

State of Florida

Florida Assessment of Student Thinking (FAST), Benchmarks for Excellent Student Thinking (B.E.S.T.), and Science & Social Studies Statewide Assessments Technical Report

2023–2024

Volume 2 Test Development

ACKNOWLEDGMENTS

This technical report was produced on behalf of the Florida Department of Education (FDOE). Requests for additional information concerning this technical report or the associated appendices should be directed to Dr. Salih Binici at FDOE (Salih.Binici@fldoe.org).

Major contributors to this technical report include the following staff from Cambium Assessment, Inc. (CAI): Dr. Yuan Hong, Dr. Myvan Bui, Dr. Sherry Li, Dr. Peter Diao, Matt Gordon, Zoe Dai, Oliver Brown, and Daniel Lam. The major contributors from FDOE are as follows: Dr. Salih Binici, Vince Verges, Susie Lee, Catherine Altmaier, Racquel Harrell, Sally Donnelly, Dr. Shakia Johnson, Dr. Stacy Skinner, Leah Glass, Kristina Lamb, Dr. Wenyi Li, Jieling Ming, Saeyan Yun, and Yiting Yao.

TABLE OF CONTENTS

1. INTRODUCTION	1
2. TEST SPECIFICATIONS	3
2.1 Blueprint Development Process	3
2.1.1 Target Blueprints	4
2.1.2 Cognitive Complexity	9
2.2 Content-Level and Psychometric Considerations	10
3. ITEM DEVELOPMENT PROCEDURES	12
3.1 Summary of Item Sources	13
3.2 Item Types	14
3.3 Cognitive Laboratories	15
3.4 Item Translations to Braille Format	17
3.5 Development and Review Process for New Items	18
3.5.1 Development of New Items	18
3.5.2 Rubric Validation	21
3.6 Development and Maintenance of the Item Pool	23
3.7 Alignment Process for Existing Items and Results from Alignment Studies	24
4. TEST CONSTRUCTION	26
4.1 Overview	26
4.2 Item Selection Algorithm	26
4.3 Test Construction Summary Materials	30
4.3.1 Item Cards	30
4.3.2 Bookmaps for Accommodated Forms	31
4.4 Accommodation Form Construction	32
5. REFERENCES	34

LIST OF APPENDICES

Appendix A1: Science Reporting Categories Descriptors
Appendix A2: ELA Reporting Categories Descriptors
Appendix A3: NGSSS EOC Reporting Categories Descriptors
Appendix A4: Mathematics and EOC Reporting Categories Descriptors
Appendix B1: ELA Blueprints
Appendix B2: Mathematics and EOC Blueprints
Appendix B3: U.S. History Blueprints
Appendix B4: Civics Blueprints
Appendix B5: Biology and Science Blueprints
Appendix C: Example Item Types
Appendix D: Spring 2024 Operational Item Blueprint Match
Appendix E: ELA and Mathematics Alignment Study Proposal

Appendix F: Spring 2024 Simulation Results
Appendix G: Cognitive Lab Final Report

LIST OF TABLES

Table 1: Blueprint Test Length by Grade and Subject or Course.....	5
Table 2: Number of Items Available in the Spring 2024 Item Pool by Grade and Subject or Course.....	5
Table 3: Blueprint Percentage of Test Items Assessing Each Reporting Category in ELA Reading	6
Table 4: Blueprint Percentage of Test Items Assessing Each Reporting Category in Mathematics	7
Table 5: Reporting Categories Used in Mathematics	7
Table 6: Blueprint Percentage of Test Items Assessing Each Reporting Category in Mathematics EOC	7
Table 7: Reporting Categories Used in EOC	8
Table 8: Blueprint Percentage of Test Items Assessing Each Reporting Category in Science and Social Studies	8
Table 9: Reporting Categories Used in Science and Social Studies	8
Table 10: Blueprint Percentage of Reading Passage Types by Grade.....	9
Table 11: Blueprint Percentage of Items by Cognitive Complexity.....	9
Table 12: Item Bank Observed Percentage of Items by Cognitive Complexity.....	10
Table 13: ELA Reading Item Types and Descriptions	14
Table 14: Mathematics and Mathematics EOC Item Types and Descriptions	14
Table 15: Science and Social Studies Item Types and Descriptions	15
Table 16: Word Counts and Readabilities of Reading Passages in FAST ELA Reading	19
Table 17: Number of ELA Reading Field-Test Items by Type	23
Table 18: Number of Mathematics and EOC Field-Test Items by Type.....	23
Table 19: Number of Science and Social Studies Field-Test Items by Type	24
Table 20: Number of ELA Writing Field-Test Prompts by Grade in 2024.....	24
Table 21: Observed Spring 2024 Percentage of ELA Reading Passage Types by Grade	30

LIST OF FIGURES

Figure 1: Item Selection Process.....	29
Figure 2: Example Item Card.....	31

1. INTRODUCTION

By statute, all Florida public school students are required to participate in statewide assessments. Beginning with the 2022–2023 school year, Florida’s statewide, standardized assessments in English Language Arts (ELA) Reading, ELA Writing, and Mathematics were aligned with the Benchmarks for Excellent Student Thinking (B.E.S.T.). Assessments for Science and Social Studies remain aligned to Florida’s state academic standards, which were adopted in 2008. These Science and Social Studies standards were previously referred to as Next Generation Sunshine State Standards (NGSSS).

In December 2008, the Florida State Board of Education (SBE) adopted NGSSS for Social Studies. These standards were used to develop the U.S. History End-of-Course (EOC) Assessment. The 2010 Florida Legislature authorized the Florida EOC assessments. The grade 7 Civics and Government strand of the NGSSS was used to develop the initial Civics EOC Assessment. In July 2021, the SBE adopted new state academic standards for Civics and Government (CG). Beginning with the spring 2024 test administration, these standards will be assessed by the Civics EOC Assessment.

The State of Florida implemented new online computer-adaptive tests (CATs) for operational use beginning with the 2022–2023 school year for ELA Writing and Mathematics, and in 2023–2024 for Science and Social Studies. Before this, they were fixed-form online and paper tests.

The assessment program for ELA and mathematics, referred to as the Florida Assessment of Student Thinking (FAST), replaced the Florida Standards Assessments (FSA) in ELA Reading and Mathematics. The FAST assessments are progress monitoring (PM) assessments administered three times a year. Grades 4–10 Writing and Mathematics EOC Algebra 1 and Geometry are considered B.E.S.T. assessments and are not part of the PM FAST assessments. The FAST and B.E.S.T. were first administered to students during fall 2022. Writing was administered in spring 2023 as a standalone field test administered to a representative sample of Florida students. Beginning with the 2023–2024 school year, Writing will be administered during each spring administration. In spring testing windows, all the CATs are given to students as summative assessments.

Beginning in spring 2024, the summative Statewide Science Assessment in grades 5 and 8, as well as the Biology 1, Civics, and U.S. History EOC assessments, were delivered in a computer-adaptive format that allows for immediate reporting. While the core content for these tests did not change, some administration details (e.g., reduced test length) and blueprint specifications (e.g., number of items each student will see) have been updated. The fall and winter 2023 administrations of the Science and Social Studies EOC assessments were computer-based, fixed-form tests, and results were available for all students after the testing window as in previous years.

The online versions of the Science and Social Studies assessments include multiple-choice items, and the ELA and mathematics assessments include (but are not limited to) the use of several technology-enhanced item types. For all online assessments, accommodated versions are available to students whose Individualized Education Plans (IEPs) or Section 504 Plans indicate such a need. IEPs or Section 504 Plans indicating the need for versions other than online are addressed accordingly.

The interpretation, usage, and validity of test scores rely heavily on the test development process itself. This volume details that process and how it contributes to the validity of the test scores. Specifically, this volume provides evidence to support the following:

- The test design summary/blueprint stipulated the range of operational items from each reporting category that were required on each form. This document guided item selection and test construction for the assessments.
 - The test design summaries for both Mathematics and ELA Reading were updated during the 2022–2023 school year to represent the reduced test length and new reporting categories. Content Advisory Committee meetings were conducted with educators so they could provide feedback on the overall test length, number of reporting categories, and benchmarks included in those reporting categories. The design summary now specifically states that the ELA Reading and ELA Writing components are tested and reported separately. All Mathematics and ELA tests are also administered in one session, in one day.
 - The test design summaries for Science and Social Studies were updated in the 2023–2024 school year. The core content for these tests did not change; however, some items were reduced in length and the number of items each student will see per blueprint component was updated.
- The test item specifications provided detailed guidance for item writers and reviewers to ensure that the items were aligned to the standards they were intended to measure. The Florida Department of Education (FDOE) and committees of experienced Florida educators developed and approved the specifications that define the content and format of the tests and test items. The test item specifications for ELA and Mathematics were revised in 2021 and 2022 after the adoption of the B.E.S.T. Standards and the decision to move to CATs, and for Science and Social Studies in 2023–2024 when those assessments were moved to CATs. The item specifications are also updated each year as needed to document any necessary changes or clarifications that arise throughout a development cycle.
- The item development procedures employed were consistent with industry standards.
- The development and maintenance of the item-pool plan established an item bank in which test items cover the range of measured standards, grade-level difficulties, and cognitive complexity using both selected-response (SR) keyed items and constructed-response (CR) machine-scored or handscored item types.
- The thorough test development process contributed to the comparability of the online and accommodated tests.

2. TEST SPECIFICATIONS

Following the adoption and integration of the Florida State Academic and B.E.S.T. Standards into the school curriculum, items and test item specifications were developed to ensure that the tests and their items were aligned to the benchmarks and grade-level expectations that they were intended to measure. FDOE and content specialists developed test item specifications. Educator committees also reviewed and provided feedback for the initial development of the B.E.S.T. test item specifications.

The assessments are based on their relevant standards, course descriptions, and test item specifications. FAST and B.E.S.T. test item specifications are based on the B.E.S.T. Standards, the state’s science assessments are based on Florida’s Science standards, while U.S. History and Civics are based on their respective standards. In July 2021, the new state academic standards for Civics were adopted. Beginning with the spring 2024 test administration, these standards will be assessed by the Civics EOC Assessment. The specifications are a resource that defines the content and format for the test and test items for item writers and reviewers. Each grade-level and course specifications document indicates the alignment of items with the Florida Standards and informs all stakeholders about the scope and function of the assessments. In addition to these general guidelines, specifications for FAST ELA Reading and ELA Writing also include guidelines for developing reading and writing passages and prompts, such as length, type, and complexity.

2.1 BLUEPRINT DEVELOPMENT PROCESS

A test design summary/blueprint for each assessment specifies the number of items, item types, and reporting categories. The blueprint construction for the assessments are evidenced by the test design summary documents found at <https://www.fldoe.org/accountability/assessments/k-12-student-assessment>. These documents were created using Florida’s course descriptions as the basis for the design. The course descriptions can be found on the CPALMS website at <http://www.cpalms.org/Public/search/Course>.

After the decision was made to switch ELA Reading and Mathematics to a computer-adaptive test (CAT), Content Advisory Committee (CAC) meetings were held with educators to propose and approve the revised blueprints. An in-person CAC meeting was held to discuss the blueprints for ELA Reading and Mathematics grades 3–8 in April 2022. There was a virtual CAC meeting held for Mathematics and EOC assessments in September 2022. In December 2022, a virtual meeting was held for both subjects to discuss potentially shortening the blueprints even further after they had been shortened during the move from FSA to B.E.S.T. (although the decision was to not do so). The blueprints were also discussed at the Technical Advisory Committee (TAC) meetings in summer and November 2022.

The reporting categories for ELA Reading were derived from the applicable “cluster” naming convention in the Florida B.E.S.T. Standards, and the percentages of the reporting categories within the tests were derived from the number, complexity, and breadth of the standards to be assessed. Vocabulary standards were folded in with the Reading Across Genres Standards to create the Reading Across Genres & Vocabulary reporting category. Guidelines for the weight of each reporting category for FAST ELA Reading were determined by Florida’s TAC. The TAC advised FDOE that to avoid “statistical noise” generated from the items scored in a small reporting

category, a minimum of 15% of the total raw score points should be derived from each reporting category.

The reporting categories for Mathematics were also derived from the “domain” naming convention in the Florida B.E.S.T. Standards. As with ELA Reading, if a Mathematics domain has too few standards, two or more domains might be combined to make the reporting category 15% of the raw score points of that grade’s assessment.

The benchmark information provides benchmark clarification statements, assessment limits, stimulus attributes, response attributes, prior knowledge, and a sample item for each benchmark that could be assessed.

The Science, U.S. History, and Civics reporting categories assessed are defined based on the 2007 and 2008 NGSSS adoption. The NGSSS are divided into benchmarks that identify what a student should be able to do following the completion of each course. The test item specifications documents for NGSSS Science, Civics, and U.S. History contain benchmark-specific information.

Detailed descriptions for the constructs of the reporting categories are presented in Appendices A1–A4 for all assessments.

2.1.1 Target Blueprints

Test blueprints provided the following guidelines:

- Length of the test (duration and number of items)
- Content areas to be covered and the acceptable range of items in each content area or reporting category
- Acceptable range of item difficulty for the specified grade level
- Approximate number of field-test items, if applicable
- Descriptions of test item types

This section provides only a summary of the blueprints. Detailed blueprints for each content level are presented in Appendices B1–B5 for all subjects.

In all grades and subjects, the assessments are administered as CATs. Grades 3–10 ELA Reading, grades 3–8 Mathematics, Mathematics EOC assessments (Algebra 1 and Geometry), grades 5 and 8 Science, EOC assessments in Science and Social Studies (Biology 1, U.S. History, and Civics) are administered online. Additionally, ELA Writing is administered online for grades 4–10. In spring 2023, typed written response accommodations were provided for students taking ELA Writing assessments in grades 4–10; therefore, responses from these students were collected online. For grades and subjects testing online, accommodations are provided if indicated by a student’s IEP or Section 504 Plan.

Table 1 displays the blueprint for total test length by grade and subject or course. Each year, approximately 6–10 items on all tests are field-test items and are not used to calculate a student’s score. Table 2 displays the number of operational and field-test items available in the spring 2024 item pool. ELA Writing items are not included in the item counts listed for ELA Reading tests.

Table 1: Blueprint Test Length by Grade and Subject or Course

Subject/Course	Grade	Total Number of Items
ELA Reading	3	36–40
	4	36–40
	5	36–40
	6	36–40
	7	36–40
	8	36–40
	9	36–40
	10	36–40
Mathematics	3	35
	4	35
	5	35
	6	36
	7	36
	8	36
Algebra 1	9	45
Geometry	9	45
Biology 1	10	45
Grade 5 Science	5	45
Grade 8 Science	8	45
Civics	7	44
U.S. History	9	46

Table 2: Number of Items Available in the Spring 2024 Item Pool by Grade and Subject or Course

Subject/Course	Grade	Number of Operational Items	Total Item Counts in the Field-Test Pool	Total Item Counts in the Field Test and Operational Pool
ELA Reading	3	356	349	705
	4	432	233	665
	5	442	229	671
	6	357	231	588
	7	374	249	623
	8	331	293	624
	9	389	263	652
	10	391	290	681

Subject/Course	Grade	Number of Operational Items	Total Item Counts in the Field-Test Pool	Total Item Counts in the Field Test and Operational Pool
Mathematics	3	444	234	678
	4	330	349	679
	5	450	234	684
	6	450	230	680
	7	367	253	620
	8	287	350	637
Algebra 1	9	318	279	597
Geometry	9	365	152	517
Biology 1	10	838	183	1,021
Grade 5 Science	5	802	337	1,139
Grade 8 Science	8	694	257	951
Civics	7	546	320	866
U.S. History	9	615	249	764

Reporting categories were used to more narrowly define the topics assessed in each content area. Individual scores on reporting categories provide information to help identify areas in which a student may have had difficulty. Table 3 and Table 4 provide the percentage of operational items required in the blueprints by content strands, or reporting categories, for each grade level or course. The percentages shown represent an acceptable range of item counts. As many of these items in the ELA Reading component were associated with passages, flexibility was necessary for test construction for practical reasons. The ELA Writing component prompt was not included in these blueprints. Table 5 provides the reporting categories for Mathematics grades 3–8, Table 6 provides the percentage of operational items required in the blueprints by content strands, or reporting categories, for Mathematics EOC, and Table 7 provides the reporting categories for Mathematics EOC. Table 8 provides the percentage of operational items required in the blueprints by content strands, or reporting categories, for Science and Social Studies. Table 9 provides the reporting categories for Science and Social Studies.

Table 3: Blueprint Percentage of Test Items Assessing Each Reporting Category in ELA Reading

Grade	Reading Prose and Poetry	Reading Informational Text	Reading Across Genres & Vocabulary
3	25%–35%	25%–35%	35%–50%
4	25%–35%	25%–35%	35%–50%
5	25%–35%	25%–35%	35%–50%
6	25%–35%	25%–35%	35%–50%
7	25%–35%	25%–35%	35%–50%
8	25%–35%	25%–35%	35%–50%

Grade	Reading Prose and Poetry	Reading Informational Text	Reading Across Genres & Vocabulary
9	25%–35%	25%–35%	35%–50%
10	25%–35%	25%–35%	35%–50%

Table 4: Blueprint Percentage of Test Items Assessing Each Reporting Category in Mathematics

Grade	1*	2*	3*	4*
3	23%–29%	23%–29%	23%–29%	23%–29%
4	31%–37%	31%–37%	31%–37%	
5	23%–29%	23%–29%	23%–29%	23%–29%
6	33%–42%	25%–36%	25%–36%	
7	25%–31%	22%–31%	22%–28%	22%–28%
8	22%–28%	22%–28%	25%–31%	22%–28%

*See Table 5 for the reporting category names.

Table 5: Reporting Categories Used in Mathematics

Grade	Reporting Category
3	Number Sense and Additive Reasoning Number Sense and Multiplicative Reasoning Fractional Reasoning Geometric Reasoning, Measurement, and Data Analysis and Probability
4	Number Sense and Operations with Whole Numbers Number Sense and Operations with Fractions and Decimals Geometric Reasoning, Measurement, and Data Analysis and Probability
5	Number Sense and Operations with Whole Numbers Number Sense and Operations with Fractions and Decimals Algebraic Reasoning Geometric Reasoning, Measurement, and Data Analysis and Probability
6	Number Sense and Operations Algebraic Reasoning Geometric Reasoning, Data Analysis, and Probability
7	Number Sense and Operations and Algebraic Reasoning Proportional Reasoning and Relationships Geometric Reasoning Data Analysis and Probability
8	Number Sense and Operations and Probability Algebraic Reasoning Linear Relationships, Data Analysis, and Functions Geometric Reasoning

Table 6: Blueprint Percentage of Test Items Assessing Each Reporting Category in Mathematics EOC

Course	1*	2*	3*
Algebra 1	31%–38%	31%–38%	31%–38%

Course	1*	2*	3*
Geometry	33%–40%	27%–33%	33%–40%

*See Table 7 for reporting category names.

Table 7: Reporting Categories Used in EOC

Course	Reporting Category
Algebra 1	Expressions, Functions and Data Analysis Linear Relationships Non-Linear Relationships
Geometry	Logic, Relationships, and Theorems Congruence, Similarity, and Constructions Measurement and Coordinate Geometry

Table 8: Blueprint Percentage of Test Items Assessing Each Reporting Category in Science and Social Studies

Grade/Course	1*	2*	3*	4*
Biology 1	35%	25%	40%	
Grade 5 Science	17%	29%	29%	25%
Grade 8 Science	19%	27%	27%	27%
Civics	25%–30%	15%–20%	20%–25%	20%–25%
U.S. History	33%	34%	33%	

*See Table 9 for the reporting category names.

Table 9: Reporting Categories Used in Science and Social Studies

Grade/Course	Reporting Category
Biology 1	Molecular and Cellular Biology Classification, Heredity, and Evolution Organisms, Populations, and Ecosystems
Grade 5 Science	Nature of Science Earth and Space Science Physical Science Life Science
Grade 8 Science	Nature of Science Earth and Space Science Physical Science Life Science
Civics	Origins and Purposes of Law and Government Roles, Rights, and Responsibilities of Citizens Government Policies and Political Processes Organization and Function of Government
U.S. History	Late Nineteenth and Early Twentieth Century, 1860–1910

Grade/Course	Reporting Category
	Global Military, Political, and Economic Challenges, 1890–1940
	The United States and the Defense of the International Peace, 1940–Present

The FAST ELA Reading blueprint also included specifications for the genres of text presented in the passages. Two main types of text were used: literary and informational. Table 10 provides target percentages of the test passages assessing each type of text.

Table 10: Blueprint Percentage of Reading Passage Types by Grade

Grades	Informational	Literary
3–5	50%	50%
6–8	50%	50%
9–10	50%	50%

2.1.2 Cognitive Complexity

Cognitive complexity refers to the cognitive demand associated with a test item. The cognitive classification system implemented by FDOE is based on Dr. Norman L. Webb’s Depth of Knowledge (DOK) levels (Webb, 2002). The rationale for classifying a test item by its cognitive complexity focuses on the expectations made of the test item, not on the ability of the student. When classifying a test item’s demands on thinking (i.e., what the test item requires the student to recall, understand, analyze, and do), it is assumed that the student is familiar with the basic concepts of the task. Test items are chosen for the Statewide Science Assessment based on the standards and their grade level appropriateness, but the complexity of the test items remains independent of the particular curriculum a student has experienced. On any given assessment, the cognitive complexity of a multiple-choice (MC) item may be affected by the distractors (answer options). The cognitive complexity of a test item depends on the grade level of the assessment; a test item that has a high level of cognitive complexity at one grade may not be as complex at a higher grade. The categories—low, moderate, and high complexity—form an ordered description of the demands a test item may make on a student. For example, low-complexity test items may require a student to solve a one-step problem. Moderate-complexity test items may require multiple steps. However, the number of steps is not always indicative of cognitive level. High-complexity test items may require a student to analyze and synthesize information. The distinctions made in terms of complexity ensure that test items will assess the depth of student knowledge at each benchmark. The intent of the item writer weighs heavily in determining the complexity of a test item. Table 11 presents the target range of the percentage of items by cognitive complexity on Statewide Science and Social Studies Assessments at the CAT item bank level. Table 12 presents the actual percentages in the item banks.

Table 11: Blueprint Percentage of Items by Cognitive Complexity

Grade/Course	Low	Moderate	High
Grade 5 Science	10–20%	60–80%	10–20%
Grade 8 Science	10–20%	60–80%	10–20%

Grade/Course	Low	Moderate	High
Biology 1	10–20%	60–80%	10–20%
U.S. History	20–30%	45–65%	15–25%
Civics	15–25%	45–65%	15–25%

Table 12: Item Bank Observed Percentage of Items by Cognitive Complexity

Grade/Course	Low	Moderate	High
Grade 5 Science	15.5%	73.2%	11.3%
Grade 8 Science	18.2%	71.0%	10.8%
Biology 1	14.3%	72.1%	13.6%
U.S. History	24.7%	55.1%	20.2%
Civics	26.6%	49.6%	23.8%

Cognitive complexity targets are not currently a required component of FAST/B.E.S.T. blueprints. However, item development plans target a variety of intended cognitive complexity to explore the different components of the benchmarks and the variety of rigor they offer within a grade level.

2.2 CONTENT-LEVEL AND PSYCHOMETRIC CONSIDERATIONS

In addition to the test blueprints, several content-level and psychometric considerations were used in the development of Florida’s K–12 statewide student assessment. Content-level considerations included the following:

- Selected items addressed a variety of topics.
- Identified correct answer or key was correct.
- Each item had only one correct response (some technology-enhanced items did, in fact, have more than one correct answer, and these items were reviewed to confirm that the number of correct answers matched the number asked for in the item itself).
- Identified item content or reporting category was correct.
- Items were free from typographical, spelling, punctuation, or grammatical errors.
- Items were free of any bias concerns and did not include topics that stakeholders might find offensive.
- Items fulfilled style specifications (e.g., italics, boldface).
- Items marked as do-not-use (DNU) were not selected.

Psychometric considerations included the following:

- A reasonable range of item difficulties was included.
- p -values for MC and constructed-response (CR) items were reasonable and within specified bounds.
- Corrected point-biserial correlations were reasonable and within specified bounds.
- No items with negative corrected point-biserial correlations were used.

- Item response theory (IRT) a -parameters for all items were reasonable and greater than 0.50.
- IRT b -parameters for all items were reasonable and between -2 and 3 .
- For MC items, IRT c -parameters were less than 0.40.
- Few items with model fit flags were used.
- Few items with differential item functioning (DIF) flags were used.

More information about p -values, corrected point-biserial correlations, IRT parameters, and DIF calculations can be found in Volume 1, Annual Technical Report, of this report. The spring 2023 FAST and B.E.S.T. tests were calibrated and equated to the IRT-calibrated item pool. More details about calibration, equating, and scoring can be found in Volume 1 of this report.

3. ITEM DEVELOPMENT PROCEDURES

The item development procedures employed by Cambium Assessment, Inc. (CAI) for the Florida assessments were consistent with industry practice. Just as the development of Florida’s content and performance standards was an open, consensus-driven process, the development of test items and stimuli to measure those constructs was grounded in a similar philosophy.

Item development began with the following guidelines: the test item specifications; the Florida Standards; language accessibility, bias, and sensitivity guidelines; editorial style guidelines; and the principles of universal design (UD). These guidelines ensured that each aspect of a Florida item was relevant to the measured construct and was unlikely to distract or confuse test takers. In addition, these guidelines helped ensure that the wording, required background knowledge, and other aspects of the item were familiar across identifiable groups.

The principles of UD of assessments mandate that tests are designed to minimize the impact of construct-irrelevant factors in the assessment of student achievement, removing barriers to access for the widest range of students possible. The following seven principles of UD, as clearly defined by Thompson, Johnstone, & Thurlow (2002), were applied to the assessments’ development:

1. Inclusive assessment population
2. Precisely defined constructs
3. Accessible, non-biased items
4. Amenable to accommodations
5. Simple, clear, and intuitive instructions and procedures
6. Maximum readability and comprehensibility
7. Maximum legibility

CAI applied these UD principles in the development of all test materials, including tasks, items, and manipulatives. Test development specialists receive extensive training in item development. At every step of the review process, adherence to the principles of UD was confirmed.

The application of UD principles as defined by Thompson, Johnstone, & Thurlow (2002) helps develop assessments that are usable to the greatest number of test takers, including Students with Disabilities (SWDs) and English language learners (ELLs).

As documented in this technical report, the item development procedures implemented for the Florida tests are consistent with industry practice. Specifically, Florida implements the UD principles throughout every stage of the assessment development process (i.e., initial design, item development, field testing, and implementation) to minimize the need for individual accommodations. As noted by Shaftel et al. (2015), under UD principles, accessibility is integral to the item development processes, thus minimizing access barriers associated with the tests themselves to the greatest extent possible for all students, including SWDs and ELLs.

Test development specialists receive extensive training in item development, including instruction on the UD principles and guidance on designing accessible content. Adherence to the UD principles is confirmed at every step of the review process so that the test maximizes readability, legibility, and compatibility with accommodations. Checklists that align to the Council of Chief

State School Officers (CCSSO) Principles for High-Quality Summative Assessment are used at each phase of the development cycle. As described in the Statewide Assessments Guide (FDOE, 2024), the processes of item development and test construction are carefully guided and include many quality-control measures such as the following:

- Item content on accommodated forms matches item content as administered online to the extent possible (e.g., wording, graphics, paragraph breaks, option order) via multiple rounds of content reviews. Note that some interactions will have accommodated form-specific language, such as equation and table match items. This additional language is needed to guide students on how to appropriately answer some items on accommodated forms.
- The student sees two-page items on an even then odd-numbered page simultaneously, just as they would see the entire item on one screen. Appropriate language is used for directives on the accommodated forms.

In terms of software that supports the item development process, CAI's Item Tracking System (ITS) served as the technology platform to efficiently carry out any item and test development process. ITS facilitated the creation of the item banks, item writing and revision, cataloging of changes and comments, and export of documents (items and passages). ITS enforced a structured review process, ensuring that every item that was written or imported underwent the appropriate sequence of reviews and signoffs; ITS archived every version of each item along with reviewer comments added throughout the process. ITS also provided sophisticated pool management features that increased item quality by providing real-time, detailed item inventories and item use histories. Because ITS could be configured to import items in multiple formats (e.g., Microsoft Word and Excel, XML), CAI was able to import items from multiple sources. To support online test delivery, ITS had a unique Web Preview feature that displayed items exactly as they were also presented to students, using the same program code used in CAI's Test Delivery System (TDS). An online test does not have a blueline (print approval) process like a paper-based test (PBT), and this feature provides an item-by-item blueline capability.

Before test administration, a series of user acceptance testing (UAT) is performed on all approved platforms to ensure that items are rendered as expected and have similar appearance across platforms to minimize potential device effects.

Rigorous review is in place to ensure that item content on accommodated forms matches item content as administered online (e.g., wording, graphics, paragraph breaks, option order).

The next section describes the item sources, and the subsequent sections outline the procedure used for the development and review of new items and the alignment of existing items.

3.1 SUMMARY OF ITEM SOURCES

Items for the spring 2024 assessments came from multiple sources as outlined here.

New Items Written by CAI

New field-test items were included in the spring 2024 item pool, and these items will be used on future test forms. The newly developed field-test items were written for the Florida-specific item bank. Mathematics and ELA items were written by CAI content experts or by trained partners. Pearson contracted item writers to create new items for Science and Social Studies. All items underwent a rigorous process of preliminary, editorial, and senior review by CAI and FDOE’s Test Development Center (TDC) content teams, who followed appropriate alignment, content, and style specifications. All items were also reviewed by panels of Florida educators and citizens for content accuracy, and to ensure that the test items were fair, unbiased, and included topics acceptable to the Florida public.

3.2 ITEM TYPES

One of the important features of online assessments is the administration of technology-enhanced items. Generally referred to as Machine-Scored Constructed-Response (MSCR) items, these include a wide range of item types. MSCR items require students to interact with the test content to select, construct, and/or support their answers.

Table 13–Table 15 list the item types and provide a brief description of each. For accommodated forms, some of these items must be modified or replaced with other items that assess the same standard and can be scanned and scored electronically. Please see the test design summary/blueprint documents or the test item specifications for specific details. Examples of various item types can be found in Appendix C, Example Item Types.

Table 13: ELA Reading Item Types and Descriptions

Response Type	Description
Multiple-Choice (MC)	Student selects one correct answer from a number of options.
Multiple-Select (MS)	Student selects all correct answers from a number of options.
Table Match (MI)	Student checks a box to indicate if information from a column header matches information from a row. On accommodated forms, the student fills in a bubble to indicate if information from a column header matches information from a row.
Hot-Text (HT)	Student is directed to either select or use the drag-and-drop feature to use text to support an analysis or make an inference. On accommodated forms, the student fills in bubbles to indicate which sentences are correct.
Multiple-Choice, Hot-Text Selectable (Two-part HT)	Student selects the correct answers from Part A and Part B. Part A is an MC or an MS item, and Part B is a selectable HT item.
Evidence-Based Selected-Response (EBSR)	Student selects the correct answers from Part A and Part B. Part A often asks the student to make an analysis or inference, and Part B requires the student to use text to support Part A.

Table 14: Mathematics and Mathematics EOC Item Types and Descriptions

Response Type	Description
Multiple-Choice (MC)	Student selects one correct answer from four options.

Response Type	Description
Multiple-Select (MS)	Student selects all correct answers from a number of options.
Edit Task Inline Choice (ETIC)	Student chooses the replacement for an incorrect number, word, phrase, or blank from a number of options. This includes items with one or more ETIC interactions. On accommodated forms, the student fills in a bubble to indicate the correct number, word, or phrase that should fill in the blank.
Grid (GI)	Student selects numbers, words, phrases, or images and uses the drag-and-drop feature to place them into a graphic. This item type may also require the student to use the point, line, or arrow tools to create a response on a graph.
Hot-Text (HT)	Student is directed to select text to support an analysis or make an inference. On accommodated forms, the student fills in bubbles to indicate which sentences are correct.
Equation (EQ)	Student uses a keypad with a variety of mathematical symbols to create a response. On accommodated forms, the student uses an empty response box to write in their answer.
Table Match (MI)	Student checks a box to indicate if information from a column header matches information from a row. On accommodated forms, the student is directed to fill in a bubble that matches a correct option from a column with a correct option from a row.
Multi-Interaction (MULTI)	This is an item that contains more than one response type. It could contain more than one of the same interaction type (except for multiple combinations of ETIC), or a combination of interaction types.

Table 15: Science and Social Studies Item Types and Descriptions

Response Type	Description
Multiple-Choice (MC)	Student selects one correct answer from four options.

3.3 COGNITIVE LABORATORIES

In a United States Department of Education (ED)-funded grant report investigating the accessibility of computerized assessments, Shaftel et al. (2015) point out that technology-enhanced items (TEIs) present greater accessibility barriers than traditional item types on accommodated tests, and that they should be examined to ensure that no construct-irrelevant variance is introduced. If some aspect of the technology impedes, or advantages, students in their responses to items, this could affect item responses and inferences regarding abilities on the measured construct.

Florida assessments are delivered by the same test delivery system as the Smarter Balanced Assessment Consortium (SBAC), therefore, research evidence on the SBAC platform can also be generalizable to Florida assessments. Two types of research were conducted for SBAC: (1) usability studies on system tools and features; and (2) cognitive lab studies evaluating the validity of various item types. Findings show that (1) various aspects of the test delivery system (e.g., tools, navigation, directions) provide students equitable access to the assessed content; and (2) TEI types do not introduce construct-irrelevant variance into scores. The full research report is provided in Volume 7, Special Studies, of the *Florida Standards Assessments 2014–2015 Technical Report*, which was included in an earlier submission for peer review.

FAST, B.E.S.T., and Science cognitive lab were conducted to examine the response processes of test takers for grades 3, 7, and 10 ELA, grades 3 and 7 Mathematics, Algebra 1, grades 5 and 8 Science, and Biology 1. These grades/courses were selected because they represent the item types, share similar blueprints (including the same content categories), and have the same test development procedures as the non-selected grades/courses. The assessments are all based on the same content standards and benchmarks, along with extensive content limits that define what is to be assessed. For all grades/courses, committees of educators collaborate with item development experts, assessment experts, and FDOE staff annually to review new and field-test items so that each test adequately samples the relevant domain of material the test is intended to cover. These committees review and verify the alignment of the test items with the content standards and measurement specifications so that the items measure the appropriate content. Given these commonalities between the selected and non-selected grades/courses, results from cognitive lab studies from the selected grades/courses are generalizable to non-selected grades/courses and non-selected item types.

In the studies, students worked through sample items. Eight students responded to each item, and their thinking processes were elicited through a combination of concurrent think-aloud (thinking out loud while reading and responding to an item) and focused probes that were tailored based on the anticipated solution path for a given item.

The cognitive lab interviews used recorded audio, and the students' responses to the test items were captured by the TDS. Following the cognitive lab, the interviewer reviewed all relevant information and files in a report that included, for each item attempted by the student, a detailed record of the student's think-aloud and responses to probes, as well as a record of the student's test item response.

These reports were evaluated by content experts to determine whether the evidence for any given item met the following criteria:

1. Students who receive full credit on an item display—through their think-aloud and responses to probes—defensible evidence that they based their response on the combination of skills and knowledge that make up the “intended construct.”
2. Students who do not receive full credit on an item display—through their think-aloud and responses to probes—defensible evidence that they understood (at a general level) what the item was asking them to do, and they were unable to provide a full-credit response as a result of deficiencies in one or more aspects of the skills or knowledge that make up the “intended construct.” For example, they lacked the necessary procedural knowledge for manipulating fractions or they were unable to apply the reasoning skills required by the item.

The studies were delayed due to the COVID-19 pandemic and school closings in 2020–2021. They were finally completed in 2024. In comparison to the intended cognitive complexity, it was found that the enacted cognitive complexity either met or exceeded the intended cognitive complexity in 58%–88% of the items. Evidence of linguistic complexity that was construct irrelevant was not found; however, students had significant difficulty reading algebra equations accurately, suggesting a focal point teachers should consider targeting during instruction. Study findings generalized across sampled grades. This study provides response process validity evidence that

assessment items measure the intended cognitive processes represented in the State’s academic content standards. The full findings can be found in the cognitive laboratory report in Appendix G.

3.4 ITEM TRANSLATIONS TO BRAILLE FORMAT

As is noted in Allman (2009), it is common that portions of a test may need to be modified to be translatable to braille format. Modifications may include substituting words, reformatting the layout of the item, and replacing untranslatable items with others of equal weight, content, and difficulty. As Winter (2010) acknowledges, this can pose a challenge to comparability, but this accommodation is needed for students with disabilities to properly demonstrate the knowledge, skills, and abilities the construct represents.

Florida uses a rigorous process, outlined in the *Florida Statewide Assessments Production Specifications* when creating the braille translations of its summative assessments and works with the Florida Instructional Materials Center for the Visually Impaired (FIMC-VI) and the American Printing House for the Blind (APH), both of which are leaders in the industry. Both FIMC-VI and APH follow practices determined by the Braille Authority of North America (BANA).

When forms are translated into braille, our contractors ensure that the braille forms match the regular print forms and make exceptions only when modifications for the braille reader are necessary. For instance, sometimes the item directions need to be modified for the braille reader instructing them to write in the letter instead of filling in the bubble. We also provide both Unified English Braille (UEB)-Nemeth and UEB-Technical versions of Mathematics and Science tests, and for all tests we provide both contracted and uncontracted versions to ensure that visually impaired students have the type of braille they read available. This means that in some cases, four braille transcriptions are made for each grade and subject: *UEB-Nemeth Uncontracted*, *UEB-Nemeth Contracted*, *UEB-Technical Uncontracted*, *UEB-Technical Contracted*. We ensure that the students who read braille are tested and challenged at the same level as their sighted peers. By working with FIMC-VI and APH, we ensure that all tests are reviewed and proofread by certified braille transcribers/proofreaders and teachers of the visually impaired who have vast experience and knowledge regarding students in this demographic. If modifications are made, a subject content specialist must approve any suggestions made by FIMC-VI and APH. Our content team ensures that the information vital to the item is retained in the braille format and that the student who reads braille is not given either an advantage or a disadvantage.

When transcribing pictures, cartoons, and graphics, images are either described or made in a tactual format for the braille reader or, with permission from content specialists, are sometimes omitted from the test if they do not provide any additional information. If graphics are described, we often use the descriptions currently created for text-to-speech, which all students have access to. If tactile graphics are created, they are kept as true to the original as possible. When deviation is needed, we comply with best practices in the field. Examples are as follows:

- Extraneous details such as decorative pictures, icons, or sections of a map that are not needed for the item are sometimes omitted—as the amount of information that can be interpreted through fingers is less than the amount of information the eye can process.

- Occasionally, especially with three-dimensional figures represented as two-dimensional drawings, graphics are too complex to be created tactfully and description alone either would not provide enough information or would give away the answer. In situations such as this, we develop manipulatives of the three-dimensional figures with specific directions to the Test Administrator on how to present them.

3.5 DEVELOPMENT AND REVIEW PROCESS FOR NEW ITEMS

3.5.1 Development of New Items

CAI developed field-test items to be embedded in the Florida Statewide Assessments operational tests. As part of the standard test development process, item writers followed the guidelines in FDOE’s approved test item specifications and the test design summary/blueprint.

CAI staff used the test item specifications to train qualified item writers, each of whom had prior item-writing experience. The item writers were trained at CAI item-writing workshops or had previous training on writing multiple-choice (MC) and constructed-response (CR) items. CAI content area assessment specialists worked with TDC content specialists to review measurement practices in item writing and interpret the meaning of the Florida Standards and benchmarks as illustrated by the test item specifications documents. This information, along with the purpose of the assessment, was explained to the item writers. Sample item stems that are included in the specification documents serve as models for the writers to use in creating items to match the Standards. To ensure that the items tapped the range of difficulty and taxonomic levels required, item writers use a method based on Webb’s cognitive demands (Webb, 2002) and DOK levels.

Item writing and passage selection were guided by the following principles for each of the item types. When writing items, item writers were trained to develop items that

- have an appropriate number of correct response options or combinations;
- contain plausible distractors that represent feasible misunderstandings of the content;
- represent the range of cognitive complexities and include challenging items for students performing at all levels;
- are appropriate for students in the assigned grade in terms of reading level, vocabulary, interest, and experience;
- are embedded in a real-world context, where indicated;
- do not provide answers or hints to other items in the set or test;
- are in the form of questions or directions for task completion;
- use clear language and avoid negative constructions unless doing so provides substantial advantages; and
- are free of ethnic, gender, political, socioeconomic, and religious bias.

Similarly, reading passages should

- represent literary (fiction), informational (nonfiction), and practical selections (e.g., nontraditional pieces, including tables, charts, glossaries, indices);

- provide students with the opportunity to interact with complex, authentic texts that may employ a variety of different structures;
- include multimedia elements when appropriate;
- be of high interest and appropriate readability for the grade level;
- be of appropriate length for the grade level;
- include topics that are in alignment with sensitivity guidelines;
- be free of ethnic, gender, political, and religious bias;
- not provide answers or hints to other items in the test; and
- include real-world texts (e.g., consumer or workplace documents, public documents such as letters to the editor, newspaper and magazine articles, thesaurus entries) to the extent possible.

When selecting passages, word count, readability, and text complexity are used in conjunction with other aspects of the passages (level of interest, accessibility of the topic, thematic elements) to determine appropriateness for a particular grade level. Table 16 provides the guidelines used in FAST ELA Reading.

Table 16: Word Counts and Readabilities of Reading Passages in FAST ELA Reading

Grade	Word Count (approximate)	Lexile Range (approximate)
3	100–700	450–900
4	100–900	770–1,050
5	200–1,000	770–1,050
6	200–1,100	955–1,200
7	300–1,100	955–1,200
8	350–1,200	955–1,200
9	350–1,300	1080–1,400
10	350–1,350	1080–1,400

In FAST ELA Reading, the texts are categorized as either informational or literary. *Informational texts* inform the reader and include the following types of publications:

- Exposition: informational trade books, news articles, historical documents, and essays
- Persuasive texts: speeches, essays, letters to the editor, and informational trade books
- Procedural texts and documents: directions, recipes, manuals, and contracts

Literary texts enable the reader to explore other people’s experiences or to simply read for pleasure and include the following genres:

- Narrative fiction: historical and contemporary fiction, science fiction, folktales, legends, and myths and fables
- Literary nonfiction: personal essays, biographies/autobiographies, memoirs, and speeches
- Poetry: lyrical, narrative, and epic works; sonnets, odes, and ballads

Department Item Review and Approval

After an internal review, the sets of items were reviewed by content specialists at the TDC. If needed, CAI, Pearson, and TDC content staff discussed requested revisions, ensuring that all items appropriately measured the Florida Standards. The items were then revised by CAI (for ELA and Mathematics) or Pearson (for Science and Social Studies) and brought to Florida bias, sensitivity, and content committees for review. After any final adjustments were made to the items, including an editorial review conducted by the TDC, the TDC provided a decision for each item: *Accept as Appears*, *Accept as Revised*, or *Reject*. Items that were approved by the TDC were subsequently web-approved and placed on field-test forms.

Committee Review of New Items

All items generated for use on Florida’s assessments were required to pass a series of rigorous reviews before they could appear as field-test items on operational test forms. The items were reviewed by three committees—the Bias Committee, the Community Sensitivity Committee, and the Content Item Review Committee.

The bias and sensitivity committees reviewed items for potential bias and controversial content. These committees consisted of Florida reviewers who were selected to ensure geographic and ethnic diversity. These committees ensure that items

- present racial, ethnic, and cultural groups in a positive light;
- do not contain controversial, offensive, or potentially upsetting content;
- avoid content familiar only to specific groups of students because of race or ethnicity, class, or geographic location;
- aid in the elimination of stereotypes; and
- avoid words or phrases that have multiple meanings.

The TDC and CAI reviewed the bias and sensitivity committees’ feedback and conveyed any issues to the attention of the Content Item Review Committee.

The Content Item Review Committee consisted of Florida classroom teachers or content specialists by grade for each subject area. The primary responsibility of the committee members was to review all new items to ensure that they were free from such flaws as (a) inappropriate readability level, (b) ambiguity, (c) incorrect or multiple answer keys (although some item types may include multiple answer keys by design), (d) unclear instructions, and (e) factual inaccuracy. These items were approved, approved with modifications, or rejected. Only approved items were added to the item pool for the field-test stage.

In addition, Science convenes an Expert Review Panel to confirm item accuracy and longevity and that best practices are represented. This review meeting takes place after bias, sensitivity, and content committees.

3.5.2 Rubric Validation

After items were field-tested, the rubric used for scoring MSCR items was validated by a team of grade-level Florida educators for ELA and Mathematics. These individuals reviewed the machine-assigned scores for CR items based on the scoring rubrics and either approved the scoring rubric as it appeared on the field test or suggested revisions to the scoring based on their interpretation of the item task and the rubric. The rubric validation meeting occurred in May 2024 in person in Tallahassee, Florida.

Similar to the items field-tested in previous years, rubrics were reviewed in one of two ways: items with simpler rubrics were reviewed via frequency tables of all student responses, while items with more complex rubrics were reviewed in 45-response samples.

Items with complex rubrics include grid (GI) items, hot-text (HT) draggable items, equation (EQ) items with full keypads, text entry natural language (NL) items, and multi-interaction (MULTI) items containing at least one of the preceding response types.

Items with simple rubrics include edit task choice and edit task inline choice (ETIC) items, HT selectable items, matching (MI) items, EQ items with simple numeric keypads, MC items, HT selectable (two-part HT) items, and any MULTI items comprised entirely of the preceding response types.

MC items, multiple-select (MS) items, and Evidence-Based Selected-Response (EBSR) items do not go through rubric validation.

Before the rubric validation meeting, CAI staff selected a sample of 45 student responses for each item with complex rubrics. The sample consisted of the following data:

- Fifteen responses from students who performed as expected on the item given their overall performance
- Fifteen responses from students who were predicted to perform well on the item given their overall performance, but instead performed poorly on the item
- Fifteen responses from students who were predicted to perform poorly on the item given their overall performance, but instead performed well on the item

For items with simple rubrics, CAI staff generated frequency tables that contained all student responses for each item. Frequency tables were generated out of CAI's Database of Record (DOR).

The Rubric Validation Committee reviewed 45 responses for every item with a complex rubric, having the option to approve the score or suggest a different score based on the committee's understanding of the rubric. For items with simple rubrics, the committee members were shown each item, along with the correct response and the most frequently selected incorrect responses. TDC and CAI staff ensured that the committee was scoring consistently. The committee meetings used the following procedures:

- All committee members were given a laptop allowing them to respond to the items the way a student would be able to respond in a live test.
- Each item was displayed with a projector.
- The committee discussed how to answer the item and how each point was earned.
- For items with complex rubrics, each of the 45 student responses and machine-assigned scores were displayed with a projector.
- For items with simple rubrics, the item was displayed with a projector, along with the correct response and the most frequently selected incorrect responses.
- If the committee members reached a consensus that a score was incorrect, the committee proposed modifications to the rubric.
- CAI rescored the responses using the revised rubric.
- CAI reviewed the responses that received changed scores to determine if they were correctly scored.
- The TDC reviewed the rescored responses and approved the rubric.

If any scores changed based on the Rubric Validation Committee review, CAI staff revised the machine rubric and rescored the item. After the item was rescored, CAI staff reviewed at least 10% of responses for which the score changed. This review ensured that committee suggestions were honored, that the item was scored consistently, and that no unintended changes in scoring occurred because of the revision to the machine rubric. CAI staff reviewed changes with TDC staff, and TDC staff had one final opportunity to revise the rubric or approve or reject the item.

At the end of the testing window, CAI conducted classical item analysis on these field-test items to ensure that the items functioned as intended with respect to the underlying scales. CAI's analysis program computed the required item and test statistics for each MC and CR item to check the integrity of the item and to verify the appropriateness of the difficulty level of the item. Key statistical analyses included item discrimination, distractor analysis, item difficulty analysis, and fit analysis. Details of these analyses are presented in Section 5, Item Analyses Overview, of Volume 1.

Field test items that survived statistical analysis were reviewed at the data review meeting. At the start of the data review meeting, CAI (ELA and Mathematics) and Pearson (Science and Social Studies) staff lead panels of educators in a formal training to familiarize them with the item development process, the purpose of data review, the meanings of statistical flags, and how a data review committee operates. The training included a review of the item cards, which detail specific item attributes (including grade level and alignment to the standards), the content and rubric of the item, and the various item statistics. Participants reviewed the items on laptops to interact with the item types as students do, and to view all statistics associated with each item. Data review was conducted using CAI's Content Rater system in the Item Tracking System (ITS). The Content Rater also enables educators to see the online "item card" stored for each item in the ITS. Item cards allow a viewer to see information regarding IRT statistics, fairness statistics, percent and average score of students in each point category, biserial correlations, and more. CAI can customize the Content Rater questions to the data review process and then use the system to capture individual educator decisions about the items. CAI facilitated group discussions of the item data once individual reviews were completed. CAI recorded notes from the group discussions of the

item data directly in ITS. Final decisions made by the committee were noted in ITS, and CAI and Pearson worked with FDOE to discuss any items on which the committee could not reach a consensus.

Additionally, FDOE reviewed the items based on the IRT statistics flagging criteria. The lists of items dropped from both the data review with educators and through discussions with FDOE were combined, and ITS was updated to disqualify these items from being added to the operational pool. Items that do pass both reviews are added to the operational pool for the next PM2/Winter administration for ELA and Mathematics, and during Spring for Science and Social Studies.

3.6 DEVELOPMENT AND MAINTENANCE OF THE ITEM POOL

As described earlier, new items are developed each year to be added to the operational item pool after being field-tested. Several factors determine the development of new items. The item development team conducts a gap analysis for distributions of items across multiple dimensions, such as item counts, item types, item difficulty, and numbers in each reporting category.

In spring 2024, field-test items were embedded in online forms in grades 3–10 ELA Reading, grades 3–8 Mathematics, Mathematics EOC tests, Science and Social Studies. An independent field test for ELA Writing in grades 4–10 was also conducted in spring 2023. All assessments were computer-adaptive tests (CATs) with a predetermined number and location of field-test items. Table 17–Table 19 provide the number of field-test items by type for ELA Reading, Mathematics, and Science and Social Studies. Table 20 provides the number of writing prompts for each grade.

Table 17: Number of ELA Reading Field-Test Items by Type

Item Type	3	4	5	6	7	8	9	10
EBSR	29	20	19	17	22	25	18	31
HT	20	17	6	7	7	7	14	15
MC	236	151	153	162	189	225	198	205
MI	30	12	13	14	18	17	17	19
MS	32	31	33	14	13	18	16	19
Two-Part HT	2	2	5	2	0	1	0	1

Table 18: Number of Mathematics and EOC Field-Test Items by Type

Item Type	3	4	5	6	7	8	Algebra 1	Geometry
EQ	60	128	93	102	109	113	66	50
ETIC	22	25	21	11	14	36	43	22
GI	7	4	2	6	3	12	8	7
HT	0	0	1	0	1	0	0	4
MC	79	128	82	72	99	140	79	42
MI	14	22	8	3	4	16	14	5
MS	39	25	14	21	5	14	42	11
Multi	9	12	10	11	17	17	26	12

Table 19: Number of Science and Social Studies Field-Test Items by Type

Item Type	Biology 1	Grade 5 Science	Grade 8 Science	Civics	U.S. History
MC	337	257	183	320	249

Table 20: Number of ELA Writing Field-Test Prompts by Grade in 2024

Grade	Number of Prompts
4	10
5	10
6	11
7	10
8	10
9	12
10	16

3.7 ALIGNMENT PROCESS FOR EXISTING ITEMS AND RESULTS FROM ALIGNMENT STUDIES

A third-party, independent alignment study was conducted in February 2016. This report can be found in Volume 4, Evidence of Reliability and Validity, Appendix D, FSA Alignment Report, of the *Florida Standards Assessments 2015–2016 Technical Report*.

A new third-party, independent alignment study for the new B.E.S.T. Standards is planned for 2025. For details see this volume’s Appendix E, ELA and Mathematics Alignment Study Proposal.

The new study will be designed to yield evidence pertaining to fulfillment of requirements as stated in federal statute related to the content alignment of statewide assessments with corresponding academic standards. Four main research questions will guide the work:

- 1. Framework Analysis:** To what extent do the CAT algorithms, test blueprints, and other relevant test specifications and documentation reflect structure and design that support the capacity of alignment of test events with corresponding grade-level academic standards?
- 2. Aggregate Data Review:** To what extent do the available aggregate data for test events administered in spring 2023 provide evidence that the algorithm and blueprints are yielding test forms as expected?
- 3. Validation of Internal Metadata:** To what extent is independent coding of assessment targets reasonably consistent with the assessment targets identified within internal (vendor) item metadata?

4. Test Form-Level Alignment: What is the degree of alignment of actual test events, sampled from below satisfactory, on grade level, and above satisfactory/mastery with corresponding Florida B.E.S.T. Standards, based on agreed-on criteria and minimum cutoffs?

The study will yield multiple lines of evidence that will support a validity argument that would extend across all test events generated by a computer-adaptive assessment program. Beyond the content alignment evidence for individual test events, it is important to provide additional evidence that can help extend findings across all test events generated by a particular testing program. Because CAT form assembly relies on internal metadata to meet blueprint specifications, validation of the internal metadata (based on independent item-level content analysis) allows for greater confidence that an assessment program can generate test forms that include content consistent with blueprint intent and, therefore, that test form-level findings can be reasonably generalized across all test forms generated by the assessment program. By drawing on multiple lines of evidence, the overall study design allows for the potential to craft a logic argument for the capacity for alignment of all test events generated by the FAST and EOC assessment programs included in the study with the corresponding Florida B.E.S.T. Standards, as appropriate, based on results.

The resulting logic argument, stated in the affirmative, would be as follows:

- If relevant test specifications and documentation reflect a structure and design to support the capacity of alignment of test events with corresponding grade-level academic standards, and
- if test events (sampled from below satisfactory, on grade level, and above satisfactory) meet minimum alignment criteria (based on agreed-on cutoffs for Categorical Concurrence, DOK/Cognitive Complexity Consistency, Range of Knowledge Correspondence, and Balance of Representation), and
- if the test blueprints and algorithm are generating test events as intended (based on data from all administered test events), and
- if validation of internal metadata supports generalizability of alignment findings across all test forms generated by the assessment programs,
- then it is possible to argue for the capacity for alignment for all test events resulting from Florida FAST assessments for ELA grades 3–10, FAST assessments for Mathematics grades 3–8, and B.E.S.T. EOC exams for Algebra I and Geometry with corresponding Florida B.E.S.T. Standards.

4. TEST CONSTRUCTION

4.1 OVERVIEW

During the 2022–2023 school year, the Florida Department of Education (FDOE) began transitioning from the fixed-form Florida Standards Assessment (FSA) to the computer-adaptive Florida Assessment of Student Thinking (FAST). In spring 2022, the first set of FAST items developed to align with the Benchmarks for Excellent Student Thinking (B.E.S.T.) Standards were field-tested. In summer 2022, field-test items were calibrated and placed on the FSA scale. Consistent with the PM1 and PM2 administrations, the spring 2023 FAST summative PM3 administration (as well as EOC Algebra 1 and Geometry) utilized CAI’s adaptive algorithm to administer tests using these pre-equated items on the FSA scale. During this transition year, scores were reported to students on the FSA scale.

Subsequently, calibrations in summer 2023 placed items in ELA Reading at grades 3–10, and Mathematics at grades 3–8, on a common vertical FAST scale via a linking design. EOC Algebra 1 and Geometry were placed on the B.E.S.T. scale. Standard settings were conducted for all grades in ELA Reading, Mathematics, ELA Writing, Algebra 1, and Geometry. In the 2023–2024 school year and beyond, FDOE reported scores on the new FAST scale.

In spring 2023–2024, Science and Social Studies also transitioned to CAT delivery. Calibrations in summer 2024 updated all item parameters, which were then linked back and placed on the original scale, utilizing existing cuts and performance levels.

In addition to the online CAT, Florida also has accommodated forms. Accommodated forms were administered to students instead of the online forms if such a need was indicated on their Individualized Education Program (IEP) or Section 504 Plan. For the Mathematics EOC assessments, Algebra 1 and Geometry, only one accommodated form was given. Accommodated forms used online parameters for scoring purposes and no calibrations were done on the accommodated forms.

4.2 ITEM SELECTION ALGORITHM

CAI’s adaptive algorithm takes as input two sources of information: an item pool and a test blueprint. The adaptive algorithm is then configured to execute maximally adaptive test administrations under the constraint of blueprint match. Configuration of the adaptive algorithm is critical because the composition of the item pool, which changes from administration to administration, interacts with the blueprint to influence the performance of the adaptive algorithm.

Item Pool

CAI’s ability to administer various state item pools is proven. For example, CAI administered items from the Smarter Balanced item bank during the 2013 pilot test and the 2014 field test. CAI designed and built the item renderers shared by the open-source version of the test delivery engine and CAI’s version of the item-rendering software. These renderers ensure that the items appear to students exactly as they did in the field test.

Test Blueprint

Test blueprints may contain specifications from the content hierarchy (strand, benchmark, standard, etc.) and other constraints, such as item type, or any other test item attribute that may be stored. CAI’s adaptive engine supports blueprints that meet the following conditions (which have been advocated by the Consortium for Citizens with Disabilities, an umbrella group encompassing most national advocacy groups for students with disabilities and other exceptional students):

1. Every student is tested on the full range of grade-level content, with no discernible differences in the content assessed.
2. Every student is tested on items measuring the same mix of cognitively complex skills, with no discernible difference—regardless of student proficiency.
3. Every student is tested on items reflecting the full range of other aspects of the grade-level curriculum as may be appropriate for the grade and subject.
4. Students are tested on items that provide the best measurement possible within these constraints.

These four principles ensure that every student can accurately demonstrate his or her academic skills and knowledge across the entire grade-level curriculum. CAI’s adaptive algorithm supports blueprints that align with these principles.

Item Selection

The adaptive algorithm, built on our partnerships with client states over the years, ensures that each student will receive a test that (1) matches the blueprint and (2) contains the items that best match students’ performance level, as defined by the blueprint. To accomplish this goal, the algorithm implements a highly parameterized multiple-objective utility function that includes

- a measure of the content match to the blueprint;
- a measure of overall test information; and
- measures of test information for each reporting category on the test.

We define an objective function that measures an item’s contribution to each of these objectives, weighting them to achieve the desired balance among them. The following equation sketches this objective function for a single item.

$$f_{ijt} = w_2 \left(\frac{\sum_{r=1}^R s_{rit} p_r d_{rj}}{\sum_{r=1}^R d_{rj}} \right) + w_1 \sum_{k=1}^K q_k h_{1k}(v_{kijt}, V_{knt}, t_k) + w_0 h_0(u_{ijt}, U_{it}, t_0)$$

Where the w terms represent user-supplied weights that assign relative importance to meeting each of the objectives, d_{rj} indicates whether item j has the blueprint-specified feature r , and p_r is the user-supplied priority weight for feature r . The term s_{rit} is an adaptive control parameter that is described in this section. In general, s_{rit} increases for features that have not met their designated minimum as the end of the test approaches.

The remainder of the terms represent an item’s contribution to measurement precision:

- v_{kijt} is the value of item j toward reducing the measurement error for reporting category k for test taker i at time of selection t ; and
- u_{ijt} is the value of item j in terms of reducing the overall measurement error for test taker i at time of selection t .

The terms U_{it} and V_{kit} represent the total information overall and on reporting category k , respectively.

The term q_k is a user-supplied priority weight associated with the precision of the score estimate for reporting category k . The t terms represent precision targets for the overall score (t_0) and each score reporting category score.

The functions $h(\cdot)$ are given by:

$$h_0(u_{ijt}, U_{it}, t_0) = \begin{cases} au_{ijt} & \text{if } U_{it} < t_0 \\ bu_{ijt} & \text{otherwise} \end{cases}$$
$$h_{1k}(v_{kijt}, V_{kit}, t_k) = \begin{cases} c_k v_{kijt} & \text{if } V_{kit} < t_k \\ d_k v_{kijt} & \text{otherwise} \end{cases}$$

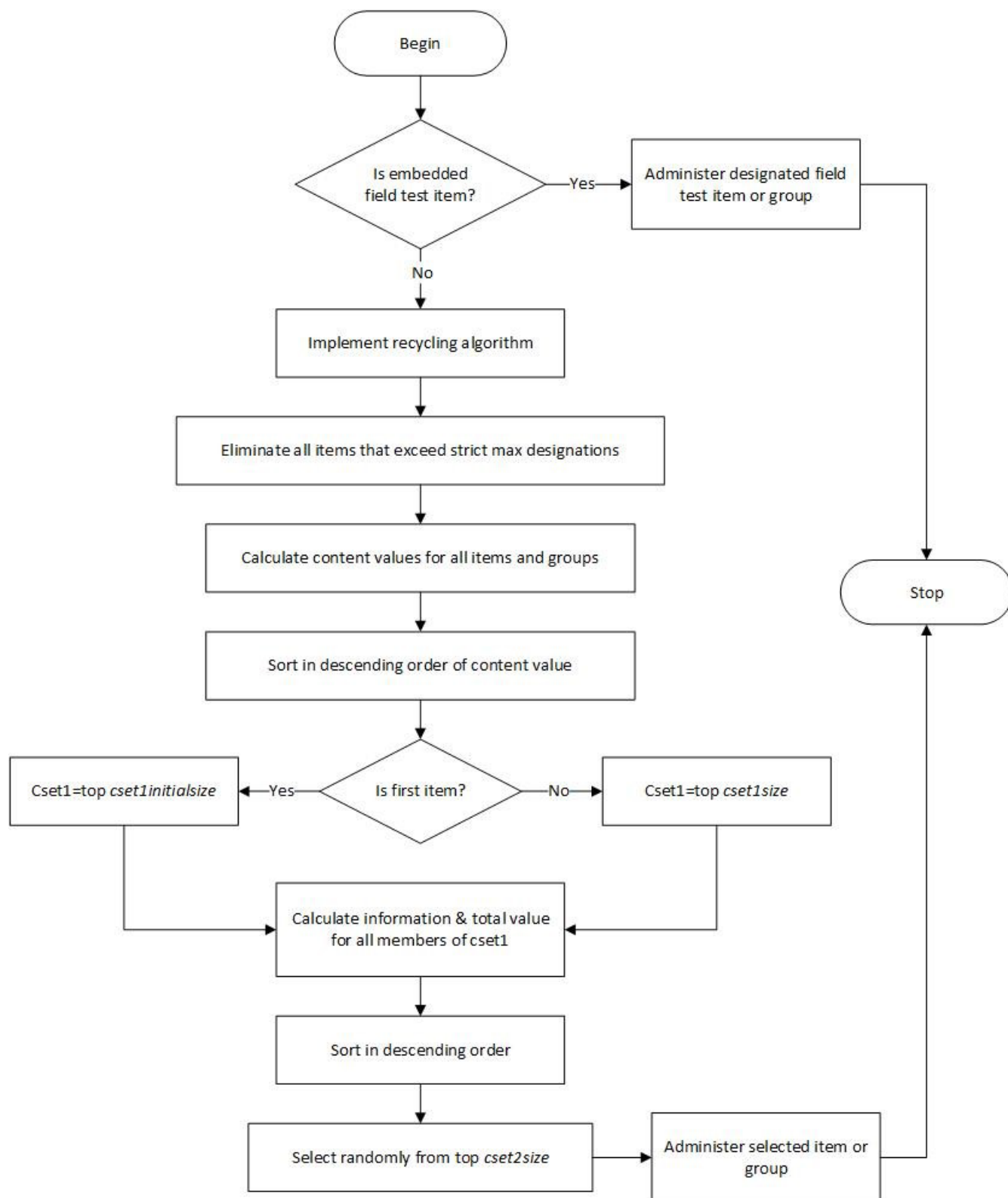
Items can be selected to maximize the value of this function. This objective function can be manipulated to produce a pure, standards-free adaptive algorithm by setting w_2 to zero or to produce a completely blueprint-driven test by setting $w_1 = w_0 = 0$. Adjusting the weights to optimize performance for a given item pool will enable users to maximize information subject to the constraint that the blueprint is virtually always met. We note that the computations of the content values and information values generate values on very different scales and that the scale of the content value varies as the test progresses. Therefore, we normalize both the information and content values before computing the value of the equation.

This normalization is given by $x = \begin{cases} 1 & \text{if } \min = \max \\ \frac{x - \min}{\max - \min} & \text{otherwise} \end{cases}$, where min and max represent the minimum and maximum, respectively, of the metric computed over the current set of items or item groups.

Items (or groups of items in the case of the ELA tests) are sorted by their “content value,” their value toward meeting the content constraints in the blueprint. Information measures are added to the content measures, and the items are sorted based on their overall value for the objective function. The final item selection is made based on a random selection from among the small subset of items that have the highest combined content and information value.

Figure 1 summarizes the item selection process. If the item position has been designated for a field-test item, then a field-test item is administered. Otherwise, the adaptive algorithm is triggered.

Figure 1: Item Selection Process



Blueprint Match

Configuration of the adaptive algorithm for the spring 2024 test administration was designed to administer tests meeting blueprint specifications while also maximizing test information to student ability for ELA and Mathematics tests. In the adaptive item-selection algorithm, item selection takes place in two discrete stages: blueprint satisfaction and match-to-ability. Due to the operational field-test design for the Science and Social Studies calibrations (for further details see Volume 1), the Science and Social Studies tests were only configured to meet the blueprint specifications.

While the simulation results described in the spring 2024 Simulation Summary Report (see Appendix F, Spring 2024 Simulation Results) indicated that the configuration resulted in the test administrations meeting all blueprint match requirements, it is also important to evaluate the blueprint match rate for the actual test administrations.

Appendix D, Spring 2024 Operational Item Blueprint Match, contains the operational item blueprint-match results for the spring 2024 grades 3–8 Mathematics, grades 3–10 ELA Reading, EOC Algebra 1 and Geometry, grades 5 and 8 Science, EOC Biology 1, U.S. History, and Civics. For the blueprint match analysis, only students who completed all parts of the test were included. As can be seen in Appendix D, in all assessments, all reporting categories met the blueprint or blueprint range. In addition to blueprint match, the observed percentage of reading passage types by grade is documented in Table 21.

In 2024 100% blueprint match was achieved for all reporting categories, including passage maximums for all grades. For some tests, item recycling in the adaptive algorithm was necessary. Eventually, with a large enough item pool, item recycling will not be necessary.

Table 21: Observed Spring 2024 Percentage of ELA Reading Passage Types by Grade

Grade	Informational	Literary
3	50%	50%
4	50%	50%
5	50%	50%
6	50%	50%
7	50%	50%
8	50%	50%
9	50%	50%
10	50%	50%

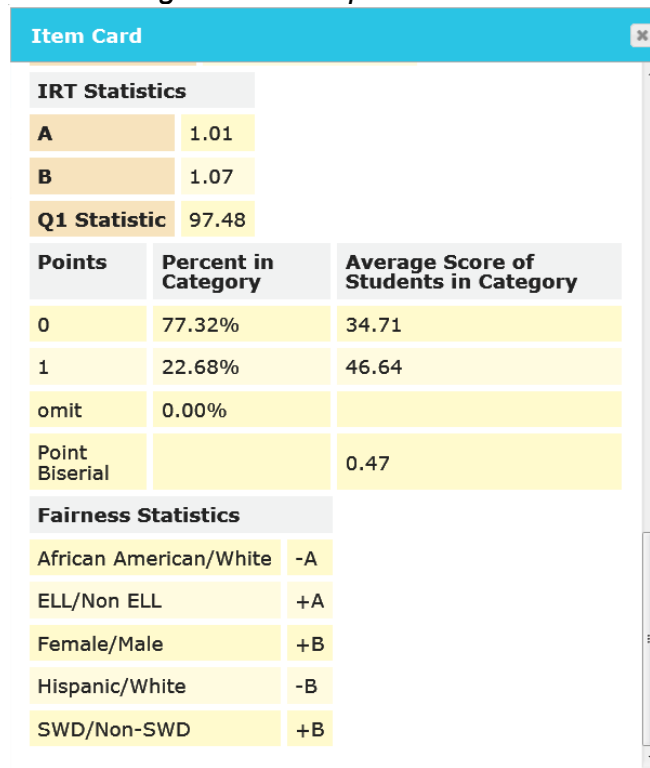
4.3 TEST CONSTRUCTION SUMMARY MATERIALS

4.3.1 Item Cards

Item cards, generated within CAI’s Item Tracking System (ITS), contained statistical information about an individual item. Item cards contained classical item statistics, item response theory (IRT)

statistics, and differential item functioning (DIF) statistics. When possible, item cards also contained a screen capture of the item. This was not possible in the case of some technology-enhanced items. In these instances, the items were viewed directly in ITS. Item cards were typically used to determine the viability of an individual field-test item for operational use in the next administration. Figure 2 shows one example of an item card.

Figure 2: Example Item Card



Item Card		
IRT Statistics		
A	1.01	
B	1.07	
Q1 Statistic	97.48	
Points	Percent in Category	Average Score of Students in Category
0	77.32%	34.71
1	22.68%	46.64
omit	0.00%	
Point Biserial		0.47
Fairness Statistics		
African American/White	-A	
ELL/Non ELL	+A	
Female/Male	+B	
Hispanic/White	-B	
SWD/Non-SWD	+B	

4.3.2 Bookmaps for Accommodated Forms

Bookmaps were only provided for accommodated forms. A bookmap is a spreadsheet that lists the characteristics of all items on a form. Bookmaps contain information such as the following:

- Item ID
- Item position
- Form
- Grade
- Role (e.g., operational or field-test)
- Item format (e.g., MC)
- Point value
- Answer key
- Reporting category
- Cognitive complexity

Bookmaps are used as an accessible resource by both content specialists and psychometricians to find information about a test form. Bookmaps differ from item cards in that there are no statistical summaries in a bookmap. Bookmaps contain useful information regarding the forms that are built in ITS.

4.4 ACCOMMODATION FORM CONSTRUCTION

Student scores should not depend on the mode of administration or type of test form. Because Florida’s tests were administered in an online test system, scores obtained via alternate modes of administration must be established as comparable to scores obtained through online testing. During test development, forms across all modes were required to adhere to the same test blueprints and content-level considerations. This section outlines the overall test development plans that ensured the comparability of online tests and accommodated tests.

To create the spring 2024 accommodated forms, CAI’s automated form-building tool inside CAI’s ITS, was used. ITS is a web-based software application that, in conjunction with the Item Authoring Tool (IAT), which can be accessed through ITS, enables users to create items and stimuli for testing purposes. Once the items and stimuli have been created, users can review, approve, and publish these items in ITS so that they can be administered to students through the Test Delivery System (TDS). ITS serves as an item repository and comprises several databases, referred to as item banks, that contain items specific to a client or project. It is an integrated system that supports all phases of test development. It facilitates direct item entry by item writers, online review and editing by reviewers, automated reporting for clients, and management and production of test forms.

Psychometric targets were set for each item and test form by CAI psychometricians in conjunction with FDOE, using the automated form-building tool inside ITS. Florida-accommodated fixed-form assessments have two distinct psychometric targets: (1) item-level targets and (2) test-level targets. Item-level targets are related to item statistics such as desired item difficulty, discrimination, parameter ranges (such as possibility of guessing), fit, and content bias. Test-level targets are related to test-level summaries and aggregated information, including desired average test difficulty, average item difficulty, target test information, standard error of measurement (SEM), and test characteristic curves. The nature of the item-level targets is predefined and does not change from one test administration to the next; however, the overall test targets might be updated with respect to policy needs and scale drift. All test forms must always meet blueprint targets. Items that failed the psychometric targets are flagged for review and removed if it is still possible to meet the blueprint without them.

For the test-level targets, the most important consideration is the SEM curve that shows the level of error of measurement expected at each ability level. The SEM is calculated as the reciprocal of the square root of the test information curve, and thus the SEM is lowest when information is highest. Ability estimates in the middle of the distribution often appear more reliable than the ability estimates at the high and low ends of the scale. The test construction always aims to minimize error at the score points at which relevant decisions (e.g., pass/fail) may be made.

Content specialists reviewed the forms and made any necessary item replacements, taking into account suitability for inclusion in an accommodated form and psychometric feedback. To build accommodated forms, content specialists began with the selected forms and removed any

technology-enhanced items (TEIs) that could not be rendered on accommodated forms or machine-scored. These items were then replaced with either MC items or other TEIs that could be rendered on accommodated forms from the same reporting category. In some instances, it was necessary to select replacement items from a different reporting category to satisfy statistical expectations; however, all parties ensured that each reporting category was still appropriately represented in the final test forms. Two of the forms with the best statistics were selected to be sent to FDOE for evaluation and selection of a final form.

5. REFERENCES

- Allman, C. (2009). *Test access: Making tests accessible for students with visual impairments: A guide for test publishers, test developers, and state assessment personnel* (4th ed.). American Printing House for the Blind. <https://sites.aph.org/wp-content/uploads/2017/09/Test-Access-Making-Tests-Accessible-2009.pdf>
- Florida Department of Education. (2024, August). *Florida Statewide Assessments Guide*. <https://www.fldoe.org/core/fileparse.php/20102/urlt/K12SAG.pdf>
- Shaftel, J., Benz, S., Boeth, E., Gahm, J., He, D., Loughran, J., Mellen, M., Meyer, E., Minor, E., & Overland, E. (2015). *Accessibility for Technology-Enhanced Assessments (ATEA): Report of project activities*. University of Kansas. <https://ateassessments.atlas4learning.org/sites/default/files/ATEA%20Project%20Report%20111615.pdf>
- Thompson, S. J., Johnstone, C. J., & Thurlow, M. L. (2002). *Universal design applied to large scale assessments* (Synthesis Report 44). University of Minnesota, National Center on Educational Outcomes. <https://nceo.umn.edu/docs/onlinepubs/synth44.pdf>
- Webb, N. L. (2002). *Alignment study in language arts, mathematics, science, and social studies of state standards and assessments for four states*. Council of Chief State School Officers.
- Winter, P. (2010). Comparability and test variations. In P. Winter (Ed.), *Evaluating the comparability of scores from achievement test variations* (pp. 1–11). Council of Chief State School Officers. <https://files.eric.ed.gov/fulltext/ED543067.pdf>