mathematics. Student results are reported according to academic achievement level descriptors utilizing cut scores established through embedded standard setting for each of the four achievement levels: *Well Below State Expectations*, *Below State Expectations*, *At State Expectations*, *Above State Expectations*.

1.3. Required Assessment and Policies for Including All Students

All students in grades 3–8 and the second year of high school enrolled in Maine's public schools, Special Purpose Private Schools (SPPS), regional programs, charter schools, or private schools with at least 60% publicly funded students are required to participate in the MTYA. These students are eligible for and required to participate in Maine's state assessment program at state expense. Students with disabilities and multilingual learners may participate in the MTYA with accommodations.

Exceptions to participation in the MTYA would occur in cases involving students with the most significant cognitive disabilities who have been found eligible for alternate assessments via the IEP Team Process. Only about 1% of all publicly funded Maine students in grades eligible for assessment participate in an alternate assessment; the rest of the student population (approximately 99%) participate in the MTYA.

1.4. Meaningful Consultation

1.4.1. Schedule of Major Events

Table 1.1 presents the major events that occurred for the 2025 Maine Through Year Assessment.

Table 1.1. Schedule of Major Events for the Spring 2025 Administration

Event	Date(s)
High School Content and Bias Review	October 22–24, 2024
Data Review	October–November 2024
Technical Advisory Committee (TAC) Meeting	December 9 and 10, 2024
Test Administration Training Module ^a	Spring 2025
Operational Test Window	April 14–May 30, 2025

^a Test Administration Training slides are included in Appendix A.

This list provides more details about the events presented in the table.

- **Content and Bias Review:** a meeting with Maine educators to review all high school items authored for the program by NWEA
- **Data Review:** a review/analysis of field test items that were flagged for item performance. NWEA shares/discusses with Maine DOE the results of this review, and decisions are made regarding the next steps for the flagged items.
- **Technical Advisory Committee (TAC) Meeting:** a meeting with selected and designated assessment experts to review, discuss, and advise Maine’s assessment program. Additional TAC member information and meeting topics can be found in Appendix H.
- **Test Administration Training:** training to prepare District Assessment Coordinators, School Assessment Coordinators, and proctors. Topics covered include Maine Through Year Assessment Overview, Technology Readiness, Assessment Management in Acacia, Accessibility & Not-Tested Codes, Preparing & Monitoring the Assessment, Proctor & Student Experience, Operational Reports, Data & Reporting, Preparation, Resources, & Tips, and Communication & Support.
- **Operational Test Window:** the time period during which Maine students take the summative assessments

Below is a list of topics from the TAC meeting leading up to the Spring 2025 Through Year Assessment administration:

- December 9, 2024
 - Operational Leadup to Spring 2025
 - Overview of Maine Through Year Assessment
 - Simulation Report from Spring 2024
 - RIT Scale Adjustment (Spring 2024)
 - Maine DOE Narrative: How Did We Get Here?
 - Descriptions, Rationales, and Impacts of the Process
 - Technical Report (Spring 2024 administration)
- December 10, 2024
 - Peer Review Evidence (Spring 2023 administration)
 - December 2023 Submission
 - Peer Review Notes Received & Additional Evidence

Section 2: Test Design and Content Development

This section describes the test design and content development processes for the Spring 2025 Maine Through Year Assessment.

2.1. Test Design & Development

The Maine Through Year Assessment (MTYA) is designed to measure Maine’s accountability standards, the Common Core State Standards (CCSS), in reading and mathematics. In Spring 2025, Maine administered computer adaptive assessments in reading and mathematics for grades 3–8 and the second year of high school (HS). The reading and mathematics HS assessments were fixed forms in Spring 2023; the Spring 2024 assessments were adaptive tests with shallow item pools and included a focus on field testing, which continued in Spring 2025 as part of an effort to expand the operational pool and support increased adaptability in the future. The summative items in the grades 3–8 and HS assessments were licensed from NWEA; the grades 3–8 items came from existing NWEA item banks, and the HS items were part of a bank created by NWEA in collaboration with the Maine DOE and Maine educators. All items on the assessment were aligned to the CCSS and underwent a rigorous item-development process.

Table 2.1 summarizes the versions of the assessments. The computer adaptive test (CAT) design allows for the estimation of a student’s ability with greater precision using fewer items (see Section 3.3 for a description of the CAT constraint-based engine). Only summative items are used in the determination of a student’s summative Maine scale score. Diagnostic items appear on the spring assessment to help determine subscores. The paper forms were available in three formats (standard print, large print, or braille) and did not include field test items.

Table 2.1. Summary of Assessments by Content Area & Grade

Grade	Summative Items	Approximate Summative Points	Field Test Items	Diagnostic Items	Total Items	Paper Form?
Reading						
3	27	30–31	5	14	46	Yes
4	27	30–31	5	14	46	Yes
5	27	30–31	5	14	46	Yes
6	27	30–31	5	14	46	Yes
7	27	30–31	5	14	46	Yes
8	27	30–31	5	14	46	Yes
HS	30	33–34	7	12	49	Yes
Mathematics						
3	27	30–31	5	18	50	Yes
4	27	30–31	5	18	50	Yes
5	27	30–31	5	18	50	Yes
6	27	30–31	5	18	50	Yes
7	27	30–31	5	18	50	Yes
8	27	30–31	5	18	50	Yes
HS	30	33–34	5	17	52	Yes

2.2. Achievement Level Descriptors

An achievement level is a range of scores that defines a specific level of student achievement, as articulated in the Achievement Level Descriptors (ALDs). Maine’s policy ALDs were adopted in Spring 2023 to broadly define the characteristics of student performance on the state assessment at each of the four achievement levels: *Well Below*, *Below*, *At*, and *Above State Expectations*. Table 2.2 provides detailed explanations of each policy ALD.

Table 2.2. Maine’s Policy Achievement Level Descriptors

<i>Well Below State Expectations</i>	<i>Below State Expectations</i>	<i>At State Expectations</i>	<i>Above State Expectations</i>
On this assessment, students at this achievement level demonstrate limited understanding of the knowledge and skills necessary at this grade level, as specified in the Common Core State Standards. The students <i>need substantial academic support</i> to be prepared for the next grade level and to be on track for college and career readiness.	On this assessment, students at this achievement level demonstrate partial understanding of the knowledge and skills necessary at this grade level, as specified in the Common Core State Standards. The students <i>need additional academic support</i> to be prepared for the next grade level and to be on track for college and career readiness.	On this assessment, students at this achievement level demonstrate the knowledge and skills necessary at this grade level, as specified in the Common Core State Standards. The students <i>are prepared</i> for the next grade level and are on track for college and career readiness.	On this assessment, students at this achievement level demonstrate advanced understanding of the knowledge and skills necessary at this grade level, as specified in the Common Core State Standards. The students <i>are well prepared</i> for the next grade level and are well prepared for college and career readiness.

While the policy ALDs define broad, statewide expectations of student performance, the range ALDs translate those expectations into more detailed, grade- and standard-specific descriptions that guide item development, test design, and interpretation. The range ALDs show a progression of skills within a standard over multiple achievement levels. Range ALDs describe what a student should likely be able to do at a particular achievement level regarding on-grade content based on the broader policy ALDs. For each assessed standard, the ALDs show the range of on-grade content, highlighting the progression in the complexity and sophistication of skills across achievement levels. They do this by focusing on differentiating factors within each standard that represent the progression of student knowledge and understanding of the specified skill. The ALDs also strive to preserve differentiation between the skills as they progress across grades. The intent is that the ALDs, when viewed as a whole, provide a wide range of knowledge, skills, and abilities students can demonstrate over the course of the year while also considering the work from the previous grade and the upcoming work in the next grade.

Some content may appear in multiple places in the standards, but the ALDs are written to minimize overlap between grades. For example, CCSS mathematics standards 3.NBT.1 and

4.NBT.3 both assess rounding whole numbers. The ALDs for these standards use grade-level content limits to ensure that an item assessing rounding will only align to one grade. Range ALDs allow students at various levels to demonstrate their knowledge and skills. By helping to describe a student's current level of understanding, range ALDs support stakeholders in pinpointing areas of strength and areas of growth. Range ALDs are also used to guide NWEA content specialists in writing items for assessments and, by doing so, help develop a deeper item bank that can better serve the needs of each individual student.

NWEA content specialists wrote the initial draft of the Maine range ALDs and then held a workshop with Maine educators in September 2022 to review and revise the ALDs. Maine educators were asked to review these NWEA ALDs in relation to the Common Core State Standards used in Maine. Each participant reviewed range ALDs for grades 3–8 and HS in either reading or mathematics. The review's purpose was to allow Maine educators to study the ALDs and share their feedback with NWEA content specialists.

The number of committee members for each content area was limited, and for this reason, educators with expertise in all grade levels were recruited to participate. The four selected participants represented three different regions of the state, including Southern Maine, Southern-Central Maine, and Down East Maine, and one educator represented a virtual academy. All participants had experience working in schools with a high number of economically disadvantaged students, and some participants had experience working with special education students, English language learners, and gifted and talented students.

Both the reading and mathematics ALDs had progressions updated based on feedback from the Maine educators. These updates included reassigning ALD statements to another level within the progression, removing ALD statements, revising ALD statements, and crafting new ALD statements.

The complete set of range ALD statements utilized for the Maine Through Year Assessment is publicly available within the [Achievement Level Explorer tool](#).

2.3. Test Blueprints

All items on the end-of-year summative portion of the MTYA are aligned to a Common Core State Standard and to an achievement level descriptor specific to the standard. To ensure coverage of grade-level academic standards, the Maine Through Year summative test blueprints for reading and mathematics outline the overall structure of the assessments. For Maine, the summative blueprints are structured around instructional areas for reporting based on content categories within the CCSS.

The assessment blueprints for each grade level in reading and mathematics in Appendix G list:

- the grade-level content standards included within each instructional area,
- the item-count targets for each instructional area,
- the approximate points targets and approximate percentage-of-overall-points targets on the assessment for each instructional area,
- the percentage of on-grade items, and
- the achievement level percentage targets.

Appendix B provides more detailed information about standard coverage at each grade level within the blueprints based on empirical data from the Spring 2025 administration.

2.3.1. Cognitive Complexity Blueprint Considerations

Cognitive complexity is considered part of the guidelines for constructing the test blueprints. Range ALDs for each standard clarify the level of mastery and cognitive processes expected of students performing at each achievement level, with a focus on increased complexity and sophistication of skills across the achievement levels. For the MTYA, each item in the reading and mathematics pools was either written for or aligned to a specific content standard and an associated ALD. Complete range ALDs for each MTYA subject and grade can be accessed via the [Achievement Level Explorer tool](#) and then selecting “State of Maine” as the assessment partner.

The MTYA blueprints include targets specifying boundaries for the percentages of items that should be selected from the different achievement levels, with at least 60% or more of the items coming from *At State Expectations* and *Above State Expectations* levels. Because the MTYAs in reading and mathematics are adaptive, the exact distribution of Achievement Level Descriptors and ALD levels for any given test event will vary based on individual student achievement and other blueprint constraints.

To ensure that the assessments include a deep pool of items that span a full range of skills and cognitive complexities, both the standards and the ALD distributions are important factors when determining item pool needs and item-development plans.

In addition, Table 2.3 shows that all items are aligned to a Depth of Knowledge (DOK) level based on Webb’s 2009 framework (see Appendix J), further supporting the diversity of the item pool.

Table 2.3. Item Counts per DOK Level

Subject	Depth of Knowledge		
	DOK Level 1	DOK Level 2	DOK Level 3
Reading	287	1,713	393
Mathematics	1,088	1,416	50

2.3.2. Reading Summative Blueprint Considerations

In reading, the instructional areas shown in Table 2.4 are closely connected to the CCSS Reading Strands.

Table 2.4. Instructional Areas for Maine Reading Summative Blueprints

Instructional Areas for Grades 3–8 & HS
Literary Text
Informational Text
Vocabulary

When creating the reading blueprints for Maine, focus was given to the weight and breadth of the reading standards designed to assess literary and informational texts and vocabulary skills. The blueprints were influenced by the Priority Instructional Content guidance from Student

Achievement Partners ([Achieve the Core](#)) and reflect the belief that not all content standards are emphasized equally in the classroom. The reading assessments are designed to keep the text at the center and use text-based questions. These assessment items highlight close reading skills, text analysis, textual evidence, and academic vocabulary.

Table 2.5 shows the approximate percentages for the instructional areas for each grade. Additional information about how these percentages are represented in the assessments can be found in Appendix G.

Table 2.5. Approximate Summative Blueprint Percentages by Instructional Area: Reading, Grades 3–8 & HS

Instructional Area	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	HS
Literary Text	45–50%	40–45%	35–40%	35–40%	30–35%	30–35%	30–35%
Informational Text	30–35%	35–40%	35–40%	40–45%	45–50%	45–50%	45–50%
Vocabulary	20–25%	20–25%	20–25%	20–25%	20–25%	20–25%	20–25%

It is important that reading assessments provide a balance of content between literary and informational texts and represent a range of text complexity. According to the Common Core State Standards, students are expected to demonstrate an understanding of increasingly complex texts as a result of grade-level and discipline-specific content expectations.

Reading text content is classified as either literary or informational. The balance of percentages shifts from more literary content to more informational content as the grade level increases. These percentages originated for grade bands with the Common Core State Standards and have been extrapolated to be grade-specific for Maine. Table 2.6 shows the percentages of literary and informational text by grade.

Table 2.6. Approximate Summative Blueprint Percentages by Text Type: Reading, Grades 3–8 & HS

Text Type	Grade 3	Grade 4	Grade 5	Grade 6	Grade 7	Grade 8	HS
Literary	55–60%	55–60%	50%	45–50%	40–45%	40–45%	40–45%
Informational	40–45%	40–45%	50%	50–55%	55–60%	55–60%	55–60%

Text complexity is a measure of how challenging a text is to read and understand. Many factors may make a text complex, so a text complexity measurement is the process of evaluating a text for quantitative data, qualitative data, and the considerations for the reader and task. The texts in the reading item bank should include those that cover a range of text complexity within a grade level, including minimally complex, moderately complex, and highly complex.

Quantitative data includes concrete measures such as word length or frequency, sentence length, text cohesion, and vocabulary. These are communicated through readability measures including Lexile, word count, and Flesch-Kincaid. Quantitative measures are only a guide; exceptions can be made if the qualitative measures and/or grade-level alignments are appropriate. Table 2.7 shows acceptable Lexile ranges for each grade.

Table 2.7. Approximate Reading Lexile Ranges, Grades 3–8 & HS

Grade(s)	Lexile Range
3	450L–790L
4–5	745L–980L
6–8	925L–1155L
HS	960L–1305L

Note. These Lexile bands reflect the adaptive nature of the assessments and the need to include a slightly larger range of readabilities than outlined in the [CCSS](#).

Table 2.8 provides acceptable word count ranges for each grade. For paired passages, each individual passage should fall within the word count range.

Table 2.8. Approximate Reading Word Count Ranges, Grades 3–8 & HS

Grade	Word Count Range
3	200–700
4	200–900
5	300–1000
6	400–1100
7	400–1100
8	400–1200
HS	600–1400

Qualitative data includes the following dimensions: meaning/purpose, structure, language, and knowledge demands. Additionally, considerations regarding the reader and their interaction with a passage and the items they will answer for each passage help acknowledge students' role in the assessment. NWEA conducts a review using a Passage Quality Checklist (included in Appendix J) that documents the complexity and suitability of each passage for assessment. For more information about text complexity, see <https://achievethecore.org/page/2725/text-complexity>.

2.3.3. Mathematics Summative Blueprint Considerations

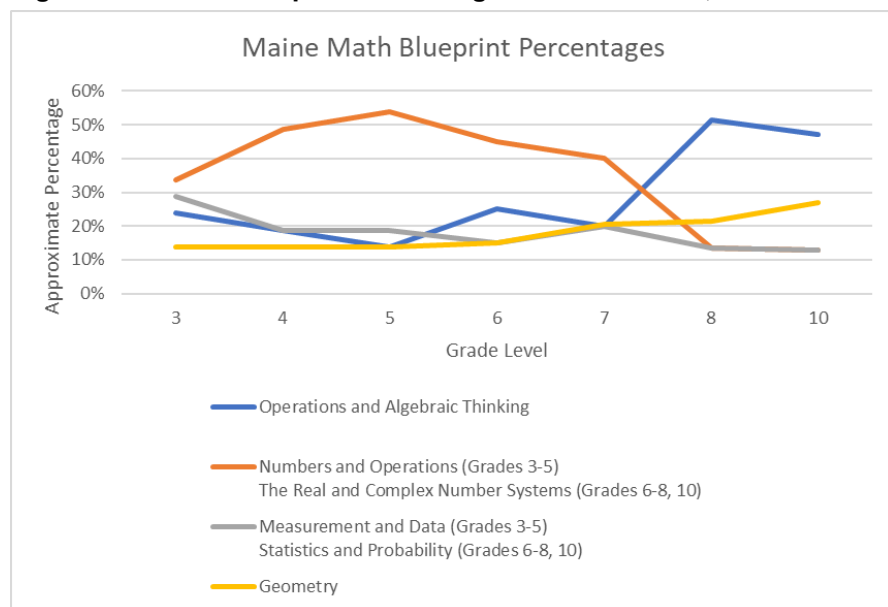
For mathematics, the instructional areas are closely connected to the CCSS mathematics domains, as shown in Table 2.9.

Table 2.9. Instructional Areas for Maine Mathematics Summative Blueprints

Instructional Areas for Grades 3 to 5
Operations and Algebraic Thinking
Numbers and Operations
Measurement and Data
Geometry
Instructional Areas for Grades 6 to 8 and HS
Operations and Algebraic Thinking
The Real and Complex Number Systems
Geometry
Statistics and Probability

The mathematics blueprints reflect the instructional emphasis of the content at each grade. For example, Geometry receives more instructional time as the grades progress, which is reflected in how the percentage increases from 13% in grade 3 to 30% in grade 10 (HS). The blueprints were also influenced by Student Achievement Partners' Focus Areas in Mathematics ([Achieve the Core](#)), which calls out the major work of each grade and guides educators on how to focus instructional time. Figure 2.1 shows how the percentage for each reporting category shifts across the grades.

Figure 2.1. Maine Blueprint Percentages—Mathematics, Grades 3–8 & 10 (HS)



The graph shows that students' skills in the Numbers and Operations instructional area are emphasized as they work with whole numbers less than 1,000 and fractions with a limited set of denominators in grade 3 before moving on to decimals and a larger set of fractions in grade 5. After students grasp these skills, the significance of the instructional area (which shifts to The Real and Complex Number Systems starting in grade 6) gradually lessens as students work with the set of rational numbers in grade 6 before moving on to the set of irrational numbers in high school.

Conversely, the emphasis on students' skills in the Operations and Algebraic Thinking instructional area steadily increases as students solve simple two-step problems in context in grade 3 and move to working with linear and quadratic functions in high school.

For the content category Measurement and Data in grades 3–5 or Statistics and Probability in grades 6–8 and HS, the emphasis remains relatively constant and ranges from 13% to 30%. Students' skills gradually progress from working with picture graphs in grade 3 to scatter plots in high school.

For the Geometry instructional area, the emphasis gradually increases from around 14% in grade 3 as students work with area and perimeter to around 28% in high school as students work with more complex figures and geometric proofs.

In grade 7, three content categories are each assessed at approximately 20% in the blueprint because it is the point at which the Operations and Algebraic Thinking instructional area and the Geometry instructional area continue to increase while the third instructional area (Measurement and Data in grades 3–5 or Statistics and Probability in grades 6–8 and HS) remains relatively constant near 20%.

Table 2.10 and Table 2.11 show the approximate percentages for the instructional areas for each grade. Additional information about how these percentages are represented in the assessments can be found in Appendix G.

Table 2.10. Approximate Summative Blueprint Percentages: Mathematics, Grades 3–5

Instructional Area	Grade 3	Grade 4	Grade 5
Operations and Algebraic Thinking	23–25%	18–20%	13–15%
Numbers and Operations	33–35%	48–50%	53–55%
Measurement and Data	28–30%	18–20%	18–20%
Geometry	13–15%	13–15%	13–15%

Table 2.11. Approximate Summative Blueprint Percentages: Mathematics, Grades 6–8 & HS

Instructional Area	Grade 6	Grade 7	Grade 8	HS
Operations and Algebraic Thinking	25%	20%	48–53%	46–50%
The Real and Complex Number Systems	45%	40%	13–15%	13–15%
Geometry	15%	20%	21–23%	26–30%
Statistics and Probability	15%	20%	13–15%	13–15%

2.3.4. Fixed Forms

Content specialists select items for the fixed forms from the MTYA item bank to be consistent with the approximate summative blueprint percentages for the online adaptive forms (as described in Section 2.3) alongside other considerations such as choosing only items that are appropriate for the print formats offered in Maine. Test forms are also reviewed through a psychometrics lens to ensure difficulty targets are adhered to and that the blueprints are sufficiently represented, allowing for discrepancies that may exist due to rounding and the nature of approximations for the instructional area ranges. The fixed form assessment blueprints for each grade level and content area can be found in Appendix G following the computer adaptive assessment blueprints.

2.4. Item Types

The Maine Through Year Assessment consists of several item types, as outlined in Table 2.12.

Table 2.12. Online Item Types

Item Type	Description
Multiple-Choice (Choice)	Students select one response from multiple options.
Multi-Select (Choice Multiple)	Students select two or more responses from multiple options. Some multi-select items are also two-point items for which students can earn partial credit.

Item Type	Description
Composite	Students interact with multiple interaction types included within a single item. Students may receive partial credit for composite items.
Gap Match	A type of drag-and-drop item in which students select one or more answer options from the item toolbox and populate a defined area, or “gap.” Some gap match items are also two-point items for which students can earn partial credit.
Graphic Gap Match	A type of drag-and-drop item in which students move one or more answer options from the toolbox and populate a defined area, or “gap,” that has been embedded within an image in the item response area. Some graphic gap match items are also two-point items for which students can earn partial credit.
Hot Text	Students select a response from within a piece of text or a table of information (e.g., word, section of a passage, number, symbol, or equation) that is highlighted in the selected text. Some hot text items are also two-point items for which students can earn partial credit.
Text Entry	Students input numeric answers using a keyboard.

Table 2.13 and Table 2.14 outline the percentages of item types by content area and grade level in the available Spring 2025 summative pools.

Table 2.13. Item Type Percentages by Grade—Reading Summative Pools, Spring 2025

Grade	Item Type				
	Multiple-Choice	Multi-Select	Composite	Gap Match	Hot Text
3	79%	8%	6%	6%	1%
4	82%	8%	5%	5%	0%
5	84%	9%	2%	4%	2%
6	83%	6%	4%	4%	2%
7	76%	11%	6%	6%	1%
8	83%	6%	4%	5%	2%
HS	68%	12%	13%	3%	3%

Note. Due to rounding of individual percentages, some totals may not equal 100% exactly.

Table 2.14. Item Type Percentages by Grade—Mathematics Summative Pools, Spring 2025

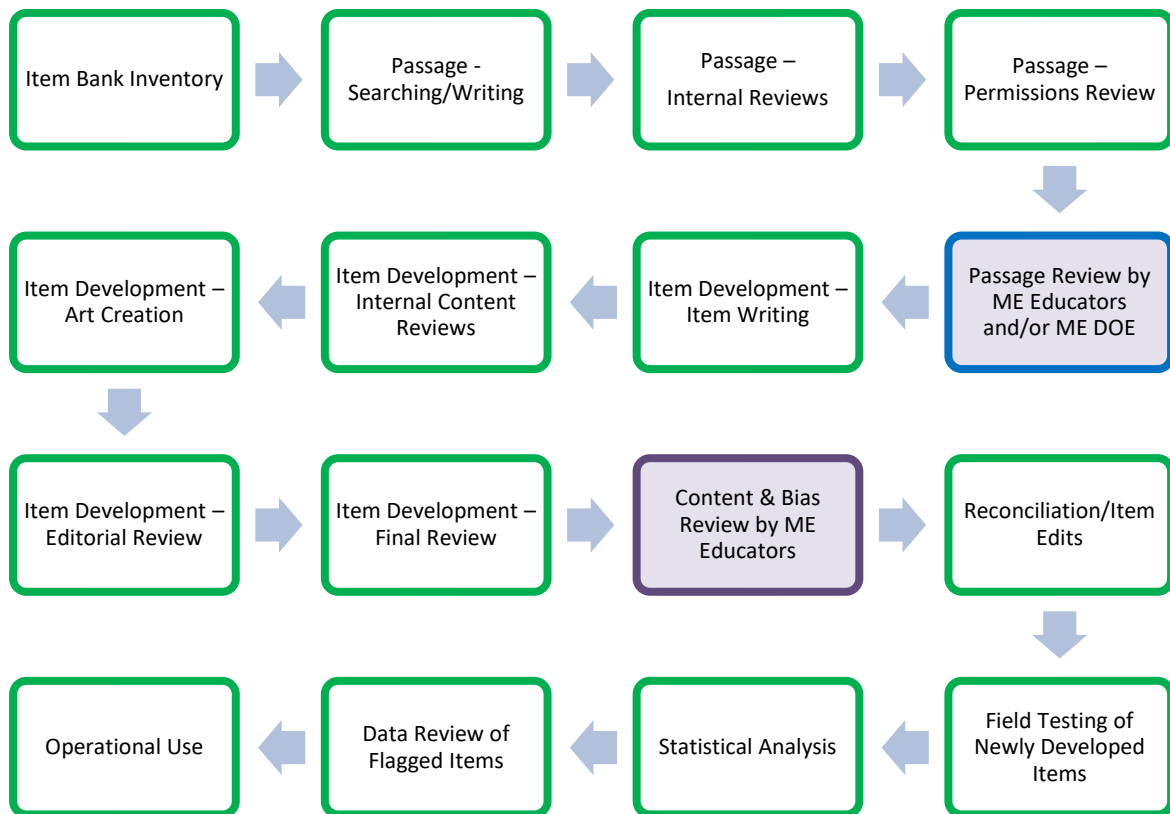
Grade	Item Type						
	Multiple-Choice	Multi-Select	Composite	Gap Match	Graphic Gap Match	Hot Text	Text Entry
3	53%	10%	6%	10%	6%	3%	12%
4	52%	11%	4%	8%	9%	7%	9%
5	53%	8%	7%	7%	8%	4%	12%
6	55%	7%	6%	9%	3%	6%	14%
7	61%	9%	3%	7%	2%	7%	11%
8	51%	8%	5%	7%	4%	11%	15%
HS	25%	21%	14%	18%	1%	17%	4%

Note. Due to rounding of individual percentages, some totals may not equal 100% exactly.

2.5. Content Development

New content development for Maine was focused on reading and mathematics HS items; however, the general content development process also applies to the grades 3–8 items that were selected for use in Maine from existing NWEA item banks. Items are developed in accordance with Universal Design for Learning principles and are each aligned to a standard and an ALD. Items are rigorously reviewed and edited during internal reviews, including reviews for bias/sensitivity at each stage of development. For the newly developed HS items, feedback from an external content and bias/sensitivity review involving Maine educators was incorporated into the items prior to field testing. Items that pass the review stages are field tested and subsequently reviewed based on their item statistics, which may include re-examining the item content. Items that pass the data review become operational in the item pool. Figure 2.2 provides an example of an item development workflow for high school reading; the high school mathematics workflow is similar but excludes the passage steps. At the conclusion of a test window, the process begins again, using the most recent pool and test simulation data to determine the areas of focus for future item development, with a particular focus on standard and ALD coverage.

Figure 2.2. High School Reading Development Workflow



2.5.1. Item Development and Guidelines

Item development begins with a review and inventory of the existing item bank while also determining areas of focus based on information derived from simulations. From this, an item development plan can be created that may include item-writing specifications and targets such as:

- Specific standards to be targeted to strengthen and add more depth to the pool
- Specific cognitive complexity targets (in the form of ALDs) to strengthen and add more depth to the pool
- General item-writing guidelines in terms of overall content, item stems, item responses, style, and scoring rules
- Guidelines for using technology-enhanced items (TEIs)

The reading and mathematics HS items are written internally by NWEA content specialists or external professional item writers. Grades 3–8 include items written by NWEA content specialists, external professional item writers, and educators trained at item writing workshops. Regardless of the source, all items undergo a rigorous review and editing process, including bias/sensitivity reviews at multiple stages of development, and are developed in accordance with Universal Design for Learning principles. Following best practices, including style, ensures that items are accurately measuring student knowledge at each level by focusing the items on construct-relevant information and presentation and ensuring that items are accessible and fair to all students. The subsequent field-testing process, statistical analysis of item performance, and data review help further ensure the quality of items that become operational. All summative operational items in the MTYA, regardless of their original source, go through Maine-specific field testing to ensure they are effective items for Maine students.

For reading items, this process also includes writing or identifying passages and the identification or development of passage resources. Passage resources may include sources from the public domain, copyright works that are permissioned for use, and commissioned works. For HS item development, all passages are commissioned or taken from the public domain. Passages, like items, undergo a rigorous review process, including bias/sensitivity reviews.

Passages are developed or selected to:

- offer appropriate content, length (emphasis on word counts), and text complexity
- provide engaging reading opportunities for students as they take the test
- include ample variation to appeal to a wide range of student audiences
- contain the characteristics required for the development of items that target a range of standards and ALDs

Items developed for the Maine assessments are tightly aligned to either a part of or an entire standard as well as to a corresponding ALD for the standard, providing additional information around the cognitive demand and rigor of each item. Additional data points may also be tracked, such as Webb's depth of knowledge level, to support overall alignment decisions. Examples of the item and passage review checklists used during the content development process can be found in Appendix J.

2.5.2. Universal Design

Ensuring that assessments are accessible to students with a variety of needs, including those with disabilities, is a critical part of item development. With a strong foundation in Universal Design for Learning (UDL; Rose & Meyer, 2006), the assessments become engaging and accessible for all students. The NWEA content team ensures that each item is created with the principles of UDL in mind. These principles provide a framework for developing flexible items to

support many kinds of learners and maximize options for assessments that provide multiple means of representation, action and expression, and engagement. Considerations NWEA takes into account when developing items include:

- Items are free of unnecessary linguistic complexity.
- Information presented in items is clear, concise, and relevant to the standard being assessed.
- Context and language are fair and familiar to students at their grade level and do not give advantages or disadvantages to subgroups.
- Items are free of stereotypes and potential disrespect regarding age, gender, race, ethnicity, language, religion, sexual orientation, social economic status, disability, or geographic region.
- Items do not challenge personal beliefs or values and avoid emotionally charged topics.
- Names and gender are avoided unless necessary. If names must be used, a variety of genders and ethnicities are represented.
- Graphics are intentional and not merely decorative.
- Graphics are not color dependent.
- MathML uses equation tags compatible with text-to-speech and screen readers.
- Art is tagged to be compatible with screen readers where possible.

Applying UDL principles to assessments helps minimize the amount of background knowledge needed to correctly respond to an item and ensures that items do not contain sources of construct-irrelevant variance so that assessments can more appropriately capture what each student knows. Including UDL principles in item development also helps ensure that there will be available items for the creation of accommodated forms such as large print and braille.

2.5.3. *Sensitivity and Fairness*

NWEA takes seriously the task of creating items that are free from bias and sensitivity issues and are fair to all students. Items are revised to eliminate bias, sensitivity, and fairness issues—or rejected if an issue cannot be remedied through the revision process. Items are reviewed for sensitivity and fairness multiple times throughout the development process, with some items also being reviewed in a collaborative effort with Maine educators (see Section 2.6).

- **Bias:** This is defined as item content, unrelated to the concept or skill being assessed, that may unfairly influence a student's performance or an item construct that does not have equivalent meaning for all students.
- **Sensitivity:** This can result if the experience of taking a test differs from the classroom experience in that students do not have the opportunity to discuss the material with a teacher or their peers. Sensitive content risks drawing students out of the testing experience by provoking negative emotional responses.
- **Fairness:** This is defined as the equitable treatment of all students during the assessment process. To make a test fair, test developers must work to eliminate any barriers that prevent students from understanding and interacting with item content in a manner that accurately demonstrates what they know or are able to do.

A successful item is free of bias and sensitivity issues and is accessible to all students. An item should NOT:

- distract, upset, or confuse in any way
- contain inappropriate or offensive topics
- require construct-irrelevant knowledge or specialized knowledge
- favor students from certain language communities
- favor students from certain cultural backgrounds
- favor students based on gender
- favor students based on social economic issues
- employ idiomatic or regional phrases and expressions
- stereotype certain groups of people or behaviors
- favor students from certain geographic regions
- favor students who have no visual impairments
- use height, weight, test scores, or homework scores as content or data in an item

There is no hard and fast “list” of material that is potentially distracting or upsetting, but some topics are seldom appropriate for K–12 assessments, such as sexuality, illegal substances, illegal activities, excessive violence, discriminatory descriptions, death, grieving, catastrophes, animal neglect or abuse, and loss of a family member.

2.6. Content and Bias Review Meeting

The purpose of the Content and Bias Review (CBR) meeting is to have Maine educators evaluate new test items developed for the field test item bank. Educators review content, alignment to standards, and each key to gain actionable feedback on all items. Only the high school items went through a CBR meeting in October 2024, as this was the only grade for which items were developed specifically for Maine. Grades 3–8 items, while not part of the Maine CBR process, would have undergone similar reviews either internally or externally during their initial development.

Educators are asked to review the items in advance of the virtual CBR and decide if they feel the items should be accepted, accepted with revisions, or rejected. Training slides can be found in Appendix I. The CBR meeting begins with a general session in which participants are given an overview of the purpose of the meeting and the process to be followed. Following the general session, participants report to either the reading or mathematics breakout room, where reminders about the criteria by which items should be reviewed are provided.

Each breakout room includes an NWEA facilitator who leads a discussion regarding any items that have been flagged as accept with revisions or rejected by any of the reviewers with the goal of coming to a final decision for the item. If needed, requests are reconciled with Maine DOE in the days following the meeting and then revisions can be applied.

Educators review items and provide comments based on the following criteria that is provided on the checklists:

- Item aligns to the standards.
- Item is clearly worded.
- Item type is appropriate for the content/standard.
- Item has one and only one best correct answer.
- Item distractors are plausible.

- Item art is clear and necessary.
- Item is mathematically correct.
- Item is factually correct.

PDF copies of the Achievement Level Descriptors and the item review criteria checklist are available for the educators to use during their review.

Table 2.15 outlines the total numbers of items taken to the Content and Bias Review meeting, as well as the numbers of items accepted, accepted with revisions, and rejected.

Table 2.15. October 2024 Content and Bias Review Results

Content Area	Total Items Reviewed	Accepted	Accepted with Revisions	Rejected
Reading	168	128	40	0
Mathematics	200	161	38	1

2.7. Field Testing and Data Review

Data review is the process of reviewing field-tested items for quality and appropriateness based on the results of statistical analysis of student responses. In performing this data review, NWEA adheres to the *Standards for Educational and Psychological Testing* (AERA et al., 2014) by implementing quality control procedures to ensure accurate information about student learning. The requirements regarding test administration, scoring, and reporting are as follows:

- Standard 4.8: The test review process should include empirical analyses and/or the use of expert judges to review items and scoring criteria.
- Standard 6.0: Adherence to the established procedures should be monitored, and any material errors should be documented and, if possible, corrected.
- Standard 6.9: 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 (AERA et al., 2014).

A data review took place in August 2025. Field test items were flagged based on statistical criteria. NWEA assessment specialists then conducted a close examination of the items based on the flags. As a result, some items were removed from the pool, some were deemed appropriate to remain in the pool and changed to an operational status, and some were revised and will be re-field tested in Spring 2026. Table 2.16 presents the criteria for flagging items that were field tested in Spring 2025, and Table 2.17 provides the results of the data review process.

Table 2.16. Data Review Flagging Criteria

Flag Name	Flag Text	Meaning	Implication for Data Review
Pvalue_LOW	P value less than 0.2	Less than 20% of students got this item correct. This item seems very difficult.	Does it seem reasonable that this item appears to be very difficult?
Pvalue_Dis	Distractor p value higher than key	Distractor percentages $> p$ value, meaning more	Is the answer key correct?

Flag Name	Flag Text	Meaning	Implication for Data Review
		students chose a distractor than the key	
Pbis_LOW	Item-total correlation less than 0.2	Item-total correlation < 0.20, meaning the item does not differentiate between high- and low-performing students	Is the answer key correct?
Pbis_Dis	Positive distractor correlation or higher than key	Item-total correlation for distractors > 0.05, meaning the item does not differentiate between high- and low- performing students. Some high-performing students chose a distractor over the key.	Is the answer key correct? Is there a reason why high- performing students select a distractor as an answer?
Score_0_Pbis	Positive Score 0 correlation	Item-total correlation for score of 0 > 0.0, meaning score of 0 on the item does not differentiate achievement levels as expected.	Is there a reason earning 0 points is happening more often for high-performing students than low-performing students?
Score_2_Pbis	Score 2 correlation less than 0.2	Item-total correlation for score of 2 < 0.2, meaning the score of 2 on the item does not differentiate achievement levels as expected.	Is there a reason earning 2 points happens more often for low-performing students than high-performing students?
Scores_0Vs1_Pbis	Score 0 correlation higher than Score 1	Item-total correlation for score of 0 > item-total correlation for score of 1, meaning a score of 0 on the item is better at differentiating achievement levels than a score of 1.	Is there anything that could cause the item to perform the opposite of what is expected for high- vs. low-performing students who received a score of 0 vs. 1?
Scores_1Vs2_Pbis	Score 1 correlation higher than Score 2	Item-total correlation for score of 1 > item-total correlation for score of 2, meaning score of 1 on the item is better at differentiating achievement levels than a score of 2.	Is there anything that could cause the item to perform the opposite of what is expected for high- vs. low-performing students who received a score of 1 vs. 2?
DIF of subgroups	A+ or A- through C+ or C-	Item is flagged for potential bias toward a certain group of students.	Is there anything that could trigger bias toward certain groups of students? Anything at a C+ or C- should likely be retired.

Table 2.17. Data Review Results

Grade	Accepted	Revise and Re-Field Test	Reject	Total
Total	409	20	64	493
Reading				
3	15			15
4	18		1	19
5	16		2	18
6	17		1	18
7	14		1	15
8	15		3	18
HS	115	6	32	153
Total Reading	210	6	40	256
Mathematics				
3	14	1		15
4	12	3		15
5	13	2		15
6	11	3		14
7	12	3		15
8	12	1		13
HS	125	1	24	150
Total Mathematics	199	14	24	237

Please note that item calibrations will estimate the item parameters of the field test items labeled “Accepted” on both the Maine summative and interim scales. This may result in some accepted items being dropped due to flagging.

Section 3: Administration and Security

3.1. Administration

District and School Assessment Coordinators are primarily responsible for ensuring a uniform assessment administration, including scheduling logistics, training and supervision of proctors, and maintaining assessment security. *The Maine Through Year Assessment Coordinator Guide* provides clear guidance on preparing for, monitoring, and concluding the administration of the Maine Through Year Assessment. *The Maine Through Year Assessment Administration Guide* contains explicit directions and proctor scripts for consistency of administration across different schools and School Administrative Units (SAUs).

3.2. Spring 2025 Administration

This section provides an overview of the observed demographics of participating students, their estimated ability distributions, and descriptions of the item pool.

3.2.1. Student Population

Table 3.1–Table 3.4 display demographic information and ability distributions for Maine’s general student population.

Table 3.1. Demographic Information—Reading

Grade	Type	Total	Gender		Ethnicity						
			Female	Male	Hispanic/ Latino	Am. Indian or Alaska Native	Asian	Black or African American	Native HI or Pacific Islander	White	Two or More Races
3	N	12,515	6,074	6,439	506	85	151	635	11	10,607	520
	%	100	48.53	51.45	4.04	0.68	1.21	5.07	0.09	84.75	4.16
4	N	11,822	5,724	6,097	425	89	145	637	12	10,040	474
	%	100	48.42	51.57	3.59	0.75	1.23	5.39	0.10	84.93	4.01
5	N	12,378	6,026	6,348	443	87	161	635	17	10,537	498
	%	100	48.68	51.28	3.58	0.70	1.30	5.13	0.14	85.13	4.02
6	N	12,235	5,886	6,346	438	101	176	633	14	10,413	460
	%	100	48.11	51.87	3.58	0.83	1.44	5.17	0.11	85.11	3.76
7	N	11,978	5,805	6,173	416	99	166	583	11	10,224	479
	%	100	48.46	51.54	3.47	0.83	1.39	4.87	0.09	85.36	4.00
8	N	12,279	6,014	6,260	418	87	170	654	8	10,485	457
	%	100	48.98	50.98	3.40	0.71	1.38	5.33	0.07	85.39	3.72
HS	N	12,344	5,918	6,405	451	114	252	616	17	10,476	418
	%	100	47.94	51.89	3.65	0.92	2.04	4.99	0.14	84.87	3.39

Table 3.2. Demographic Information—Mathematics

Grade	Type	Total	Gender		Ethnicity						
			Female	Male	Hispanic/ Latino	Am. Indian or Alaska Native	Asian	Black or African American	Native HI or Pacific Islander	White	Two or More Races
3	N	12,531	6,083	6,444	510	85	153	657	11	10,597	518
	%	100	48.54	51.42	4.07	0.68	1.22	5.24	0.09	84.57	4.13
4	N	11,846	5,745	6,100	429	89	152	657	12	10,035	472
	%	100	48.50	51.49	3.62	0.75	1.28	5.55	0.10	84.71	3.98
5	N	12,429	6,058	6,367	456	87	164	659	17	10,548	498
	%	100	48.74	51.23	3.67	0.70	1.32	5.30	0.14	84.87	4.01
6	N	12,263	5,897	6,363	447	102	184	640	14	10,415	461

Grade	Type	Total	Gender		Ethnicity						
			Female	Male	Hispanic/ Latino	Am. Indian or Alaska Native	Asian	Black or African American	Native HI or Pacific Islander	White	Two or More Races
	%	100	48.09	51.89	3.65	0.83	1.50	5.22	0.11	84.93	3.76
7	N	12,019	5,829	6,190	426	99	174	602	13	10,228	477
	%	100	48.50	51.50	3.54	0.82	1.45	5.01	0.11	85.10	3.97
8	N	12,272	6,025	6,242	427	89	173	662	8	10,460	453
	%	100	49.10	50.86	3.48	0.73	1.41	5.39	0.07	85.23	3.69
HS	N	12,358	5,929	6,408	460	113	250	633	17	10,468	417
	%	100	47.98	51.85	3.72	0.91	2.02	5.12	0.14	84.71	3.37

Table 3.3. Ability Distribution—Summative Scale Scores

Grade	Summative Scale Score			
	Reading		Mathematics	
	Mean	SD	Mean	SD
3	1,504	17	1,505	20
4	1,505	17	1,503	18
5	1,506	17	1,501	18
6	1,507	16	1,497	19
7	1,506	17	1,497	19
8	1,503	18	1,497	17
HS	1,503	15	1,501	20

Table 3.4. Ability Distribution—Summative Theta

Grade	Summative Theta			
	Reading		Mathematics	
	Mean	SD	Mean	SD
3	-0.24	1.41	-0.53	1.57
4	0.11	1.36	-0.42	1.76
5	0.06	1.32	-0.06	1.70
6	0.13	1.21	-0.35	1.68
7	0.27	1.31	-0.82	1.72
8	0.23	1.40	-0.68	1.54
HS	0.07	0.94	-1.22	1.15

3.2.2. Item Pool Characteristics

To ensure the adequacy of the item pool for administering a computer adaptive test (CAT), Table 3.5 details the numbers of items of various types and levels in the item pool for Maine by instructional area in the summative item pools for reading and mathematics.

Table 3.5. Numbers of Items by Content and Instructional Areas

Content Area	Instructional Area	Grade						
		3	4	5	6	7	8	HS
Reading	Informational Text	470	402	442	516	525	503	140
	Literary Text	400	460	436	430	438	314	118
	Vocabulary	261	416	380	303	325	352	36
	Total	1,131	1,278	1,258	1,249	1,288	1,169	294
Mathematics	Geometry	57	154	152	161	250	370	138
	Measurement and Data	346	21	30	21	6	–	–
	Numbers and Operations	290	54	63	63	48	–	–
	Operations and Algebraic Thinking	367	214	176	438	278	525	176
	Statistics and Probability	137	235	94	135	350	178	44
	The Real and Complex Number Systems	70	704	844	505	356	58	32
	Total	1,267	1,382	1,359	1,323	1,288	1,131	390

Beyond the instructional areas, the lower standard levels were also examined by assessing the number of items available at each standard. The percentages of students who received at least one item from each standard are shown in Appendix B.

Table 3.6 and Table 3.7 represent the ALD distribution across instructional areas (IAs) by grade in reading and mathematics, respectively. Table 3.8 represents the median item response time by ALD, and Table 3.9 shows on- and off-grade item counts for the summative items. The data indicate that the engine worked as intended across subjects and grades.

Table 3.6. ALD Distribution Across Instructional Areas—Reading

Grade	ALD	In the Pool			Administered		
		IA1	IA2	IA3	IA1	IA2	IA3
3	1	31	31	15	28	30	9
	2	64	64	24	59	53	17
	3	22	42	24	17	41	18
	4	24	24	30	17	22	26
4	1	19	23	23	15	19	20
	2	47	39	31	44	32	25
	3	39	27	47	35	25	42
	4	24	28	21	22	24	19
5	1	21	47	16	16	41	13
	2	34	53	31	29	48	25

Grade	ALD	In the Pool			Administered		
		IA1	IA2	IA3	IA1	IA2	IA3
	3	42	15	41	35	13	35
	4	25	11	16	17	11	13
6	1	15	37	19	14	32	14
	2	28	36	17	26	33	17
	3	35	36	31	31	35	25
	4	40	34	17	34	33	16
7	1	33	32	3	31	30	2
	2	48	78	13	40	70	13
	3	24	23	55	22	19	50
	4	25	14	19	15	9	14
8	1	29	50	5	27	48	4
	2	45	61	17	38	55	17
	3	21	40	69	17	36	61
	4	20	32	31	16	29	27
HS	1	9	9	1	9	9	1
	2	14	31	5	14	31	5
	3	30	24	5	30	24	5
	4	9	8	1	9	8	1

Note. IA1 = Literary Text, IA2 = Informational Text, IA3 = Vocabulary;
 ALD1 = Well Below State Expectations, ALD2 = Below State Expectations, ALD3 = At State Expectations, and
 ALD4 = Above State Expectations

Table 3.7. ALD Distribution Across Instructional Areas—Mathematics

Grade	ALD	In the Pool				Administered			
		IA1	IA2	IA3	IA4	IA1	IA2	IA3	IA4
3	1	38	30	25	1	37	30	24	1
	2	39	38	50	8	39	38	50	8
	3	31	39	43	10	30	39	42	10
	4	23	19	55	3	23	18	49	3
4	1	11	34	13	7	11	33	13	7
	2	18	60	21	19	18	58	20	19
	3	21	94	30	9	21	90	30	9
	4	11	34	11	5	11	33	11	5
5	1	8	37	5	6	8	37	5	6
	2	12	70	16	8	12	69	16	8
	3	23	80	38	24	23	77	38	24
	4	8	36	12	7	8	33	12	7
6	1	36	40	7	17	36	40	7	17
	2	34	46	19	4	34	46	19	4
	3	37	54	12	20	37	52	12	20

Grade	ALD	In the Pool				Administered			
		IA1	IA2	IA3	IA4	IA1	IA2	IA3	IA4
	4	11	30	4	6	11	28	4	6
7	1	12	16	5	14	12	16	5	13
	2	29	28	14	52	29	28	14	52
	3	30	64	4	47	30	63	4	47
	4	3	8	4	22	3	8	4	22
8	1	61	5	15	28	61	5	13	21
	2	62	7	26	18	60	7	24	14
	3	33	7	39	33	31	7	39	33
	4	32	1	19	16	32	1	19	16
HS	1	15	3	11	3	15	3	11	3
	2	34	6	23	7	32	6	23	7
	3	30	3	23	5	29	3	22	5
	4	6	4	9	1	6	4	8	1

Note. For grades 3–5: IA1 = Operations and Algebraic Thinking, IA2 = Number and Operations, IA3 = Measurement and Data, IA4 = Geometry;
 For grades 6–8 & HS: IA1 = Operations and Algebraic Thinking, IA2 = The Real and Complex Number Systems, IA3 = Geometry, IA4 = Statistics and Probability;
 ALD1 = *Well Below State Expectations*, ALD2 = *Below State Expectations*, ALD3 = *At State Expectations*, and ALD4 = *Above State Expectations*

Table 3.8 shows the median item response times by item ALD. Each item was matched to an associated ALD during the standard setting meeting. The median response times for items matched to each ALD were calculated following test administration. Results show that items in ALD1 had shorter average response times than items in ALD2, items in ALD2 had shorter average response times than items in ALD3, and items in ALD3 had shorter average response times than items in ALD4 across all grades in reading and mathematics, except for reading grades 3, 7, and HS. The results shown in the table indicate that items in higher ALDs are more complex and require more cognitive load than items in lower ALDs.

Table 3.8. Median Item Response Time by ALD

Grade	Median Response Time (in seconds)			
	ALD1	ALD2	ALD3	ALD4
Reading				
3	40	58	56	62
4	44	52	57	69
5	52	55	57	68
6	52	54	56	59
7	49	54	53	66
8	43	53	58	59
HS	43	40	46	41
Mathematics				
3	35	41	47	59

Grade	Median Response Time (in seconds)			
	ALD1	ALD2	ALD3	ALD4
4	32	43	47	67
5	40	47	50	74
6	43	46	51	81
7	45	48	53	72
8	35	50	56	66
HS	30	40	45	46

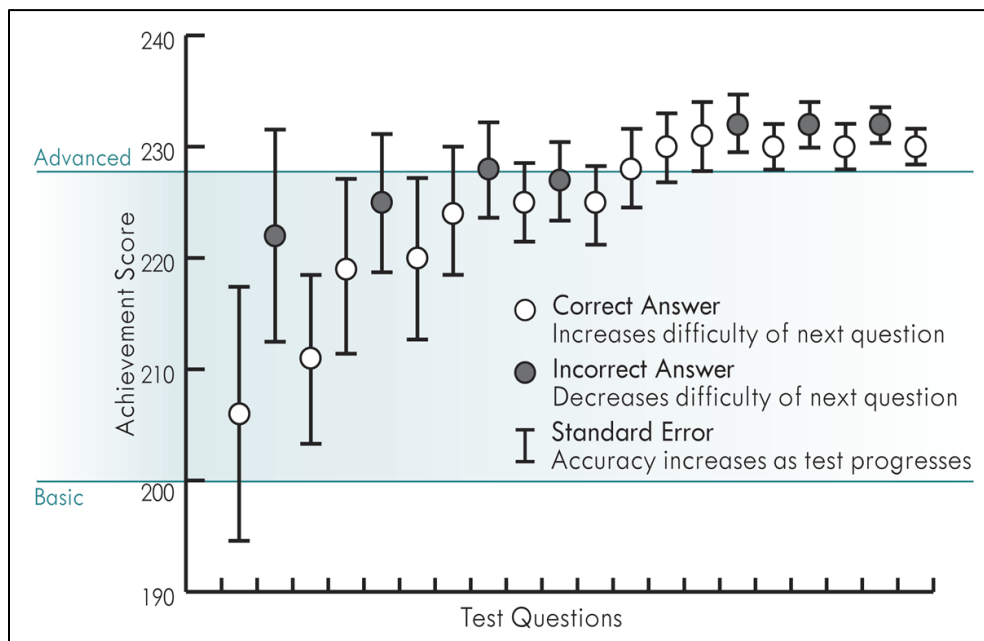
Note. ALDs matched to an associated item during standard setting are defined as: ALD1 = *Well Below State Expectations*, ALD2= *Below State Expectations*, ALD3 = *At State Expectations*, & ALD4 = *Above State Expectations*

3.3. Constraint-Based Engine Adaptive Test Administration

A CAT administers items to match the ability level of the students: different students receive different items based on item difficulty and their ability levels. For example, students with lower ability levels (based on their answers to previous items) receive easier items compared with students with higher ability levels who receive harder items as the test progresses.

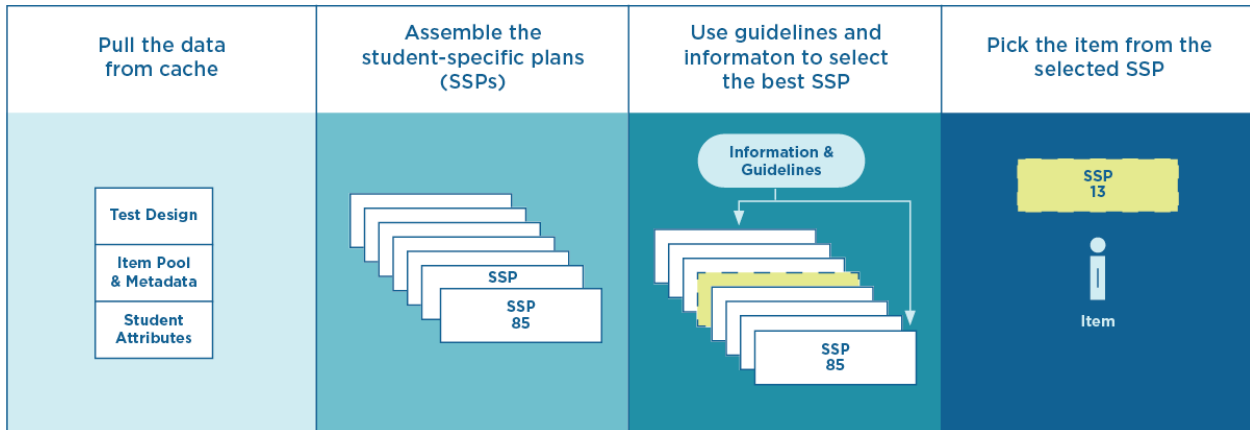
The constraint-based engine (CBE) uses the blueprint and a student’s momentary theta (θ) to drive item selection, as shown in Figure 3.1. Momentary theta is the ability estimate of the student that is recalculated and updated after answering each item. Items are selected based on item difficulty. The goal of the constraint-based engine’s item selection is to provide a test that meets “must-have” constraints and “nice-to-have” guidelines. For example, a constraint of the summative portion is that the engine must deliver 70% on-grade items, while the remaining 30% may adapt by one grade level below or above. The CBE has two stages of consideration as it selects the items necessary to conform to the test blueprint while providing the maximum information about the student based on the student’s momentary ability estimate.

Figure 3.1. Adaptive Engine Overview



The student-specific plan (SSP), similar to the shadow test approach (Van der Linden & Reese, 1998), selects items based on the required aspects of the test blueprint and the student’s momentary theta, as shown in Figure 3.2. Item selection for the SSP occurs through a process of choosing multiple feasible SSPs and then choosing the complete SSP that best maximizes guideline adherence and information. Only after the best SSP has been chosen are items ordered (NWEA, 2020).

Figure 3.2. Student-Specific Plan Approach



Note. Selections are based on the similar shadow test approach.

3.3.1. Engine Evaluation

NWEA checks the adaptive engine at two points: pre-administration simulations and a post-administration evaluation. These two studies are important evidence, along with post-administration analyses, for confirming interpretation and test-score use arguments regarding student proficiency with the state standards.

Pre-administration simulations are conducted prior to the operational testing window to evaluate the CBE’s item-selection algorithm and estimation of student ability based on the test blueprints and adaptive specifications. The simulation tool uses the operational CBE, thereby providing results with the same properties and functionality as what would be seen operationally. Detailed information regarding the simulation study can be found in Appendix C. After the testing window closes, a post-administration evaluation study is conducted to determine whether the CBE performed as expected. The results of the post-administration evaluation study are presented in this section.

In order to deliver a quality test, various constraints and guidelines are set up in the CBE to specify details of the test requirements. While constraints are rules that must be followed, weights are used to differentiate the importance of different guidelines. One constraint is meeting the requirements of the test blueprint. Because the adaptive test selects items according to individual student abilities in order to provide reliable scores, score precision and item-exposure rates are also important factors. Results for blueprint constraint accuracy, item-exposure rates, and score precision and accuracy are presented below.

3.3.2. Blueprint Constraint Accuracy

Table 3.9 presents the blueprint constraint results at the reporting category level for the spring administration. This analysis exclusively focused on students who completed the maximum/full-length test for each test event, and, in all cases, it yielded a perfect match for the number of items at the reporting category level.

Table 3.9. Blueprint Constraint Accuracy by Reporting Category

Grade	Summative Content Across Instructional Areas	#Items Intended		#Items Administered			%Match
		Min	Max	Average	Min	Max	
Reading							
3	Literary Text	12	14	12	12	14	100
	Informational Text	8	9	9	8	9	100
	Vocabulary	5	7	6	5	7	100
4	Literary Text	11	12	11	11	12	100
	Informational Text	9	11	9	9	11	100
	Vocabulary	5	7	6	5	7	100
5	Literary Text	9	11	11	9	11	100
	Informational Text	9	11	10	9	11	100
	Vocabulary	5	7	7	5	7	100
6	Literary Text	9	11	10	9	11	100
	Informational Text	11	12	11	10	12	100
	Vocabulary	5	7	6	5	7	100
7	Literary Text	8	9	9	8	9	100
	Informational Text	12	14	12	12	14	100
	Vocabulary	5	7	6	5	7	100
8	Literary Text	8	9	9	8	9	100
	Informational Text	12	14	12	12	14	100
	Vocabulary	5	7	6	5	7	100
HS	Literary Text	8	9	8	8	8	100
	Informational Text	12	14	14	14	14	100
	Vocabulary	5	8	8	8	8	100
Mathematics							
3	Operations and Algebraic Thinking	6	6	6	6	6	100
	Numbers and Operations	9	9	9	9	9	100
	Measurement and Data	8	8	8	8	8	100
	Geometry	4	4	4	4	4	100
4	Operations and Algebraic Thinking	5	5	5	5	5	100
	Numbers and Operations	13	13	13	13	13	100
	Measurement and Data	5	5	5	5	5	100
	Geometry	4	4	4	4	4	100
5	Operations and Algebraic Thinking	4	4	4	4	4	100
	Numbers and Operations	14	14	14	14	14	100
	Measurement and Data	5	5	5	5	5	100
	Geometry	4	4	4	4	4	100

Grade	Summative Content Across Instructional Areas	#Items Intended		#Items Administered			%Match
		Min	Max	Average	Min	Max	
6	Operations and Algebraic Thinking	7	7	7	7	7	100
	The Real and Complex Number Systems	12	12	12	12	12	100
	Geometry	4	4	4	4	4	100
	Statistics and Probability	4	4	4	4	4	100
7	Operations and Algebraic Thinking	5	5	5	5	5	100
	The Real and Complex Number Systems	11	11	11	11	11	100
	Geometry	6	6	6	6	6	100
	Statistics and Probability	5	5	5	5	5	100
8	Operations and Algebraic Thinking	13	13	13	13	13	100
	The Real and Complex Number Systems	4	4	4	4	4	100
	Geometry	6	6	6	6	6	100
	Statistics and Probability	4	4	4	4	4	100
HS	Operations and Algebraic Thinking	14	14	14	14	14	100
	The Real and Complex Number Systems	4	4	4	4	4	100
	Geometry	8	8	8	8	8	100
	Statistics and Probability	4	4	4	4	4	100

3.3.3. Item Exposure Rates

Because different students receive different items based on blueprint constraints and their ability during an adaptive administration, it is ideal to have a low exposure rate. The exposure rate for each operational item is calculated as the percentage of students who received that item, as shown in Table 3.10. For example, if Item 1 was administered to 500 out of 1,000 students, the exposure rate would be 50%. “Total” is the total number of items in the operational item pool.

Table 3.10. Operational Item Exposure Rates

Grade	#Items				Item Exposure Rate											
					0–20%		21–40%		41–60%		61–80%		81–99%		100%	
	Total	Used	Unused	Unused %	N	%	N	%	N	%	N	%	N	%	N	%
Reading																
3	666	499	167	25.08	460	92.18	32	6.41	5	1.00	1	0.20	0	0.00	0	0.00
4	986	590	396	40.16	558	94.58	31	5.25	1	0.17	0	0.00	0	0.00	0	0.00
5	901	603	298	33.07	566	93.86	33	5.47	4	0.66	0	0.00	0	0.00	0	0.00
6	894	524	370	41.39	502	95.80	18	3.44	4	0.76	0	0.00	0	0.00	0	0.00
7	966	516	450	46.58	477	92.44	39	7.56	0	0.00	0	0.00	0	0.00	0	0.00
8	691	540	151	21.85	511	94.63	24	4.44	5	0.93	0	0.00	0	0.00	0	0.00
HS	146	146	0	0.00	85	58.22	36	24.66	21	14.38	4	2.74	0	0.00	0	0.00
Mathematics																
3	745	513	232	31.14	502	97.86	8	1.56	0	0.00	1	0.19	2	0.39	0	0.00
4	1049	525	524	49.95	518	98.67	7	1.33	0	0.00	0	0.00	0	0.00	0	0.00
5	965	511	454	47.05	498	97.46	11	2.15	2	0.39	0	0.00	0	0.00	0	0.00
6	960	444	516	53.75	422	95.05	20	4.50	2	0.45	0	0.00	0	0.00	0	0.00
7	966	456	510	52.80	429	94.08	25	5.48	1	0.22	1	0.22	0	0.00	0	0.00
8	694	441	253	36.46	420	95.24	17	3.85	4	0.91	0	0.00	0	0.00	0	0.00
HS	183	178	5	2.73	134	75.28	12	6.74	17	9.55	13	7.30	2	1.12	0	0.00

A number of field test items were embedded in the Spring 2025 test for possible operational use in future test administrations. Field test items were distributed using target demographic characteristics of the Maine student population. For example, each item should be administered to approximately 50% female and 50% male students if the Maine student population has a 50/50 gender proportion. The results presented in Table 3.11 show that all field test items were appropriately administered to each demographic subgroup.

Table 3.11. Field Test Item Exposure Rates

Grade	FT Items	Mean	Female	Male	Hispanic/Latino	Am. Indian or Alaska Native	Asian	Black or African American	Native HI or Pacific Islander	White	Two or More Races
Reading											
3	15	4,181	48	51	4	1	1	5	0	84	4
4	19	3,119	48	52	4	1	1	6	0	84	4
5	19	3,262	49	50	4	1	1	5	0	86	4
6	18	3,401	48	52	4	1	1	5	0	85	4
7	15	4,001	48	51	3	1	1	5	0	85	4
8	18	3,420	47	51	3	1	1	6	0	84	3
HS	153	561	48	52	4	1	2	5	0	85	3
Mathematics											
3	15	4,181	48	51	4	1	1	5	0	84	4
4	15	3,952	48	51	4	1	1	6	0	85	4
5	15	4,107	49	51	4	1	1	5	0	85	4
6	14	4,100	48	52	4	1	1	5	0	85	4
7	15	3,993	48	51	4	1	1	5	0	85	4
8	13	4,726	49	51	3	1	1	5	0	85	4
HS	150	574	48	52	4	1	2	5	0	85	3

3.3.5. Item Sequence

The distribution of items that each student receives is not inherently subject to a predefined sequence or grouping for the summative, diagnostic, and field test items. In the absence of specific preferences, the adaptive engine arranges the items based on the individual student's test performance. An exception to this rule pertains to items that are part of a set with a common reading passage or paired passages; in such cases, the engine ensures that these items are delivered as a cohesive group rather than being dispersed. NWEA's evaluation reveals that items were allocated based on their performance without adhering to any predefined sequence or grouping, except for the designated locations for the field test items. In the reading tests, the actual placement of field test items varied due to the arrangement of reading passage sets and the engine's design to avoid introducing unrelated items in the midst of a reading passage set.

3.4. Paper Form Administration

For Spring 2025, the majority of Maine's students participated through the computer adaptive assessment. Students with an IEP or 504 Plan could request an alternate, accommodated paper-based form in standard print, large print, or braille. A fixed test form was built for each grade and content area to fulfill the needs of the three accommodated test forms. Braille and large print forms were prepared in advance according to registration data, and the required materials were packed and shipped to the requesting schools. Standard paper-based forms were available via print on demand. These materials were sent to School Assessment Coordinators via NWEA's secure SFTP site.

Table 3.12 presents the numbers of summative operational items needed for the spring fixed forms.

- All items are on grade level.
- There are no anchor or linking items on the paper forms.

Table 3.12. Paper Form Summative Item Totals by Content and Grade

Content	Grade	Summative Operational
Reading	3–8	27
Mathematics	3–8	27
Reading	HS	30
Mathematics	HS	30

3.4.1. Receiving and Taking Inventory of School Materials

The quantity of materials shipped to each school is based on data collected during the rostering process. School Assessment Coordinators are required to open packages containing braille or large print forms immediately upon receipt to inventory the contents. School Assessment Coordinators are responsible for the printing and secure handling of standard paper-based forms, as well as for providing secure assessment materials to proctors. All standard assessment booklets are provided as single materials. School Assessment Coordinators do not distribute any assessment materials, except the *Maine Through Year Assessment Proctor User Guide* and *The Maine Through Year Assessment Administration Guide*, until the day of each session.

On the day of the assessment, the School Assessment Coordinator distributes the correct assessment booklets needed for that day's assessment to each proctor. Assessment booklets are distributed to proctors early enough on the day of the assessment to give them ample time to review the directions prior to the assessment. After each day of the assessment is complete, all assessment materials are returned to the School Assessment Coordinator for secure storage as soon as possible. All materials, including used and unused booklets and scratch paper, are returned at the end of each day of testing.

3.4.2. Score Transcription

During or immediately following assessment administration, student responses for paper-based accommodated assessments are transcribed into the online assessment engine. To transcribe responses requires the proctor or other designated and authorized district or school personnel to log in to the NWEA State Solutions Secure Browser using the student's test ticket. The required steps for the proctor to transcribe student answers are as follows:

1. Obtain the student's test ticket from the School Assessment Coordinator.
2. After the student has completed the paper accommodated assessment, use a device that has the NWEA State Solutions Secure Browser software installed and use the student's test ticket to log in to the student's assessment.
3. For security reasons, Maine DOE recommends, when feasible, that a second trained staff member be present to verify all transcriptions.
4. Once transcribing student responses is complete, the assessment is submitted. The proctor should then return all printed assessment materials to the School Assessment Coordinator.

Transcribing is the process of moving the student's assessment response to another medium by a district employee. The process should be as faithfully completed as possible and follow the qualifications and procedures as outlined:

1. The transcriber must be a current employee of the school district.
2. The transcriber must be trained in assessment administration and have signed the Assessment Security and Data Privacy Agreement.
3. Transcription must take place in a secure location.
4. The assessment must be transcribed exactly as the student answered the assessment items.

Local SAU policy will determine whether School Assessment Coordinators should securely destroy test tickets, scrap paper, and accommodated paper forms on-site or if all materials should be sent to the district office to be securely destroyed by the District Assessment Coordinator. If shipping to the district office, security and record-keeping guidance must be followed.

3.5. Assessment Security

In a centralized assessment process, it is critical that equity of opportunity, standardization of procedures, and fairness to students is maintained. Therefore, Maine DOE requires that all assessment administrators and proctors review the information in the *Maine Assessment Security Handbook*.

The Maine DOE recommends that assessment administrators (or proctors) report any potential irregularities to the School Assessment Coordinator. This is especially important for any irregularities that may:

- (1) involve a breach of assessment item security
- (2) lead to assessment invalidation
- (3) involve student misbehavior
- (4) involve educator misbehavior

The School Assessment Coordinator, or other administrator, should report irregularities to Krista Averill, Maine DOE Assessment Coordinator, at Krista.Averill@maine.gov or 1-207-215-6528. See the *Maine Assessment Security Handbook* for more details on this process.

3.5.1. Assessment Ethics and Appropriate Practice

All teachers need to be familiar with appropriate assessment ethics and security practices related to assessments. Proctors are expected to actively monitor student participation during the assessment to ensure students remain on-task. Professionalism, common sense, and practical procedures provide the right framework for assessment ethics. The *Maine Assessment Security Handbook* outlines clear practices for appropriate security.

3.5.2. Online Security

Student test tickets contain student-level password information for accessing the assessment and must be kept secure. Proctors should print or be given the student test tickets prior to assessment administration, allowing them ample time to review and organize the tickets for distribution before the assessment begins. Once an assessment session is started, only the student taking the assessment is allowed to view the student's screen. No one is allowed to view or copy assessment content while a student is taking the assessment.

The Maine Through Year Assessment Coordinator Guide, as well as other manuals and guides available online, are not considered secure assessment materials.

3.5.3. Student Assessment Security

Students should look only at their individual computers. For further security, folders may be set up around each computer screen to eliminate any possibility of students looking at other computer screens. For larger groups, it is advisable to have a sufficient number of proctors to monitor the room.

3.5.4. Returning or Destroying Secure Materials

Proctors should collect all student test tickets, scratch paper, and assessment booklets (where applicable) from students after the assessment so that those materials can be securely destroyed.

3.6. Systems for Protecting Data Integrity and Privacy

School Assessment Coordinators, assessment administrators, and proctors are required to complete and sign the MEA Assessment Security and Data Privacy Agreement. Signed copies should be filed and kept on-site, available for delivery to the Maine DOE if requested.

NWEA maintains the following protocols to ensure that the sensitive data that are captured are protected and secure from unauthorized use, hacks, or other forms of compromise.

Test Content Security

NWEA encrypts test data both prior to transmission and in-transit and then delivers the data through a secure downloadable browser that is only accessible through 256-bit TLS user authentication and proctor-provided usernames and passwords. NWEA's test system also saves students' work at frequent intervals, and assessment packages are encrypted while on students' workstations.

Data Protection

- Data at rest are protected across a wide range of Amazon Web Services (AWS) and state applications.
- Encryption is enabled for all network traffic, including Transport Layer Security for web-based network infrastructure
- Policies and procedures to protect personally identifiable information (PII) data are strictly enforced.

Secure Identity and Access Management

- A centralized identity provider is used to manage account access, restricting access to authorized personnel only.
- A least privilege model is used to ensure operational staff have only those privileges needed to complete their tasks.
- Multi-factor authentication and other account-level controls are enabled.
- Passwords and other credentials are securely stored using AWS tools that handle encryption, rotation, and access control.

Infrastructure Protection

- Operating systems, middleware, applications, and code are patched on a regular basis.
- Distributed Denial of Service (DDoS) protection layers are used for all internet-facing applications.
- Intrusion detection/prevention services are utilized.
- Inbound and outbound traffic is controlled and monitored based on established rules.

Detection and Monitoring

- AWS are leveraged to comprehensively monitor all layers.
- Application and system-level logs are analyzed periodically to gain insights into the information contained within them.
- An incident management process is maintained for security events that may affect the confidentiality, integrity, or availability of systems or data.
- Monitoring and alerts are configured and investigated regularly for any unexpected events, including hacking attempts and attacks.

Section 4: Item Statistics, Calibration, and Scaling

This section presents item statistics and the methods and process of establishing the Maine scale.

4.1. Classical Item Statistics

4.1.1. Expected P Value

Item difficulty is measured by a p value that represents the proportion of students who answered an item correctly and ranges from 0 to 1. A high p value indicates an easy item, with a high percentage of students answering it correctly, whereas a low p value indicates a difficult item. For example, a p value of 0.79 indicates that 79% of students answered the item correctly. In the case of polytomous items, the p value is calculated as the average item score divided by the number of possible score points on the item.

Table 4.1 and Table 4.2 present the summary statistics for the p values across operational and field test items, respectively, and the count of items falling within different p -value ranges (e.g., less than or equal to 0.1, 0.2, etc.). The data include adaptive items for all grades. For adaptive items that were administered without a representative student sample, their expected p values are provided. An expected p value represents the proportion of correct responses if the item was administered to a representative student sample. Appendix D provides the summary p -value statistics by item type.

Table 4.1. Summary of P Values—Operational Items

Grade	N	P Value Summary					P Value Counts									
		Mean	Median	SD	Min	Max	≤ 0.1	≤ 0.2	≤ 0.3	≤ 0.4	≤ 0.5	≤ 0.6	≤ 0.7	≤ 0.8	≤ 0.9	> 0.9
Reading																
3	331	0.47	0.46	0.11	0.15	0.92	1	13	56	151	79	22	4	4	0	1
4	406	0.49	0.49	0.11	0.14	0.87	5	15	57	143	133	41	7	5	0	0
5	411	0.49	0.49	0.12	0.06	0.92	4	13	44	179	103	43	13	8	2	2
6	384	0.50	0.49	0.12	0.13	0.96	3	9	54	150	108	39	17	3	0	1
7	369	0.48	0.47	0.11	0.09	0.86	0	10	56	159	96	35	10	2	1	0
8	378	0.50	0.49	0.12	0.12	0.93	3	8	43	155	117	31	10	9	0	2
HS	140	0.48	0.48	0.11	0.24	0.77	0	8	21	54	38	16	3	0	0	0
Mathematics																
3	451	0.49	0.49	0.07	0.17	0.75	2	6	25	252	146	12	8	0	0	0
4	453	0.49	0.49	0.09	0.14	0.90	4	10	28	235	139	26	8	3	0	0
5	455	0.48	0.49	0.07	0.04	0.83	2	6	29	278	123	13	2	1	1	0
6	404	0.48	0.48	0.09	0.16	0.84	2	10	33	215	120	12	10	2	0	0
7	389	0.48	0.48	0.07	0.10	0.75	3	5	21	251	97	10	2	0	0	0
8	399	0.46	0.47	0.07	0.10	0.85	4	6	43	261	78	6	0	1	0	0
HS	171	0.41	0.42	0.07	0.15	0.66	1	7	50	107	4	2	0	0	0	0

Table 4.2. Summary of P Values—Field Test Items

Grade	N	P Value Summary					P Value Counts									
		Mean	Median	SD	Min	Max	≤ 0.1	≤ 0.2	≤ 0.3	≤ 0.4	≤ 0.5	≤ 0.6	≤ 0.7	≤ 0.8	≤ 0.9	> 0.9
Reading																
3	5	0.50	0.53	0.16	0.27	0.71	0	2	3	2	2	2	4	0	0	0
4	5	0.46	0.45	0.16	0.13	0.71	1	2	5	4	2	1	4	0	0	0
5	7	0.48	0.44	0.15	0.15	0.80	1	0	3	9	1	1	2	0	1	0
6	5	0.44	0.45	0.12	0.09	0.59	0	0	6	5	6	0	0	1	0	0
7	5	0.53	0.54	0.19	0.18	0.79	1	1	2	2	2	2	5	0	0	0
8	12	0.43	0.37	0.16	0.24	0.75	0	2	8	3	1	2	2	0	0	0
HS	117	0.46	0.46	0.16	0.09	0.83	5	21	28	32	32	10	21	1	3	0
Mathematics																
3	6	0.39	0.34	0.19	0.11	0.76	2	4	2	3	2	2	0	0	0	0
4	11	0.42	0.44	0.19	0.10	0.75	2	0	3	5	2	1	1	1	0	0
5	5	0.42	0.42	0.17	0.04	0.66	1	1	2	4	5	0	1	1	0	0
6	18	0.37	0.38	0.18	0.05	0.62	2	2	3	2	2	0	2	1	0	0
7	19	0.38	0.37	0.21	0.09	0.70	2	4	2	1	3	0	2	1	0	0
8	13	0.35	0.30	0.21	0.05	0.73	0	5	0	3	1	1	1	2	0	0
HS	145	0.33	0.33	0.16	0.05	0.80	20	29	45	22	11	3	7	13	0	0

4.1.2. Item Discrimination (Item-Total Correlation)

Item-total correlation describes the relationship between performance on an item and performance on the entire test (test scaled score). Students who perform well on a test are expected to have a higher probability of selecting the right answer to any given item, and students who perform poorly are more likely to select the wrong answer. This means that for a highly discriminating item, students who get the item correct will have a higher test score than students who get the item incorrect. The item-total correlation coefficient ranges between -1.0 and $+1.0$. An item with a high positive item-total correlation discriminates between low-performing and high-performing students better than an item with an item-total correlation near zero. A negative item-total correlation indicates that lower-performing students did better on that item than higher-performing students. However, if an item is either very difficult or very easy, there will be little variation in student responses, as most students would either respond incorrectly or correctly. The resulting item-total correlation for such items is typically low.

Table 4.3 and Table 4.4 present the summary statistics for the item-total correlations across operational and field items, respectively. Instead of using the number-correct raw score, the estimated final scaled score was used to compute the item-total correlations because number-correct scores would not provide much insight into student performance on an adaptive test. For items administered adaptively in grades 3–8, their item-total correlations tend to be lower because these adaptive items were seen by students within a restricted ability range. Additionally, most of the items displaying negative item-total correlations had very few responses (less than 10 student responses). Appendix E provides the summary item-total correlation statistics by item type.

Table 4.3. Summary of Item-Total Correlations—Operational Items

Grade	N	Item-Total Correlation Summary					Item-Total Correlation Counts									
		Mean	Median	SD	Min	Max	≤ 0.1	≤ 0.2	≤ 0.3	≤ 0.4	≤ 0.5	≤ 0.6	≤ 0.7	≤ 0.8	≤ 0.9	> 0.9
Reading																
3	331	0.35	0.35	0.12	-0.19	0.62	21	76	124	81	18	4	0	7	0	0
4	406	0.35	0.35	0.11	-0.14	0.86	18	91	156	92	36	1	1	10	1	0
5	411	0.34	0.35	0.10	-0.04	0.64	22	103	169	80	28	2	0	7	0	0
6	384	0.35	0.35	0.11	-0.16	0.71	26	69	172	88	21	1	1	6	0	0
7	369	0.34	0.34	0.10	-0.03	0.66	21	99	148	81	14	3	0	3	0	0
8	378	0.34	0.35	0.12	-0.29	0.74	24	82	138	95	22	3	1	13	0	0
HS	140	0.37	0.37	0.12	-0.03	0.73	6	34	45	32	18	2	1	2	0	0
Mathematics																
3	451	0.35	0.34	0.09	-0.02	0.71	12	92	222	105	13	3	1	3	0	0
4	453	0.37	0.37	0.09	-0.04	0.78	12	71	219	114	28	4	1	4	0	0
5	455	0.36	0.35	0.09	0.03	0.86	10	96	224	98	24	1	0	1	1	0
6	404	0.36	0.36	0.09	-0.15	0.65	14	76	181	103	24	3	0	3	0	0
7	389	0.35	0.35	0.09	0.01	0.60	17	91	178	83	16	1	0	3	0	0
8	399	0.36	0.35	0.09	0.10	0.74	10	85	186	93	17	7	1	0	0	0
HS	171	0.33	0.32	0.10	0.08	0.64	14	55	67	25	6	3	0	1	0	0

Table 4.4. Summary of Item-Total Correlations—Field Test Items

Grade	N	Item-Total Correlation Summary					Item-Total Correlation Counts										
		Mean	Median	SD	Min	Max	≤ 0.1	≤ 0.2	≤ 0.3	≤ 0.4	≤ 0.5	≤ 0.6	≤ 0.7	≤ 0.8	≤ 0.9	> 0.9	
Reading																	
3	15	0.36	0.39	0.12	0.07	0.50	1	2	5	6	1	0	0	0	0	0	
4	19	0.36	0.39	0.13	0.01	0.54	0	6	3	6	1	3	0	0	0	0	
5	18	0.31	0.33	0.11	0.04	0.47	1	5	7	4	1	0	0	0	0	0	
6	18	0.31	0.31	0.10	0.14	0.50	2	7	6	2	0	1	0	0	0	0	
7	15	0.35	0.36	0.10	0.15	0.50	2	2	6	4	0	1	0	0	0	0	
8	18	0.27	0.30	0.18	-0.19	0.50	3	4	5	4	2	0	0	0	0	0	
HS	153	0.26	0.28	0.15	-0.07	0.57	17	39	39	23	32	3	1	0	0	0	
Mathematics																	
3	15	0.32	0.36	0.14	-0.04	0.50	2	1	6	5	1	0	0	0	0	0	
4	15	0.42	0.43	0.09	0.30	0.61	0	0	7	5	0	2	1	0	0	0	
5	15	0.39	0.41	0.17	0.07	0.59	1	2	3	3	1	5	0	0	0	0	
6	14	0.34	0.38	0.18	-0.07	0.52	0	1	5	4	2	2	0	0	0	0	
7	15	0.30	0.34	0.15	-0.01	0.53	2	2	5	3	2	1	0	0	0	0	
8	13	0.35	0.38	0.13	0.04	0.53	0	3	4	3	1	2	0	0	0	0	
HS	150	0.29	0.31	0.16	-0.11	0.55	16	29	49	29	20	7	0	0	0	0	

4.2. IRT Calibration

When Maine’s scale was established in 2023, the first step was to calibrate items to a standardized scale and then use the calibrated items to derive student scores. The Rasch model (Rasch, 1960, 1980; Wright, 1977) for dichotomous items and the partial-credit model (PCM; Masters, 1982) for polytomous items were used to calibrate items and create the Maine scale. These two models have had a long-standing presence in applied testing programs. For all content areas, item parameter estimations were implemented using WINSTEPS 3.90.2.0 (Linacre, 2015) that used joint maximum likelihood estimation (MLE), as described by Wright (1977) and Masters (1982).

Under the Rasch model, the probability of a student with ability θ responding correctly to item i is as follows, where θ_j and b_i are the person and item parameters, respectively:

$$P(u_{ij} = 1|\theta_j, b_i) = \frac{e^{(\theta_j - b_i)}}{1 + e^{(\theta_j - b_i)}}$$

Under the PCM, the probability of a student with ability θ having a score at the k th level of item i is:

$$P(u_{ij} = k|\theta_i) = \frac{e^{[\sum_{u=1}^k(\theta_j - b_i + d_{iu})]}}{\sum_{v=1}^{m_i} e^{[\sum_{u=1}^k(\theta_j - b_i + d_{iu})]}}$$

where k is the score on the item, m_i is the total number of score categories for the item, d_{iu} is the threshold parameter for the threshold between scores u and $u - 1$, and θ_j and b_i are the person and item parameters, respectively.

The free calibration¹ method was used to derive item parameters for the summative items by subject and grade. Table 4.5 presents the summary of IRT item statistics across all operational items.

Table 4.5. Summary of IRT Item Statistics—Operational Items

Grade	#Items	Mean	Median	SD	Min	Max	Range (Max–Min)
Reading							
3	142	-0.11	-0.19	1.34	-3.04	3.59	6.63
4	183	-0.08	-0.10	1.42	-3.11	4.12	7.23
5	184	-0.09	-0.08	1.37	-2.92	2.97	5.89
6	197	-0.12	-0.20	1.25	-2.68	2.98	5.66
7	186	0.10	0.10	1.36	-2.49	3.34	5.82
8	185	0.22	0.02	1.26	-2.15	3.23	5.38
HS	108	-0.02	-0.08	0.87	-2.00	2.96	4.96
Mathematics							
3	242	-0.40	-0.49	1.81	-4.89	3.99	8.88

¹ Calibration can be done by itself or combined with equating. The former is referred to as free calibration, and the latter is the anchor/fixed parameter method.

Grade	#Items	Mean	Median	SD	Min	Max	Range (Max–Min)
4	252	-0.17	-0.45	2.03	-4.32	4.44	8.75
5	233	-0.47	-0.56	1.92	-3.91	4.08	7.99
6	226	-0.83	-0.67	1.96	-5.00	4.41	9.41
7	208	-0.84	-0.84	1.70	-4.64	3.72	8.36
8	218	-0.71	-0.83	1.73	-4.39	3.44	7.83
HS	118	-0.35	-0.53	1.45	-3.36	3.19	6.55

4.3. IRT Model Assumptions

Being one of the item response theory models (IRT), Rasch and PCM models have the same assumptions as other IRT models: local independency, model fit, and unidimensionality (Hambleton & Swaminathan, 1985). These three assumptions are checked to evaluate the appropriateness of using the Rasch and PCM models for the assessment.

4.3.1. Local Independence

Local independence refers to a response to an item that is not affected by other items after removing the contribution of ability measures. The IRT model assumes that the response to an item is only affected by the item’s difficulty and student’s ability. Local dependence violates this assumption by introducing factors irrelevant to those two factors. Examples of local independence violation are:

- The response to an item depends on the response to a prior item—such as, derive a value from Item A, then use Item A’s response to solve Item B’s equation. If Item A is answered incorrectly, then the response to Item B must be wrong. Scores on Item B are affected by the answer to Item A, a factor other than item difficulty and student ability.
- Other items on the test give away the answer to Item A—this is referred to as “clueing” in test development.

When constructing items, each item has a complete concept in itself and does not rely on other items. When selecting items for an adaptive test, item enemy information is incorporated to avoid clueing.

4.3.2. Model Fit

Model fit refers to how well an item fits the calibration model. It is usually a statistical chi-square, representing the difference between the observed score (i.e., actual student responses to items) and the expected score (i.e., what the model predicts students with a certain ability should be getting on items). Individual item fit is evaluated using infit and outfit statistics:

- **Infit:** an information-weighted fit statistic that is more sensitive to unexpected student behavior affecting responses to items near the student’s ability level
- **Outfit:** an outlier-sensitive fit statistic that is more sensitive to unexpected student behavior on items far from the student’s ability level

Both infit and outfit provide mean-square fit (MNSQ) statistics. The expected value of MNSQ is 1.0. Summary statistics for the infit and outfit MNSQ statistics are presented in Table 4.6. The fit statistics were computed using response data from on-grade items with a minimum of 500

responses to ensure statistical stability. A cutoff of greater than 2.0 is used for item-fit flagging. The review process pays more attention to the infit than to the outfit because infit is the more stable statistic.

The table shows that all average infit and outfit values are close to 1.0, indicating that items fit well at their intended grade level. Infits are very stable across grades, and maximum values are all less than 2.0, as they reflect how well an individual item fits the overall measurement model. This stability ensures that results are not significantly affected by grade-level differences, making it particularly useful for longitudinal or multi-grade assessments. While some grades have cases of item outfit values greater than 2.0, the majority of such values are within the value of 3.0. These items have less impact on the measurement system because “outfit problems are less of a threat to measurement than infit ones” (Linacre, 2002). The results from the model fit analyses and item statistics will be used to inform future item development. For instance, if items with model fit statistics that fall outside of the acceptable range are found to be relatively easy or difficult, they will be replaced during item development to ensure proper coverage of the student ability scale.

Table 4.6. Summary of Mean-Square Infit and Outfit Statistics

Grade	N	Infit				Outfit			
		Mean	SD	Min	Max	Mean	SD	Min	Max
Reading									
3	142	0.98	0.07	0.86	1.42	0.99	0.12	0.84	1.94
4	183	0.99	0.10	0.85	1.85	1.03	0.54	0.81	8.21
5	184	0.98	0.09	0.84	1.72	0.99	0.18	0.78	2.88
6	197	0.99	0.08	0.82	1.41	1.01	0.13	0.76	1.67
7	186	0.99	0.08	0.81	1.31	1.00	0.13	0.63	1.65
8	185	0.98	0.08	0.76	1.35	0.99	0.14	0.65	1.84
HS	108	0.99	0.09	0.85	1.31	1.00	0.14	0.82	1.67
Mathematics									
3	242	0.97	0.07	0.84	1.40	1.00	0.21	0.73	2.97
4	252	0.99	0.10	0.86	1.80	1.04	0.21	0.84	2.18
5	233	0.98	0.08	0.84	1.52	1.01	0.23	0.78	3.41
6	226	0.98	0.08	0.83	1.35	1.00	0.16	0.76	1.92
7	208	0.98	0.08	0.86	1.67	1.03	0.33	0.68	5.19
8	218	0.98	0.09	0.79	1.88	1.01	0.21	0.73	3.41
HS	118	0.99	0.08	0.79	1.47	1.00	0.14	0.69	1.83

Correlations among instructional area (IA) scores are presented in Table 4.7 and Table 4.8 for reading and mathematics, respectively. The correlations range from 0.94 to 0.98 in reading and 0.85 to 0.97 in mathematics, indicating strong correlations between IAs across subjects and grades. Note that the correlation coefficients in Tables 4.7 and 4.8 are disattenuated. The correlations have been fully corrected for measurement error in both variables (instructional areas).

Table 4.7. Correlations Among Instructional Area Scores—Reading

Grade		IA1	IA2	IA3
3	IA1	1.00	0.98	0.95
	IA2		1.00	0.97
	IA3			1.00
4	IA1	1.00	0.97	0.96
	IA2		1.00	0.97
	IA3			1.00
5	IA1	1.00	0.96	0.94
	IA2		1.00	0.96
	IA3			1.00
6	IA1	1.00	0.97	0.94
	IA2		1.00	0.96
	IA3			1.00
7	IA1	1.00	0.98	0.95
	IA2		1.00	0.97
	IA3			1.00
8	IA1	1.00	0.97	0.95
	IA2		1.00	0.97
	IA3			1.00
HS	IA1	1.00	0.96	0.94
	IA2		1.00	0.95
	IA3			1.00

Note. The correlations provided are disattenuated. IA1 = Literary Text, IA2 = Informational Text, IA3 = Vocabulary

Table 4.8. Correlations Among Instructional Area Scores—Mathematics

Grade		IA1	IA2	IA3	IA4
3	IA1	1.00	0.92	0.95	0.91
	IA2		1.00	0.94	0.90
	IA3			1.00	0.93
	IA4				1.00
4	IA1	1.00	0.95	0.92	0.85
	IA2		1.00	0.96	0.87
	IA3			1.00	0.89
	IA4				1.00
5	IA1	1.00	0.93	0.90	0.89
	IA2		1.00	0.93	0.89
	IA3			1.00	0.88
	IA4				1.00
6	IA1	1.00	0.95	0.88	0.90
	IA2		1.00	0.92	0.93
	IA3			1.00	0.90

Grade		IA1	IA2	IA3	IA4
	IA4				1.00
7	IA1	1.00	0.97	0.94	0.95
	IA2		1.00	0.95	0.95
	IA3			1.00	0.92
	IA4				1.00
8	IA1	1.00	0.96	0.94	0.96
	IA2		1.00	0.92	0.93
	IA3			1.00	0.92
	IA4				1.00
HS	IA1	1.00	0.97	0.95	0.91
	IA2		1.00	0.95	0.93
	IA3			1.00	0.94
	IA4				1.00

Note. The correlations provided are disattenuated. For grades 3–5: IA1 = Operations and Algebraic Thinking, IA2 = Number and Operations, IA3 = Measurement and Data, IA4 = Geometry; For grades 6–8 & HS: IA1 = Operations and Algebraic Thinking, IA2 = The Real and Complex Number Systems, IA3 = Geometry, IA4 = Statistics and Probability

4.3.3. Unidimensionality

The unidimensionality assumption is that items on the test measured only one latent trait. It can be assessed by examining the model fit. Essentially, if the model fit is not adequate, then the unidimensional assumption is not tenable. The specific steps taken and criteria to assess model fit are discussed in detail in the previous section. The results indicate that the unidimensionality assumption holds for most tests.

4.4. Scaling

A scale can be established through different methods (Kolen & Brennan, 2004). The fix two cut score method was selected because it eases the use and interpretation of score and achievement levels. This list shows the steps for implementing this method:

1. Maine DOE determines:
 - a. the number of achievement levels,
 - b. the initial scale score range, and
 - c. two fixed cut scores across grades and content areas.
2. Cut scores are obtained from the standard setting meeting. Note that the recommended cut scores are approved by the Commissioner of Education.
3. The equations below are used to derive equating constants.
4. The lowest and highest obtainable scores (LOSS & HOSS) of the scale are finalized.

Puhan & Dorans (2018) was consulted when determining the scale properties. Relevant key points considered were:

1. The mean score centers around the midpoint of the scale in order to maximize the longevity of the scale.
 - a. Because the fix two cut scores method is used, the *At State Expectations* level cut score should be centered around the midpoint of the scale.

2. The range of scores is wide enough to accommodate population shift. In other words, the number of score units preserves the score differentiation but does not yield unjustified differentiation.
 - a. Puhan and Dorans (2018) recommends that the number of scale units is similar to the raw score points. However, empirical data shows that this approach may cause many scale scores to be rounded to the same values or truncated to LOSS/HOSS.
 - b. Instead, the number of theta values within (-10, 10) one decimal point is used to estimate the number of scale points needed. This method yields 200 score units.

There are four achievement levels defined for the Maine scale: *Well Below State Expectations*, *Below State Expectations*, *At State Expectations*, and *Above State Expectations*. The two fixed cuts are set at the *At State Expectations* and *Above State Expectations* levels. Table 4.9 presents the scaling constants, scale score cuts, and LOSS/HOSS. It is worth noting that only summative items are included in the calculation of state summative scale scores. The summative item counts are 27 for grades 3–8 and 30 for the second year of high school. Table 4.10 and Table 4.11 show the frequency of the summative scale scores for reading and mathematics, respectively.

Table 4.9. Maine Grade-Level Scale Properties

Grade	Scaling Constants		Scale Score Cuts			Range	
	Intercept	Slope	<i>Below State Expectations</i>	<i>At State Expectations</i>	<i>Above State Expectations</i>	LOSS	HOSS
Reading							
3	1507.14	11.90	1483	1500	1525	1400	1600
4	1503.75	12.50	1486	1500	1525	1400	1600
5	1505.26	13.16	1487	1500	1525	1400	1600
6	1505.26	13.16	1486	1500	1525	1400	1600
7	1502.63	13.16	1483	1500	1525	1400	1600
8	1500.00	12.50	1484	1500	1525	1400	1600
HS	1501.56	15.63	1489	1500	1525	1400	1600
Mathematics							
3	1511.25	12.50	1486	1500	1525	1400	1600
4	1507.00	10.00	1488	1500	1525	1400	1600
5	1502.08	10.42	1484	1500	1525	1400	1600
6	1501.14	11.36	1481	1500	1525	1400	1600
7	1505.44	10.87	1482	1500	1525	1400	1600
8	1504.35	10.87	1484	1500	1525	1400	1600
HS	1523.22	17.86	1489	1500	1525	1400	1600

Table 4.10. Summative Scale Score Frequency Table—Reading

Scale Score	Grade						
	3	4	5	6	7	8	HS
1400							
...							
1426		1					

Scale Score	Grade						
	3	4	5	6	7	8	HS
1428			1				
1436						1	
1438		1					
1444							1
1447			1				
1451			1			2	
1452	1						
1453		4	1		1	3	
1454		1			3	2	
1455		1				1	
1456		3				12	
1457			9	1	3		
1458	3	10	2	3	4	9	
1459		5	7	2	1	5	1
1460		5				14	
1461	2	20	9	3	7	31	2
1462	14		10	3	6		4
1463	1	16	23	12	6	15	
1464	5	23	15	8	25	48	3
1465	41	26				30	
1466		29	39	16	21	51	8
1467	2		34	19	41		18
1468	24	53	63	25	31	49	
1469	63	26				80	31
1470	35	66	50	34	64	56	36
1471	98	53	78	40	69	126	
1472			70	53	96		52
1473	30	75				96	95
1474	139	67	93	61	92	148	
1475	62	99	97	65	125	126	101
1476	175	100	117	77	113	125	
1477	90						111
1478		104	109	73	127	135	170
1479	213	111	145	99	134	153	
1480	163	145	162	117	143	142	210
1481	194	111				170	214
1482	203		152	142	147		
1483	164	140	166	132	152	148	236
1484		151	192	172	156	184	267
1485	194	181				202	

Scale Score	Grade						
	3	4	5	6	7	8	HS
1486	194	150	215	181	157	238	328
1487	186		194	190	173		308
1488	204	195	221	226	206	235	
1489	231	198	226	181	208	282	378
1490	227	209				328	
1491		232	208	205	233	318	345
1492	243		206	203	243		382
1493	244	250	210	223	206	296	
1494	259	234				259	395
1495	258	280	242	260	197	237	421
1496	280	291	256	272	207	221	
1497			287	293	220		488
1498	306	312				236	484
1499	274	301	274	299	238	218	
1500	312	314	270	396	241	236	525
1501	302	321	314	374	291	266	
1502	319						506
1503		296	329	404	387	283	509
1504	335	337	318	429	412	290	
1505	334	344	370	408	418	332	469
1506	320	304				318	472
1507	329		354	444	424		
1508	354	330	356	405	396	387	494
1509		308	323	360	422	378	472
1510	326	288				377	
1511	314	265	371	354	335	341	466
1512	295		333	371	358		
1513	270	309	403	397	366	369	440
1514	248	325	392	356	341	347	369
1515	260	349				277	
1516		324	381	352	334	290	350
1517	257		366	361	297		377
1518	300	340	320	348	347	259	
1519	256	324				225	304
1520	254	273	341	351	294	224	242
1521	245	273	306	320	294	237	
1522			307	292	249		207
1523	280	280				214	165
1524	237	243	272	266	226	182	
1525	227	172	245	235	202	225	160

Scale Score	Grade						
	3	4	5	6	7	8	HS
1526	216	164	254	172	200	163	
1527	192						155
1528		160	199	183	177	152	134
1529	174	128	187	140	154	183	
1530	150	125	142	115	145	119	75
1531	125	102				121	77
1532	81		101	109	124		
1533	98	98	88	81	129	89	72
1534		72	89	68	101	94	50
1535	61	70				50	
1536	48	50	87	89	84	71	70
1537	35		92	60	80		
1538	39	63	40	64	71	33	31
1539	24	29	60	49	55	41	32
1540	20	39				22	
1541		21	33	33	45	28	19
1542	22		29	35	34		5
1543	12	34	29	27	21	6	.
1544	4	16				16	16
1545	15	16	11	30	20	4	
1546	4	9	28	16	14		
1547			11	15	15		
1548	7	7				8	11
1549	1	11	12	8	9		
1550	2		6	13	3	1	
1551	1	7	12	5	6	1	
1552	1						
1553		4	3	8	4	2	5
1554			7	1	5		
1555		1		4	2		
1556		4					1
1557	1		3		3		
1558	1	1	1				
1559		1	1	1			1
1561			3	1			
1562			.	.	3		
1563			1	2	1		
1564		2	1	.			
1566				1			
1567				1			

Scale Score	Grade						
	3	4	5	6	7	8	HS
1571			1				
1576				1			
...							
1600							

Table 4.11. Summative Scale Score Frequency Table—Mathematics

Scale Score	Grade						
	3	4	5	6	7	8	HS
1400							
...							
1432							1
1436							1
1438	1						
1440				2			
1441						1	
1442				.	3		
1443	2			2			
1444	1			5			
1445				2			
1446	2			.			
1447				6	4		
1448	5			2	2		11
1449			1	1		3	
1450	5			11		1	1
1451	15			17	9		
1452		1		19	3		2
1453	1		5	15	2		
1454	6	2			17	6	
1455	20	4		16	9	1	8
1456	13	2	4	36			
1457		8	5	34	18		5
1458	36	4	7	36	8	26	
1459	22	3	11	33	24	6	3
1460	33	7	12	52	30	1	
1461	20	14	11	60	30	31	27
1462		21	20		52	9	40
1463	49	18		66	58	20	
1464	37	22	25	73	78	42	12
1465	57	32	42	69	84	20	
1466	43	39	37	106	107	94	121

Scale Score	Grade						
	3	4	5	6	7	8	HS
1467		42	43	106	107	39	
1468	61	28	47	124	102	93	39
1469	60	49	64	132			
1470	85	40	67	123	113	57	182
1471	74	50	89		137	125	36
1472		70	102	133	141	117	
1473	95	68	81	123	146	114	322
1474	112	70	84	150	152	148	
1475	115	81	110	148	139	150	44
1476	123	75	87	168	159	161	
1477		99	118	180	155	159	382
1478	147	118	134	145	187	198	
1479	137	85	152		158	197	104
1480	147	111	166	150	174	196	373
1481	141	132	174	169			
1482		129	200	201	176	237	477
1483	141	134	180	208	176	253	
1484	151	152	175	239	207	306	144
1485	160	182	196	239	195	330	
1486	188	180	198	191	225	307	483
1487		198	204		283	304	
1488	192	239		254	305	334	483
1489	206	215	194	244	353	328	494
1490	189	222	210	251	372	333	
1491	214	258	235	228	396	325	451
1492		200	227	233	360	315	
1493	202	235	233	273	338	294	490
1494	207	231	256	288			
1495	256	251	264	322	321	303	567
1496	269	252	332		321	315	563
1497		300	349	360	322	307	
1498	298	298	347	289	266	300	553
1499	298	296	296	338	285	324	
1500	337	290	314	346	221	285	536
1501	358	265	298	335	237	295	
1502		276	345	333	216	290	490
1503	411	273	297	281	226	283	
1504	369	259	310		235	249	456
1505	369	276	266	264	198	241	430
1506	387	278	287	283			

Scale Score	Grade						
	3	4	5	6	7	8	HS
1507		248	222	239	217	289	386
1508	362	269	235	238	220	247	
1509	344	223	243	232	196	211	333
1510	315	213	266	234	159	209	
1511	320	195	263	238	171	210	293
1512		168			136	171	
1513	333	173	228	197	159	182	244
1514	290	154	222	217	161	136	220
1515	274	160	207	176	146	127	
1516	297	149	215	167	140	144	219
1517		175	189	181	121	128	
1518	253	210	170	158	126	134	176
1519	281	191	164	145			
1520	231	187	161	125	99	123	174
1521	215	169	189		108	117	166
1522		158	137	121	93	84	
1523	196	126	131	109	82	85	166
1524	168	145	127	99	70	92	
1525	150	110	103	77	73	76	153
1526	142	111	87	75	68	56	
1527		105	66	60	81	62	147
1528	130	102	57	63	78	46	
1529	110	84	66		78	48	123
1530	112	75	57	59	63	53	110
1531	110	82	67	58			
1532		69	64	60	61	39	93
1533	101	62	56	33	47	38	
1534	101	41	56	42	38	30	92
1535	96	44	45	41	41	33	
1536	83	56	44	23	37	24	97
1537		41		27	35	28	
1538	96	59	41		31	22	97
1539	59	31	27	21	24	15	103
1540	75	25	28	17	34	19	
1541	47	40	30	25	22	6	94
1542		38	23	15	17	17	
1543	30	25	21	21	18	15	85
1544	44	24	29	23			
1545	31	14	25	14	11	12	55
1546	35	24	25		9	8	49

Scale Score	Grade						
	3	4	5	6	7	8	HS
1547		16	17	19	15	12	
1548	26	15	15	20	14	7	50
1549	23	14	24	14	11	13	
1550	24	4	13	19	11	9	41
1551	16	6	15	9	8	11	
1552		8	11	4	10	1	33
1553	22	6	12	9	7	5	
1554	19	2	7		8	2	32
1555	16	5	2	6	6	4	28
1556	19	4	3	6			
1557		7	3	5	8	2	22
1558	15		3	1	4	3	
1559	10	3	5	3	1	4	22
1560	9	1		1	4		
1561	7			3	1	1	20
1562	.	2			3		
1563	11	1	2		1		20
1564	5		1	1	2	1	14
1565	1				2		
1566	5	1	1		2	2	14
1567			3	2		1	
1568	3			1	1		12
1569	2	1					
1570	2				1		11
1571	1		1		1		14
1572				1	1		
1573				1			6
1575		1				1	6
1576						1	
1577							10
1579	2						6
1580	2			1			4
1581	1						
1582							4
1584		1					
1585	1						
1586					1		2
1588							1
1589				1			2
1591							2

Scale Score	Grade						
	3	4	5	6	7	8	HS
1595							1
1596							1
1600	1						5

The field test calibration results are shown in Table 4.12. There were 257 reading items and 239 mathematics items. After removing items due to either small sample size, not passing content review, or not passing psychometric review, 201 reading items and 157 mathematics items were promoted to become next year's operational items. In total, 358 items will be available for next spring's operational testing administration.

Table 4.12. Field Test Calibration Results

Grade	Promote to OP	Removal Stage			Total
		Sample Size	Content Review	Psychometric Review	
Total	358	3	84	51	496
Reading					
3	13			2	15
4	18		1		19
5	15	1	2	1	19
6	15		1	2	18
7	13		1	1	15
8	13		3	2	18
HS	114		38	1	153
Sub Total	201	1	46	9	257
Mathematics					
3	13		1	1	15
4	12		3		15
5	11		2	2	15
6	11	1	3		15
7	9		3	3	15
8	9	1	1	3	14
HS	92		25	33	150
Sub Total	157	2	38	42	239

Section 5: Technical Quality—Validity

Validity is defined by the *Standards for Educational and Psychological Testing* as “the degree to which evidence and theory support the interpretations of test scores for proposed uses of tests. Validity is, therefore, the most fundamental consideration in developing tests and evaluating tests” (AERA et al., 2014, p. 11). Validating a test score interpretation is not a quantifiable property but an ongoing process, beginning at initial conceptualization of the construct and continuing throughout the entire process of assessment development and implementation. Every aspect of an assessment development and administration provides evidence in support of (or a challenge to) the validity of the intended inferences about what students know based on their score, including design, content specifications, item development, test constraints, psychometric quality, standard setting, and administration.

As this technical report has progressed, it has covered the different phases of the testing cycle and provided different pieces of technical quality evidence. It provides relevant evidence and a rationale in support of test score interpretations and intended uses based on the *Standards*, which are considered to be “the most authoritative statement of professional consensus regarding the development and evaluation of educational and psychological tests” (Linn, 2006, p. 27). The validity argument begins with a statement of the assessment’s intended purposes, followed by the evidentiary framework, where available validity evidence is provided to support the argument that the test actually measures what it purports to measure (SBAC, 2016).

First, the Through Year Assessment went through psychometric analyses—such as test reliability, classification accuracy, conditional standard error of measurement (CSEM), test information, differential item function (DIF), and convergent validity check—and the results so far strongly support the reliability and validity claims of this assessment. In addition, the test-development process ensures validity of the intended test score interpretations provided through the scale score. Last but not least, this assessment is aligned to grade-level content, and test scores are suitable for use in accountability systems as a result of a robust development process to determine the test blueprint, passage and item specifications, and ALDs.

5.1. Intended Purposes and Validity Evidence Framework

The *Standards* describes validation as a process of constructing and evaluating arguments for the intended interpretation and use of test scores:

“A sound validity argument integrates various strands of evidence into a coherent account of the degree to which existing evidence and theory support the intended interpretation of test scores for specific uses. . .

Ultimately, the validity of an intended interpretation of test scores relies on all the available evidence relevant to the technical quality of a testing system” (AERA et al., 2014, pp. 21–22).

The *Standards* (AERA et al., 2014, pp. 13–19) outlines the following five main sources of validity evidence:

- Evidence based on test content
- Evidence based on response processes

- Evidence based on internal structure
- Evidence based on relations to other variables
- Evidence based on validity and consequences of testing

Evidence based on test design refers to traditional forms of content validity or content-related evidence. Evidence based on response processes refers to the cognitive process engaged in by students when answering test items, or the “evidence concerning the fit between the construct and the detailed nature of performance or response actually engaged in by test takers” (AERA et al., 2014, p. 15). Evidence based on internal structure refers to the psychometric analyses of “the degree to which the relationships among test items and test components conform to the construct on which the proposed test score interpretations are based” (AERA et al., 2014, p. 16). Evidence based on relations to other variables refers to traditional forms of criterion-related validity evidence, such as predictive and concurrent validity. Evidence based on validity and consequences of testing refers to the evaluation of the intended and unintended consequences associated with a testing program.

Table 5.1 lists the intended purposes of the test and presents an overview of the validity components evidenced in this technical report.

Table 5.1. Intended Test Purposes and Sources of Validity Evidence

Test Purpose	Sources of Validity Evidence			
	Test Content	Response Processes	Internal Structure	Relation to Other Variables
1. To report individual student achievement relative to the state-adopted content standards in reading and mathematics	✓	✓	✓	✓
2. To provide information to the public about school performance through the state’s Every Student Succeeds Act (ESSA) reporting system, the ESSA Dashboard (https://www.maine.gov/doe/dashboard)	✓	✓	✓	
3. To support school identification within the state’s ESSA compliant system of school identification and support	✓	✓	✓	
4. To provide a source of information for ongoing local program evaluation	✓	✓	✓	

5.2. Purposes and Evidence

5.2.1. Test Purpose 1

Purpose: To report individual student achievement relative to the state-adopted content standards in reading and mathematics

Sources of Validity Evidence Based on Test Content:

- Test blueprint, content specifications, and item specifications are aligned to the full breadth and depth of grade-level content, process skills, and associated cognitive complexity.

- Blueprint specifications are evaluated for each test event for regular and accommodated populations. The evaluations are performed prior to test administration by simulation and then again following test administration.
- Tests are linked to the Maine Learning Results by the incorporation of the CCSS into item- and test-development specifications.
- Bias is minimized through Universal Design and accessibility resources.
- The item pool and item-selection procedures adequately support the test design.
- Operational computer adaptive test events meet all blueprint constraints, both for the general student population and for students taking accommodated test forms.
- Relevant sections within this report: 2, 3, 7, 8

Sources of Validity Evidence Based on Response Processes:

- Item-development and quality-control processes include screening and reviewing field test items for potential construct-irrelevant difficulty due to bias against demographic groups.
- The item types used in the assessment require response processes specified in the CCSS.
- The standard setting process relies on stakeholder judgments about proficiency based on student responses to, and the response processes elicited by, test items.
- Statistics regarding item response times show that, in general, more time is required as the associated achievement level of the items increases.
- Relevant sections within this report: 2, 7, 8

Sources of Validity Evidence Based on Internal Structure:

- ALDs were developed in consultation with committees of Maine educators with a secondary goal of providing information to all Maine educators.
- Achievement levels were set consistent with best practices through the embedded standard setting procedures.
- The assessment supports precise measurement and consistent classification to support analysis and reporting of scores.
- Item-calibration results support good fit of the test model and intended internal structure of the measurement construct.
- Differential item functioning (DIF) analysis was completed for all items across identifiable subgroups of students.
- Tests reliably measure on a scale that is established by achievement levels at every grade and reliably classify students into the achievement levels.
- Relevant sections within this report: 2, 3, 4, 6, 8

Sources of Validity Evidence Based on Relations to Other Variables:

- Claims for evidence based on relations to other variables are often supported by correlations. The correlations of reading summative scores and reading interim scores, reading summative scores and science summative scores, mathematics summative scores and mathematics interim scores, and mathematics summative scores and science summative scores indicate the strength and direction of the relationship between two variables and supports the validity of the test score as a measure of that construct.
- Convergent validity is assessed by comparing the agreement between two measures of the same construct, while discriminant validity is assessed by comparing the agreement between two measures that claim to measure different constructs. These types of

evidence help researchers and scientists make informed claims about the relationships between variables and the validity of their tests and measurements.

- Relevant section within this report: 8

5.2.2. Test Purpose 2

Purpose: To provide information to the public about school performance through the state's Every Student Succeeds Act (ESSA, 2015) reporting system, the ESSA Dashboard

Sources of Validity Evidence Based on Test Content:

- Test content is aligned with the reporting requirements of Maine's ESSA Dashboard.
- Bias is minimized through Universal Design and accessibility resources.
- Blueprint, passage specifications, and item specifications are aligned to grade-level content, process skills, and associated cognitive complexity.
- The item pool and item-selection procedures adequately support the test design.
- Reporting categories align with the structure of Maine's standards to support the interpretation of the test results.
- Relevant sections within this report: 2, 5, 7, 8

Sources of Validity Evidence Based on Response Process:

- Blueprint, passage specifications, and item specifications are aligned to grade-level content, process skills, and associated cognitive complexity.
- Achievement levels were set consistent with best practices.
- Relevant sections within this report: 2, 4, 7

Sources of Validity Evidence Based on Internal Structure:

- The assessment supports precise measurement and consistent classification to support analysis and reporting of longitudinal data.
- Reporting categories align with the structure of Maine's standards to support the interpretation of the test results.
- Achievement levels were set consistent with best practices.
- Item-calibration results support good fit of the test model and intended internal structure of the measurement construct.
- Differential item functioning (DIF) analysis was completed for all items across identifiable subgroups of students.
- Relevant sections within this report: 2, 3, 4, 6

5.2.3. Test Purpose 3

Purpose: To support school identification within the state's ESSA compliant system of school identification and support

Sources of Validity Evidence Based on Test Content:

- Maine's model of school support emphasizes the importance of measurement for academic achievement and progress of English language arts and mathematics.
- Blueprint, passage specifications, and item specifications are aligned to grade-level content, process skills, and associated cognitive complexity.
- Reporting categories align with the structure of Maine's standards to support the interpretation of the test results.
- Relevant sections within this report: 2, 7

Sources of Validity Evidence Based on Response Process:

- Bias is minimized through Universal Design and accessibility resources.
- Blueprint, passage specifications, and item specifications are aligned to grade-level content, process skills, and associated cognitive complexity.
- Achievement levels are vertically articulated.
- Relevant sections within this report: 2, 3, 4, 6, 7

Sources of Validity Evidence Based on Internal Structure:

- The assessment supports precise measurement and consistent classification to support analysis and reporting of longitudinal data.
- Achievement levels are vertically articulated.
- Item-calibration results support good fit of the test model and intended internal structure of the measurement construct.
- Differential item functioning (DIF) analysis was completed for all items across identifiable subgroups of students.
- Relevant sections within this report: 2, 3, 4, 6

5.2.4. Test Purpose 4

Purpose: To provide a source of information for ongoing local program evaluation

Sources of Validity Evidence Based on Test Content:

- Reporting categories align with the structure of Maine’s standards to support the interpretation of test results.
- Relevant sections within this report: 2, 8

Sources of Validity Evidence Based on Response Process:

- Bias is minimized through Universal Design and accessibility resources.
- Blueprint, passage specifications, and item specifications are aligned to grade-level content, process skills, and associated cognitive complexity.
- ALDs were developed in consultation with committees of Maine educators with a secondary goal of providing information to all Maine educators.
- Relevant sections within this report: 2, 3, 4, 6, 7, 8

Sources of Validity Evidence Based on Internal Structure:

- The assessment supports precise measurement and consistent classification for all students.
- ALDs were developed in consultation with committees of Maine educators with a secondary goal of providing information to all Maine educators.
- Scale is vertically articulated and supports longitudinal tracking of students’ academic progress.
- Achievement levels are vertically articulated.
- Relevant sections within this report: 2, 3, 4, 6

5.3. Interpretive Argument Claims

The test scores for the spring administration support their intended purposes. Claims to support this are documented in this technical report, as shown in Table 5.2.

Table 5.2. Interpretive Argument Claims—Evidence to Support the Essential Validity Elements

Argument	Tech Report Section(s)	Evidence
Tests and items were carefully developed to ensure that the test measured the Maine content standards.	2. Test Design and Content Development	Description of the development and review process for items, passages, and tests
Test score interpretations are comparable across students.	3.3. Constraint-Based Adaptive Test Engine 4. Item Statistics, Calibration, and Scaling 6. Technical Quality—Other	Simulations, analysis of test information, conditional standard errors of measurement, classification accuracy, and reliability estimates; blueprint comparability across students; item analysis, calibration, and scaling procedures
Test administrations were secure and standardized.	3. Administration and Security	Test administration procedures, including administration training, test accommodations, test security, and availability of help desk during testing window
Scoring was standardized and accurate.	6.4. Scoring 8.3. Reporting	Scoring rules and procedures; quality control of operational scoring
Achievement standards were rigorous and technically sound.	8. Achievement Standards and Reporting	Documentation of standard-setting procedures, including the methodology, identification of workshop participants, implementation process, and ALD development and validation
Assessments were accessible to all students and fair across student subgroups.	2. Test Design and Content Development 3. Administration and Security 6. Technical Quality—Other 7. Inclusion of All Students	Accommodation policy and implementation, sensitivity review, availability of translations, and DIF analyses
Assessments are positively correlated to other test scores.	8. Achievement Standards and Reporting	Correlations of reading summative scores and reading interim scores, reading summative scores and science summative scores, mathematics summative scores and mathematics interim scores, and mathematics summative scores and science summative scores

5.4. Summary of Validity Arguments

The Through Year Assessment is designed to align with the Common Core State Standards (CCSS), Maine’s accountability standards. Test blueprints, item specifications, and passage specifications ensure comprehensive coverage of grade-level content and cognitive complexity. Universal Design principles and accessibility resources minimize bias, and blueprint constraints are verified through simulations and operational test events. Reporting categories reflect the CCSS, supporting meaningful interpretation of results for individual achievement and accountability reporting.

Item development and quality-control procedures include reviews for construct-irrelevant difficulty and bias. The assessment uses item types that elicit response processes consistent with CCSS expectations. Standard-setting activities incorporate stakeholder judgments based on student responses, ensuring that the cognitive processes engaged during testing align with the intended constructs. Additionally, response time analyses show patterns consistent with item difficulty, further supporting validity.

Psychometric analyses confirm that the internal structure of the test supports the intended measurement model. Item calibration results demonstrate good fit, and achievement levels were established through rigorous standard-setting procedures. The assessment provides precise measurement and reliable classification across achievement levels with vertically articulated scales for longitudinal tracking. Differential item functioning (DIF) analyses ensure fairness across student subgroups, and reliability estimates, classification accuracy statistics, and conditional standard errors of measurement further support structural validity.

The assessment demonstrates expected relationships with other measures of academic achievement, including positive correlations of reading summative scores and reading interim scores, reading summative scores and science summative scores, mathematics summative scores and mathematics interim scores, and mathematics summative scores and science summative scores (see Section 8.3.4). These findings support convergent validity and reinforce the argument that the test measures constructs consistent with other established indicators of student performance. Discriminant validity checks also confirm that the test does not overlap with unrelated constructs, strengthening its interpretive claims.

The assessment fulfills its intended purposes, including individual achievement reporting, public accountability under the Every Student Succeeds Act (ESSA, 2015), school identification for support systems, and local program evaluation. Secure and standardized administration procedures, accurate scoring, and fair accessibility for all students are documented. Achievement standards were set using best practices, and the consequences of testing—such as informing school improvement and policy decisions—are aligned with educational goals. Collectively, these elements provide a coherent validity argument that integrates evidence across all phases of test development and implementation.

Section 6: Technical Quality—Other

The *Standards for Educational and Psychological Testing* refers to reliability as the “consistency of scores across replications of a testing procedure” (AERA et al., 2014, p. 33). The level of reliability/precision of scores has implications for validity. In other words, scores must be consistent and precise enough to be useful for their intended purposes. If scores are to be meaningful, tests should produce stable scores if the same group of students were to take the same test repeatedly without any fatigue or memory of the test. In addition, the range of certainty around the scores should be small enough to support educational decisions. The reliability/precision of the assessment was examined through analyses of measurement error under simulated and operational conditions, as follows:

- Marginal reliability for adaptive tests
- Cronbach’s alpha and standard error of measurement (SEM) for fixed forms
- Classification accuracy

Combined, these data provide several ways of looking at the reliability of student scores on a test. Classification accuracy provides important information related to achievement level classifications. These are of particular interest in the context of state accountability requirements.

6.1. Reliability

6.1.1. Marginal Reliability for Adaptive Tests

Traditional reliability coefficients from classical test theory consider individual items and depend on all test takers to take common items; however, in a CAT, different students receive different items. Therefore, the marginal reliability coefficient for the CAT administration was calculated. Samejima (1994) recommends the marginal reliability coefficient because it uses test information (e.g., variance of estimated theta and SEM) to estimate the reliability of student scores:

$$\text{Marginal Reliability} = \frac{\text{var}(\hat{\theta}) - \sigma^2}{\text{var}(\hat{\theta})}$$

where σ is defined as:

$$\sigma = E\{[I(\theta)]^{-1/2}\}$$

Table 6.1 and Table 6.2 present the overall error of estimated theta and test reliability for the grades 3–8 and second year of high school adaptive tests. Each table includes the average number of items administered, the standard deviation (SD) of the estimated theta, the mean conditional standard error of measurement (CSEM), and the marginal reliability coefficient. The SD of estimated theta and mean SEM are relatively small, and the marginal reliability of the overall scores is 0.88 or higher for reading and 0.90 or higher for mathematics. These results indicate that, overall, the score precision is reasonable: the overall mean SEM values were approximately 0.40, while the reliability estimates are consistent with the guidelines for reliability in a graduation test (Phillips & Camara, 2006). Additional reliability data by student group is provided in Appendix K.

Table 6.1. Reliability Statistics—Reading

Grade	Average # Items	SD of Estimated Theta	Mean SEM	Reliability
3	27	1.41	0.39	0.92
4	27	1.36	0.38	0.92
5	27	1.32	0.39	0.91
6	27	1.21	0.38	0.90
7	27	1.31	0.36	0.93
8	27	1.40	0.39	0.92
HS	30	0.94	0.32	0.88

Table 6.2. Reliability Statistics—Mathematics

Grade	Average # Items	SD of Estimated Theta	Mean SEM	Reliability
3	27	1.57	0.40	0.94
4	27	1.76	0.39	0.95
5	27	1.70	0.40	0.94
6	27	1.68	0.40	0.94
7	27	1.72	0.40	0.95
8	27	1.54	0.40	0.93
HS	30	1.15	0.36	0.90

6.1.2. Classification Accuracy

Classification accuracy is a measure of how accurately test scores place students into reporting category levels. It refers to the agreement between the actual classifications using observed cut scores and true classifications based on known true cut scores. It is common to estimate classification accuracy by using a psychometric model to find true scores corresponding to observed scores. The likelihood of inaccurate placement depends on the amount of error associated with scores, especially those nearest cut points.

Classification accuracy was calculated as follows (SBAC, 2016):

1. For each student, a normal distribution was constructed, with means equal to the scale score estimate and standard deviation equal to the SEM as a plausible true score distribution.
2. For each student, the proportion of that normal distribution that fell within each achievement level was calculated.
3. Within the groups of students assigned to a particular achievement level (Level 4, 3, 2, or 1 for the overall score), the sums of the proportions over students were computed. This provided estimates of the number of students whose true score falls within a level for each assigned achievement level. These sums were then expressed as a proportion of the total sample (i.e., expected proportion).
4. With the table of expected proportions, correct classification rates were then defined. This is the proportion of students whose true classification agrees with the assigned level among the subset of students with that assigned level.
5. The overall classification rate is the sum of the proportions of students whose true score level agrees with the assigned level divided by the total proportion of students assigned to a level.

Table 6.3 and Table 6.4 present the respective reading and mathematics classification accuracy results by grade and achievement level for grades 3–8 and the second year of high school. Overall classification accuracy ranges from 0.76 to 0.87.

Table 6.5 presents the classification *accuracy* results by grade at each achievement level and each cut for grades 3–8 and the second year of high school. Overall classification accuracy ranges from 0.76 to 0.87. The classification accuracy at each achievement level ranges from 0.58 to 0.93, whereas the classification accuracy at each cut ranges from 0.85 to 0.97. Additional classification accuracy data by student group is provided in Appendix L.

Table 6.6 presents the classification *consistency* results by grade at each achievement level and each cut for grades 3–8 and the second year of high school. Overall classification consistency ranges from 0.69 to 0.81. The classification consistency at each achievement level ranges from 0.48 to 0.85, whereas the classification consistency at each cut ranges from 0.85 to 0.97. Note that the lower levels of classification accuracy and consistency for the second year of high school tests in Spring 2025 will be improved in Spring 2026 due to the increased item pool.

Table 6.3. Classification Accuracy by Achievement Level—Reading

Grade	Achievement Level	N	Prop.	Expected Proportion ^a				Class. Acc.	Overall Class. Acc.
				L1	L2	L3	L4		
3	<i>Well Below State Expectations</i>	1,561	0.12	0.82	0.18	0.00	0.00	0.82	0.81
	<i>Below State Expectations</i>	6,147	0.49	0.10	0.77	0.12	0.00	0.77	
	<i>At State Expectations</i>	3,264	0.26	0.00	0.09	0.85	0.06	0.85	
	<i>Above State Expectations</i>	1,558	0.12	0.00	0.00	0.23	0.77	0.77	
4	<i>Well Below State Expectations</i>	1,406	0.12	0.86	0.14	0.00	0.00	0.86	0.82
	<i>Below State Expectations</i>	6,147	0.52	0.12	0.73	0.15	0.00	0.73	
	<i>At State Expectations</i>	2,652	0.22	0.00	0.09	0.85	0.06	0.85	
	<i>Above State Expectations</i>	1,627	0.14	0.00	0.00	0.19	0.81	0.81	
5	<i>Well Below State Expectations</i>	1,776	0.14	0.87	0.12	0.00	0.00	0.87	0.82
	<i>Below State Expectations</i>	6,426	0.52	0.17	0.69	0.14	0.00	0.69	
	<i>At State Expectations</i>	2,324	0.19	0.00	0.08	0.86	0.05	0.86	
	<i>Above State Expectations</i>	1,861	0.15	0.00	0.00	0.21	0.79	0.79	
6	<i>Well Below State Expectations</i>	1,567	0.13	0.86	0.14	0.00	0.00	0.86	0.82
	<i>Below State Expectations</i>	6,988	0.57	0.16	0.71	0.14	0.00	0.71	
	<i>At State Expectations</i>	2,533	0.21	0.00	0.09	0.86	0.05	0.86	
	<i>Above State Expectations</i>	1,157	0.09	0.00	0.00	0.20	0.80	0.80	
7	<i>Well Below State Expectations</i>	1,707	0.14	0.86	0.14	0.00	0.00	0.86	0.85
	<i>Below State Expectations</i>	6,432	0.54	0.12	0.78	0.10	0.00	0.78	
	<i>At State Expectations</i>	2,596	0.22	0.00	0.08	0.88	0.04	0.88	
	<i>Above State Expectations</i>	1,259	0.10	0.00	0.00	0.18	0.82	0.82	
8	<i>Well Below State Expectations</i>	1,429	0.12	0.89	0.11	0.00	0.00	0.89	0.83
	<i>Below State Expectations</i>	5,832	0.47	0.14	0.76	0.10	0.00	0.76	
	<i>At State Expectations</i>	3,254	0.26	0.00	0.08	0.87	0.05	0.87	
	<i>Above State Expectations</i>	1,778	0.14	0.00	0.00	0.21	0.79	0.79	
HS	<i>Well Below State Expectations</i>	914	0.07	0.85	0.15	0.00	0.00	0.85	0.79
	<i>Below State Expectations</i>	6,367	0.51	0.20	0.65	0.16	0.00	0.65	
	<i>At State Expectations</i>	2,893	0.23	0.00	0.12	0.84	0.04	0.84	

Grade	Achievement Level	N	Prop.	Expected Proportion ^a				Class. Acc.	Overall Class. Acc.
				L1	L2	L3	L4		
	<i>Above State Expectations</i>	2,196	0.18	0.00	0.00	0.24	0.76	0.76	

^a Level 1 = *Well Below State Expectations*, Level 2 = *Below State Expectations*, Level 3 = *At State Expectations*, and Level 4 = *Above State Expectations*.

Table 6.4. Classification Accuracy by Achievement Level—Mathematics

Grade	Achievement Level	N	Prop.	Expected Proportion ^a				Class. Acc.	Overall Class. Acc.
				L1	L2	L3	L4		
3	<i>Well Below State Expectations</i>	1,796	0.14	0.89	0.11	0.00	0.00	0.89	0.83
	<i>Below State Expectations</i>	6,110	0.49	0.13	0.71	0.15	0.00	0.71	
	<i>At State Expectations</i>	2,519	0.20	0.00	0.11	0.85	0.05	0.85	
	<i>Above State Expectations</i>	2,117	0.17	0.00	0.00	0.13	0.87	0.87	
4	<i>Well Below State Expectations</i>	1,360	0.11	0.90	0.10	0.00	0.00	0.90	0.84
	<i>Below State Expectations</i>	5,230	0.44	0.13	0.75	0.12	0.00	0.75	
	<i>At State Expectations</i>	2,997	0.25	0.00	0.10	0.86	0.04	0.86	
	<i>Above State Expectations</i>	2,270	0.19	0.00	0.00	0.14	0.86	0.86	
5	<i>Well Below State Expectations</i>	1,155	0.09	0.89	0.11	0.00	0.00	0.89	0.85
	<i>Below State Expectations</i>	5,486	0.44	0.09	0.80	0.11	0.00	0.80	
	<i>At State Expectations</i>	3,716	0.30	0.00	0.10	0.86	0.04	0.86	
	<i>Above State Expectations</i>	2,078	0.17	0.00	0.00	0.13	0.87	0.87	
6	<i>Well Below State Expectations</i>	881	0.07	0.91	0.09	0.00	0.00	0.91	0.83
	<i>Below State Expectations</i>	4,718	0.38	0.10	0.79	0.11	0.00	0.79	
	<i>At State Expectations</i>	4,327	0.35	0.00	0.13	0.83	0.04	0.83	
	<i>Above State Expectations</i>	2,345	0.19	0.00	0.00	0.15	0.85	0.85	
7	<i>Well Below State Expectations</i>	989	0.08	0.93	0.07	0.00	0.00	0.93	0.87
	<i>Below State Expectations</i>	3,737	0.31	0.08	0.84	0.08	0.00	0.84	
	<i>At State Expectations</i>	4,901	0.41	0.00	0.11	0.86	0.03	0.86	
	<i>Above State Expectations</i>	2,408	0.20	0.00	0.00	0.13	0.87	0.87	
8	<i>Well Below State Expectations</i>	727	0.06	0.90	0.10	0.00	0.00	0.90	0.84
	<i>Below State Expectations</i>	4,332	0.35	0.12	0.79	0.09	0.00	0.79	

Grade	Achievement Level	N	Prop.	Expected Proportion ^a				Class. Acc.	Overall Class. Acc.
				L1	L2	L3	L4		
	<i>At State Expectations</i>	4,725	0.38	0.00	0.11	0.86	0.03	0.86	
	<i>Above State Expectations</i>	2,505	0.20	0.00	0.00	0.15	0.85	0.85	
HS	<i>Well Below State Expectations</i>	1,681	0.14	0.83	0.16	0.01	0.00	0.83	0.76
	<i>Below State Expectations</i>	4,289	0.35	0.24	0.58	0.19	0.00	0.58	
	<i>At State Expectations</i>	3,118	0.25	0.01	0.15	0.80	0.04	0.80	
	<i>Above State Expectations</i>	3,301	0.27	0.00	0.00	0.11	0.89	0.89	

^a Level 1 = *Well Below State Expectations*, Level 2 = *Below State Expectations*, Level 3 = *At State Expectations*, and Level 4 = *Above State Expectations*.

Table 6.5. Classification Accuracy by Achievement Level and Cut

Grade	Accuracy at AL				Accuracy at Cut			Overall
	AL1	AL2	AL3	AL4	Cut1	Cut2	Cut3	
Reading								
3	0.82	0.77	0.85	0.77	0.93	0.89	0.92	0.82
4	0.86	0.73	0.85	0.81	0.93	0.89	0.92	0.82
5	0.88	0.69	0.86	0.79	0.93	0.90	0.92	0.82
6	0.86	0.71	0.86	0.80	0.94	0.89	0.92	0.82
7	0.86	0.78	0.88	0.82	0.94	0.91	0.93	0.85
8	0.89	0.76	0.87	0.79	0.93	0.91	0.93	0.83
HS	0.85	0.65	0.84	0.76	0.89	0.86	0.95	0.79
Mathematics								
3	0.89	0.71	0.85	0.87	0.94	0.89	0.94	0.83
4	0.90	0.75	0.86	0.86	0.93	0.90	0.95	0.84
5	0.89	0.80	0.86	0.87	0.94	0.89	0.96	0.85
6	0.91	0.79	0.83	0.85	0.92	0.88	0.96	0.83
7	0.93	0.84	0.86	0.87	0.93	0.91	0.97	0.87
8	0.90	0.79	0.86	0.85	0.91	0.90	0.97	0.84
HS	0.83	0.58	0.80	0.89	0.86	0.85	0.96	0.76

Note. AL1 = *Well Below State Expectations*, AL2 = *Below State Expectations*, AL3 = *At State Expectations*, and AL4 = *Above State Expectations*.

Table 6.6. Classification Consistency by Achievement Level and Cut

Grade	Consistency at AL				Consistency at Cut			Overall	Kappa
	AL1	AL2	AL3	AL4	Cut1	Cut2	Cut3		
Reading									
3	0.72	0.67	0.81	0.68	0.93	0.89	0.92	0.74	0.62
4	0.77	0.62	0.82	0.70	0.93	0.89	0.92	0.75	0.63
5	0.79	0.57	0.82	0.72	0.93	0.90	0.92	0.75	0.63
6	0.72	0.60	0.83	0.71	0.94	0.89	0.92	0.75	0.61
7	0.76	0.67	0.85	0.76	0.94	0.91	0.93	0.79	0.67
8	0.78	0.68	0.83	0.70	0.93	0.91	0.93	0.77	0.66
HS	0.73	0.52	0.80	0.65	0.89	0.86	0.95	0.71	0.56
Mathematics									
3	0.82	0.59	0.81	0.80	0.94	0.89	0.94	0.76	0.66
4	0.82	0.66	0.83	0.80	0.93	0.90	0.95	0.78	0.69
5	0.82	0.71	0.82	0.80	0.94	0.89	0.96	0.79	0.69
6	0.82	0.71	0.79	0.76	0.92	0.88	0.96	0.77	0.66
7	0.85	0.80	0.81	0.82	0.93	0.91	0.97	0.81	0.73
8	0.80	0.73	0.81	0.78	0.91	0.90	0.97	0.78	0.68
HS	0.75	0.48	0.72	0.84	0.86	0.85	0.96	0.69	0.57

Note. AL1 = Well Below State Expectations, AL2 = Below State Expectations, AL3 = At State Expectations, and AL4 = Above State Expectations.

6.1.3. Score Precision

Conditional standard error of measurement (CSEM) quantifies the degree of measurement error in scale score units, and its calculation is contingent on the student’s ability. This means that the test exhibits varying levels of error at different positions along the ability scale. In the context of an adaptive assessment, the CSEM will vary for identical scale scores. Therefore, it is imperative to provide averages in reporting.

In the context of item response theory (IRT), CSEMs for each scale score are defined as the reciprocal of the square root of the test information function (Hambleton & Swaminathan, 1985).

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

where $CSEM(\theta)$ is the IRT CSEM for a scale score, and $I(\theta)$ is the test information function.

CSEMs are especially useful for characterizing measurement precision with respect to score thresholds employed in decision-making, such as the cut score used to determine student proficiency on an assessment. Table 6.7 presents the CSEMs for the achievement level cut scores that demark the three cut scores on the Maine Through Year Assessment. It includes data on the numbers of students within ± 10 scale score points from these thresholds, the mean CSEMs for students in proximity to the cut scores, and the standard deviation (SD) of the CSEMs. In general, CSEMs of middle-range scale scores and cut scores are smaller than those at the two ends, indicating low measurement error and high score precision.

Table 6.7. CSEMs at the Cut Scores

Content Area	Grade	Below State Expectations			At State Expectations			Above State Expectations		
		N	Mean CSEM	SD	N	Mean CSEM	SD	N	Mean CSEM	SD
Reading	3	3,155	5.07	0.26	5,011	4.69	0.47	3,412	4.99	0.15
	4	3,079	4.98	0.22	4,945	4.67	0.48	3,496	4.86	0.36
	5	3,189	5.01	0.14	4,313	4.94	0.25	3,596	4.75	0.44
	6	2,751	5.03	0.23	4,969	4.68	0.47	3,392	4.81	0.40
	7	2,607	5.01	0.22	4,531	4.44	0.51	3,269	4.46	0.50
	8	3,486	5.01	0.19	4,969	4.81	0.41	3,104	4.96	0.22
	HS	4,451	5.25	0.45	5,948	5.01	0.15	2,363	5.36	0.51
Mathematics	3	3,069	5.00	0.07	5,179	4.93	0.25	2,966	5.00	0.04
	4	4,115	3.81	0.40	5,407	3.82	0.39	2,554	4.01	0.07
	5	3,528	4.00	0.03	5,827	3.98	0.15	2,414	4.00	0.06
	6	3,418	5.00	0.09	5,364	4.93	0.38	1,849	4.94	0.36
	7	4,452	4.01	0.13	5,099	3.99	0.13	1,613	3.99	0.18
	8	4,865	4.01	0.09	5,406	3.99	0.15	1,513	4.00	0.04
	HS	5,176	6.58	0.83	5,242	5.81	0.99	1,615	5.40	0.81

Table 6.8 presents the average CSEM by score decile, including the overall student ability distribution. A decile is similar to a percentile rank, with 10 ranks corresponding to the 10th, 20th, 30th, . . . 90th, and 100th percentile ranks. A higher SEM indicates a shallower pool of items suitable for students with these abilities. For instance, results indicate that the summative reading item pool is notably limited for students with very high abilities, while the mathematics item pool is shallower for students with very low and high abilities.

Table 6.8. CSEMs by Summative Score Decile

Grade	Overall	Proficiency Score Decile									
		1st	2nd	3rd	4th	5th	6th	7th	8th	9th	10th
Reading											
3	4.92	5.43	5.02	5.01	4.84	4.44	4.48	4.91	4.99	4.97	5.04
4	4.83	5.30	5.01	4.93	4.65	4.55	4.48	4.56	4.68	4.97	5.17
5	4.93	5.39	5.01	5.00	4.99	4.88	4.74	4.60	4.56	4.86	5.22
6	4.84	5.50	5.00	5.00	4.63	4.45	4.46	4.63	4.68	4.85	5.21
7	4.59	5.23	5.01	4.89	4.44	4.16	4.15	4.26	4.29	4.46	5.02
8	4.93	5.37	5.02	5.01	4.99	4.87	4.58	4.59	4.81	5.00	5.14
HS	5.26	6.13	5.32	5.02	5.00	5.01	5.01	5.01	5.02	5.03	5.94
Mathematics											
3	5.00	5.21	5.01	5.00	4.99	4.87	4.89	4.96	5.00	5.00	5.07
4	3.93	4.12	4.00	3.74	3.66	3.77	3.95	4.00	4.00	4.01	4.07
5	4.00	4.08	4.00	4.00	4.00	3.99	3.94	3.97	4.00	4.00	4.05
6	4.95	5.09	5.01	5.00	5.00	5.00	4.97	4.99	4.59	4.83	5.04
7	4.02	4.20	4.01	4.00	4.00	4.00	4.00	4.00	3.95	3.95	4.05
8	4.04	4.33	4.01	4.00	4.00	4.01	4.00	3.94	3.98	4.00	4.07
HS	6.19	7.90	7.01	7.00	6.60	5.90	5.57	5.37	5.51	5.41	5.69

Table 6.9 and Table 6.10 show the summary of CSEMs by instructional area (IA) in reading and mathematics, respectively. The scale score range for reading IAs is 100–320, and the scale score range for mathematics IAs is 100–350. The mean CSEMs range from 4.68 to 7.04 in reading and from 4.82 to 8.54 in mathematics. Note that the maximum CSEMs for reading were 250 in grade 6 IA2, grade 8 IA2, and HS IA3, indicating that a few students in these particular grades received zero raw scores in these IAs. While CSEMs of 250 may seem unreasonably high, they provide additional information on how students perform at these particular grades and IAs.

Table 6.9. Summary of CSEMs by Instructional Area—Reading

Grade	IA	Mean	SD	Min	Max
3	1	4.84	0.69	4.02	23.05
	2	5.92	1.09	4.67	25.04
	3	6.95	1.63	5.24	29.72
4	1	4.94	0.62	4.09	19.12
	2	5.53	0.83	4.49	23.63
	3	6.90	1.70	5.25	22.98
5	1	5.15	0.56	4.47	19.37

Grade	IA	Mean	SD	Min	Max
	2	5.34	0.94	4.21	24.87
	3	6.75	1.52	5.43	26.53
6	1	5.34	0.62	4.47	19.65
	2	5.00	3.22	4.02	250.00
	3	7.04	1.83	5.31	24.80
7	1	5.55	0.77	4.55	21.02
	2	4.71	0.59	3.90	19.24
	3	6.59	1.43	4.97	26.60
8	1	5.92	1.12	4.78	23.79
	2	4.83	2.33	4.11	250.00
	3	6.73	1.71	5.06	29.72
HS	1	5.44	1.13	4.36	25.71
	2	4.68	0.90	3.83	25.58
	3	6.37	2.66	4.85	250.00

Note. IA = Instructional Area; IA1 = Literary Text, IA2 = Informational Text, IA3 = Vocabulary

Table 6.10. Summary of CSEMs by Instructional Area—Mathematics

Grade	IA	Mean	SD	Min	Max
3	1	6.90	1.61	5.74	24.16
	2	5.80	0.88	4.90	21.73
	3	6.09	0.81	5.14	25.20
	4	6.81	1.41	5.80	29.73
4	1	6.97	1.69	5.74	24.89
	2	4.94	0.46	4.29	10.63
	3	7.17	1.21	5.98	27.39
	4	7.39	2.06	5.62	29.73
5	1	7.52	2.01	6.00	26.18
	2	4.82	0.48	4.06	22.69
	3	7.26	1.64	5.62	24.88
	4	7.59	1.99	5.56	26.41
6	1	6.72	1.47	5.52	23.55
	2	5.00	0.46	4.39	19.96
	3	7.29	1.48	6.15	27.38
	4	7.28	1.99	6.05	29.72
7	1	6.90	1.55	5.33	24.42
	2	5.34	0.58	4.60	21.84
	3	6.95	1.04	6.08	22.25
	4	6.87	1.63	5.50	29.72
8	1	5.04	0.58	4.42	19.54
	2	7.37	1.84	6.14	28.10
	3	7.14	1.63	6.05	22.77

Grade	IA	Mean	SD	Min	Max
	4	7.29	1.78	5.84	23.33
HS	1	4.98	0.84	4.18	21.46
	2	8.54	2.64	6.75	29.72
	3	6.22	1.44	5.19	29.50
	4	7.55	2.09	6.28	27.34

Note. IA = Instructional Area;

For grades 3–5: IA1 = Operations and Algebraic Thinking, IA2 = Number and Operations, IA3 = Measurement and Data, IA4 = Geometry;

For grades 6–8 & HS: IA1 = Operations and Algebraic Thinking, IA2 = The Real and Complex Number Systems, IA3 = Geometry, IA4 = Statistics and Probability

6.2. Fairness and Accessibility

Assessment fairness and accessibility are addressed through multiple approaches in this report. First, Universal Design is used to design the test and items (see Section 2.5.2). Second, accommodations are provided according to special student needs during administration and through various paper forms (Section 3.4 and Section 7). Third, analyses are conducted to evaluate item fairness and accessibility. While the first two approaches are qualitative methods, the last approach is quantitative. This section addresses the methods and results of these analyses.

Differential item functioning (DIF) is a statistical procedure that flags items for potential bias. The fundamental measurement assumption of DIF is that the probability of a correct response to a test item is a function of the item’s difficulty and the student’s ability. This function is expected to remain invariant to other characteristics unrelated to ability, such as gender and ethnicity. Therefore, if two students with the same ability respond to the same item, they are assumed to have an equal probability of answering the item correctly. To test this assumption, responses to items by students sharing an aspect of a characteristic (e.g., gender) are compared with responses to the same items by other students who share a different aspect of the same characteristic (e.g., males vs. females). The group representing students in a specific demographic group is referred to as the *focal* group. The group comprised of students from outside this group is referred to as the *reference* group.

When DIF is detected and the fundamental measurement assumption does not hold (i.e., students with the same ability in different groups of interest have different probabilities of correctly answering an item), the item is said to be functioning differently for the two groups. The presence of DIF in an item suggests that the item is functioning unexpectedly regarding the groups included in the comparison. The cause of the unexpected functioning is not revealed in a DIF analysis. It may be that item content is inadvertently providing an advantage or disadvantage to members of one of the two groups. Content experts who have special knowledge of the groups involved can often identify a cause of this type. DIF may also result from differential instruction closely associated with group membership.

Because fairness is a fundamental validity issue, it is essential that items be reviewed and assessed for DIF. Many methods for assessing DIF have been used and compared in conventional paper-pencil tests; however, DIF detection may be more important for a CAT than it is for traditional paper-pencil tests for two reasons (Zwick et al., 1994): First, items with DIF may be more consequential for the examinees because fewer items are administered in a CAT. Second, several potential sources of DIF may be introduced, such as differential computer

familiarity, facility, and anxiety. The difficulty of DIF analysis in a CAT is introduced by the fact that different sets of items are administered to different examinees. Therefore, the logistic regression (LR) procedure was applied to items that were administered in this CAT.

6.2.1. Logistic Regression (LR) DIF Method

The LR DIF procedure models item responses (for both dichotomous and polytomous items) as a function of group memberships, ability estimates, and their interaction. Testing for the presence of DIF based on logistic regression provides a model-based approach to identify uniform and nonuniform DIF. DIF is classified as uniform if the effect is constant; that is, uniform DIF exists when the difference in the probabilities of a correct answer for the two groups is the same at all ability levels. DIF is classified as nonuniform if the effect varies conditional on the ability level; that is, nonuniform DIF exists if the interaction between item-response function and group membership is disordinal.

The LR procedure compares the following three models (Fu & Monfils, 2016; Swaminathan & Rogers, 1990; Zumbo, 1999):

$$\text{Model 1: } \text{logit}(P) = \beta_0 + \beta_1 X + \beta_2 E$$

$$\text{Model 2: } \text{logit}(P) = \beta_0 + \beta_1 X + \beta_2 G + \beta_3 E$$

$$\text{Model 3: } \text{logit}(P) = \beta_0 + \beta_1 X + \beta_2 G + \beta_3 XG + \beta_4 E$$

where:

- P is the probability of a test taker answering an item incorrectly (for a dichotomous item) and the probability of getting an item score or lower (for a polytomous item).
- X is the criterion variable (typically an ability estimate).
- G is the group membership.
- E is a vector, including additional explanatory variables.
- β are the associated regression parameters for model k .

For both dichotomous and polytomous items, Models 1, 2, and 3 are also referred to as a no-DIF model, a uniform DIF model, and a nonuniform DIF model, respectively. The group estimates (β_2) are related to uniform DIF, and the interaction estimates (β_3) are associated with nonuniform DIF. Note that for a dichotomously scored item, the target probability that the LR estimates is the probability of answering an item incorrectly, which is different from the probability of answering an item correctly that many people may be accustomed to. Similarly, the target probability in the regression model for a polytomously scored item is the probability of obtaining an item score or below, to be consistent with that for a dichotomously scored item.

The item shows DIF if the modeled fit statistic is improved when group and interaction are added to the model, in order. To test the presence of nonuniform DIF, Model 2 and Model 3 are compared, using the likelihood ratio test with 1 degree of freedom (df) in chi-square distribution:

$$\chi^2 = [-2 \ln L(\text{Model2})] - [-2 \ln L(\text{Model3})]$$

Similarly, to test the presence of uniform DIF, Model 1 and Model 2 are compared, using the likelihood ratio test with 1 df:

$$\chi^2 = [-2 \ln L(\text{Model1})] - [-2 \ln L(\text{Model2})]$$

To test overall DIF (uniform DIF or nonuniform DIF), Model 1 and Model 3 are compared, using the likelihood ratio test with 2 df:

$$\chi^2 = [-2 \ln L(\text{Model1})] - [-2 \ln L(\text{Model3})]$$

The effect size is also used to avoid practically trivial but statistically significant results (French & Miller, 1996). Effect size is indicated by the difference of the Nagelkerke R^2 between two models (Gómez-Benito et al., 2009). Table 6.11 presents the DIF classification rules for the LR DIF procedure. These rules were confirmed to be consistent to the Mantel-Haenszel DIF classification rule for dichotomous items used by ETS (Fu & Monfils, 2016).

Table 6.11. LR DIF Categories

DIF Category	Level of DIF	Definition
A	Negligible	χ^2 test is not significant at 0.05 level or $\Delta R^2 < 0.035$.
B	Moderate	χ^2 test is significant at 0.05 level and $0.035 \leq \Delta R^2 < 0.070$.
C	Strong	χ^2 test is significant at 0.05 level and $\Delta R^2 \geq 0.070$.

Note. ΔR^2 is the Nagelkerke R^2 difference between two models.

6.2.2. DIF Results

DIF analysis is performed between a pair of demographic subgroups, typically defined by gender or ethnicity. For gender, male was used for the reference group, and female was used for the focal group; for ethnicity, white was used for the reference group, and a different minority subgroup was used for the focal group. More than 80% of students are white for the spring test. The large discrepancy in counts between reference group and focal group may cause statistical bias in estimates. Therefore, DIF was not conducted if the sample size for either group was less than 100. There are reduced counts of adaptive items meeting the minimum student counts required for DIF analyses due to the nature of adaptive item selection, while field test items were controlled to have required student counts and to be distributed across demographic groups.

Table 6.12 and Table 6.13 present the numbers of items identified for DIF for operational items and field test items, respectively. Considering that the Rasch model is applied (i.e., the same slope is assumed for all items), uniform DIF results are reported. The “+” sign next to the DIF category indicates that the item is in favor of the reference group, and the “-” sign indicates that the item is in favor of the focal group. As shown in the tables, most items were classified into Category A DIF, indicating negligible differential item functioning. Among the items eligible for DIF screening, the maximum proportion of items displaying Category B DIF did not exceed 1.5% per grade. Typically, item review is focused on items classified as exhibiting Category C DIF; a few C DIF items were found per grade in the item pool.

Table 6.12. DIF Analysis Results—Operational Items

Grade	Focal Group	Item Count by DIF Category					
		Total	A	B+	B-	C+	C-
Reading							
3	Female	190	189	1	–	–	–
	Black or African American	144	54	–	–	–	–
	Hispanic/Latino	144	43	–	–	–	–

Grade	Focal Group	Item Count by DIF Category					
		Total	A	B+	B-	C+	C-
	Two or More Races	144	45	-	-	-	-
	FRL	177	177	-	-	-	-
	ELL	53	53	-	-	-	-
	IEP	141	140	1	-	-	-
4	Female	228	228	-	-	-	-
	Black or African American	226	217	3	2	2	2
	Hispanic/Latino	228	215	-	4	4	5
	Two or More Races	-	-	-	-	-	-
	FRL	-	-	-	-	-	-
	ELL	-	-	-	-	-	-
	IEP	228	226	1	-	-	1
5	Female	233	233	-	-	-	-
	Black or African American	233	229	1	1	-	2
	Hispanic/Latino	233	225	1	1	5	1
	Two or More Races	-	-	-	-	-	-
	FRL	-	-	-	-	-	-
	ELL	-	-	-	-	-	-
	IEP	233	230	1	-	2	-
6	Female	202	202	-	-	-	-
	Black or African American	199	186	1	5	3	4
	Hispanic/Latino	202	185	2	5	2	8
	Two or More Races	-	-	-	-	-	-
	FRL	-	-	-	-	-	-
	ELL	-	-	-	-	-	-
	IEP	202	197	1	-	3	1
7	Female	208	208	-	-	-	-
	Black or African American	208	200	2	-	2	4
	Hispanic/Latino	208	201	1	1	2	3
	Two or More Races	-	-	-	-	-	-
	FRL	-	-	-	-	-	-
	ELL	-	-	-	-	-	-
	IEP	208	206	-	-	2	-
8	Female	236	236	-	-	-	-
	Black or African American	233	213	3	7	3	7
	Hispanic/Latino	234	212	2	8	4	8
	Two or More Races	-	-	-	-	-	-
	FRL	-	-	-	-	-	-
	ELL	-	-	-	-	-	-
	IEP	234	228	3	-	2	1
HS	Female	47	47	-	-	-	-
	Black or African American	47	47	-	-	-	-
	Hispanic/Latino	47	47	-	-	-	-

Grade	Focal Group	Item Count by DIF Category					
		Total	A	B+	B-	C+	C-
	Two or More Races	-	-	-	-	-	-
	FRL	-	-	-	-	-	-
	ELL	-	-	-	-	-	-
	IEP	47	47	-	-	-	-
Mathematics							
3	Female	341	334	4	3	-	-
	Black or African American	46	46	-	-	-	-
	Hispanic/Latino	12	12	-	-	-	-
	Two or More Races	11	11	-	-	-	-
	FRL	315	315	-	-	-	-
	ELL	53	53	-	-	-	-
	IEP	192	192	-	-	-	-
4	Female	361	354	4	3	-	-
	Black or African American	45	45	-	-	-	-
	Hispanic/Latino	7	7	-	-	-	-
	Two or More Races	9	9	-	-	-	-
	FRL	336	333	3	-	-	-
	ELL	66	66	-	-	-	-
	IEP	199	196	2	1	-	-
5	Female	352	339	6	6	1	-
	Black or African American	47	47	-	-	-	-
	Hispanic/Latino	13	13	-	-	-	-
	Two or More Races	13	13	-	-	-	-
	FRL	313	311	2	-	-	-
	ELL	57	57	-	-	-	-
	IEP	180	179	1	-	-	-
6	Female	336	322	7	5	1	1
	Black or African American	61	61	-	-	-	-
	Hispanic/Latino	16	16	-	-	-	-
	Two or More Races	17	17	-	-	-	-
	FRL	276	276	-	-	-	-
	ELL	70	70	-	-	-	-
	IEP	158	158	-	-	-	-
7	Female	311	303	4	3	-	1
	Black or African American	44	44	-	-	-	-
	Hispanic/Latino	22	22	-	-	-	-
	Two or More Races	24	24	-	-	-	-
	FRL	264	264	-	-	-	-
	ELL	58	58	-	-	-	-
	IEP	131	131	-	-	-	-
8	Female	307	302	2	3	-	-
	Black or African American	51	51	-	-	-	-

Grade	Focal Group	Item Count by DIF Category					
		Total	A	B+	B-	C+	C-
	Hispanic/Latino	24	24	-	-	-	-
	Two or More Races	22	22	-	-	-	-
	FRL	264	264	-	-	-	-
	ELL	55	55	-	-	-	-
	IEP	114	114	-	-	-	-
	HS	Female	158	154	3	1	-
	Black or African American	40	40	-	-	-	-
	Hispanic/Latino	40	40	-	-	-	-
	Two or More Races	40	40	-	-	-	-
	FRL	118	118	-	-	-	-
	ELL	38	38	-	-	-	-
	IEP	48	48	-	-	-	-

Note. FRL = Free and Reduced Lunch, ELL = English Language Learner, IEP = Individualized Education Plan

Table 6.13. DIF Analysis Results—Field Test Items

Grade	Focal Group	Item Count by DIF Category					
		Total	A	B+	B-	C+	C-
Reading							
3	Female	15	15	-	-	-	-
	Black or African American	15	15	-	-	-	-
	Hispanic/Latino	15	15	-	-	-	-
	Two or More Races	15	15	-	-	-	-
	FRL	15	15	-	-	-	-
	ELL	15	15	-	-	-	-
	IEP	15	15	-	-	-	-
4	Female	18	18	-	-	-	-
	Black or African American	16	16	-	-	-	-
	Hispanic/Latino	15	15	-	-	-	-
	Two or More Races	16	16	-	-	-	-
	FRL	18	18	-	-	-	-
	ELL	16	16	-	-	-	-
	IEP	18	18	-	-	-	-
5	Female	16	16	-	-	-	-
	Black or African American	16	16	-	-	-	-
	Hispanic/Latino	16	16	-	-	-	-
	Two or More Races	16	16	-	-	-	-
	FRL	16	16	-	-	-	-
	ELL	16	16	-	-	-	-
	IEP	16	16	-	-	-	-
6	Female	18	18	-	-	-	-
	Black or African American	17	17	-	-	-	-
	Hispanic/Latino	16	16	-	-	-	-

Grade	Focal Group	Item Count by DIF Category					
		Total	A	B+	B-	C+	C-
	Two or More Races	17	17	-	-	-	-
	FRL	18	18	-	-	-	-
	ELL	17	17	-	-	-	-
	IEP	18	18	-	-	-	-
7	Female	15	15	-	-	-	-
	Black or African American	15	15	-	-	-	-
	Hispanic/Latino	15	15	-	-	-	-
	Two or More Races	15	15	-	-	-	-
	FRL	15	15	-	-	-	-
	ELL	15	15	-	-	-	-
	IEP	15	15	-	-	-	-
8	Female	15	15	-	-	-	-
	Black or African American	14	14	-	-	-	-
	Hispanic/Latino	14	14	-	-	-	-
	Two or More Races	14	14	-	-	-	-
	FRL	15	15	-	-	-	-
	ELL	14	14	-	-	-	-
	IEP	15	15	-	-	-	-
HS	Female	153	152	-	1	-	-
	Black or African American	-	-	-	-	-	-
	Hispanic/Latino	-	-	-	-	-	-
	Two or More Races	-	-	-	-	-	-
	FRL	153	153	-	-	-	-
	ELL	-	-	-	-	-	-
	IEP	50	50	-	-	-	-
Mathematics							
3	Female	15	15	-	-	-	-
	Black or African American	14	14	-	-	-	-
	Hispanic/Latino	14	14	-	-	-	-
	Two or More Races	14	14	-	-	-	-
	FRL	15	15	-	-	-	-
	ELL	14	14	-	-	-	-
	IEP	15	15	-	-	-	-
4	Female	15	15	-	-	-	-
	Black or African American	15	15	-	-	-	-
	Hispanic/Latino	15	15	-	-	-	-
	Two or More Races	15	15	-	-	-	-
	FRL	15	15	-	-	-	-
	ELL	15	15	-	-	-	-
	IEP	15	14	1	-	-	-
5	Female	15	15	-	-	-	-
	Black or African American	15	15	-	-	-	-

Grade	Focal Group	Item Count by DIF Category					
		Total	A	B+	B-	C+	C-
	Hispanic/Latino	15	15	–	–	–	–
	Two or More Races	15	15	–	–	–	–
	FRL	15	15	–	–	–	–
	ELL	15	15	–	–	–	–
	IEP	15	15	–	–	–	–
6	Female	14	14	–	–	–	–
	Black or African American	14	14	–	–	–	–
	Hispanic/Latino	14	14	–	–	–	–
	Two or More Races	14	14	–	–	–	–
	FRL	14	14	–	–	–	–
	ELL	14	14	–	–	–	–
	IEP	14	14	–	–	–	–
7	Female	15	15	–	–	–	–
	Black or African American	15	15	–	–	–	–
	Hispanic/Latino	15	15	–	–	–	–
	Two or More Races	15	15	–	–	–	–
	FRL	15	15	–	–	–	–
	ELL	15	15	–	–	–	–
	IEP	15	15	–	–	–	–
8	Female	13	13	–	–	–	–
	Black or African American	13	13	–	–	–	–
	Hispanic/Latino	13	13	–	–	–	–
	Two or More Races	13	13	–	–	–	–
	FRL	13	13	–	–	–	–
	ELL	13	13	–	–	–	–
	IEP	13	13	–	–	–	–
HS	Female	150	149	1	–	–	–
	Black or African American	–	–	–	–	–	–
	Hispanic/Latino	–	–	–	–	–	–
	Two or More Races	–	–	–	–	–	–
	FRL	150	147	2	1	–	–
	ELL	–	–	–	–	–	–
	IEP	40	38	2	–	–	–

Note. FRL = Free and Reduced Lunch, ELL = English Language Learner, IEP = Individualized Education Plan

6.3. Full Achievement Continuum

It is important for an assessment to cover the full achievement continuum in order to provide reliable scores of the entire score range, or at least at the cut scores to provide higher classification accuracy. The summative item bank covers a wide range of difficulties, as shown in Table 4.5. This enables the summative assessment to effectively differentiate between lower- and higher-performing students. Most importantly, it increases accuracy in classifying students' achievement levels, especially for students just above or below the cut scores. The evidence on

CSEMs from Section 6.1.3 indicates the tests can accurately estimate ability across the full ability scale, especially at the middle range of the scale and around the cut scores.

6.4. Scoring

There are two scoring approaches to estimate student scores: number correct and pattern scoring. The number correct method uses student responses to determine student scores: correct vs. incorrect for dichotomous items and earned score points for polytomous items. This method yields a one-to-one correspondence between raw scores and scale scores. Pattern scoring not only considers student responses but also item difficulty in score decisions. Answering a difficult item correctly will yield a higher score than answering an easier item correctly; thus, when two students earn the same raw scores through different item sets, their scale scores may differ because of the difference in difficulty between the two sets of items. Consequently, pattern scoring yields multiple correspondences between raw scores and scale scores.

The goal of computer adaptive testing is to reach a desirable score precision across the student's ability range. Student ability estimates (thetas) are computed during test administration to select subsequent items that assist in obtaining reliable scores. Pattern scoring helps attain stable student ability estimates quicker than the number correct method because of the inclusion of item difficulty in estimation. Thus, it is typically used for an adaptive test.

6.4.1. Constructing the Maine Scale

Rationales and procedures for constructing the Maine Through Year Assessment are described in Section 4.4. Both literature and practical considerations play important parts in the procedures. The rationales and procedures were discussed with the TAC members. The TAC's feedback was also considered when determining the scale properties. Achievement levels established on the Maine scale score are determined by the standard setting meeting and approved by the Commissioner of Education (see Section 8).

6.4.2. Machine-Scored Items

The Maine Through Year Assessment has only machine-scored items. The item pool included technology-enhanced items and constructed-response items; however, those items typically have multiple correct answer keys. The keys have been evaluated, checked, and then hard coded into the database for scoring purposes. Calibration and validation of test item parameters were described in Section 4.2 and Section 4.3. Note that technology-enhanced items were excluded when constructing paper forms (including large print and braille forms) due to the limitations of the media.

6.4.3. Attemptedness Rule and Not-Tested Codes

Attemptedness for the Maine Through Year Assessment is defined as answering at least 25% of the summative items. With different test lengths across grades and content areas, a fixed value (7 items) is selected for all tests. Besides this attemptedness rule, there are also situations that could invalidate student test scores. Different Not-Tested Codes (NTC) are assigned to pinpoint different causes of score invalidation. Table 6.14 lists the various NTC codes. A student's Maine scale score and achievement level are not reported when the attemptedness threshold is not met or an NTC code is present.

Table 6.14. Available Not-Tested Codes

NTC Code	Description
INV	Invalid: Student's assessment was invalidated, such as due to a security breach.
EMW	Emergency Medical Waiver: Student was not assessed because of an approved emergency medical waiver.

6.5. Multiple Assessment Forms

An adaptive test has a large item pool in comparison with the number of items used in a fixed-form test. Items administered to individual students are selected according to the students' responses to prior items. Each student may have received a different set of items by the completion of the test. In other words, an adaptive test has multiple test forms by nature.

6.6. Multiple Versions of an Assessment

The Maine Through Year Assessment is mostly an adaptive test, but various paper accommodation forms are built for students with special needs. The number of students taking paper forms is not large enough for calibration. Instead, item parameters are derived from the adaptive test. The parameters are then used to derive scores for students who took paper forms. This approach makes the scores of the adaptive test and paper forms comparable.

6.7. Technical Analysis and Ongoing Maintenance

When planning the Spring 2025 assessment, test blueprint, test design, item development, specifications for CBE setup, and various psychometric analyses were considered. The test design, procedures, and methods documented in this report were applied to the Spring 2025 administration and will continue be used as guidelines for maintaining test consistency across administrations.

Section 7: Inclusion of All Students

Multiple guides were created for the Maine Through Year Assessment to explain the target population, supports, and accommodations for all students or specific populations, as well as guidance for test coordination and administration. The guides provided include:

1. *Maine Through Year Assessment Checklist Spring 2025 Administration*
2. *The Maine Through Year Assessment Coordinator Guide*
3. *The Maine Through Year Assessment Administration Guide*
4. *The Maine Through Year Assessment Proctor User Guide*
5. *The Maine Through Year Manage Online Testing Guide*
6. *The Maine Through Year Assessment User and Student Management Guide*
7. *The Maine Through Year Assessment Accessibility Guide*
8. *NWEA State Solutions: NWEA System and Technology Guide*

7.1. Testing Population

The Maine Through Year Assessment Coordinator Guide states that the Maine Through Year Assessment is designed for students in grades 3–8 and their second year of high school, with the exception of students with the most significant cognitive disabilities who have been found eligible for alternate assessments via the IEP Team Process. It is expected that approximately 99% of the student population participates in the Maine Through Year Assessment. The Every Student Succeeds Act (ESSA, 2015) requires that all students (who are eligible to test) participate in the state assessments.

7.2. Procedures for Including Students Who Utilize Accessibility Features

The Maine Through Year Assessment Coordinator Guide states that “All students are expected to participate in state assessments. No student, including students with disabilities, may be excluded from the state assessment and accountability system” (p. 13).

Three tiers of accessibility features have been developed to support the inclusion of all students, such as students with disabilities (SWDs): universal tools, designated supports, and accommodations (as described in Section 7.4).

7.3. Procedures for Including Multilingual Learners

In compliance with the Every Student Succeeds Act (ESSA, 2015) and state law on the inclusion of Multilingual Learners (MLs), *The Maine Through Year Assessment Coordinator Guide* states that “[School Administrative Units] should carefully consider the tools and resources utilized by MLs on a routine basis to access classroom instruction. These should be implemented as designated supports for the student during the assessment experience” (p. 14). Guidelines for the participation of newly arrived multilingual learners are also addressed in the guide.

7.4. Accommodations

Accommodations increase accessibility to a test by removing barriers without affecting the test construct. Accessibility is an important part of score validity, as student scores should represent the knowledge, skills, and abilities of the student. If a student cannot fully access the test, then the score cannot properly represent the individual’s achievement. Accessibility to the test was considered at different stages of test development and administration.

At the development stage: Universal Design was used to guide item development and style (see Section 2.5.2 for more details). Content and Bias Review and Data Review meetings checked for potential item bias through qualitative and quantitative methods.

In addition to the adaptive test, fixed-form standard print, large print, and braille forms were created for students with a documented need in an IEP or 504 Plan. During paper-based form creation, items were hand selected to ensure the blueprints were met at each grade level for each content area. Items were carefully sequenced and reviewed to avoid clueing within a grade level. The item types selected for the paper-based forms include multiple choice, multi-select, and composite (which uses elements of both multiple choice and multi-select).

Additionally, items do not include any art that is inappropriate for the visually impaired population. As a back-up, the braille vendor will reach out to NWEA if something cannot be brailled, which did not happen this year. The psychometric team provided statistical targets to the content team and reviewed and approved all selections to ensure that items on the paper forms were of similar difficulty, complexity, and compatibility to those selected by the constraint-based engine for the adaptive tests.

At the administration stage: Universal tools were provided within the test platform and accessible by all students. Students have the choice to use any of the available tools. Some of the universal tools are embedded in the online secure browser and do not require activation, such as answer eliminator, zoom, guideline, calculator for select math items, etc. Scrap/scratch paper is a nonembedded universal tool required to be provided to all students by the proctor. Information on the use of universal tools is not recorded.

Another tier of accessibility features is designated supports. Designated supports can be provided to students who meet the following two criteria:

1. An educational team with knowledge of the student's achievement has determined that the support is appropriate for the student.
2. The support is consistent with the student's routine instruction and assessment.

Text-to-Speech (TTS) is available as an embedded designated support that needs to be assigned within the assessment platform. Table 7.1 provides the numbers of students who were assigned TTS. Other designated supports that cannot be embedded in the online system are made available by the test administrator/proctor, such as small group or individual setting, mathematical supports, and bilingual word glossary.

In addition to the paper-based form accommodations, other accommodations include human reader, American sign language, scribe, calculator, and human reader for reading passages.

Refer to *The Maine Through Year Assessment Accessibility Guide* for more details regarding universal tools, designated supports, and accommodations.

Table 7.1. Numbers of Students Who Were Assigned TTS

Grade	Content Area	Number of Students
3	Reading	3,145
	Mathematics	3,609
4	Reading	2,864
	Mathematics	3,209
5	Reading	2,906
	Mathematics	3,157
6	Reading	2,148
	Mathematics	2,400
7	Reading	2,032
	Mathematics	2,241
8	Reading	1,911
	Mathematics	2,087
HS	Reading	534
	Mathematics	581




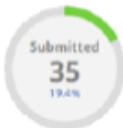
7.5. Monitoring Test Administration for Special Populations

Monitoring of the test administration is conducted in two ways: through the assessment administration and management system and through Maine DOE site visits.

7.5.1. Monitoring in Acacia

The Acacia system provides multiple pieces of information related to monitoring test status both during and after assessment. During the testing window, a testing status icon can be used to help proctors monitor student testing status with ease (Figure 7.1). After the testing window, the testing time marks at the item level can be analyzed to help understand the total test duration, time spent on each item, and any student test behavior related to testing time.

Figure 7.1. Monitoring Testing Status in Acacia

Icon	Assessment Status Icon Description
	<p>The Ready to Test icon displays the number and percentage of students who are enrolled and ready to take the assessment. It includes assessments in the Registered, Enrolled, and Ready to Test statuses. All assessments remaining in these statuses at the end of the assessment window are changed to Expired.</p>
 	<p>The In Progress icon displays the number and percentage of students actively testing. It includes assessments in the In Progress status only.</p> <p>The Alerts icon displays the number and percentage of students who have logged out and have not completed an assessment or have an enrollment hold. These students need test ticket login information to log back in and complete an assessment. This count includes assessments in the Inactive and Enrollment Hold statuses.</p> <p>Note: If any assessment registrations are in the Enrollment Hold status during the week before the assessment starts, contact NWEA Partner Support to resolve the hold.</p>
	<p>The Submitted icon displays the number and percentage of students who completed and submitted assessments. It includes assessments in the Submitted status only.</p>

7.5.2. Maine DOE Administration Monitoring and Support

The Maine DOE Assessment Team regularly provides three levels of technical assistance for all administrations of the Maine educational assessments. Level 1 technical assistance is available upon request to all Maine School Administrative Units (SAUs) and can be provided in the form of virtual/in-person meetings, on-site training, and support with developing resources. Level 1 technical assistance is provided year-round and is in high demand by the educator workforce. Level 2 technical assistance is assigned based on a 10-year, randomly selected cohort cycle to monitor the reliable and valid administration of state assessments across Maine. Level 3 technical assistance is the highest level of support for Maine SAUs and schools. The focus of Level 3 technical assistance is participation of eligible students and the secure and uniform administration of state assessments.

In spring 2025, 14 SAUs were randomly selected for Level 2 technical assistance. As part of the Level 2 technical assistance monitoring plan, schools participate in an initial virtual meeting with the Maine DOE Assessment Team during which SAU leadership receives the Assessment Observation Form shown in Figure 7.2. SAU leadership then conducts at least one local observation across the Maine Educational Assessments and submits notes to the Maine DOE Assessment Team via the Checklist Tool. Following the submission of the tool, the Maine DOE

Assessment Team conducts a follow-up meeting with the school to review the observation form as well as challenges and lessons learned from the administration. Of the observation forms completed by SAUs this year, eight were for the Maine Through Year Assessment across grades 3, 4, 5, 7, 8, and high school. Observation forms completed by SAU leadership indicated acceptable or exemplary conditions across all aspects of the assessment administration.

In spring 2025, 11 schools were identified for Level 3 technical assistance, the most intensive level of technical assistance, for the Maine Through Year Assessment. Two schools were identified due to assessment security concerns during the 2023–2024 school year, eight schools were identified due to low participation rates during the 2023–2024 school year, and one school was identified for both assessment security concerns and low participation rates during the 2023–2024 school year.

For all schools receiving Level 3 technical assistance, action plans are developed by the Maine DOE Maine Through Year Assessment Coordinator. Action plans include multiple elements, such as close monitoring of the completion of required assessment coordinator and proctor trainings, strategies for improving assessment participation based on local challenges, and on-site visits. On-site visits were conducted for three schools during the required fall administration of the Maine Through Year Assessment, and the Maine Educational Assessment Observation Form (shown in Figure 7.2) was completed. Two of the three observed schools were identified due to assessment security concerns. Afterward, a meeting was conducted with school leadership to review the findings of the observation and (for those schools identified for needing technical assistance due to security concerns) to determine the next steps for ensuring that local assessment security procedures align with state security expectations as outlined in the *Maine Assessment Security Handbook*.

Figure 7.2. 2025 Maine Educational Assessment Observation Form

School Name:	
Assessment Administrator:	Proctor/TA/AA(s):
Observer:	Subject:
Date of Observation:	Grade:

	Item	Code*	Comments
1	Instructional materials that may provide clues or answers are not visible in the room.		
2	The desks/tables are arranged with enough space between them to minimize opportunities to review each other's work.		
3	Desks/tables are clear of all materials except what is allowed in the assessment administrator manual.		
4	Electronic devices were collected or otherwise stored away and unavailable for student use.		

5	The Assessment Administrator read directions clearly, loudly, and exactly as printed in the Assessment Administration Manual.		
6	Students worked independently of each other.		
7	The assessment room was free of disruptions (talking, fire drills, intercom announcements).		
8	Booklets/tickets were distributed to and collected from the students individually by the Assessment Administrator/Proctor(s) and not passed by students.		
9	The Assessment Administrator answered only questions related to the directions.		
10	Students were provided a break individually, (where applicable) during an assessment session with close supervision.		
11	Students worked on appropriate sections of the assessment and did not return to or go forward to other sections.		
12	All students remained quiet as everyone completed the assessment session.		
13	Assessment tickets/booklets, answer documents, and scrap paper were never left unattended.		
14	The assessment room was supervised at all times.		
15	The Assessment Administrator/Proctor(s) were actively monitoring the room at all times.		
16	Assessment signs were posted on room doors (e.g., Do Not Disturb, Electronic Devices Not Allowed, Quiet Please Assessments in Progress).		

* Use Codes: NA = Not Applicable; 1 = Exemplary; 2 = Acceptable; 3 = Minor Issue; 4 = Major Issue; UO = Unable to Observe

Is this the TA's first time administering the assessment?

Yes

No

TA's level of confidence administering the assessment.

- High
- Neutral
- Low

Does the proctor/TA/AA feel they received sufficient training and support to administer the assessment?

- Yes
- No

If no, please explain.

Did you observe any students or did the specifically observed student complete the entire assessment?

- Yes
- No

If no, please provide a reason why the student or students did not complete the assessment. Please check all that apply.

- Student became ill and left the room
- Student became overwhelmed
- Student was dismissed
- Student left the room and did not return
- Student has an accommodation that allows taking breaks
- Student was administered the assessment administration over multiple days
- Student refused to complete the assessment
- Environmental disruption resulted in student not completing the assessment

Other reason, please describe.

Was the student(s) provided an opportunity to participate in a practice session?

- All students were provided the opportunity
- Some students were provided the opportunity
- None of the students were provided the opportunity

Were any of the students or the specifically observed student observed choosing the same answer repeatedly?

- Yes
- No

If yes, was it related to any of the following?

- Test content
- Test preparation
- Student characteristic
- TA/Proctor/AA behavior
- Environment
- Unknown

Were any of the students or the specifically observed student observed hurrying through the assessment?

Yes

No

If yes, was it related to any of the following?

Test content

Test preparation

Student characteristic

TA/Proctor/AA behavior

Environment

Unknown

Were any of the students observed using the universal tools provided in the assessment?

Yes

No

If yes, how did the student appear to be using the tool(s)?

Appropriately utilizing the tools

Trying the tool out

Playing around (tool appeared to be a distraction)

Other, please describe.

List any observed accommodations provided to students.

Please provide any insight, including specific topics for additional assessment training offered by the Maine Department of Education.

Did the assessment platform function as expected?

Yes

No

If no, please describe and include what type of device was used (e.g., iPad, Chromebook, Windows).

Section 8: Achievement Standards and Reporting

Achievement standards describe student performance across four levels: *Well Below State Expectations*, *Below State Expectations*, *At State Expectations*, and *Above State Expectations*. This section describes the procedures for defining achievement standards, setting achievement standards, and reporting.

8.1. State Adoption of Achievement Standards

The Maine Through Year Assessment (MTYA) program is Maine’s statewide system of summative assessments in reading and mathematics in grades 3–8 and the second year of high school that was first administered in Spring 2023. The Maine Department of Education (DOE) contracted with NWEA to design and develop the MTYA, and NWEA contracted with edCount LLC and Creative Measurement Solutions LLC to design and implement the alignment study and standard setting.

The MTYA standard setting design was a systematic approach grounded in principled assessment design (PAD). Under this design, the Achievement Level Descriptors (ALDs) shown in Table 2.16 were developed early in the test-development lifecycle to support domain definition (e.g., explication of the construct of interest), item development, and standard setting.

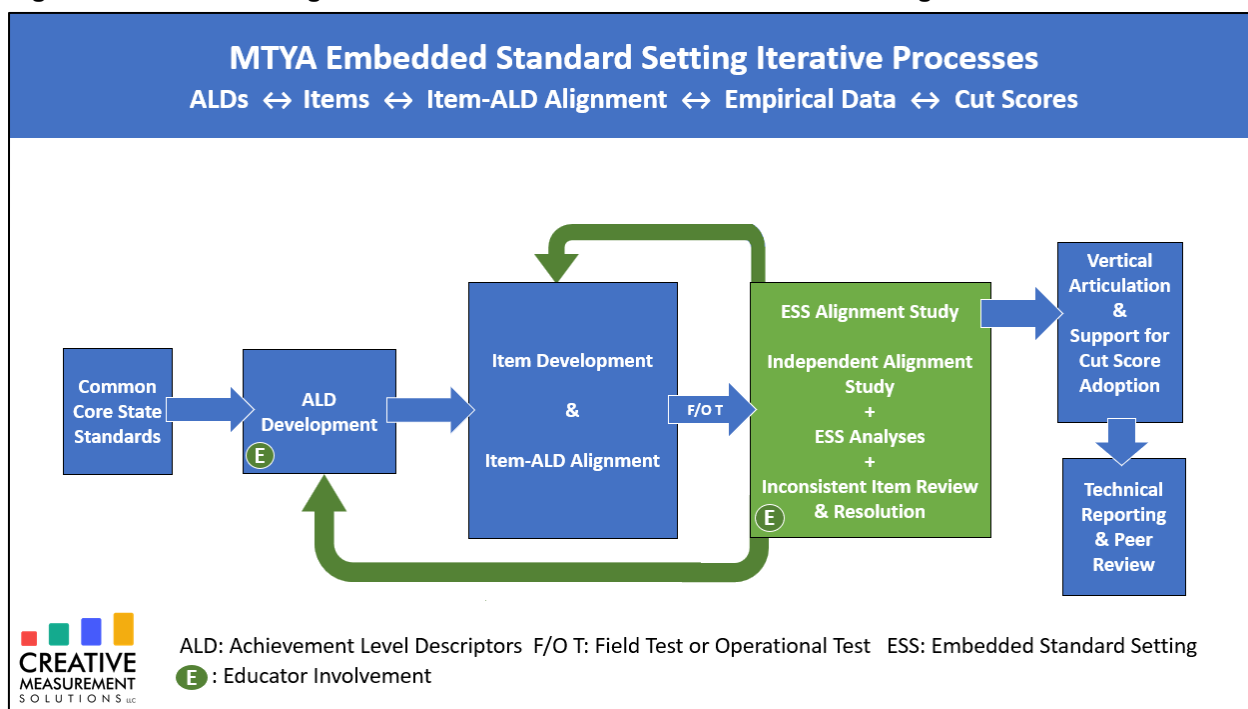
Three cut scores were adopted, defining the four levels of achievement:

- The *Below State Expectations* cut score separates the *Well Below State Expectations* and *Below State Expectations* levels.
- The *At State Expectations* cut score separates the *Below State Expectations* and *At State Expectations* levels.
- The *Above State Expectations* cut score separates the *At State Expectations* and *Above State Expectations* levels.

8.2. Achievement Standard Setting

Embedded Standard Setting (ESS) was employed to establish the MTYA achievement level cut scores. The ESS methodology was selected because it is the natural extension of principled assessment design to standard setting (Lewis & Cook, 2020). It transforms standard setting from a standalone workshop to a set of processes actively integrated throughout the assessment-development lifecycle, as illustrated in Figure 8.1. The iterative nature of the ESS processes (represented by the green feedback arrows in the figure) supports the coherence of various assessment components and artifacts, including Achievement Level Descriptors (ALDs), item development, item-ALD alignment, empirical data, and cut scores (and, therefore, score interpretation). Thus, adherence to these iterative processes supports validity of the assessments and score interpretation. The standard setting technical report is provided in Appendix N.

Figure 8.1. Maine Through Year Assessment Embedded Standard Setting Iterative Processes



ESS processes directly contribute to the valid interpretation and use of test scores and improve test quality and the strength of validity arguments by maintaining a consistent focus on optimizing the evidentiary relationship between test items and the Common Core State Standards (CCSS), as reflected by the associated ALDs. ESS processes include:

- **Achievement Level Descriptor development:** This is an articulation of the intended interpretations of the Maine Through Year Assessment across the achievement levels.
- **The ESS Alignment Study:** This is a review of a representative sampling of MTYA items by Maine educators in which they provide independent alignments of these items to the Common Core State Standards and Maine achievement levels and review and resolve items with alignments that are inconsistent with the data.
- **ESS analyses and the estimation of cut scores:** Educators' alignments of items to the Maine achievement levels are employed to identify optimal cut scores.
- **Post-ESS Alignment Study workshop:** These activities lead to the adoption of cut scores, including cut score refinement to support an integrated, vertically articulated system of cross-grade cut scores meeting workshop panelists' and other stakeholders' expectations and in consideration of Maine DOE policy goals.
- **Documentation of validity evidence supporting Maine's adopted cut scores:** This includes those forms of evidence commonly cited in the measurement literature and those used to satisfy federal peer review requirements.

Findings from each of these activities provide evidence that the ESS processes work together to promote the coherence of the assessment. Specifically:

- Range ALDs were developed to align to the CCSS; final ALDs were reviewed and refined by Maine educators.

- Results from the ESS Alignment Study demonstrated the efficacy of panelists’ consensus regarding the alignment of items to the ALDs; high correlations with empirical difficulty, weighted kappa values, and panelist agreement rates demonstrated a strong panelist understanding of their role and judgment tasks.
- ESS analyses produced cut scores that optimally reflect the panelists’ judgments by minimizing inconsistencies between those judgments and empirical data.
- Results from the Review and Resolution workshop showed iterative improvement in the consensus regarding item-ALD alignments and associated efficacy measures, including correlations, kappa values, and agreement rates, as expected of a consensus-building activity.
- Post-workshop vertical articulation produced a well-articulated, cross-grade system of cut scores in reading and mathematics that reflect the panelists’ and other stakeholders’ expectations for impact data, using methods supported by MTYA Technical Advisory Committee members.
- Thorough documentation of validity evidence supporting the MTYA adopted cut scores demonstrated strong adherence to principles of test-score validation, as articulated in the measurement literature and in the guidelines for federal peer review.

Together, these findings support the validity of the MTYA program’s adopted cut scores. Linkages from ALDs to test scores are consistent with the tenets of Principled Assessment Design, support intended score interpretations, and inform decision-making.

In support of this iterative process, NWEA reviewed the results of the July 2023 alignment study and identified a recommended plan of action that was initially presented to the Maine DOE in April 2024. This plan included action items to add additional information into the assessment blueprint to ensure clarity, to undertake a focused review of the ALD language based on feedback provided in the alignment study, and to identify/develop additional items for Spring 2025 field testing based on an item bank review and analysis. NWEA will consult with the Maine DOE as this work progresses.

For reading and mathematics, the adopted cut scores were presented to the Commissioner of Education and were approved on August 28, 2023. Table 8.1–Table 8.4 present the final approved cut scores that were used for scoring and the associated impact data.

Table 8.1. Final Approved Cut Scores—Reading

Grade	Cut Scores		
	<i>Well Below State Expectations</i>	<i>Below State Expectations</i>	<i>At State Expectations</i>
3	1483	1500	1525
4	1486	1500	1525
5	1487	1500	1525
6	1486	1500	1525
7	1483	1500	1525
8	1484	1500	1525
HS	1489	1500	1525

Table 8.2. Impact Data Associated with Cut Scores—Reading

Grade	Percent at Level			
	<i>Well Below State Expectations</i>	<i>Below State Expectations</i>	<i>At State Expectations</i>	<i>Above State Expectations</i>
3	12.6%	27.1%	47.3%	13.0%
4	12.2%	23.9%	48.5%	15.4%
5	12.8%	18.6%	53.0%	15.6%
6	10.4%	22.5%	53.5%	13.6%
7	11.4%	24.9%	50.4%	13.3%
8	10.1%	24.2%	53.4%	12.3%
HS	13.3%	24.7%	49.7%	12.3%

Table 8.3. Final Approved Cut Scores—Mathematics

Grade	Cut Scores		
	<i>Well Below State Expectations</i>	<i>Below State Expectations</i>	<i>At State Expectations</i>
3	1486	1500	1525
4	1488	1500	1525
5	1484	1500	1525
6	1481	1500	1525
7	1482	1500	1525
8	1484	1500	1525
HS	1489	1500	1525

Table 8.4. Impact Data Associated with Cut Scores—Mathematics

Grade	Percent at Level			
	<i>Well Below State Expectations</i>	<i>Below State Expectations</i>	<i>At State Expectations</i>	<i>Above State Expectations</i>
3	17.3%	21.1%	43.9%	17.7%
4	18.6%	24.5%	44.0%	12.9%
5	18.5%	30.7%	40.0%	10.8%
6	18.8%	36.4%	35.9%	8.9%
7	20.1%	36.0%	35.4%	8.5%
8	20.5%	39.1%	33.5%	6.9%
HS	25.0%	32.0%	35.5%	7.5%

8.3. Reporting

The Maine Through Year Assessments are administered in reading and mathematics. These assessments were developed specifically for Maine to provide teachers, students, and parents with information on student learning strengths and needs throughout the year, as well as student progress in mastering college and career-ready skills based on Maine’s accountability standards, the Common Core State Standards.

8.3.1. Achievement Level Descriptors

Achievement Level Descriptors (ALDs) are a plain-language description of what students must know as defined by each of the achievement levels established through cut scores. The ALDs firmly root the cut scores and achievement levels in the content that students are supposed to learn. In qualitative and quantitative terms, the ALDs and cut scores *together* define the difference between a student who is performing at, below, or above grade-level expectations (see Section 2.4 and Table 2.16 for more details about ALDs). The cut scores for these achievement levels were established and validated in summer 2023 by Maine educators, the Maine DOE, and the Maine Technical Advisory Committee.

8.3.2. Setting the Cut Scores

To establish the cut scores, a process called “embedded standard setting” helps determine two points along the scale score range (known as cut scores) that define the score range for each achievement level. Maine educators and stakeholders from around the state participated in the embedded standard-setting process for the MTYA, facilitated by edCount and Creative Measurement. The cut score recommendations from this statewide committee were presented to the Maine Department of Education and were approved in late August 2023.

8.3.3. Reports

For the MTYA, reports were developed and are available at the district, school, group, and individual student levels. Table 8.5 presents a description of each report. A more detailed report explanation can be found in Appendix F.

Table 8.5. Report Levels

Report Name	Aggregation Level	Summary
District Report	District (SAU)	Shows the average scale scores for schools in the district, the distribution of school average scale scores across the achievement levels, and the distribution of student scale scores in each school
School Report	School	Shows the average scale scores for students in the school, the distribution of student scale scores across the achievement levels, the average scale scores and score distributions for each group in the school, and the individual scale scores for each student in the school
Teacher Report	Group	Shows the average scale scores for students in the group, the distribution of student scale scores across the achievement levels, and the individual scale scores for each student in the group
Student Report	Individual Student	Shows all the details for an individual student’s test
Individual Student Report (Spring Only)	Individual Student	Shows all tests in all available content areas for a student in the spring; designed for parents and families
RIT Report	Varies—based on user type	Shows RIT (i.e., interim) score information for all students matching the search criteria, including RIT score, achievement percentile, and instructional area RIT

Report Name	Aggregation Level	Summary
Demographic Report	Varies—based on user type	Shows the average scale scores, average reporting category scores, and distribution of scale scores for demographic groups such as gender, ethnicity/race, and targeted group
Comparison Summary Report	School	Shows aggregate comparison of multiple organizations by grade, subject, and student demographics
Student Results File	District and State	Downloadable export of student-level data at district and state levels during the test window

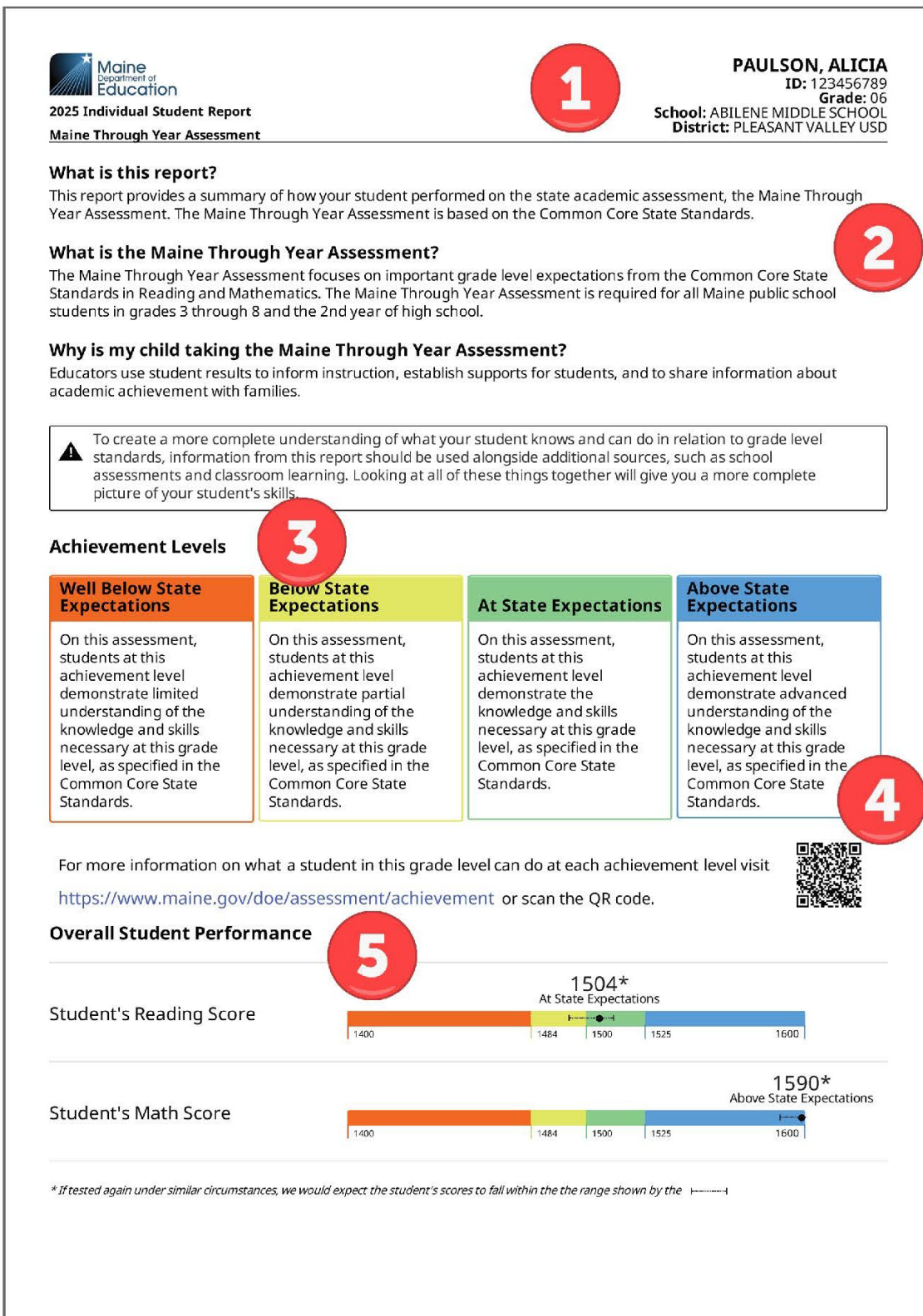
Figure 8.2 shows a mockup of the Individual Student Report (ISR). The ISR is a one-page report designed to show a student’s achievement on the Maine Through Year reading and mathematics assessments to parents and families.

In November 2023, a group of Maine educators worked collaboratively with Maine DOE to develop parent-friendly, accessible language and formatting for the Maine Through Year Assessment ISRs. In addition, this educator panel created supplemental documents to explain the reading and mathematics skills focused on at each grade level in easy-to-understand terms. These Individual Student Report supplemental pages for families have been translated into Maine’s top ten languages. An English example can be found in Appendix O.

Educators can print ISRs in batches, making them easy to distribute after testing is complete.

The ISRs are generated for the spring assessment and will not be available for the fall and winter assessments.

Figure 8.2. Individual Student Report



For more report screenshots and report explanations, please see Appendix F.

8.3.4. Relations to Other Scores

Evidence based on relations to other variables refers to traditional forms of criterion-related validity evidence, such as concurrent and predictive validity, and more comprehensive investigations of the relationships among test scores and other variables, such as multitrait-multimethod studies.

Table 8.6 shows the correlations between the Spring 2025 summative scores and the Fall 2024 interim scores for reading and mathematics, as well as the correlations between the Spring 2025 science summative scores for grades 5 and 8 (provided by Maine DOE) and the respective Spring 2025 summative scores for reading and mathematics. The correlations range from 0.74 to 0.79 for reading, 0.81 to 0.86 for mathematics, and 0.74 to 0.75 for science grades 5 and 8. Additional correlation data by student group is provided in Appendix M.

The strong correlations between spring summative and fall interim scores for reading (ranging from 0.74 to 0.79) and between spring summative and fall interim scores for mathematics (ranging from 0.81 to 0.86) indicate strong convergent validity, meaning that the same skills or traits are being measured using different methods. Meanwhile, the correlations between spring summative reading and science (0.75 for grades 5 and 8) and between spring summative mathematics and science (ranging from 0.74 to 0.75) demonstrate discriminant validity by showing the relationship between different skills or traits measured by the same assessment method (spring summative tests).

Table 8.6. Correlations Among Spring 2025 Summative Scores and Fall 2024 Interim Scores

Grade	Spring Summative/Fall Interim Correlation (Reading/Reading)	Spring Summative Correlation (Reading/Science)	Spring Summative/Fall Interim Correlation (Math/Math)	Spring Summative Correlation (Math/Science)
3	0.78	–	0.81	–
4	0.79	–	0.84	–
5	0.79	0.75	0.85	0.75
6	0.78	–	0.85	–
7	0.79	–	0.86	–
8	0.78	0.75	0.85	0.74
HS	0.74	–	0.82	–

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.
https://www.testingstandards.net/uploads/7/6/6/4/76643089/standards_2014edition.pdf
- Every Student Succeeds Act (ESSA), 20 U.S.C. § 6301 (2015).
<https://www.congress.gov/114/plaws/publ95/PLAW-114publ95.pdf>
- French, A. W., & Miller, T. R. (1996). Logistic regression and its use in detecting differential item functioning in polytomous Items. *Journal of Educational Measurement*, 33(3), 315–332.
<https://www.jstor.org/stable/1435375>
- Fu, J., & Monfils, L. (2016). *LDIF_ES: A SAS macro for logistic regression tests for differential item functioning of dichotomous and polytomous items*. (Research Memorandum ETS RM–16-17). Educational Testing Service (ETS).
<https://www.ets.org/Media/Research/pdf/RM-16-17.pdf>
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Springer. <https://doi.org/10.1007/978-94-017-1988-9>
- Kolen, M. J., & Brennan, R. L. (2004). *Test equating, scaling, and linking: Methods and practices* (2nd ed.). Springer. <https://doi.org/10.1007/978-1-4939-0317-7>
- Lewis, D., & Cook, R. (2020). Embedded standard setting: Aligning standard setting methodology with contemporary assessment design principles. *Educational Measurement: Issues and Practice*, 39(1), 8–21. <https://doi.org/10.1111/emip.12318>
- Linacre, J. M. (2015). *Winsteps®* (Version 3.90.2) [Computer software]. Winsteps.com. Available from <https://www.winsteps.com/>
- Linacre, J. M. (2002) What do infit and outfit, mean-square and standardization mean? *Archives of Rasch Measurement*, 16(2), 878. <https://www.rasch.org/rmt/rmt162f.htm>
- Linn, R. L. (2006). The standards for educational and psychological testing: Guidance in test development. In S. M. Downing & T. M. Haladyna (Eds.), *Handbook of Test Development*. Routledge.
- Masters, G. N. (1982). A Rasch model for partial credit scoring. *Psychometrika*, 47(2), 149–174.
<https://doi.org/10.1007/BF02296272>
- NWEA. (2020). *Constraint-based engine scientific approach and methodology* [Confidential Tech. Rep.].
- Phillips, S., & Camara, W. J. (2006). Legal and ethical issues. In R. L. Brennan (Ed.), *Educational measurement* (4th ed., pp. 734–755). Praeger.

- Puhan, G., & Dorans, N. (2018). *Technical considerations in scale development*. Annual Meeting of the National Council on Measurement in Education, New York, NY, United States.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. MESA Press.
- Rasch, G. (1980). *Probabilistic models for some intelligence and attainment tests* (Expanded ed.). University of Chicago Press.
- Samejima, F. (1994). Estimation of reliability coefficients using the test information function and its modifications. *Applied Psychological Measurement*, 18(3), 229–244.
<https://doi.org/10.1177/014662169401800304>
- Smarter Balanced Assessment Consortium (SBAC). (2016). *Smarter Balanced Assessment Consortium: 2014–15 technical report*.
<https://portal.smarterbalanced.org/library/en/2014-15-technical-report.pdf>.
- Swaminathan, H., & Rogers, H. J. (1990). Detecting differential item functioning using logistic regression procedures. *Journal of Educational Measurement*, 27(4), 361–370.
<https://www.jstor.org/stable/1434855>
- Van der Linden, W. J., & Reese, L. M. (1998). A model for optimal constrained adaptive testing. *Applied Psychological Measurement*, 22(3), 259–270.
<https://doi.org/10.1177/01466216980223006>
- Wright, B. D. (1977). Solving measurement problems with the Rasch model. *Journal of Educational Measurement*, 14(2), 97–116. <https://www.jstor.org/stable/1434010>
- 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., Thayer, D. T., & Wingersky, M. (1994). DIF analysis for pretest items in computer-adaptive testing. (Research Report No. RR-94-33). Educational Testing Service (ETS).
<https://doi.org/10.1002/j.2333-8504.1994.tb01606.x>