



New England Common Assessment Program

Science
2013–14
Technical Report
January 2015



100 EDUCATION WAY, DOVER, NH 03820 (800) 431-8901
WWW.MEASUREDPROGRESS.ORG

TABLE OF CONTENTS

CHAPTER 1	OVERVIEW	1
1.1	PURPOSE OF THE NEW ENGLAND COMMON ASSESSMENT PROGRAM.....	1
1.2	PURPOSE OF THIS REPORT.....	1
1.3	ORGANIZATION OF THIS REPORT.....	2
CHAPTER 2	CURRENT YEAR UPDATES.....	3
CHAPTER 3	TEST DESIGN AND DEVELOPMENT	4
3.1	SCIENCE TEST SPECIFICATIONS	4
3.1.1	Overview of Test Design	4
3.1.2	Standards (Assessment Targets)	4
3.1.3	Item Types	6
3.1.4	Test Design.....	6
3.1.5	Blueprint	7
3.1.6	Depth of Knowledge.....	7
3.1.7	Use of Calculators and Reference Sheets.....	8
3.2	TEST DEVELOPMENT PROCESS.....	9
3.2.1	Item Development.....	9
3.2.2	Item Reviews by Measured Progress	9
3.2.3	Item Reviews by the States	10
3.2.4	Bias and Sensitivity Review	11
3.2.5	Reviewing and Refining	11
3.2.6	Item Editing	11
3.2.7	Item Selection and Operational Test Assembly.....	11
3.2.8	Operational Test Draft Review.....	12
3.2.9	Alternative Presentations	13
CHAPTER 4	TEST ADMINISTRATION	14
4.1	RESPONSIBILITY FOR ADMINISTRATION	14
4.2	ADMINISTRATION PROCEDURES	14
4.3	PARTICIPATION REQUIREMENTS AND DOCUMENTATION	14
4.3.1	Students With Disabilities	15
4.3.2	Students With Limited English Proficiency	15
4.4	ADMINISTRATOR TRAINING.....	15
4.5	DOCUMENTATION OF ACCOMMODATIONS.....	16
4.6	TEST SECURITY.....	16
4.7	TEST AND ADMINISTRATION IRREGULARITIES	17
4.8	TEST ADMINISTRATION WINDOW	17
4.9	NECAP SERVICE CENTER.....	17
CHAPTER 5	SCORING.....	19
5.1	MACHINE-SCORED ITEMS	19

5.2	PERSON-SCORED ITEMS.....	19
5.3	INQUIRY TASK SCORING	20
5.4	SCORING LOCATION AND STAFF	20
5.4.1	Reader Recruitment and Qualifications	21
5.4.2	Reader Training	22
5.4.2.1	Anchor Set	22
5.4.2.2	Training Set.....	22
5.4.2.3	Qualifying Set.....	23
5.4.2.4	Retraining.....	23
5.4.3	QAC and SR Training	23
5.4.4	Benchmarking Meetings	24
5.5	METHODOLOGY FOR SCORING CONSTRUCTED-RESPONSE ITEMS	24
5.5.1	Monitoring of Scoring Quality Control and Consistency	26
5.5.1.1	Embedded CRRs	26
5.5.1.2	Read-Behind Procedures.....	27
5.5.1.3	Double-Blind Scoring	28
5.5.1.4	Recalibration Sets.....	28
5.5.2	Scoring Reports	29
CHAPTER 6	CLASSICAL ITEM ANALYSIS	31
6.1	CLASSICAL DIFFICULTY AND DISCRIMINATION INDICES.....	31
6.2	DIFFERENTIAL ITEM FUNCTIONING	33
6.3	DIMENSIONALITY ANALYSES.....	34
CHAPTER 7	SCALING AND EQUATING	38
7.1	ITEM RESPONSE THEORY.....	39
7.2	ITEM RESPONSE THEORY RESULTS.....	41
7.3	EQUATING.....	42
7.4	EQUATING RESULTS	43
7.5	ACHIEVEMENT STANDARDS.....	44
7.6	REPORTED SCALED SCORES	44
CHAPTER 8	RELIABILITY	47
8.1	RELIABILITY AND STANDARD ERROR OF MEASUREMENT	48
8.2	2013–14 SUBGROUP RELIABILITY.....	48
8.3	REPORTING SUBCATEGORY RELIABILITY.....	49
8.4	INTERRATER CONSISTENCY	49
8.5	RELIABILITY OF ACHIEVEMENT-LEVEL CATEGORIZATION.....	50
8.5.1	Accuracy and Consistency Results.....	51
CHAPTER 9	SCORE REPORTING	53
9.1	PRIMARY REPORTS	53
9.2	STUDENT REPORT	53
9.3	INTERACTIVE REPORTING	54
9.3.1	Item Analysis Roster Report	54
9.3.2	Achievement Level Summary	55
9.3.3	Released Items Summary Data Report	56

9.3.4	Longitudinal Data Report	56
9.4	SCHOOL, DISTRICT, AND STATE GRADE-LEVEL RESULTS REPORTS	56
9.5	DISTRICT AND STATE SUMMARY REPORTS.....	59
9.6	DECISION RULES.....	59
9.7	QUALITY ASSURANCE	59
CHAPTER 10	VALIDITY	61
10.1	QUESTIONNAIRE DATA.....	62
10.1.1	Difficulty of Assessment.....	62
10.1.2	Content.....	64
10.1.3	Homework.....	66
10.1.4	Performance in Science Class.....	67
REFERENCES	69
APPENDICES	71
APPENDIX A	GUIDELINES FOR THE DEVELOPMENT OF SCIENCE INQUIRY TASKS	
APPENDIX B	NECAP SCIENCE COMMITTEE MEMBERS	
APPENDIX C	PARTICIPATION RATES	
APPENDIX D	ACCOMMODATION FREQUENCIES	
APPENDIX E	NECAP TABLE OF STANDARD ACCOMMODATIONS	
APPENDIX F	ITEM-LEVEL CLASSICAL STATISTICS	
APPENDIX G	ITEM-LEVEL SCORE POINT DISTRIBUTIONS	
APPENDIX H	DIFFERENTIAL ITEM FUNCTIONING RESULTS	
APPENDIX I	ITEM RESPONSE THEORY CALIBRATION RESULTS	
APPENDIX J	TEST CHARACTERISTIC CURVE AND TEST INFORMATION FUNCTION CHARTS	
APPENDIX K	DELTA ANALYSES AND RESCORE ANALYSES	
APPENDIX L	α -PLOTS AND b -PLOTS	
APPENDIX M	ACHIEVEMENT LEVEL DISTRIBUTIONS	
APPENDIX N	RAW TO SCALED SCORE LOOKUP TABLES	
APPENDIX O	SCALED SCORE DISTRIBUTIONS	
APPENDIX P	CLASSICAL RELIABILITIES	
APPENDIX Q	INTERRATER CONSISTENCY	
APPENDIX R	DECISION ACCURACY AND CONSISTENCY RESULTS	
APPENDIX S	SAMPLE REPORTS	
APPENDIX T	INTERACTIVE REPORTS	
APPENDIX U	ANALYSIS AND REPORTING DECISION RULES	

CHAPTER 1 OVERVIEW

1.1 PURPOSE OF THE NEW ENGLAND COMMON ASSESSMENT PROGRAM

The New England Common Assessment Program (NECAP) is the result of collaboration among New Hampshire, Rhode Island, and Vermont to build a set of tests for grades 3 through 8 and 11 to meet the requirements of the No Child Left Behind (NCLB) Act. The specific purposes of the NECAP Science tests are (1) to provide data on student achievement in science at grades 4, 8, and 11 to meet NCLB requirements; (2) to provide information to support program evaluation and improvement; and (3) to provide information to parents and the public on the performance of students and schools. The tests are constructed to meet rigorous technical criteria, to include universal design elements and accommodations so that students can access test content, and to gather reliable student demographic information for accurate reporting. School improvement is supported by

- providing a transparent test design through the NECAP Science assessment targets, distributions of emphasis, and practice tests;
- reporting results by science domain, released items, and subgroups; and
- hosting test interpretation workshops and webinars to foster understanding of results.

Student-level results are provided to schools and families to be used as one element among all the collected evidence about progress and learning that occurred on the assessment targets for the respective grade span (K–4, 5–8, 9–11). The results are a status report of a student’s performance against the assessment targets, and they should be used cautiously in concert with local data.

1.2 PURPOSE OF THIS REPORT

The purpose of this report is to document the technical aspects of the 2013–14 NECAP Science tests. Students in grades 4, 8, and 11 participated in the seventh operational administration of NECAP Science in May 2014. This report provides evidence on the technical quality of those tests, including descriptions of the processes used to develop, administer, and score the tests, as well as the processes used to analyze the results. This report is intended to serve as a guide for replicating or improving the procedures in subsequent years.

Though some parts of this technical report may be used by educated laypeople, it is intended for experts in psychometrics and educational research. The report assumes a working knowledge of measurement concepts, such as reliability and validity, and statistical concepts, such as correlation and central tendency. In some chapters, the reader is presumed also to have basic familiarity with advanced topics in measurement and statistics.

1.3 ORGANIZATION OF THIS REPORT

The organization of this report is based on the conceptual flow of a test’s life span. The report begins with the initial test specifications and addresses all intermediate steps that lead to final score reporting. Chapters 1 through 5 give a description of NECAP Science by covering the test design and development process, the administration of the test, and scoring. Chapters 6 through 8 provide statistical and psychometric information, including chapters on scaling and equating, item analysis, and reliability. Chapter 9 is devoted to NECAP Science score reporting, and Chapter 10 is devoted to discussions on validity. Finally, the references cited throughout the report are provided, followed by the report appendices.

CHAPTER 2 CURRENT YEAR UPDATES

The assessment structure and guidelines for the 2013–14 NECAP Science administration remained consistent with previous administrations. There were no changes in assessment requirements or Test Blueprints.

CHAPTER 3 TEST DESIGN AND DEVELOPMENT

3.1 SCIENCE TEST SPECIFICATIONS

3.1.1 Overview of Test Design

The NECAP Science test is a criterion-referenced test. Items on the test are developed specifically for NECAP and are directly aligned to NECAP’s science assessment targets. These assessment targets guide the development of test items and are the basis for the reporting categories.

The 2013–14 NECAP Science test was administered in spring 2014 at grades 4, 8, and 11. The test consisted of four forms per grade. Each form included common items, equating items, and embedded field-test items. Common items are those that appear on every form of the test and are used to determine a student’s test score. Each equating item appears on one form only, and because these items have been on previous tests, they are used by psychometricians to keep the test scores on the same scale from year to year. This design provides reliable and valid results at the student level (the common items) and breadth of science coverage for school results (the common plus equating items) while minimizing testing time.

The NECAP Science test included an embedded field test. Each embedded field-test item generally appears on only one of the four forms. The field-test items were distributed equally among the forms. Since students do not know which items count for their test score, embedding field-test items into the operational test ensures that students take these items seriously. Because each field-test form is taken by approximately one-fourth of the NECAP students, the sample size is large enough to produce reliable field-test data. The embedded field test yields a pool of replacement items, which are needed due to the release of approximately twenty-five percent of the common items every year.

Each form of the test has three sessions. Physical Science, Earth Space Science, and Life Science are assessed in Sessions 1 and 2. These sessions contain common, equating, and embedded field-test items. Scientific Inquiry is assessed in Session 3 by an inquiry task. Session 3 contains only common items, as the inquiry tasks go through a separate field test (rather than an embedded field test).

3.1.2 Standards (Assessment Targets)

Although the NECAP Science assessment targets are unique for each grade, the assessment targets across the grades are classified under the same statements of enduring knowledge or broad areas of inquiry

that make up the four reporting categories of Physical Science, Earth Space Science, Life Science, and Scientific Inquiry.

Life Science

1. Survival of Organisms—All living organisms have identifiable structures and characteristics that allow for survival (organisms, populations, and species).
2. Matter and Energy in Ecosystems—Matter cycles and energy flows through an ecosystem.
3. Organisms Change over Time—Groups of organisms show evidence of change over time (structures, behaviors, and biochemistry).
4. Humans Are Similar Yet Unique—Humans are similar to other species in many ways and yet are unique among Earth’s life forms.

Physical Science

1. Properties and Structure of Matter—All living and nonliving things are composed of matter having characteristic properties that distinguish one substance from another (independent of size or amount of substance).
2. Energy—Energy is necessary for change to occur in matter. Energy can be stored, transferred, and transformed, but it cannot be destroyed.
3. Forces and Motion—The motion of an object is affected by forces.

Earth Space Science

1. Earth and Earth Materials—The Earth and its materials as we know them today have developed over long periods of time, through continual change processes.
2. Solar System—The Earth is part of a solar system made up of distinct parts that have temporal and spatial interrelationships.
3. Universe and Galaxies—The origin and evolution of galaxies and the universe demonstrate fundamental principles of physical science across vast distances and time.

Scientific Inquiry

1. Formulating Questions and Hypothesizing
2. Planning and Critiquing of Investigations
3. Conducting Investigations
4. Developing and Evaluating Explanations

3.1.3 Item Types

Since the beginning of the program, the goal of NECAP has been to measure what students know and are able to do by using a variety of test item types. The 2013–14 NECAP Science test consisted of standalone items and an inquiry task at each grade level. At grade 4 the inquiry task consisted of a hands-on investigation. At grades 8 and 11, the inquiry task consisted of a description of a scientific investigation.

The item types used and the functions of each are described below.

- **Multiple-choice (MC)** items were administered to provide breadth of coverage of the assessment targets. Because they require approximately one minute for most students to answer, these items make efficient use of limited testing time and allow coverage of a wide range of knowledge and skills. Each multiple-choice item was worth one point. Multiple-choice items were administered in Sessions 1 and 2 of the test in the Physical Science, Earth Space Science, and Life Science domains.
- **Short-answer (SA)** items assess students' skills and their abilities to work with brief, well-structured problems that have one solution or a very limited number of solutions. The advantage of this item type is that it requires students to demonstrate knowledge and skills by generating, rather than merely selecting, an answer. Each short-answer item was worth two points. Short-answer items require approximately two to five minutes for most students to answer. Short-answer items were employed in the Session 3 inquiry task.
- **Constructed-response (CR)** items typically require students to use higher-order thinking skills—evaluation, analysis, and summarization—in constructing a satisfactory response. Constructed-response items should take most students approximately five to ten minutes to complete. Four-point constructed-response items were administered in Sessions 1 and 2 of the test in the Physical Science, Earth Space Science, and Life Science domains. Three-point constructed-response items were administered in the Session 3 inquiry task.

At each grade, approximately twenty-five percent of the common standalone items and the entire inquiry task were released to the public. The released items are posted on a Web site hosted by Measured Progress and on the state agency/departments of education Web sites. Schools are encouraged to incorporate the use of released items in their instructional activities so that students will be familiar with them.

3.1.4 Test Design

Table 3-1 summarizes the numbers and types of common items that were used to determine student scores on the NECAP Science assessment. In addition, each test form contained 18 multiple-choice and 3 constructed-response standalone items that were either used for equating or field-testing.

Table 3-1. 2013–14 NECAP Science: Common Items

Grade	Session			Total			
	1	2	3	MC 1 pt	SA 2 pt	CR 3 pt	CR 4 pt
4	16 MC 2 CR	17 MC 1 CR	6 SA 2 CR	33	6	2	3
8	16 MC 2 CR	17 MC 1 CR	6 SA 2 CR	33	6	2	3
11	16 MC 2 CR	17 MC 1 CR	6 SA 2 CR	33	6	2	3

MC = multiple choice; SA = short answer; CR = constructed response

3.1.5 Blueprint

NECAP Science items are categorized into the following reporting categories: Physical Science, Earth Space Science, Life Science, and Scientific Inquiry. The item distribution across the reporting categories is shown in Table 3-2.

Table 3-2. 2013–14 NECAP Science: Distribution of Common Items

Domain	MC 1 pt	SA 2 pt	CR 3 pt	CR 4 pt
Physical Science	11			1
Earth Space Science	11			1
Life Science	11			1
Scientific Inquiry		6	2	
Total	33	6	2	3

MC = multiple choice; SA = short answer; CR = constructed response

Table 3-3 displays how the raw score points are distributed across the reporting categories.

Table 3-3. 2013–14 NECAP Science: Distribution of Score Points

Domain	Points	Percentage of Test
Physical Science	15	24%
Earth Space Science	15	24%
Life Science	15	24%
Scientific Inquiry	18	28%
Total	63	100%

3.1.6 Depth of Knowledge

Each item on the NECAP Science test is assigned a Depth of Knowledge (DOK) level. The DOK level reflects the complexity of mental processing students must use to answer an item. The DOK is not synonymous with difficulty. Each of the four DOK levels is described below.

- **Level 1 (Recall).** This level requires the recall of information such as a fact, definition, term, or simple procedure. These items require students only to demonstrate a rote response, use a well-known formula, or follow a set procedure.
- **Level 2 (Skill/Concept).** This level requires mental processing beyond that of recalling or reproducing a response. These items require students to make some decisions about how to approach the item.
- **Level 3 (Strategic Thinking).** This level requires reasoning, planning, and using evidence. These items require students to handle more complexity and abstraction than items at the previous two levels.
- **Level 4 (Extended Thinking).** This level requires planning, investigating, and complex reasoning over an extended period of time. Students are required to make several connections within and across content areas. This level may require students to design and conduct experiments. Due to the nature of this level, there are no level 4 items on the NECAP Science test.

It is important that the NECAP Science test measures a range of DOK levels. Table 3-4 lists the percentage of total score points assigned to each DOK level.

Table 3-4. 2013–14 NECAP Science: DOK Percentages

	<i>Grade 4</i>	<i>Grade 8</i>	<i>Grade 11</i>
DOK 1	11%	16%	16%
DOK 2	75%	70%	70%
DOK 3	14%	14%	14%

3.1.7 Use of Calculators and Reference Sheets

Science specialists from the New Hampshire and Rhode Island departments of education and the Vermont agency of education acknowledge that the use of calculators is a necessary and important skill. Calculators can save time and allow students to solve more sophisticated and intricate problems by reducing errors in calculations. For these reasons, it was decided that calculators should be permitted in all three sessions of the NECAP Science assessment. However, the state science specialists chose to prohibit scientific and graphing calculators in Session 3 because the inquiry task includes a graphing item.

A reference sheet is provided for the grade 8 and grade 11 tests. The grade 8 reference sheet includes solar system data, algebraic formulas used in science, the biological classification system, quantities with corresponding standard units of measure, a map of plate movements, and the periodic table of the elements. The grade 11 reference sheet includes genetic codes for amino acids, the electromagnetic spectrum, a map of plate movements, algebraic formulas used in science, and the periodic table of the elements.

3.2 TEST DEVELOPMENT PROCESS

3.2.1 Item Development

The item and inquiry task development process combined the expertise of the NECAP state science specialists, committees of NECAP educators, and Measured Progress test developers to help ensure items met the needs of the NECAP program. NECAP Science items are directly aligned to the assessment targets for each science domain. Each item addresses one assessment target. Measured Progress test developers worked with the state science specialists and teachers from the states to verify the alignment of items to the appropriate assessment target. All items used on the NECAP test were reviewed by a committee of NECAP teachers and by a NECAP bias and sensitivity committee.

The inquiry tasks were developed by Measured Progress test developers using *Guidelines for the Development of Science Inquiry Tasks*, a document created by the science specialists for each state agency/department of education. The document can be found in Appendix A of this report. The inquiry tasks were field-tested separately rather than as part of the embedded field test. Measured Progress test developers, the scoring content manager, and program managers observed field-testing of the inquiry tasks in schools located in Maine. The selected schools had varying demographics and population sizes. Each inquiry task was administered to approximately 200 students. Measured Progress test developers and program managers prepared a document titled *Inquiry Task Field Test Report* for the state science specialists to review. The state science specialists then approved the final form of each inquiry task. Due to space limitations, the *Inquiry Task Field Test Report* is not reproduced here. However, it can be obtained from any of the three NECAP states as a standalone document.

3.2.2 Item Reviews by Measured Progress

Measured Progress conducted two internal reviews of the multiple-choice and constructed-response items as well as a review of the inquiry tasks before presenting the items to the state science specialists. During these reviews, science test developers focused on three major areas.

- Item alignment to the assessment target: The reviewers considered whether the item measured the content as outlined in the assessment target and whether the content was grade appropriate. The reviewers also checked the DOK level of the item.
- Correctness of science content: The reviewers considered whether the information in the item was scientifically correct. For multiple-choice items, the keyed answer had to be the only correct answer. For constructed-response items, the scoring guide had to reflect correct science content and grade-level appropriate responses.
- Universal design: The reviewers considered item structure, clarity, possible ambiguity, and the appropriateness and relevance of graphics. For constructed-response items, the reviewers

considered whether the item adequately prompted an examinee to give a response similar to the one in the scoring guide.

3.2.3 Item Reviews by the States

The state science specialists reviewed the items. Measured Progress revised the items based on edits requested by the specialists.

Item review committees (IRC) composed of state teachers and curriculum supervisors were formed to conduct another review of the items. A list of the 2013–14 NECAP IRC participants for science in grades 4, 8, and 11 and their affiliations is included in Appendix B. On August 1–2, 2013, the NECAP Science IRC was held at the Stoweflake Mountain Resort & Spa in Stowe, Vermont. Participants' primary role was to evaluate and provide feedback on potential field-test items. For each grade level, the committee members reviewed potential multiple-choice and constructed-response field-test items as well as potential inquiry tasks. During the meeting, committee members were asked to evaluate the items for the following criteria:

- Assessment target alignment:
 - Is the test item aligned to the identified assessment target?
- DOK:
 - Are the items coded to the appropriate DOK level?
- Scientific correctness:
 - Are the items and distractors correct with respect to content and grade-level appropriateness?
 - Are the scoring guides consistent with the item and do they provide grade-level appropriate responses?
- Universal design:
 - Is the item language clear and grade appropriate?
 - Is the item language accurate (syntax, grammar, conventions)?
 - Is there an appropriate use of simplified language (is language that interferes with the assessment target avoided)?
 - Are charts, tables, and diagrams easy to read and understandable?
 - Are charts, tables, and diagrams necessary to the item?
 - Are instructions easy to follow?
 - Is the item amenable to accommodations—read aloud, signed, or in Braille?

3.2.4 Bias and Sensitivity Review

Bias review is an essential component of the development process. During the bias review process, NECAP Science items were reviewed by a committee of general education teachers, English language learner (ELL) specialists, special education teachers, and other educators and members of major constituency groups who represent the interests of legally protected and/or educationally disadvantaged groups. A list of bias and sensitivity review committee participants and affiliations is included in Appendix B. Items were examined for issues that might offend or dismay students, teachers, or parents. Including such educators in the development of assessment items and materials can avoid many unduly controversial issues, and unfounded concerns can be allayed before the test forms are produced.

3.2.5 Reviewing and Refining

After the IRC and bias and sensitivity review committee meetings, Measured Progress test developers and the state science specialists met to review the committees' feedback. The specialists decided what edits should be made to the items.

3.2.6 Item Editing

Measured Progress editors then reviewed and edited the items to ensure uniform style (based on *The Chicago Manual of Style*, 15th edition) and adherence to sound testing principles. These principles included the stipulation that items were

- correct with regard to grammar, punctuation, usage, and spelling;
- written in a clear, concise style;
- written at a reading level that allows the student to demonstrate his or her knowledge of science, regardless of reading ability;
- written in a way that did not cue the correct answer (for multiple-choice options); and
- free of potentially sensitive content.

3.2.7 Item Selection and Operational Test Assembly

In preparation for the item selection meeting with the state science specialists, test developers and psychometricians at Measured Progress considered the following when selecting sets of items to propose for the common (including items for release) and the embedded field tests:

- **Content coverage/match to test design.** The test design stipulates a specific number of multiple-choice and constructed-response items from each content area. Item selection for the embedded field test was based on the number of items in the existing pool of items eligible for the common item set.

- **Item difficulty and complexity.** Item statistics drawn from the data analysis of previously tested items were used to ensure similar levels of difficulty and complexity from year to year as well as quality psychometric characteristics.
- **“Cueing” items.** Items were reviewed for any information that might “cue” or provide information that would help to answer another item.

At the item selection meeting, the state specialists reviewed the proposed sets of items and made the final selection of items for the common item set, including which items would be released after the test was administered. The state specialists also made the final selection of items for the embedded field test and approved the final wording of these items.

During assembly of the test forms, the following criteria were considered:

- **Option balance.** Items were balanced among the forms so that each form contained a fairly equal distribution of keys (correct answers).
- **Key patterns.** The sequence of keys was reviewed to ensure that key order appeared random.
- **Page fit.** Item placement was modified to ensure the best fit and arrangement of items on any given page.
- **Facing-page issues.** For multiple items associated with a single stimulus (inquiry task) and multiple-choice items with large graphics, consideration was given to whether those items needed to begin on a left- or right-hand page and to the nature and amount of material that needed to be placed on facing pages. These considerations serve to minimize the amount of page flipping required of students.
- **Relationship between forms.** Although equating and field-test items differ across forms, these items must take up the same number of pages in each form so that sessions begin on the same page in every form. Therefore, the number of pages needed for the longest form often determines the layout of each form.
- **Visual appeal.** The visual accessibility of each page of the form was taken into consideration, including such aspects as the amount of white space, the density of the text, and the number of graphics.

3.2.8 Operational Test Draft Review

After the forms were laid out as they would appear in the final test booklets, they were again thoroughly reviewed by Measured Progress editors and test developers to ensure that the items appeared exactly as the state science specialists had requested. Finally, all the forms were reviewed by the state science specialists for their final approval.

3.2.9 Alternative Presentations

Common items for grades 4, 8, and 11 were translated into Braille by a subcontractor that specializes in test materials for students who are blind or visually impaired. In addition, Form 1 for each grade was also adapted into a large-print version.

CHAPTER 4 TEST ADMINISTRATION

4.1 RESPONSIBILITY FOR ADMINISTRATION

The 2014 NECAP Science *Principal/Test Coordinator Manual* indicated that principals and/or their designated NECAP test coordinators were responsible for the proper administration of the NECAP Science test. The *Test Administrator Manual*, which contained explicit directions and read-aloud scripts, was used to ensure the uniformity of administration procedures from school to school.

4.2 ADMINISTRATION PROCEDURES

Principals and/or their schools' designated NECAP coordinators were instructed to read the *Principal/Test Coordinator Manual* before testing and familiarize themselves with the instructions provided in the *Test Administrator Manual*. The *Principal/Test Coordinator Manual* provided each school with checklists to help them prepare for testing. The checklists outlined tasks to be performed by school staff before, during, and after test administration. Besides these checklists, the *Principal/Test Coordinator Manual* described the testing material being sent to each school and how to inventory the material, track it during administration, and return it after testing was complete. The *Test Administrator Manual* included checklists for the administrators to ready themselves, their classrooms, and the students for the administration of the test. It also contained sections detailing the procedures to be followed for each test session and instructions for preparing the material before its return to Measured Progress.

4.3 PARTICIPATION REQUIREMENTS AND DOCUMENTATION

For New Hampshire, Rhode Island, and Vermont, the intent of No Child Left Behind (NCLB) legislation is for *all* students in grades 4, 8, and 11 to participate in science testing through standard administration, administration with accommodations, or alternate assessment. Furthermore, any student who is absent during any session of the NECAP Science test is expected to make up the missed sessions within the three-week testing window.

Schools were required to return a student answer booklet for every enrolled student in the grade level. On those occasions when it was deemed impossible to test a particular student, school personnel were required to inform their state agency/department of education. The states included a grid on the student answer booklets listing the approved reasons why a student answer booklet could be returned blank for one or more sessions of the science test.

- **Student withdrew from school after May 5, 2014.** If a student withdrew after May 5, 2014, but before completing all of the test sessions, school personnel were instructed to code this reason on the student’s answer booklet.
- **Student enrolled in school after May 5, 2014.** If a student enrolled after May 5, 2014, and was unable to complete all of the test sessions before the end of the testing administration window, school personnel were instructed to code this reason on the student’s answer booklet.
- **State-approved special consideration.** Each state agency/department of education had a process for documenting and approving circumstances that made it impossible or not advisable for a student to participate in testing. Schools were required to obtain state approval before beginning testing.
- **Student was enrolled on May 5, 2014, and did not complete test for reasons other than those listed above.** If a student was not tested for a different reason, school personnel were instructed to code this reason on the student’s answer booklet. These “other” categories were considered not state approved.

Appendix C lists the science participation rates of the three states combined.

4.3.1 Students With Disabilities

All students were expected to participate in the 2014 NECAP Science tests, unless they completed the state-specific alternate assessment in New Hampshire, Rhode Island, or Vermont during the 2013–14 school year.

Form 1 of the grades 4, 8, and 11 science tests was enlarged to 20-point font for students with visual impairments. At all three grades, only the common items were translated into Braille by American Printing House for the Blind, a subcontractor that specializes in test materials for students who are blind or have visual impairments.

4.3.2 Students With Limited English Proficiency

Students who were new to the United States after October 1, 2013, and were designated as limited English proficient (LEP) were required to take the 2014 NECAP Science test.

4.4 ADMINISTRATOR TRAINING

In addition to distributing the *Principal/Test Coordinator Manual* and *Test Administrator Manual*, the New Hampshire and Rhode Island departments of education and Vermont agency of education, along with Measured Progress, conducted test administration workshops or a webinar to inform school personnel about the NECAP Science test and to provide training on the policies and procedures regarding administration.

4.5 DOCUMENTATION OF ACCOMMODATIONS

Though every effort was made to provide a test that would be as accessible as possible, a need still remained to allow some students to take the test with accommodations. An operating principle employed during the development of the accommodations protocols and policy development was to allow only accommodations that would not change the construct of what was being measured by the item.

The *Principal/Test Coordinator Manual* and *Test Administrator Manual* provided directions for coding the information related to accommodations and modifications on page 2 of the student answer booklet. All accommodations used during any test session were required to be coded by authorized school personnel—not students—after testing was completed.

The *NECAP Accommodations Guide* also provides detailed information on planning and implementing accommodations. This guide can be found on each state’s agency/department of education Web site. The states collectively made the decision that accommodations be made available to all students based on individual need, regardless of disability status. Decisions regarding accommodations were made by the students’ educational teams on an individual basis and were consistent with those used during the students’ regular classroom instruction. Making accommodations decisions on an entire group basis rather than on an individual basis was not permitted. If the decision made by a student’s educational team required an accommodation not listed in the state-approved NECAP Table of Standard Accommodations, schools were instructed to contact their agency/department of education in advance of testing for specific instructions for coding the “Other Accommodations (O)” and/or “Modifications (M)” sections.

Table 4-1 shows the accommodations observed for the May 2014 NECAP Science administration. Appendix D shows the breakdown of students who tested with which accommodation. The accommodation codes are defined in the NECAP Table of Standard Accommodations, found in Appendix E.

Table 4-1. 2013–14 NECAP Science: Number of Students Tested With and Without Accommodations by Grade

Grade	Number of Students Tested	
	Without Accommodations	With Accommodations
4	24,452	6,391
8	25,902	4,867
11	26,337	3,176

4.6 TEST SECURITY

Maintaining test security is critical to the success of NECAP and the continued partnership among the three states. The *Principal/Test Coordinator Manual* and *Test Administrator Manual* explain in detail all test security measures and test administration procedures. School personnel were informed that any concerns about breaches in test security were to be reported to the test coordinator and principal immediately. The test

coordinator and/or principal were responsible for immediately reporting the concern to the district superintendent and the state director of testing at the agency/department of education. Test security was strongly emphasized in the test administration workshops and webinars conducted for all three states. The states required the principal of each school that participated in testing to log on to a secure Web site to complete the Principal's Certification of Proper Test Administration form for each grade level tested. The principal was required to provide the number of secure tests received from Measured Progress, the number of tests administered to students, and the number of secure test materials being returned to Measured Progress. The principal was then required to enter his or her name in the online form as an electronic signature. By signing the form, the principal was certifying that the tests were administered according to the procedures outlined in the *Principal/Test Coordinator Manual* and *Test Administrator Manual*, that he or she maintained the security of the test materials, that no secure material was duplicated or in any way retained in the school, and that all test materials had been accounted for and scheduled for return to Measured Progress.

4.7 TEST AND ADMINISTRATION IRREGULARITIES

During the preliminary review of results by districts and schools, several Rhode Island districts questioned the scoring of Inquiry Task items on which the majority of students received a score of 0, including students who performed at the Proficient level or higher on the test. The questions regarding scoring were often associated with general questions regarding the accuracy of reported declines in test performance from 2013 to 2014. Prior to the release of results, the Rhode Island Department of Education coordinated an internal and external review of both the equating procedures and results and the scoring of those Inquiry Task items. The review found no systematic errors in either the equating or scoring that would impact the 2014 NECAP Science Results. Results were released without any changes. As a result of the review, Measured Progress and the states discussed potential adjustments to the item development and benchmarking process to better distinguish between student responses that are totally incorrect or irrelevant and those that demonstrate some relevant knowledge understanding.

4.8 TEST ADMINISTRATION WINDOW

The test administration window was May 5–22, 2014.

4.9 NECAP SERVICE CENTER

To provide additional support to schools before, during, and after testing, Measured Progress established the NECAP Service Center. The additional support that the service center provided was an essential element to the successful administration of the three-state test program. Individuals in the field could call the centralized location using a toll-free number and ask questions or report any problems they were experiencing.

The service center was staffed based on call volume and was available from 8:00 a.m. to 4:00 p.m. beginning two weeks before the start of testing and ending two weeks after testing. The representatives were responsible for receiving, responding to, and tracking calls and then routing issues to the appropriate person(s) for resolution.

CHAPTER 5 SCORING

Upon receipt of used NECAP Science answer booklets following testing, Measured Progress scanned all student responses, along with student identification and demographic information. Imaged data for multiple-choice items were machine scored. Images of constructed-response items were processed and organized by iScore, secure server-to-server electronic scoring software designed by Measured Progress, for hand-scoring.

Student responses that could not be physically scanned (e.g., answer documents damaged during shipping) were physically reviewed and scored on an individual basis by trained, qualified readers. These scores were linked to the student’s demographic data and merged with the student’s scoring file by Measured Progress’s data processing department.

5.1 MACHINE-SCORED ITEMS

Multiple-choice responses were compared to scoring keys using item analysis software. Correct answers were assigned a score of 1 point; incorrect answers were assigned a score of 0 points. Student responses with multiple marks or blank responses were also assigned 0 points.

The hardware elements of the scanners monitored themselves continuously for correct reads, and the software driving these scanners monitored the correct data reads. Standard checks included recognition of a sheet that did not belong, was upside down, or was backward; identification of missing critical data, including a student ID number or test form that was out of range or missing; and identification of page/document sequence errors. When a problem was detected, the scanner stopped and displayed an error message directing the operator to investigate and correct the situation.

5.2 PERSON-SCORED ITEMS

The images of student responses to constructed-response items were hand-scored through the iScore system. Using iScore minimized the need for readers to physically handle actual answer booklets and related scoring materials. Student confidentiality was easily maintained, as all NECAP Science scoring was blind (i.e., district, school, and student names were not visible to readers). The iScore system maintained the link between the student response images and their associated test booklet numbers.

Through iScore, qualified readers accessed electronically scanned images of student responses at computer terminals. The readers evaluated each response and recorded each student’s score via keypad or mouse entry through the iScore system. When a reader finished one response, the next response immediately appeared on the computer screen.

Imaged responses from all answer booklets were sorted into item-specific groups for scoring purposes. Readers reviewed responses from only one item at a time; however, when necessary, imaged responses from a student’s entire booklet were available for viewing, and the physical booklet was also available to the on-site chief reader.

The use of iScore also helped ensure that access to student response images was limited to only those who were scoring or who were working for Measured Progress in a scoring management capacity.

5.3 INQUIRY TASK SCORING

Of special interest during this cycle of scoring the 2013–14 NECAP Science test was implementing the scoring requirements associated with inquiry task items. These items were unique in that students conducted a single scientific experiment and then answered the eight assigned questions about that experiment. The questions were designed to stand alone, meaning that each one could be scored separately, instead of as part of a set of several combined questions. This maximized the number of readers that could be assigned to score responses for each student.

5.4 SCORING LOCATION AND STAFF

Scoring Location

The iScore database, its operation, and its administrative controls are all based in Dover, New Hampshire; in addition, all 2013–14 NECAP Science responses were scored in Longmont, Colorado (grade 8 inquiry task items) and Dover (all operational grades as well as grades 4 and 11 inquiry task items).

The iScore system monitored accuracy, reliability, and consistency across the scoring site. Constant daily communication and coordination were accomplished through email, telephone, and secure Web sites, to ensure that critical information and scoring modifications were shared and implemented throughout the scoring site.

Staff Positions

The following staff members were involved with scoring the 2013–14 NECAP Science responses:

- The NECAP Science scoring project manager—an employee of Measured Progress based in Dover, New Hampshire—oversaw the communication and coordination of scoring constructed-response items.
- The iScore operational manager and iScore administrators—employees of Measured Progress based in Dover, New Hampshire—coordinated technical communication pertaining to the scoring of constructed-response items.
- A chief reader in science ensured the consistency of scoring across the scoring site for all grades tested. The chief reader—an employee of Measured Progress based in Dover, New Hampshire—also provided read-behind activities for quality assurance coordinators (QACs).

- Numerous QACs, selected from a pool of experienced senior readers (SRs) for their ability to score accurately and their ability to instruct and train readers, participated in benchmarking activities for each grade. QACs provided read-behind activities for SRs. The ratio of QACs and SRs to readers was approximately 1 to 11.
- Numerous SRs, selected from a pool of skilled and experienced readers, provided read-behind activities for the readers at their scoring tables (2 to 12 readers at each table).
- Readers at the scoring sites scored the 2013–14 NECAP Science operational and field-test student responses.

5.4.1 Reader Recruitment and Qualifications

For scoring of the 2013–14 NECAP Science test, Measured Progress actively sought a diverse scoring pool that was representative of the population of the three participating NECAP Science states. The broad range of readers included scientists, editors, business professionals, authors, teachers, graduate school students, and retired educators. Demographic information for readers (e.g., gender, race, educational background) was electronically captured and reported.

Although a four-year college degree or higher was preferred for all readers, readers of the NECAP Science test responses of grades 4, 8, and 11 students were required to have successfully completed at least two years of college and to have a demonstrated knowledge of science. This permitted the recruitment of readers who were currently enrolled in a college program, a sector of the population that had relatively recent exposure to classroom practices and current trends in their field of study. In all cases, potential readers submitted documentation (e.g., resume and/or transcripts) of their qualifications.

Table 5-1 summarizes the qualifications of the 2013–14 NECAP Science scoring leadership (QACs and SRs) and readers.

Table 5-1. 2013–14 NECAP Science: Qualifications of Scoring Leadership and Readers

<i>Scoring Responsibility</i>	<i>Spring 2013 Administration Educational Credentials</i>				<i>Total</i>
	<i>Doctorate</i>	<i>Master's</i>	<i>Bachelor's</i>	<i>Other</i>	
Scoring leadership	3.5%	37.9%	48.3%	10.3%*	100.0%
Readers	5.2%	27.7%	58.0%	9.1%**	100.0%

* Indicates the 1 QAC/SR with an associate's degree and the 2 QAC/SRs with at least 48 college credits.

** Indicates the 8 readers with associate's degrees and the 13 readers with at least 48 college credits.

Readers were either temporary Measured Progress employees or were secured through the services of one or more temporary employment agencies. All readers signed a nondisclosure/confidentiality agreement.

5.4.2 Reader Training

Reader training began with an introduction of on-site scoring staff and an overview of the NECAP Science program’s purpose and goals, including a discussion about the security, confidentiality, and proprietary nature of testing, scoring materials, and procedures.

Next, readers thoroughly reviewed and discussed the scoring guide for the item to be scored. Each item-specific scoring guide included the item itself and score point descriptions.

Following review of the item-specific scoring guide for any constructed-response item, readers began reviewing or scoring response sets organized for specific training purposes:

- Anchor set
- Training set
- Qualifying set

During training, readers were able to highlight or mark hard copies of the anchor and training sets, even if all or part of the sets was also presented online via computer.

5.4.2.1 Anchor Set

Readers first reviewed an anchor set of exemplary responses, approved by the state science specialists, for the item to be scored. Responses in anchor sets were typical rather than unusual or uncommon; solid rather than controversial or borderline; and true, meaning that they had scores that could not be changed by anyone other than the NECAP Science client and Measured Progress test development staff.

For constructed-response items, each item-specific anchor set contained, for each respective score point, a client-approved sample response that was to be considered a midrange example of its respective score point. When necessary, a second sample response was included to illustrate a substantially alternate way to achieve that score point.

Responses were read aloud to the room of readers and presented in descending score order. Trainers then announced the true score of each anchor response and facilitated a group discussion of the response in relation to the score point descriptions to allow readers to internalize typical characteristics of each score point.

This anchor set served as a reference for readers as they continued with calibration, scoring, and recalibration activities for that item.

5.4.2.2 Training Set

Next, readers practiced applying the scoring guide and anchors to responses in the training set. The training set typically included 10 to 15 student responses designed to help establish the score point range and the range of responses within each score point. The training set often represented unusual responses that were

less clear or solid (e.g., were shorter than normal, employed atypical approaches, contained both very low and very high attributes, or included difficult handwriting). Responses in the training set were presented in randomized score point order.

After readers had independently read and scored a training set response, trainers polled readers or used online training system reports to record the initial range of scores. Then they led a group discussion of one or two responses, directing reader attention to scoring issues that were particularly relevant to the specific scoring group, such as the line between two score points. Trainers modeled for readers how to discuss scores by referring to the anchor set and scoring guides.

5.4.2.3 Qualifying Set

After the training set had been completed for an item, readers were required to measurably demonstrate their ability to accurately and reliably score all responses in the item's qualifying set, according to the appropriate anchor set in concert with its scoring rubric, by scoring the qualifying set. The qualifying set consisted of 10 responses selected from an array of responses that clearly illustrated the range of score points for that item. The set was chosen in accordance with the responses reviewed and approved by the state specialists.

To be eligible to score operational 2013–14 NECAP Science responses, readers were required to demonstrate scoring accuracy rates of at least 80 percent exact agreement and at least 90 percent exact or adjacent agreement across all qualifying set responses. In other words, exact scores were required on at least eight of the qualifying set responses and either exact or adjacent scores were required on at least nine. Readers were allowed one discrepant score as long as they had at least eight exact scores.

5.4.2.4 Retraining

Readers who did not pass the first qualifying set were, at the discretion of scoring leadership, retrained as a group by reviewing their performance with scoring leadership and then scoring a second qualifying set of responses. If they achieved a minimum scoring accuracy rate of 80 percent exact and 90 percent exact or adjacent agreement on this second set, they were allowed to score operational responses.

If readers did not achieve the required scoring accuracy rates on the second qualifying set, they were not allowed to score responses for that item. Instead, they were either trained on a different item or dismissed from scoring for that day.

5.4.3 QAC and SR Training

QACs and select SRs were trained in a separate training session that occurred immediately prior to reader training. In addition to discussing the items and their responses, QAC and SR training included

emphasis on the states’ rationale behind the score points. This rationale was discussed in greater detail with QACs and SRs than with regular readers to better equip leadership to handle questions from the readers.

5.4.4 Benchmarking Meetings

In preparation for implementing NECAP Science guidelines for the scoring of field-test responses, Measured Progress scoring staff prepared and facilitated benchmarking meetings held with the NECAP state science specialists. The purpose of the meetings was to establish item-specific guidelines for scoring each NECAP Science item for the current field-test scoring session and for future operational scoring sessions.

Prior to these meetings, scoring staff collected a set of several dozen student responses that chief readers identified as being illustrative midrange examples of their respective score points. The chief readers and science specialists worked collaboratively during benchmarking meetings to finalize an authoritative set of score point exemplars for each field-test item. As a matter of practice, each of these authoritative sets is included as part of the scoring training materials and used to train readers each time that item is scored—both as a field-test item and as part of a future NECAP Science administration.

This repeated use of approved sets of midrange score point exemplars helps ensure that, each time a particular NECAP Science item is scored, readers follow the guidelines established by the state science specialists.

5.5 METHODOLOGY FOR SCORING CONSTRUCTED-RESPONSE ITEMS

Constructed-response items were scored based on possible score points and scoring procedures, as shown in Table 5-2.

Table 5-2. 2013–14 NECAP Science: Possible Score Points for Open-Response Items

<i>Item Type</i>	<i>Possible</i>	
	<i>Score Points</i>	<i>Highest Score</i>
Constructed-response	0–4	4
Inquiry task—constructed-response	0–3	3
Inquiry task—short-answer	0–2	2
Nonscorable	0	0

Nonscorable Items

Readers could designate a response as nonscorable for any of the following reasons:

- Response was blank (no attempt to respond to the question).
- Response was unreadable (illegible, too faint to see, or only partially legible/visible).

- Response was written in the wrong location (seemed to be a legitimate answer to a different question).¹
- Response was written in a language other than English.
- Response was completely off task or off topic.
- Response included an insufficient amount of material to make scoring possible.
- Response was an exact copy of the assignment.
- Response was incomprehensible.
- Student made a statement refusing to write a response to the question.

Scoring Procedures

Scoring procedures for constructed-response items included both single scoring and double scoring. Single-scored items were scored by one reader. Double-scored items were scored independently by two readers whose scores were tracked for agreement (known as interrater agreement). For further discussion of double scoring and interrater agreement, see Subsection 5.5.1.3.

Table 5-3 shows the method(s) by which common and equating constructed-response items for each operational test were scored.

Table 5-3. 2013–14 NECAP Science: Methods of Scoring Common and Equating Constructed-Response Items by Grade and Test

<i>Grade</i>	<i>Test/Field Test Name</i>	<i>Responses</i>	
		<i>Single Scored*</i>	<i>Double Scored*</i>
4	Science	100%	2% randomly
8	Science	100%	2% randomly
11	Science	100%	2% randomly
All	Unreadable responses	100%	100%
All	Blank responses	100%	100%

* Per grade and test/field test.

For each field-test item, 1,500 student responses were scored.

¹ Unreadable and wrong-location responses were eventually resolved, whenever possible, by researching the actual answer document (electronic copy or hard copy, as needed) to identify the correct location or to more closely examine the response and then assign a score.

5.5.1 Monitoring of Scoring Quality Control and Consistency

Readers were monitored for continued accuracy rates and scoring consistency throughout the scoring process, using the following methods and tools:

- Embedded committee-reviewed responses (CRRs)
- Read-behind procedures
- Double-blind scoring
- Recalibration sets
- Scoring reports

If readers met or exceeded the expected accuracy rate, they continued scoring operational responses. Any reader who fell below the expected accuracy rate for the particular item and monitoring method was retrained on that item and, upon approval by the QAC or chief reader as appropriate, was allowed to resume scoring.

It is important to note the difference between the accuracy rate each reader must have achieved to qualify for scoring live responses and the accuracy rate each reader must have maintained to continue scoring live responses. Specifically, the qualification accuracy rate was stricter than the live scoring accuracy rate. The reason for this difference is that an “exact score” in double-blind statistics requires that two readers both identify the same score for a response; an exact score during qualification requires that an individual reader match the score predefined by leadership. Thus, the latter is dependent on matching an expert, not a peer.

During live scoring, reader accuracy rates are monitored using an array of techniques, thereby providing a more complete picture of a reader’s performance than would be the case if relying on just one technique. These techniques are described in the next subsections.

5.5.1.1 *Embedded CRRs*

Previously scored CRRs were selected and loaded into iScore for blind distribution to readers as a way to monitor accuracy. Embedded CRRs, either chosen before scoring had begun or selected by leadership during scoring, were inserted into the scoring queue so as to be indistinguishable from all other live student responses.

Between 5 and 30 embedded CRRs were distributed at random points throughout the first full day of scoring an item to ensure that readers were sufficiently calibrated at the beginning of the scoring period. Individual readers often received up to 20 embedded CRRs within the first 100 responses scored, and up to 10 CRRs within the next 100 responses scored on the first day of scoring that item.

If any reader fell below the required live scoring accuracy rate, he or she was retrained before being allowed by the QAC to continue. Once the reader was allowed to resume scoring, leadership carefully monitored him or her by increasing the number of read-behinds.

5.5.1.2 Read-Behind Procedures

Read-behind scoring refers to the practice of scoring leadership, usually an SR, scoring a response after a reader has already scored it.

Responses to be placed into the read-behind queue were randomly selected by scoring leadership; readers were not made aware as to which of their responses would be reviewed by their SR. The iScore system allowed one, two, or three responses per reader at a time to be placed into the read-behind queue.

The SR entered his or her score into iScore before being allowed to see the score assigned by the reader for whom the read-behind was being performed. The SR then compared the two scores, and the ultimate reported score was determined as follows:

- If there was exact agreement between the scores, no action was taken; the regular reader’s score remained.
- If the scores were adjacent (i.e., the difference was not greater than 1), the SR’s score became the score of record. If there was a significant number of adjacent scores for this reader across items, an individual scoring consultation was held with the reader, and the QAC determined whether or when the reader could resume scoring.
- If there was a discrepant difference between the scores (greater than 1 point), the SR’s score became the score of record. An individual consultation was held with the reader, with the QAC determining whether or when the reader could resume scoring.

These three scenarios are illustrated in Table 5-4.

Table 5-4. 2013–14 NECAP Science: Examples of Read-Behind Scoring Resolutions

<i>Reader</i>	<i>QAC/SR Resolution</i>	<i>Final*</i>
4	4	4
4	3	3
4	2	2

* QAC/SR score is score of record.

Approximately 2.5 percent of all student responses were reviewed by QACs and SRs as read-behinds. In cases where a reader’s scoring rate fell below the required accuracy percentage, QACs and SRs conducted additional read-behinds for that reader.

In addition to the daily read-behinds, scoring leadership could choose to read behind any reader at any point during the scoring process and thereby take an immediate, real-time snapshot of a reader’s accuracy.

5.5.1.3 Double-Blind Scoring

Double-blind scoring refers to the practice of two readers independently scoring a response, each without knowing the response had already been or soon would be scored by another reader. Table 5-3 provides information about the proportion of responses that were double scored.

If there was a discrepancy (a difference greater than 1) between scores, the response was placed in an arbitration queue. Arbitration responses were reviewed by scoring leadership (SR or QAC) without any background knowledge of the scores assigned by the two previous readers.

Scoring leadership consulted individually with any reader whose scoring rates on the different monitoring methods fell below the required accuracy percentage, and the QAC determined whether or when the reader could resume scoring. Once the reader was allowed to resume scoring, leadership carefully monitored him or her by increasing the frequency of read-behinds.

5.5.1.4 Recalibration Sets

To determine whether readers were still calibrated to the scoring standard, readers were required to take an online recalibration set at the start and midpoint of the shift of their resumption of scoring.

Each recalibration set consisted of five responses representing the entire range of possible scores, including some with a score point of 0.

- Readers who were discrepant on two of five responses of the first recalibration set, or exact on two or fewer, were not permitted to score on that item that day and were either assigned to a different item or dismissed for the day.
- Readers, who were discrepant on only one of five responses of the first recalibration set, and/or exact on three, were retrained by their SR by discussing the recalibration set responses in terms of the score point descriptions and the original anchor set. After this retraining, such readers began scoring operational responses under the proviso that the reader's scores for that day and that item would be kept only if the reader was exact on all five of five responses of the second recalibration set administered at the shift midpoint. The QAC determined whether or when these readers had received enough retraining to resume scoring operational responses. Scoring leadership also carefully monitored the accuracy of such readers by significantly increasing the number of their read-behinds.
- Readers who were not discrepant on any response of the first recalibration set, and exact on at least four, were allowed to begin scoring operational responses immediately, under the proviso that this recalibration performance would be combined with that of the second recalibration set administered at the shift midpoint.

The results of both recalibration sets were combined with the expectation that readers would have achieved an overall 80 percent-exact and 90 percent-adjacent standard for that item for that day.

The scoring project manager voided all scores posted on that item for that day by readers who did not meet the accuracy requirement. Responses associated with voided scores were reset and redistributed to readers with demonstrated accuracy for that item.

Recalibration sets were employed for all constructed-response items. They were not used for 2-point short-answer items, for which read-behind and double-blind techniques are more informative and cost effective.

5.5.2 Scoring Reports

Measured Progress’s electronic scoring software, iScore, generated multiple reports that were used by scoring leadership to measure and monitor readers for scoring accuracy, consistency, and productivity.

Reports Generated During Scoring

Because the 2013–14 NECAP Science test administration was complex, computer-generated reports were necessary to ensure all of the following:

- Overall group-level accuracy, consistency, and reliability of scoring
- Immediate, real-time individual reader data availability for early reader intervention when necessary
- Scoring schedule maintenance

The following reports were produced by iScore:

- The **Read-Behind Summary** showed the total number of read-behind responses for each reader and noted the numbers and percentages of scores that were exact, adjacent, and discrepant between that reader and the SR or QAC. Scoring leadership could choose to generate this report by selecting options such as “Today,” “Past Week,” or “Cumulative” from a pull-down menu. The report could also be filtered to display data for a particular item or across all items. This report was used in conjunction with other reports to determine whether a reader’s scores would be voided (i.e., sent back out to the floor to be rescored by other readers). The benefit of this report is that it measures the degree to which individual readers agree with their QAC or SR on how to best score live responses.
- The **Double-Blind Summary** showed the total number of double-score responses scored by each reader and noted the numbers and percentages of scores that were exact, adjacent, and discrepant between that reader and the second reader. This report was used in conjunction with other reports to determine whether a reader’s scores would be voided. The benefit of this report is that it reveals the degree to which readers are in agreement with each other about how to best score live responses.
- The **Accuracy Summary** combined read-behind and double-score data, showing the total number of read-behind and double-score responses scored for each reader and noting his or her accuracy percentages and score point distributions.
- The **Embedded CRR Summary** showed, for each reader and for either a particular item or across all items, the total number of responses scored, the number of embedded CRRs scored, and the numbers and percentages of scores that were exact, adjacent, and discrepant between the reader and the chief reader (by virtue of the chief reader’s approval of the prescored embedded CRRs). This report was used in conjunction with other reports to determine

whether a reader's scores would be voided. The benefit of this report is that it measures the degree to which individual readers agree with their chief reader on how to best score live responses—and since embedded responses are administered during the first hours of scoring, this report provides an early indication of agreement between readers and their chief reader.

- The **Qualification Statistics Report** listed all readers by name and ID number, identifying which qualifying set(s) they did and did not take and, for the ones they did take, whether they passed or failed. The total number of qualifications passed and failed was noted for each reader, as was the total number of individuals passing or failing a particular qualifying set. The QAC could use this report to determine how the readers within his or her specific scoring group performed on a specific qualifying set.
- The **Summary Report** showed the total number of student responses for an item and identified, for the time at which the report was generated, (1) the number of single and double scorings that had been performed and (2) the number of single and double scorings yet to be performed.

CHAPTER 6 CLASSICAL ITEM ANALYSIS

As noted in Brown (1983), “A test is only as good as the items it contains.” A complete evaluation of the quality of a test must include an evaluation of each item. Both *Standards for Educational and Psychological Testing* (AERA et al., 2014) and *Code of Fair Testing Practices in Education* (2004) include standards for identifying quality items. Items should assess only knowledge or skills that are identified as part of the domain being tested and should avoid assessing irrelevant factors. Items should also be unambiguous and free of grammatical errors, potentially insensitive content or language, and other confounding characteristics. In addition, items must not unfairly disadvantage students, in particular racial, ethnic, or gender groups.

Both qualitative and quantitative analyses are conducted to ensure that NECAP Science test items meet these standards. Qualitative analyses are described in earlier chapters of this report; this chapter focuses on quantitative evaluations. Statistical evaluations are presented in four parts: (1) difficulty indices, (2) item-test correlations, (3) differential item functioning (DIF) statistics, and (4) dimensionality analyses. The item analyses presented here are based on the statewide administration of the NECAP Science test in spring 2014. Note that the information presented in this chapter is based on the items common to all forms, since those are the items on which student scores are calculated. (Item analyses are also performed for field-test items, and the statistics are then used during the item review process and form assembly for future administrations.)

6.1 CLASSICAL DIFFICULTY AND DISCRIMINATION INDICES

All multiple-choice and constructed-response items are evaluated in terms of item difficulty according to standard classical test theory practices. *Difficulty* is defined as the average proportion of points achieved on an item and is measured by obtaining the average score on an item and dividing it by the maximum possible score for the item. Multiple-choice items are scored dichotomously (correct versus incorrect); so, for these items, the difficulty index is simply the proportion of students who correctly answered the item. Polytomously scored items include short-answer, for which students can receive scores of 0, 1, or 2, and constructed-response, which are worth either 3 or 4 points total. By computing the difficulty index as the average proportion of points achieved, the indices for the different item types are placed on a similar scale, ranging from 0.0 to 1.0 regardless of the item type. Although this index is traditionally described as a measure of difficulty, it is properly interpreted as an *easiness* index, because larger values indicate easier items. An index of 0.0 indicates that all students received no credit for the item, and an index of 1.0 indicates that all students received full credit for the item.

Items that are answered correctly by almost all students provide little information about differences in student abilities, but they do indicate knowledge or skills that have been mastered by most students. Similarly, items that are correctly answered by very few students provide little information about differences in student

abilities, but they may indicate knowledge or skills that have not yet been mastered by most students. In general, to provide the best measurement, difficulty indices should range from near-chance performance (0.25 for four-option multiple-choice items or essentially 0 for constructed-response items) to 0.90, with the majority of items generally falling between around 0.4 and 0.7. However, on a standards-referenced assessment such as NECAP Science, it may be appropriate to include some items with very low or very high item difficulty values to ensure sufficient content coverage.

A desirable characteristic of an item is for higher-ability students to perform better on the item than lower-ability students do. The correlation between student performance on a single item and total test score is a commonly used measure of this characteristic of the item. Within classical test theory, the item-test correlation is referred to as the item’s discrimination, because it indicates the extent to which successful performance on an item discriminates between high and low scores on the test. For constructed-response items, the item discrimination index used was the Pearson product-moment correlation; for multiple-choice items, the corresponding statistic is commonly referred to as a point-biserial correlation. The theoretical range of these statistics is 0.0 to 1.0, with a typical observed range from 0.2 to 0.6.

Discrimination indices can be thought of as measures of how closely an item assesses the same knowledge and skills assessed by other items contributing to the criterion total score. That is, the discrimination index can be thought of as a measure of construct consistency.

A summary of the item difficulty and item discrimination statistics for each grade is presented in Table 6-1. Note that the statistics are presented for all items as well as by item type (multiple-choice and constructed-response). The mean difficulty and discrimination values shown in the table are within generally acceptable and expected ranges.

Table 6-1. 2013–14 NECAP Science: Summary of Item Difficulty and Discrimination Statistics by Grade

Grade	Item Type	Number of Items	p-Value		Discrimination	
			Mean	Standard Deviation	Mean	Standard Deviation
4	ALL	44	0.60	0.20	0.35	0.08
	CR	5	0.36	0.11	0.49	0.07
	MC	33	0.69	0.14	0.33	0.07
	SA	6	0.32	0.12	0.35	0.06
8	ALL	44	0.59	0.18	0.38	0.11
	CR	5	0.40	0.19	0.54	0.10
	MC	33	0.63	0.15	0.35	0.08
	SA	6	0.50	0.22	0.45	0.08
11	ALL	44	0.53	0.17	0.38	0.14
	CR	5	0.36	0.06	0.60	0.02
	MC	33	0.57	0.17	0.32	0.09
	SA	6	0.43	0.06	0.54	0.05

A comparison of indices across grade levels is complicated because these indices are population dependent. Direct comparisons would require that either the items or students were common across groups. Since that is not the case, it cannot be determined whether differences in performance across grade levels are because of differences in student abilities, differences in item difficulties, or both. With this caveat in mind, it appears generally that students in grades 8 and 11 found their items more difficult than students in grade 4.

Comparing the difficulty indices of multiple-choice items and open-response (short-answer or constructed-response) items is inappropriate because multiple-choice items can be answered correctly by guessing. Thus, it is not surprising that the difficulty indices for multiple-choice items tend to be higher (indicating that students performed better on these items) than the difficulty indices for open-response items. Similarly, discrimination indices for the open-response items were larger than those for the dichotomous items because of the greater variability of the former (i.e., the partial credit these items allow) and the tendency for correlation coefficients to be higher given greater variances of the correlates.

In addition to the item difficulty and discrimination summaries presented above, item-level classical statistics and item-level score point distributions were also calculated. Item-level classical statistics are provided in Appendix F; item difficulty and discrimination values are presented for each item. The item difficulty and discrimination indices are within generally acceptable and expected ranges. Very few items were answered correctly at near-chance or near-perfect rates. Similarly, the positive discrimination indices indicate that students who performed well on individual items tended to perform well overall. There was a small number of items with low discrimination indices, but none was negative. While it is not inappropriate to include items with low discrimination values or with very high or very low item difficulty values to ensure that content is appropriately covered, there were very few such cases on the NECAP Science test. Item-level score point distributions are provided for open-response items in Appendix G; for each item, the percentage of students who received each score point is presented.

6.2 DIFFERENTIAL ITEM FUNCTIONING

Code of Fair Testing Practices in Education (2004) explicitly states that subgroup differences in performance should be examined when sample sizes permit and that actions should be taken to ensure that differences in performance are because of construct-relevant, rather than irrelevant, factors. *Standards for Educational and Psychological Testing* (AERA et al., 2014) includes similar guidelines. As part of the effort to identify such problems, NECAP Science items were evaluated in terms of DIF statistics.

For NECAP Science, the standardization DIF procedure (Dorans & Kulick, 1986) was employed to evaluate subgroup differences. The standardization DIF procedure is designed to identify items for which subgroups of interest perform differently, beyond the impact of differences in overall achievement. The DIF procedure calculates the difference in item performance for two groups of students (at a time) matched for achievement on the total test. Specifically, average item performance is calculated for students at every total

score. Then an overall average is calculated, weighting the total score distribution so that it is the same for the two groups.

When differential performance between two groups occurs on an item (i.e., a DIF index in the low or high categories, explained below), it may or may not be indicative of item bias. Course-taking patterns or differences in school curricula can lead to DIF, but for construct-relevant reasons. On the other hand, if subgroup differences in performance could be traced to differential experience (such as geographical living conditions or access to technology), the inclusion of such items should be reconsidered.

Computed DIF indices have a theoretical range from -1.0 to 1.0 for multiple-choice items, and the index is adjusted to the same scale for open-response items. Dorans and Holland (1993) suggested that index values between -0.05 and 0.05 should be considered negligible. The preponderance of NECAP Science items fell within this range. Dorans and Holland further stated that items with values between -0.10 and -0.05 and between 0.05 and 0.10 (i.e., “low” DIF) should be inspected to ensure that no possible effect is overlooked and that items with values outside the -0.10 to 0.10 range (i.e., “high” DIF) are more unusual and should be examined very carefully.²

For the 2013–14 NECAP Science tests, seven subgroup comparisons were evaluated for DIF:

- Male versus female
- No disability versus disability
- Noneconomically disadvantaged versus economically disadvantaged
- NonLEP versus LEP
- White versus Asian
- White versus Black
- White versus Hispanic

The tables in Appendix H present the number of items classified as either low or high DIF, overall and by group favored.

6.3 DIMENSIONALITY ANALYSES

The NECAP Science tests were each designed to measure and report a single score on science achievement using a unidimensional scale. Thus, each of these tests is said to be measuring a single dimension, and the term “unidimensionality” is used to describe it.

Because each test is constructed with multiple content area subcategories and item types, and their associated knowledge and skills, the subtests associated with each of these could potentially result in a large number of secondary dimensions being invoked beyond the primary dimension that all the items on a test

² It should be pointed out here that DIF for items is evaluated initially at the time of field-testing. If an item displays high DIF, it is flagged for review by a Measured Progress content specialist. The content specialist consults with the agency/department of education to determine whether to include the flagged item in a future operational test administration.

have in common. Generally, the scores on such subtests are highly correlated with each other; therefore, the primary dimension they share typically explains an overwhelming majority of variance in test scores. In fact, the presence of just such a dominant primary dimension is the psychometric assumption that provides the foundation for the unidimensional IRT models that were used for calibrating, linking, scaling, and equating the 2013–14 NECAP Science test forms.

The purpose of dimensionality analysis is to investigate whether violation of the assumption of test unidimensionality is statistically detectable and, if so, (a) the degree to which unidimensionality is violated and (b) the nature of the multidimensionality. Findings from dimensionality analyses performed on the 2013–14 NECAP Science common items for grades 4, 8, and 11 are reported below. (Note: only common items were analyzed since they are used for score reporting.)

Dimensionality analyses were conducted using the nonparametric IRT-based methods DIMTEST (Stout, 1987; Stout, Froelich, & Gao, 2001) and DETECT (Zhang & Stout, 1999). Nonparametric techniques were preferred for this analysis because such techniques avoid strong parametric modeling assumptions while still adhering to the fundamental principles of item response theory. Parametric techniques, such as nonlinear factor analysis, make strong assumptions that are often inappropriate for real data, such as assuming a normal distribution for ability and lower asymptotes of zero for the item characteristic curves.

Both DIMTEST and DETECT use as their basic statistical building block the estimated average conditional covariances for item pairs. A conditional covariance is the covariance between two items conditioned on expected total score for the rest of the test, and the average conditional covariance is obtained by averaging over all possible conditioning scores. When a test is strictly unidimensional, all conditional covariances are expected to take on values within random noise of zero, indicating statistically independent item responses for examinees with equal expected scores. Non-zero conditional covariances are essentially violations of the principle of local independence, and local *dependence* implies multidimensionality. Thus, non-random patterns of positive and negative conditional covariances are indicative of multidimensionality. In particular, when multiple dimensions are present, items measuring the same dimension will have positive conditional covariance with each other, whereas items measuring different dimensions will have negative conditional covariances with each other. For example, if multiple-choice (MC) items measure a different dimension from constructed-response (CR) items, we would expect MC items to have positive conditional covariances with each other, CR items to have positive conditional covariances with each other, and MC items to have negative conditional covariances with CR items.

DIMTEST is a hypothesis-testing procedure for detecting violations of local independence. For the exploratory analyses conducted for the NECAP Science tests, the data were first divided into a training sample and a cross-validation sample. Then an analysis of the conditional covariances was conducted on the training sample data to find the cluster of items that displays the greatest evidence of local dependence. The cross-validation sample was then used to test whether the conditional covariances of the selected cluster of

items displays local dependence, conditioning on total score on the non-clustered items. The DIMTEST statistic follows a standard normal distribution under the null hypothesis of unidimensionality.

DETECT is an effect-size measure of multidimensionality. For the exploratory analyses conducted for the NECAP Science tests, as with DIMTEST, the data were first randomly divided into a training sample and a cross-validation sample. (Note: The training and cross-validation samples used for the DETECT analyses were randomly drawn independently of the samples used for the DIMTEST analyses.) The training sample was then used to find a set of mutually exclusive and collectively exhaustive clusters of items that best fit a systematic pattern of positive conditional covariances for pairs of items from the same cluster and negative conditional covariances from different clusters. Next, the clusters from the training sample were used with the cross-validation sample data to average the conditional covariances: within-cluster conditional covariances were summed, from this sum the between-cluster conditional covariances were subtracted, this difference was divided by the total number of item pairs, and this average was multiplied by 100 to yield an index of the average violation of local independence for an item pair. DETECT values less than 0.2 indicate very weak multidimensionality (or near unidimensionality), values of 0.2 to 0.4 weak to moderate multidimensionality, values of 0.4 to 1.0 moderate to strong multidimensionality, and values greater than 1.0 very strong multidimensionality.

DIMTEST and DETECT were applied to the spring 2013–14 NECAP Science assessments for grades 4, 8, and 11. The data for each grade were split into a training sample and a cross-validation sample. Each grade had at least 29,500 student examinees. Because DIMTEST was limited to using 24,000 students, the training and cross-validation samples for the DIMTEST analyses used 12,000 each, randomly sampled from the total sample. DETECT, on the other hand, had an upper limit of 50,000 students, so every training sample and cross-validation sample used with DETECT had at least 14,700 students. DIMTEST was then applied to each of the science grades. DETECT was applied to each dataset for which the DIMTEST null hypothesis was rejected in order to estimate the effect size of the multidimensionality.

Because of the large sample sizes of the NECAP Science tests, DIMTEST would be sensitive to even quite small violations of unidimensionality, and the null hypothesis was strongly rejected for every dataset ($p < 0.00005$ for all three grades). These results were not surprising because strict unidimensionality is an idealization that almost never holds exactly for a given dataset. Thus, it was important to use DETECT to estimate the effect size of the violations of local independence found by DIMTEST. Table 6-2 below displays the multidimensional effect size estimates from DETECT.

As shown in Table 6-2, all of the DETECT values indicate multidimensionality is either very weak (DETECT values less than 0.20) or weak (DETECT values on the weak side of the 0.20 to 0.40 weak-to-moderate range) for every test analyzed. Also shown in Table 6-2 are the DETECT values from last year's dimensionality analysis. This year's results are seen to be similar to last year's in that both sets of results indicated weak or very weak multidimensionality for every test.

Table 6-2. 2013–14 NECAP Science: Multidimensionality Effect Sizes

Grade	Multidimensionality Effect Size	
	2013–14	2012–13
4	0.19	0.17
8	0.19	0.21
11	0.22	0.17

We also investigated how DETECT divided the tests into clusters to see if there were any discernable patterns with respect to the item types (MC as compared to CR). In all grades, the MC items and CR items tended to cluster separately from each other. This separation has occurred in all six years of NECAP Science dimensionality analyses. In 18 analyses (6 years \times 3 grade levels), only once has more than 25 percent of the CR points occurred in an MC-dominated cluster; and never has more than 10 percent of the MC points occurred in a CR-dominated cluster. Despite this evidence of multidimensionality between the MC and CR items, the multidimensional effect sizes are weak, so no changes in test design, scoring or administration are warranted.

CHAPTER 7 SCALING AND EQUATING

This chapter describes the procedures used to calibrate, equate, and scale the NECAP Science assessment. During the course of these psychometric analyses, a number of quality control procedures and checks on the processes were implemented. These procedures included evaluations of the calibration processes (e.g., checking the number of Newton cycles required for convergence for reasonableness, checking item parameters and their standard errors for reasonableness, examination of test characteristic curves [TCCs] and test information functions [TIFs] for reasonableness); evaluation of model fit; evaluation of equating items (e.g., delta analyses, rescore analyses, examination of α -plots and b -plots for reasonableness); and evaluation of the scaling results (e.g., parallel processing by the Measured Progress Psychometrics and Research and Data and Reporting Services departments, comparing lookup tables to the previous year's). An equating report, which provided complete documentation of the quality control procedures and results, was submitted to the agency/departments of education for their approval prior to production of student reports.

Table 7-1 lists items that required intervention either during item calibration or as a result of the evaluations of the equating items. For each flagged item, the table shows the reason it was flagged and what action was taken. The number of items identified for evaluation was very typical across the grades. Descriptions of the evaluations and results are included in the Item Response Theory (IRT) Results and Equating Results sections below.

Table 7-1. 2013–14 NECAP Science: Items That Required Intervention During IRT Calibration and Equating

<i>Grade</i>	<i>Item Number</i>	<i>Reasons</i>	<i>Action</i>
04	142265	c-parameter	set c = 0
	174255	c-parameter	set c = 0
	242420	c-parameter	set c = 0
	242739	c-parameter	set c = 0
	258704	a-parameter	a set to initial
	46416	delta analysis	removed from equating
	46426	c-parameter	set c = 0
	46525	c-parameter	set c = 0
	47448	b/b analysis	removed from equating
	59906	c-parameter	set c = 0
08	144589	delta analysis	removed from equating
	144589	c-parameter	set c = 0
	174976	c-parameter	set c = 0
	220376	c-parameter	set c = 0
	242842	c-parameter	set c = 0
	258720	a-parameter	a set to initial
	258721	a-parameter	a set to initial
	58352	c-parameter	set c = 0

continued

Grade	Item Number	Reasons	Action
	135344	c-parameter	set c = 0
	46140	c-parameter	set c = 0
	47917	c-parameter	set c = 0
	48115	c-parameter	set c = 0
	48214	c-parameter	set c = 0
11	49913	c-parameter	set c = 0
	62084	c-parameter	set c = 0
	89632	b/b analysis	removed from equating
	90282	c-parameter	set c = 0

7.1 ITEM RESPONSE THEORY

All NECAP Science items were calibrated using IRT. The IRT uses mathematical models to define a relationship between an unobserved measure of student performance, usually referred to as theta (θ), and the probability (p) of getting a dichotomous item correct or of getting a particular score on a polytomous item. In IRT, it is assumed that all items are independent measures of the same construct (i.e., of the same θ). Another way to think of θ is as a mathematical representation of the latent trait of interest. Several common IRT models are used to specify the relationship between θ and p (Hambleton & Swaminathan, 1985; Hambleton & van der Linden, 1997). The process of determining the specific mathematical relationship between θ and p is called item calibration. After items are calibrated, they are defined by a set of parameters that specify a nonlinear, monotonically increasing relationship between θ and p . Once the item parameters are known, an estimate of θ for each student can be calculated. This estimate, $\hat{\theta}$, is considered to be an estimate of the student's true score or a general representation of student performance. It has characteristics that may be preferable to those of raw scores for equating purposes.

For the 2013–14 NECAP Science assessment, the three-parameter logistic (3PL) model was used for dichotomous (multiple-choice) items and the graded-response model (GRM) was used for polytomous (open-response) items. The 3PL model for dichotomous items can be defined as

$$P_i(1|\theta_j, \xi_i) = c_i + (1 - c_i) \frac{\exp[Da_i(\theta_j - b_i)]}{1 + \exp[Da_i(\theta_j - b_i)]}$$

where

i indexes the items,

j indexes students,

a represents item discrimination,

b represents item difficulty,

c is the pseudo guessing parameter,

ξ_i represents the set of item parameters (a , b , and c) for item i , and

D is a normalizing constant equal to 1.701.

In the GRM for polytomous items, an item is scored in $k + 1$ graded categories, denoted as m categories, that can be viewed as a set of k dichotomies. At each point of dichotomization (i.e., at each threshold), a two-parameter model can be used. This implies that a polytomous item with m categories can be characterized by k item category threshold curves (ICTCs) of the two-parameter logistic form:

$$P_{ik}^*(1|\theta_j, a_i, b_i, d_{ik}) = \frac{\exp\left[D a_i (\theta_j - b_i + d_{ik})\right]}{1 + \exp\left[D a_i (\theta_j - b_i + d_{ik})\right]}$$

where
 i indexes the items,
 j indexes students,
 k indexes threshold,
 a represents item discrimination,
 b represents item difficulty,
 d represents threshold, and
 D is a normalizing constant equal to 1.701.

After computing k ICTCs in the GRM, $k + 1$ item category characteristic curves (ICCCs) are derived by subtracting adjacent ICTCs:

$$P_{ik}(1|\theta_j) = P_{i(k-1)}^*(1|\theta_j) - P_{ik}^*(1|\theta_j)$$

where
 P_{ik} represents the probability that the score on item i falls in category k , and
 P_{ik}^* represents the probability that the score on item i falls above the threshold k
($P_{i0}^* = 1$ and $P_{i(m+1)}^* = 0$).

The GRM is also commonly expressed as

$$P_{ik}(k|\theta_j, \xi_i) = \frac{\exp\left[D a_i (\theta_j - b_i + d_k)\right]}{1 + \exp\left[D a_i (\theta_j - b_i + d_k)\right]} - \frac{\exp\left[D a_i (\theta_j - b_i + d_{k+1})\right]}{1 + \exp\left[D a_i (\theta_j - b_i + d_{k+1})\right]}$$

where
 ξ_i represents the set of item parameters for item i .

Finally, the item characteristic curve (ICC) for polytomous items is computed as a weighted sum of ICCCs, where each ICCC is weighted by a score assigned to a corresponding category:

$$P_i(1|\theta_j) = \sum_k^{m+1} w_{ik} P_{ik}(1|\theta_j)$$

For more information about item calibration and determination, refer to Lord and Novick (1968), Hambleton and Swaminathan (1985), or Baker and Kim (2004).

7.2 ITEM RESPONSE THEORY RESULTS

The tables in Appendix I give the IRT item parameters of all common items on the 2013–14 NECAP Science tests by grade. In addition, Appendix J shows graphs of the TCCs and TIFs, which are defined below.

The TCCs display the expected (average) raw score associated with each θ_j value between -4.0 and 4.0. Mathematically, the TCC is computed by summing the ICCs of all items that contribute to the raw score. Using the notation introduced in Section 7.1, the expected raw score at a given value of θ_j is

$$E(X|\theta_j) = \sum_{i=1}^n P_i(1|\theta_j)$$

where

i indexes the items (and n is the number of items contributing to the raw score),

j indexes students (here, θ_j runs from -4 to 4), and

$E(X|\theta_j)$ is the expected raw score for a student of ability θ_j .

The expected raw score monotonically increases with θ_j , consistent with the notion that students of high ability tend to earn higher raw scores than do students of low ability. Most TCCs are S-shaped: flatter at the ends of the distribution and steeper in the middle.

The TIF displays the amount of statistical information that the test provides at each value of θ_j . Information functions depict test precision across the entire latent trait continuum. There is an inverse relationship between the information of a test and its standard error of measurement (SEM). For long tests, the SEM at a given θ_j is approximately equal to the inverse of the square root of the statistical information at θ_j (Hambleton, Swaminathan, & Rogers, 1991), as follows:

$$SEM(\theta_j) = \frac{1}{\sqrt{I(\theta_j)}}$$

Compared to the tails, TIFs are often higher near the middle of the θ_j distribution where most students are located and where most items are sensitive by design.

Table 7-1 above lists items that were flagged based on the quality control checks implemented during the calibration process. (Note that some items were flagged as a result of the evaluations of the equating items; those results are described below.) In all cases, items flagged during this step were identified because the pseudo-guessing parameter (c -parameter) was poorly estimated. Difficulty in estimating the c -parameter is not at all unusual and is well documented in the psychometric literature (see, e.g., Nering & Ostini, 2010), especially when the item's discrimination is below 0.50. In all cases, fixing the c -parameter resulted in reasonable and stable item parameter estimates and improved model fit.

The number of Newton cycles required for convergence for each grade during the IRT analysis can be found in Table 7-2. The number of cycles required fell within acceptable ranges.

Table 7-2. 2013–14 NECAP Science: Number of Newton Cycles Required for Convergence

<i>Subject</i>	<i>Grade</i>	<i>Cycles</i>
Science	4	35
	8	39
	11	44

7.3 EQUATING

The purpose of equating is to ensure that scores obtained from different forms of a test are equivalent to each other. Equating may be used if multiple test forms are administered in the same year, as well as to equate one year's forms to those given in the previous year. Equating ensures that students are not given an unfair advantage or disadvantage because the test form they took is easier or harder than those taken by other students.

The 2013–14 administration of NECAP Science used a raw-score-to-theta equating procedure in which test forms were equated to the theta scale established on the reference form (i.e., the form used in the most recent standard setting). This is accomplished through the chained linking design, in which every new form is equated back to the theta scale of the previous year's test form. It can therefore be assumed that the theta scale of every new test form is the same as the theta scale of the reference form, since this is where the chain originated.

The groups of students who took the equating items on the 2013–14 NECAP Science tests are not equivalent to the groups who took them in the reference years. IRT is particularly useful for equating scenarios that involve nonequivalent groups (Allen & Yen, 1979). Equating for NECAP Science uses the anchor-test-nonequivalent-groups design described by Petersen, Kolen, and Hoover (1989). In this equating design, no assumption is made about the equivalence of the examinee groups taking different test forms (i.e., naturally occurring groups are assumed). Comparability is instead evaluated by utilizing a set of anchor items (also called equating items). However, the equating items are designed to mirror the common test in terms of item types and distribution of emphasis. Subsets of the equating items are distributed across forms.

Item parameter estimates for 2013–14 were placed on the 2012–13 scale by using the method of Stocking and Lord (1983), which is based on the IRT principle of item parameter invariance. According to this principle, the equating items for both the 2012–13 and 2013–14 NECAP Science tests should have the same item parameters. After the item parameters for each 2013–14 test were estimated using PARSCALE (Muraki & Bock, 2003), the Stocking and Lord method was employed to find the linear transformation (slope and intercept) that adjusted the equating items' parameter estimates such that the 2013–14 TCC for the equating items was as close as possible to that of 2012–13.

7.4 EQUATING RESULTS

Prior to calculating the Stocking and Lord transformation constants, a variety of evaluations of the equating items were conducted. Items that were flagged as a result of these evaluations are listed in Table 7-1 at the beginning of this chapter. These items were scrutinized and a decision was made as to whether to include the item as an equating item or to discard it. The procedures used to evaluate the equating items are described below.

Appendix K presents the results from the delta analysis, applied to both the dichotomous and polytomous equating items. This procedure was one of several used to evaluate the adequacy of the equating items; the discard status presented in the appendix tables indicates whether or not the item was flagged as potentially inappropriate for use in equating based on the delta analysis alone. The discard status presented in the plots, however, indicates whether an item was flagged by any procedures used to evaluate the equating items. Also presented in Appendix K are the results from the rescore analysis. With this analysis, 200 random papers from the previous year were interspersed with this year's papers to evaluate scorer consistency from one year to the next. All effect sizes were well below the criterion value for excluding an item as an equating item, 0.80 in absolute value.

In addition to the delta and rescore analyses, evaluations based on analyzing \pm -plots and b -plots, which show IRT parameters for 2013–14 plotted against the values for 2012–13, were conducted and the results are presented in Appendix L. Any items that appeared as outliers in the plots were evaluated in terms of suitability for use as equating items.

Finally, \pm -plots and b -plots, which show IRT parameters for 2013–14 plotted against the values for 2012–13, are presented in Appendix L. Any items that appeared as outliers in the plots were evaluated in terms of suitability for use as equating items.

Once all flagged items had been evaluated and appropriate action taken, the Stocking and Lord method of equating was used to place the item parameters onto the previous year's scale, as described earlier. The Stocking and Lord transformation constants are presented in Table 7-3.

Table 7-3. 2013–14 NECAP Science: Stocking and Lord Transformation Constants

<i>Content Area</i>	<i>Grade</i>	<i>α-Slope</i>	<i>b-Intercept</i>
	4	0.985	-0.133
Science	8	0.990	0.077
	11	1.003	0.220

The next administration of NECAP Science (2014–15) will be scaled to the 2013–14 administration using the same equating method described above.

7.5 ACHIEVEMENT STANDARDS

Cutpoints for NECAP Science to establish the four achievement levels (Substantially Below Proficient, Partially Proficient, Proficient, and Proficient with Distinction) were set at a standard-setting meeting held in August 2008. Details of the standard-setting procedures can be found in the standard-setting report created at that time, as well as in the 2007–08 technical report. The cuts on the theta scale that were established via standard setting and used for reporting in 2013–14 are presented in Table 7-4. Also shown in the table are the cutpoints on the reporting score scale (described below). These cuts will remain fixed throughout the assessment program unless standards are reset for any reason.

Table 7-4. 2013–14 NECAP Science: Cut Scores on the Theta Metric and Reporting Scale by Grade

Grade	Theta			Scaled Score				
	Cut 1	Cut 2	Cut 3	Minimum	Cut 1	Cut 2	Cut 3	Maximum
4	-1.222	0.048	2.371	400	427	440	463	480
8	-0.612	0.751	2.578	800	829	840	855	880
11	-0.432	0.788	2.193	1,100	1,130	1,140	1,152	1,180

Table M-1 in Appendix M shows achievement level distributions by grade. Results are shown for each of the past three years.

7.6 REPORTED SCALED SCORES

Because the θ scale used in IRT calibrations is not readily understood by most stakeholders, reporting scales were developed for NECAP Science. The reporting scales are simple linear transformations of the underlying θ scale. The reporting scales are developed such that they range from $x00$ through $x80$ (where x is grade level). In other words, grade 4 scaled scores range from 400 through 480, grade 8 from 800 through 880, and grade 11 from 1100 through 1180. The lowest scaled score in the Proficient range is fixed at $x40$ for each grade level. For example, to be classified in the Proficient achievement level or above, a minimum scaled score of 440 was required at grade 4, 840 at grade 8, and 1140 at grade 11.

By providing information that is more specific about the position of a student’s results, scaled scores supplement achievement-level scores. School- and district-level scaled scores are calculated by computing the average of student-level scaled scores. Students’ raw scores (i.e., total number of points) on the 2013–14 NECAP Science tests were translated to scaled scores using a data analysis process called scaling. Scaling simply converts from one scale to another. In the same way that a given temperature can be expressed on either Fahrenheit or Celsius scales, or the same distance can be expressed in either miles or kilometers, student scores on the 2013–14 NECAP Science tests can be expressed in raw or scaled scores.

It is important to note that converting from raw scores to scaled scores does not change students’ achievement-level classifications. Given the relative simplicity of raw scores, it is fair to question why scaled

scores for NECAP Science are reported instead of raw scores. Scaled scores make consistent the reporting of results. To illustrate, standard setting typically results in different raw cut scores across grades. The raw cut score between Partially Proficient and Proficient could be, say, 38 in grade 4 and 40 in grade 8, yet both of these raw scores would be transformed to scaled scores of $\times 40$ (i.e., 440 and 840). It is this uniformity across scaled scores that facilitates the understanding of student performance. The psychometric advantage of scaled scores over raw scores comes from the fact that they are linear transformations of $\hat{\theta}$. Since the $\hat{\theta}$ scale is used for equating, scaled scores are comparable from one year to the next. Raw scores are not.

The scaled scores are obtained by a simple translation of ability estimates ($\hat{\theta}$) using the linear relationship between threshold values on the $\hat{\theta}$ metric and their equivalent values on the scaled score metric. Students' ability estimates are based on their raw scores and are found by mapping through the TCC. Scaled scores are calculated using the linear equation

$$SS = m\hat{\theta} + b$$

where
 m is the slope and
 b is the intercept.

A separate linear transformation is used for each grade level. For NECAP Science, the transformation function is determined by fixing the Partially Proficient/Proficient cut score and the bottom of the scale—that is, the $\times 40$ and the $\times 00$ values (e.g., 440 and 400 for grade 4). The $\times 00$ location on the $\hat{\theta}$ scale is beyond (i.e., below) the scaling of all items. To determine this location, a chance score (approximately equal to a student's expected performance by guessing) is mapped to a value of -4.0 on the $\hat{\theta}$ scale. A raw score of 0 is also assigned a scaled score of $\times 00$. The maximum possible raw score is assigned a scaled score of $\times 80$ (e.g., 480 in the case of grade 4). Because only two points within the $\hat{\theta}$ scaled score space are fixed, the scaled score cuts between Substantially Below Proficient and Partially Proficient and between Proficient and Proficient with Distinction are free to vary across grades.

Table 7-5 shows the slope and intercept terms used to calculate the scaled scores for each grade. Note that the values in Table 7-5 will not change unless the standards are reset.

Table 7-5. 2013–14 NECAP Science: Scaled Score Slope and Intercept by Subject and Grade

<i>Content Area</i>	<i>Grade</i>	<i>m-Slope</i>	<i>b-Intercept</i>
Science	4	9.881	439.524
	8	8.420	833.678
	11	8.354	1,133.414

Appendix N contains raw score to scaled score lookup tables for the 2013–14 NECAP Science tests. These are the actual tables used to determine student scaled scores, error bands, and achievement levels.

Appendix O contains scaled score distribution graphs for each grade. These distributions were calculated using the sparse data matrix files that were used in the IRT calibrations. As the curves move to the right, they represent an increase in performance.

CHAPTER 8 RELIABILITY

Although an individual item's performance is an important focus for evaluation, a complete evaluation of an assessment must also address the way items function together and complement one another. Tests that function well provide a dependable assessment of the student's level of ability. Unfortunately, no test can do this perfectly. A variety of factors can contribute to a given student's score being either higher or lower than his or her true ability. For example, a student may misread an item or mistakenly fill in the wrong bubble when he or she knew the answer. Collectively, extraneous factors that impact a student's score are referred to as *measurement error*. Any assessment includes some amount of measurement error; that is, no measurement is perfect. This is true of all academic assessments—some students will receive scores that underestimate their true ability and other students will receive scores that overestimate their true ability. When tests have a high amount of measurement error, student scores are very unstable. Students with high ability may get low scores, or vice versa. Consequently, one cannot reliably measure a student's true level of ability with such a test. Assessments that have less measurement error (i.e., errors made are small on average and student scores on such a test will consistently represent their ability) are described as reliable.

There are a number of ways to estimate an assessment's reliability. One possible approach is to give the same test to the same students at two different points in time. If students receive the same scores on each test, then the extraneous factors affecting performance are small and the test is reliable. (This is referred to as *test-retest reliability*.) A potential problem with this approach is that students may remember items from the first administration or may have gained (or lost) knowledge or skills in the interim between the two administrations. A solution to the "remembering items" problem is to give a different but parallel test at the second administration. If student scores on each test correlate highly, the test is considered reliable. (This is known as *alternate forms reliability*, because an alternate form of the test is used in each administration.) This approach, however, does not address the problem that students may have gained (or lost) knowledge or skills in the interim between the two administrations. In addition, the practical challenges of developing and administering parallel forms generally preclude the use of parallel forms reliability indices. One way to address the latter two problems is to split the test in half and then correlate students' scores on the two half-tests; this in effect treats each half-test as a complete test. By doing this, the problems associated with an intervening time interval and with creating and administering two parallel forms of the test are alleviated. This is known as a *split-half estimate of reliability*. If the two half-test scores correlate highly, items on the two half-tests must be measuring very similar knowledge or skills. This is evidence that the items complement one another and function well as a group. This also suggests that measurement error will be minimal.

The split-half method requires psychometricians to select items that contribute to each half-test score. This decision may have an impact on the resulting correlation, since each different possible split of the test into halves will result in a different correlation. Another problem with the split-half method of calculating

reliability is that it underestimates reliability, because test length is cut in half. All else being equal, a shorter test is less reliable than a longer test. Cronbach (1951) provided a statistic, α (alpha), which eliminates the problem of the split-half method by comparing individual item variances to total test variance. Cronbach's α was used to assess the reliability of the 2013–14 NECAP Science test:

$$\alpha \equiv \frac{n}{n-1} \left[1 - \frac{\sum_{i=1}^n \sigma_{(Y_i)}^2}{\sigma_x^2} \right]$$

where
i indexes the item,
n is the total number of items,
 $\sigma_{(Y_i)}^2$ represents individual item variance, and
 σ_x^2 represents the total test variance.

8.1 RELIABILITY AND STANDARD ERROR OF MEASUREMENT

Table 8-1 presents descriptive statistics, Cronbach's α coefficient, and raw score standard errors of measurement (SEM) for each grade. (Statistics are based on common items only.)

Table 8-1. 2013–14 NECAP Science: Raw Score Descriptive Statistics, Cronbach's Alpha, and Standard Errors of Measurement (SEM) by Grade

Grade	Number of Students	Raw Score			Alpha	SEM
		Maximum	Mean	Standard Deviation		
4	30,843	63	32.99	10.05	0.87	3.61
8	30,769	63	33.85	11.20	0.89	3.65
11	29,513	63	30.42	11.67	0.90	3.74

Because different grades have different test designs (e.g., the number of items varies by test), it is inappropriate to make inferences about the quality of one test by comparing its reliability to that of another test from a different grade.

8.2 2013–14 SUBGROUP RELIABILITY

The reliability coefficients discussed in the previous section were based on the overall population of students who took the 2013–14 NECAP Science test. Appendix P presents reliabilities for various subgroups of interest. Subgroup Cronbach's α 's were calculated using the formula defined earlier based only on the members of the subgroup in question in the computations; values are calculated only for subgroups with 10 or more students.

For several reasons, the results of this section should be interpreted with caution. First, inherent differences between grades preclude making valid inferences about the quality of a test based on statistical

comparisons with other tests. Second, reliabilities are dependent not only on the measurement properties of a test but also on the statistical distribution of the studied subgroup. For example, it can be readily seen in Appendix P that subgroup sample sizes may vary considerably, which results in natural variation in reliability coefficients. Or α , which is a type of correlation coefficient, may be artificially depressed for subgroups with little variability (Draper & Smith, 1998). Third, there is no industry standard to interpret the strength of a reliability coefficient, and this is particularly true when the population of interest is a single subgroup.

8.3 REPORTING SUBCATEGORY RELIABILITY

Of even more interest are reliabilities for the reporting subcategories within NECAP Science content areas, described in Chapter 3. Cronbach's α coefficients for subcategories were calculated via the same formula defined previously using just the items of a given subcategory in the computations. Results are presented in Appendix P. Once again as expected, because they are based on a subset of items rather than the full test, computed subcategory reliabilities were lower (sometimes substantially so) than were overall test reliabilities, and interpretations should take this into account. The subcategory reliabilities were lower than those based on the total test, approximately to the degree one would expect based on classical test theory. Qualitative differences between grades once again preclude valid inferences about the quality of the full test based on statistical comparisons among subtests.

8.4 INTERRATER CONSISTENCY

Chapter 5 of this report describes in detail the processes that were implemented to monitor the quality of the hand-scoring of student responses for constructed-response items. One of these processes was double-blind scoring: Approximately 2 percent of student responses were randomly selected and scored independently by two different scorers. Results of the double-blind scoring were used during the scoring process to identify scorers that required retraining or other intervention and are presented here as evidence of the reliability of the NECAP Science test. A summary of the interrater consistency results is presented in Table 8-2. Results in the table are collapsed across the hand-scored items by grade. The table shows number of score categories, number of included scores, percent exact agreement, percent adjacent agreement, correlation between the first two sets of scores, and percent of responses that required a third score. This same information is provided at the item level in Appendix Q.

Table 8-2. 2013–14 NECAP Science: Summary of Interrater Consistency Statistics Collapsed Across Items by Grade

Grade	Number of		Percent		Correlation	Percent of Third Scores
	Score Categories	Included Scores	Exact	Adjacent		
4	3	3,734	78.17	20.92	0.75	0.91
	4	1,250	78.16	20.08	0.85	1.76
	5	1,862	65.25	30.45	0.83	3.97
8	3	3,677	80.42	18.09	0.83	1.44
	4	1,234	78.93	19.94	0.91	1.13
	5	1,910	60.16	35.65	0.80	3.98
11	3	3,428	72.67	25.44	0.72	1.84
	4	1,148	64.72	33.01	0.76	2.26
	5	1,714	79.23	19.25	0.91	1.40

8.5 RELIABILITY OF ACHIEVEMENT-LEVEL CATEGORIZATION

While related to reliability, the accuracy and consistency of classifying students into achievement categories are even more important statistics in a standards-based reporting framework (Livingston & Lewis, 1995). After the achievement levels were specified and students were classified into those levels, empirical analyses were conducted to determine the statistical accuracy and consistency of the classifications. For NECAP Science, students are classified into one of four achievement levels: Substantially Below Proficient, Partially Proficient, Proficient, or Proficient with Distinction. This section of the report explains the methodologies used to assess the reliability of classification decisions.

Accuracy refers to the extent to which decisions based on test scores match decisions that would have been made if the scores did not contain any measurement error. Accuracy must be estimated, because errorless test scores do not exist. Consistency measures the extent to which classification decisions based on test scores match the decisions based on scores from a second, parallel form of the same test. Consistency can be evaluated directly from actual responses to test items if two complete and parallel forms of the test are given to the same group of students. In operational test programs, however, such a design is usually impractical. Instead, techniques have been developed to estimate both the accuracy and consistency of classification decisions based on a single administration of a test. The Livingston and Lewis (1995) technique was used for the 2013–14 NECAP Science assessment because it is easily adaptable to all types of testing formats, including mixed-format tests.

The accuracy and consistency estimates reported in Appendix R make use of “true scores” in the classical test theory sense. A true score is the score that would be obtained if a test had no measurement error. Of course, true scores cannot be observed and so must be estimated. In the Livingston and Lewis method, estimated true scores are used to categorize students into their “true” classifications.

For the 2013–14 NECAP Science test, after various technical adjustments (described in Livingston & Lewis, 1995), a four-by-four contingency table of accuracy was created for each grade, where cell $[i, j]$

represented the estimated proportion of students whose true score fell into classification i (where $i = 1$ to 4) and observed score into classification j (where $j = 1$ to 4). The sum of the diagonal entries (i.e., the proportion of students whose true and observed classifications matched) signified overall accuracy.

To calculate consistency, true scores were used to estimate the joint distribution of classifications on two independent, parallel test forms. Following statistical adjustments per Livingston and Lewis (1995), a new four-by-four contingency table was created for each grade and populated by the proportion of students who would be categorized into each combination of classifications according to the two (hypothetical) parallel test forms. Cell $[i, j]$ of this table represented the estimated proportion of students whose observed score on the first form would fall into classification i (where $i = 1$ to 4) and whose observed score on the second form would fall into classification j (where $j = 1$ to 4). The sum of the diagonal entries (i.e., the proportion of students categorized by the two forms into exactly the same classification) signified overall consistency.

Another way to measure consistency is to use Cohen's (1960) coefficient κ (kappa), which assesses the proportion of consistent classifications after removing the proportion of consistent classifications that would be expected by chance. It is calculated using the following formula:

$$\kappa = \frac{(\text{Observed agreement}) - (\text{Chance agreement})}{1 - (\text{Chance agreement})} = \frac{\sum_i C_{ii} - \sum_i C_{i.}C_{.i}}{1 - \sum_i C_{i.}C_{.i}}$$

where

$C_{i.}$ is the proportion of students whose observed achievement level would be Level i (where $i = 1-4$) on the first hypothetical parallel form of the test;

$C_{.i}$ is the proportion of students whose observed achievement level would be Level i (where $i = 1-4$) on the second hypothetical parallel form of the test;

C_{ii} is the proportion of students whose observed achievement level would be Level i (where $i = 1-4$) on both hypothetical parallel forms of the test.

Because κ is corrected for chance, its values are lower than are other consistency estimates.

8.5.1 Accuracy and Consistency Results

The accuracy and consistency analyses described above are provided in Table R-1 of Appendix R. The table includes overall accuracy and consistency indices, including kappa. Accuracy and consistency values conditional upon achievement level are also given. For these calculations, the denominator is the proportion of students associated with a given achievement level. For example, the conditional accuracy value is 0.85 for Substantially Below Proficient for grade 8. This figure indicates that among the students whose true scores placed them in this classification, 85 percent would be expected to be in this classification when categorized according to their observed scores. Similarly, a consistency value of 0.76 indicates that 76 percent of students with observed scores in the Substantially Below Proficient level would be expected to score in this classification again if a second, parallel test form were used.

For some testing situations, the greatest concern may be decisions around level thresholds. For example, in testing done for No Child Left Behind accountability purposes, the primary concern is distinguishing between students who are proficient and those who are not yet proficient. In this case, the accuracy of the Partially Proficient/Proficient threshold is of greatest interest. For 2013–14 NECAP Science, Table R-2 in Appendix R provides accuracy and consistency estimates at each cutpoint, as well as false positive and false negative decision rates. (A false positive is the proportion of students whose observed scores were above the cut and whose true scores were below the cut. A false negative is the proportion of students whose observed scores were below the cut and whose true scores were above the cut.)

The above indices are derived from Livingston and Lewis’s (1995) method of estimating the accuracy and consistency of classifications. It should be noted that Livingston and Lewis discuss two versions of the accuracy and consistency tables. A standard version performs calculations for forms parallel to the form taken. An “adjusted” version adjusts the results of one form to match the observed score distribution obtained in the data. The tables use the standard version for two reasons: (1) this “unadjusted” version can be considered a smoothing of the data, thereby decreasing the variability of the results; and (2) for results dealing with the consistency of two parallel forms, the unadjusted tables are symmetrical, indicating that the two parallel forms have the same statistical properties. This second reason is consistent with the notion of forms that are parallel; that is, it is more intuitive and interpretable for two parallel forms to have the same statistical distribution.

Note that, as with other methods of evaluating reliability, Decision Accuracy and Consistency (DAC) statistics calculated based on small groups can be expected to be lower than those calculated based on larger groups. For this reason, the values presented in Appendix R should be interpreted with caution. In addition, it is important to remember that it is inappropriate to compare DAC statistics between grades.

CHAPTER 9 SCORE REPORTING

The data used for the NECAP Science reports are the results of the spring 2014 administration of the NECAP Science test. NECAP Science tests are based on the NECAP Science assessment targets, which cover the grade spans K–4, 5–8, and 9–11. For example, the grade 8 NECAP Science test is based on the assessment targets of grades 5 through 8. Because the assessment targets cover grade spans, the state agency/departments of education determined that assessing science in the spring—as opposed to the fall, when mathematics, reading, and writing are assessed—would allow students and schools adequate time to cover all assessment targets through the curriculum and would also avoid a testing overload in the fall. All students who participated in NECAP Science were represented in testing year reports, because the students took the test in the school where they completed their learning of the assessment targets for their particular grade span.

9.1 PRIMARY REPORTS

Measured Progress created four primary reports for the 2013–14 NECAP Science test:

- Student Report
- Interactive Reporting
- School-, District-, and State Grade-Level Results Reports
- District and State Summary Report

With the exception of the Student Report, all reports were available for schools and districts to view or download on a password-secure Web site hosted by Measured Progress. Student-level data files were also available for districts to download. Each of these reports is described in the following subsections. Sample reports are provided in Appendix S.

9.2 STUDENT REPORT

The NECAP Science Student Report is a single-page, two-sided report printed on 8.5" by 11" paper. The front side of the report includes informational text about the design and uses of the assessment. It also describes the three corresponding sections of the reverse side of the report as well as the achievement levels. The reverse side provides a complete picture of an individual student's performance on the NECAP Science test, divided into three sections. The first section provides the student's overall performance for science. In addition to giving the student's achievement level, it presents the scaled score numerically and in a graphic

that places the score, including its standard error of measurement, within the full range of possible scaled scores demarcated into the four achievement levels.

The second section of the report displays the student's achievement level in science relative to the percentage of students at each achievement level across the school, district, and state.

The third section shows the student's performance compared to school, district, and statewide performances in each of the four tested science domains: Physical Science, Earth Space Science, Life Science, and Scientific Inquiry.

Student performance is reported in the context of possible points: average points earned for the school, district, and state; and average points earned by students who are minimally proficient on the test (scaled score of 440, 840, or 1140). The average points earned is reported as a range, because it is the average of all students who are minimally proficient, plus or minus one standard deviation.

To provide a more complete picture of the inquiry task portion of the science test (Session 3), each report includes a description of the inquiry task that was administered to all students at that grade. The grade 4 inquiry task always contains a hands-on experiment; the grade 8 inquiry task sometimes contains a hands-on experiment and sometimes contains a paper-and-pencil data analysis; and the grade 11 inquiry task always contains a paper-and-pencil data analysis.

The NECAP Student Report is confidential and should be kept secure within the school and district. The Family Educational Rights and Privacy Act (FERPA) requires that access to individual student results be restricted to the student, the student's parents/guardians, and authorized school personnel.

9.3 INTERACTIVE REPORTING

Four interactive reports were available: Item Analysis Roster, Achievement Level Summary, Released Items Summary Data, and Longitudinal Data. Each of these interactive reports is described in the following sections. Sample interactive reports are provided in Appendix T. To access these four interactive reports, the user clicked the interactive tab on the home page of the system and selected the report desired from the dropdown menu. Next, the user applied basic filtering options, such as the name of the district or school and the grade level/content area test, to open the specific report. At this point, the user had the option of printing the report for the entire grade level or applying advanced filtering options to select a subgroup of students to analyze. Advanced filtering options include gender, ethnicity, limited English proficient (LEP), individualized education program (IEP), and socioeconomic status (SES). All interactive reports, with the exception of the Longitudinal Data Report, allowed the user to provide a custom title for the report.

9.3.1 Item Analysis Roster Report

The NECAP Science Item Analysis Roster Report provides a roster of all students in a school and provides performance on the common items that are released to the public. The student names and

identification numbers are listed as row headers down the left side of the report. The items are listed as column headers in the same order they appeared in the released item document.

For each item, the following are shown:

- The Depth of Knowledge (DOK) code
- The item type
- The correct response key for multiple-choice items
- The possible points
- The content standard

For each student, multiple-choice items are marked either with a plus sign (+), indicating that the student chose the correct multiple-choice response, or a letter (from A to D), indicating the incorrect response chosen by the student. For short-answer and constructed-response items, the number of points earned is shown. All responses to released items are shown in the report, regardless of the student's participation status. The columns on the right side of the report show the Total Test results, broken into several categories. Subcategory Points Earned columns show points earned by the student in each content area subcategory relative to total possible points. A Total Points Earned column is a summary of all points earned and total possible points in the content area. The last two columns show the student's scaled score and achievement level. Students reported as Not Tested are given a code in the achievement level column to indicate the reason the student was not tested. It is important to note that not all items used to compute student scores are included in this report; only released items are included. At the bottom of the report, the average percentage correct for each multiple-choice item and average scores for the short-answer and constructed-response items are shown for the school, district, and state. When the user applies advanced filtering criteria, the School and District Percent Correct/Average Score rows at the bottom of the report are blanked out and only the Group row and the State row for the group selected will contain data. This report can be saved, printed, or exported as a PDF.

The Item Analysis Roster is confidential and should be kept secure within the school and district. FERPA requires that access to individual student results be restricted to the student, the student's parents/guardians, and authorized school personnel.

9.3.2 Achievement Level Summary

The Achievement Level Summary provides a visual display of the percentages of students in each achievement level for a selected grade. The four achievement levels are represented by various colors in a pie chart. A separate table is also included below the chart that shows the number and percentage of students in each achievement level. This report can be saved, printed, or exported as a PDF file.

9.3.3 Released Items Summary Data Report

The Released Items Summary Data Report is a school-level report that provides a summary of student responses to the released items for a selected grade/content area. The report is divided into three sections by item type (Multiple Choice, Constructed Response, and Inquiry Task). For multiple-choice items, the total number/percent of students who answered the item correctly and the number of students who chose each incorrect option or provided an invalid response are reported. An invalid response on a multiple-choice item is defined as “the item was left blank” or “the student selected more than one option for the item.” For constructed response and inquiry task items, point value and average score for the item are reported. Users are also able to view the actual released items within this report. If a user clicks on a particular magnifying glass icon next to a released item number, a pop-up box will open displaying the released item.

9.3.4 Longitudinal Data Report

The Longitudinal Data Report is a confidential student-level report that provides individual student performance data for multiple test administrations. The state-assigned student identification number is used to link students across test administrations. Student performance on future test administrations will be included on this report over time. This report can be saved, printed, or exported as a PDF file.

9.4 SCHOOL, DISTRICT, AND STATE GRADE-LEVEL RESULTS REPORTS

The NECAP School Results Report and the NECAP District Results Report consist of three parts: the grade-level summary report (page 2), the content area results (page 3), and the disaggregated content area results (page 4).

The grade-level summary report provides a summary of participation in the NECAP and a summary of NECAP results. The participation section on the top half of the page shows the number and percentage of students who were enrolled on or after May 5, 2014. The total number of students enrolled is defined as the number of students tested plus the number of students not tested.

Data are provided for the following groups of students who are considered tested in NECAP:

- **Students Tested:** This category provides the total number of students tested.
- **Students Tested with an Approved Accommodation:** Students in this category tested with an accommodation and did not have their test invalidated.
- **Current LEP Students:** Students in this category are currently receiving LEP services.
- **Current LEP Students Tested with an Approved Accommodation:** Students in this category are currently receiving LEP services, tested with an accommodation, and did not have their test invalidated.
- **IEP Students:** Students in this category have an IEP.

- **IEP Students Tested with an Approved Accommodation:** Students in this category have an IEP, tested with an accommodation, and did not have their test invalidated.

Because students who were not tested did not participate, average school scores were not affected by nontested students. These students were included in the calculation of the percentage of students participating, but not in the calculation of scores. For students who participated in some but not all sessions of the NECAP Science test, overall raw and scaled scores were reported. These reporting decisions were made to support the requirement that all students participate in the NECAP testing program.

Data are provided for the following groups of students, who may not have completed the entire NECAP Science test:

- **Alternate Assessment:** Students in this category completed an alternate assessment for the 2013–14 school year.
- **Withdrew After May 5:** Students withdrawing from a school after May 5, 2014, may have taken some sessions of the NECAP Science test prior to their withdrawal from the school.
- **Enrolled After May 5:** Students enrolling in a school after May 5, 2014, may not have had adequate time to participate fully in all sessions of the NECAP Science test.
- **Special Consideration:** Schools received state approval for special consideration for an exemption on all or part of the NECAP Science test for any student whose circumstances were not described by the previous categories but for whom the school determined that taking the NECAP Science test would not be possible.
- **Other:** Occasionally, students did not complete the NECAP Science test for reasons other than those listed. These “other” categories were considered not state approved.

The results section, on the bottom half of the page, shows the number and percentage of students performing at each achievement level in science across the school, district, and state. In addition, a mean scaled score is provided across school, district, and state levels. For the district version of this report, the school information is blank.

The content area results page provides information on performance in the four tested science domains (Physical Science, Earth Space Science, Life Science, and Scientific Inquiry). The purpose of this section is to help schools determine the extent to which their curricula are effective in helping students achieve the particular standards and benchmarks contained in the NECAP Science assessment targets. Information about the content area for school, district, and state includes

- the total number of students enrolled, not tested for a state-approved reason, not tested for another reason, and tested;
- the total number and percentage of students at each achievement level (based on the number in the Tested column); and
- the mean scaled score.

Information about each science domain includes the following:

- **The total possible points for that domain.** To provide as much information as possible for each domain, the total number of points includes both the common items used to calculate scores and additional items in each category.
- **A graphic display of the percentage of total possible points for the school, state, and district.** In this graphic display, symbols represent school, district, and state performance. In addition, a line symbolizes the standard error of measurement. This statistic indicates how much a student's score could vary if the student were examined repeatedly with the same test (assuming that no learning were to occur between test administrations).

The disaggregated content area results pages present the relationship between performance and student reporting variables in science across school, district, and state levels. The report shows the number of students categorized as enrolled, not tested for a state-approved reason, not tested for another reason, and tested. The report also provides the number and percentage of students within each of the four achievement levels and the mean scaled score by each reporting category.

The list of student reporting categories is as follows:

- All Students
- Gender
- Primary Race/Ethnicity
- LEP Status
- IEP
- SES
- Migrant
- Title I
- 504 Plan

The data for achievement levels and mean scaled score are based on the number shown in the Tested column. Reporting categories data were provided by records linked to the student labels. Because performance is being reported by categories that can contain relatively low numbers of students, school personnel are advised, under FERPA guidelines, to treat these pages confidentially.

It should be noted that for New Hampshire and Vermont, no data were reported for the 504 plan. In addition, for Vermont, no data were reported for Title I.

9.5 DISTRICT AND STATE SUMMARY REPORTS

The NECAP District Summary Report provides details on student performance for all grade levels of NECAP Science tested in the district. The purpose of the report is to help districts determine the extent to which their schools and students achieve the particular standards and benchmarks contained in the NECAP Science assessment targets. The NECAP District Summary Report contains no individual school data. The information provided includes

- the total number of students enrolled, not tested for a state-approved reason, not tested for another reason, and tested;
- the total number and percentage of students at each achievement level (based on the number in the Tested column); and
- the mean scaled score.

9.6 DECISION RULES

To ensure that the reported results for the 2013–14 NECAP Science test are accurate relative to collected data and other pertinent information, a document that delineates analysis and reporting rules was created. These decision rules were observed in the analyses of test data and in reporting the test results. Moreover, these rules served as the main reference for quality assurance checks.

The decision rules document used for reporting results of the May 2014 administration of the NECAP Science test is found in Appendix U.

The first set of rules pertains to general issues in reporting scores. Each issue is described, and pertinent variables are identified. The actual rules applied are described by the way they impact analyses and aggregations and their specific impact on each of the reports. The general rules are further grouped into issues pertaining to test items, school type, student exclusions, and number of students for aggregations.

The second set of rules pertains to reporting student participation. These rules describe which students were counted and reported for each subgroup in the student participation report.

9.7 QUALITY ASSURANCE

Quality assurance measures are embedded throughout the entire process of analysis and reporting. The data processor, data analyst, and psychometrician assigned to work on the NECAP Science assessment implement quality control checks of their respective computer programs and intermediate products. Moreover, when data are handed off to different functions within the Psychometrics and Research and Data and Reporting Services departments, the sending function verifies that the data are accurate before handoff. When a function receives a data set, the first step is to verify the data for accuracy.

Another type of quality assurance measure is parallel processing. Students' scaled scores for science are assigned by a psychometrician through a process of equating and scaling. The scaled scores are also computed by a data analyst to verify that scaled scores and corresponding achievement levels are assigned accurately. Respective scaled scores and achievement levels assigned are compared across all students for 100 percent agreement. Different exclusions assigned to students that determine whether each student receives scaled scores and/or is included in different levels of aggregation are also parallel processed. Using the decision rules document, two data analysts independently write a computer program that assigns students' exclusions. For each grade, the exclusions assigned by each data analyst are compared across all students. Only when 100 percent agreement is achieved can the rest of the data analysis be completed.

The third aspect of quality control involves the procedures implemented by the quality assurance group to check the veracity and accuracy of reported data. Using a sample of schools and districts, the quality assurance group verifies that reported information is correct. The step is conducted in two parts: (1) verify that the computed information was obtained correctly through appropriate application of different decision rules, and (2) verify that the correct data points populate each cell in the NECAP Science reports. The selection of sample schools and districts for this purpose is very specific and can affect the success of the quality control efforts. Two sets of samples are selected, though they may not be mutually exclusive.

The first set includes those that satisfy the following criteria:

- One-school district
- Two-school district
- Multischool district

The second set of samples includes districts or schools that have unique reporting situations, as indicated by decision rules. This set is necessary to check that each rule is applied correctly. The second set includes the following criteria:

- Private school
- Small school that receives no school report
- Small district that receives no district report
- District that receives a report but all schools are too small to receive a school report
- School with excluded (not tested) students
- School with homeschooled students

The quality assurance group uses a checklist to implement its procedures. After the checklist is completed, sample reports are circulated for psychometric checks and program management review. The appropriate sample reports are then presented to the client for review and signoff.

CHAPTER 10 VALIDITY

Because interpretations of test scores, and not a test itself, are evaluated for validity, the purpose of the 2013–14 NECAP Science Technical Report is to describe several technical aspects of the NECAP Science tests in support of score interpretations (AERA et al., 2014). Each chapter contributes an important component in the investigation of score validation: test development and design; test administration; scoring, scaling, and equating; item analyses; reliability; and score reporting.

Standards for Educational and Psychological Testing (AERA et al., 2014) provides a framework for describing sources of evidence that should be considered when constructing a validity argument. The evidence around test content, response processes, internal structure, relationship to other variables, and consequences of testing speaks to different aspects of validity but are not distinct types of validity. Instead, each contributes to a body of evidence about the comprehensive validity of score interpretations.

Evidence on test content validity is meant to determine how well the assessment tasks represent the curriculum and standards for each grade level. Content validation is informed by the item development process, including how the test blueprints and test items align to the curriculum and standards. Viewed through the lens provided by the standards, evidence based on test content was extensively described in Chapters 3 and 4. Item alignment with NECAP Science content standards; item bias, sensitivity, and content appropriateness review processes; adherence to the test blueprint; use of multiple item types; use of standardized administration procedures, with accommodated options for participation; and appropriate test administration training are all components of validity evidence based on test content. As discussed earlier, all NECAP Science questions are aligned by educators from the member states to specific NECAP Science content standards and undergo several rounds of review for content fidelity and appropriateness. Items are presented to students in multiple formats (constructed-response and multiple-choice). Finally, tests are administered according to state-mandated standardized procedures, with allowable accommodations, and all test coordinators and administrators are required to familiarize themselves with and adhere to all of the procedures outlined in the *NECAP Principal/Test Coordinator Manual* and *Test Administrator Manual*.

The scoring information in Chapter 5 describes the steps taken to train and monitor hand-scorers, as well as quality control procedures related to scanning and machine scoring. However, studies of student response processes would also be helpful and might include an investigation of students' cognitive methods using think-aloud protocols.

Evidence based on internal structure is presented in great detail in the discussions of item analyses, reliability, and scaling and equating in Chapters 6 through 8. Technical characteristics of the internal structure of the assessments are presented in terms of classical item statistics (item difficulty, item-test correlation), differential item functioning analyses, dimensionality analyses, reliability, standard errors of measurement, and item response theory (IRT) parameters and procedures. Each test is equated to the same grade test from

the prior year to preserve the meaning of scores over time. In general, item difficulty and discrimination indices were in acceptable and expected ranges. Very few items were answered correctly at near-chance or near-perfect rates. Similarly, the positive discrimination indices indicate that most items were assessing consistent constructs, and students who performed well on individual items tended to perform well overall.

Evidence based on the consequences of testing is addressed in the scaled score information in Chapter 7 and the reporting information in Chapter 9, as well as in the Guide to Using the 2014 NECAP Science Reports, which is a separate document. Each of these chapters speaks to the efforts undertaken to promote accurate and clear information to the public regarding test scores. Scaled scores offer the advantage of simplifying the reporting of results across grade levels and subsequent years. Achievement levels provide users with reference points for mastery at each grade level, which is another useful and simple way to interpret scores. Several different standard reports are provided to stakeholders. Additional evidence of the consequences of testing could be supplemented with broader investigation of the impact of testing on student learning.

To further support the validation of the assessment program, additional studies might be considered to provide evidence regarding the relationship of NECAP Science results to other variables, including the extent to which scores from NECAP Science converge with other measures of similar constructs and the extent to which they diverge from measures of different constructs. Relationships among measures of the same or similar constructs can sharpen the meaning of scores and appropriate interpretations by refining the definition of the construct.

10.1 QUESTIONNAIRE DATA

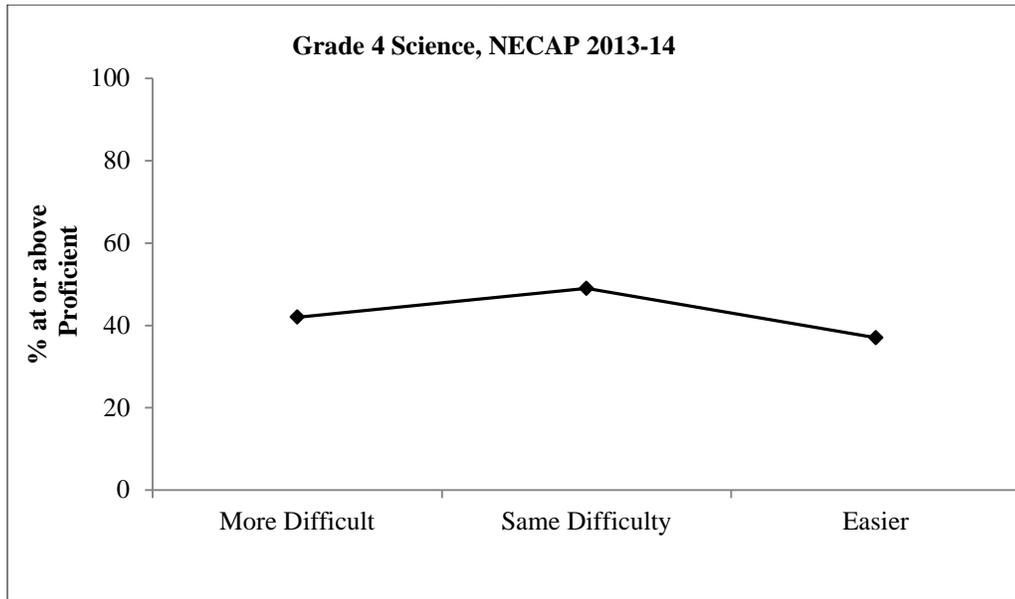
External validity of the NECAP Science assessment is conveyed by the relationship of test scores and situational variables such as time spent patterns, self-image, and attitude toward content matter. These situational variables were all based on student questionnaire data collected during the administration of the NECAP Science test. Note that no inferential statistics are included in the results presented below; however, because the numbers of students are quite large, differences in average scores may be statistically significant.

10.1.1 Difficulty of Assessment

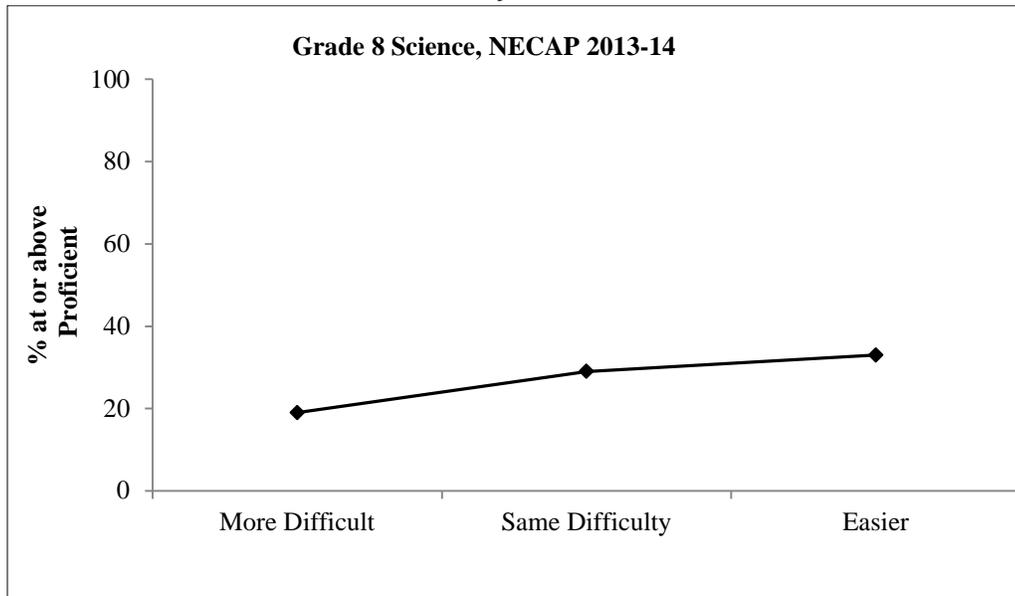
Examinees in all three grades were asked how difficult the science test was. Figures 10-1 through 10-3 show that students in grades 8 and 11 who thought the test was easier than their regular science schoolwork did better overall than those who thought it was more difficult; the trend was especially dramatic for grade 11. The pattern of responses for grade 4 was inconsistent.

Question: How difficult was this science test?

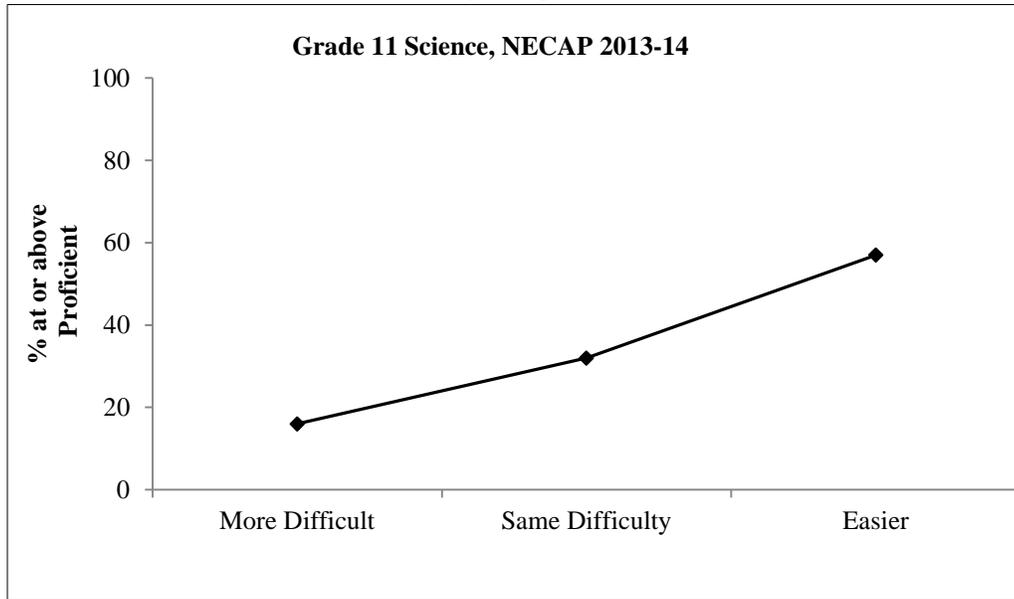
**Figure 10-1. 2013–14 NECAP Science: Questionnaire Responses
Difficulty—Grade 4**



**Figure 10-2. 2013–14 NECAP Science: Questionnaire Responses
Difficulty—Grade 8**



**Figure 10-3. 2013–14 NECAP Science: Questionnaire Responses
Difficulty—High School**

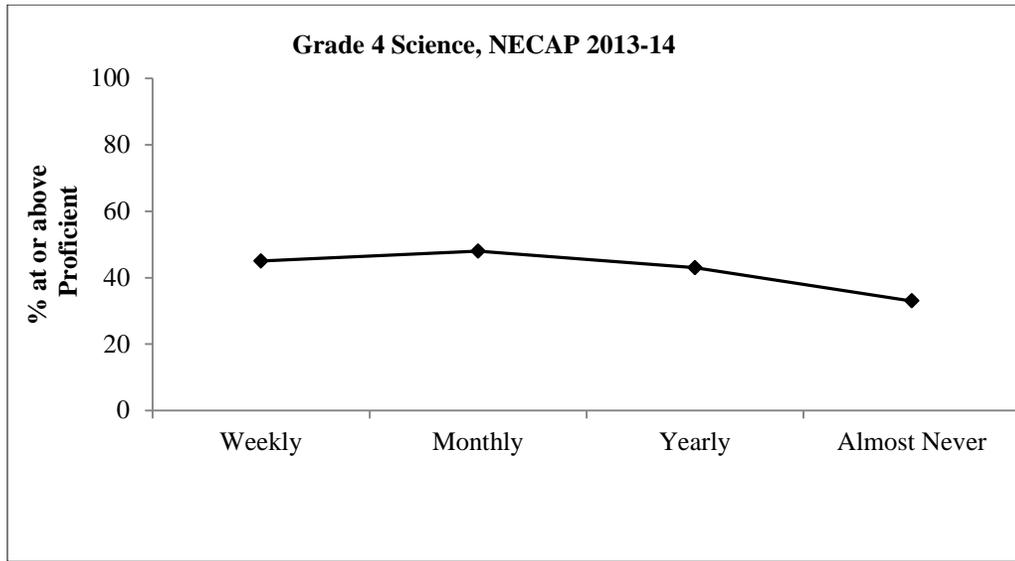


10.1.2 Content

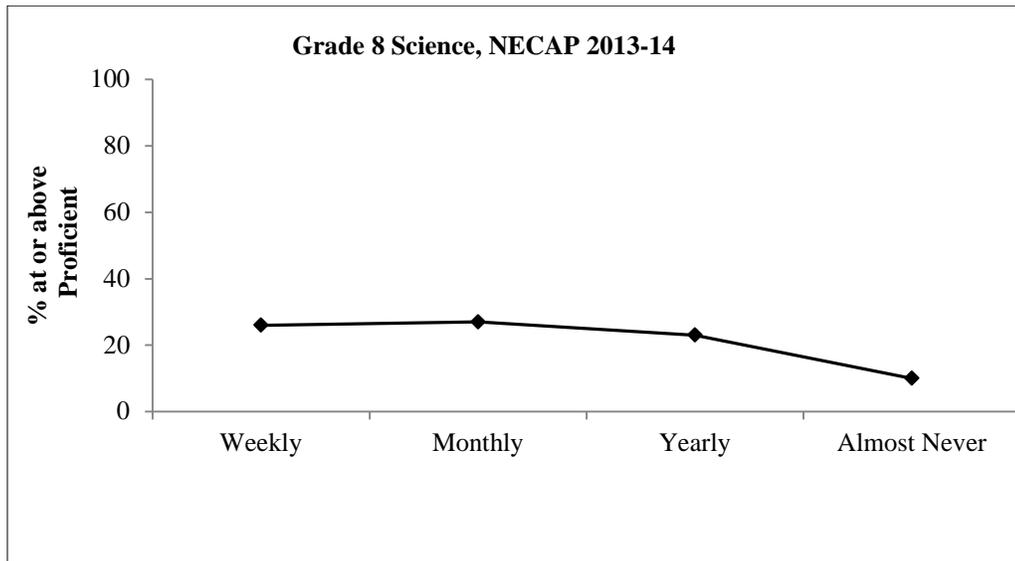
Examinees in all three grades were asked how often they do science experiments or inquiry tasks in their class. Figures 10-4 through 10-6 indicate a slight positive relationship between frequency of performing experiments or inquiry tasks and NECAP Science scores (i.e., higher scores are associated with greater frequency) for grade 8 and a relatively flat relationship for students who perform an inquiry task at least yearly in grades 4 and 11. In all cases, by excluding the “Weekly” response, the relationship is positive although the differences are slight, especially in grade 11.

Question: How often do you do science experiments or inquiry tasks in your class?

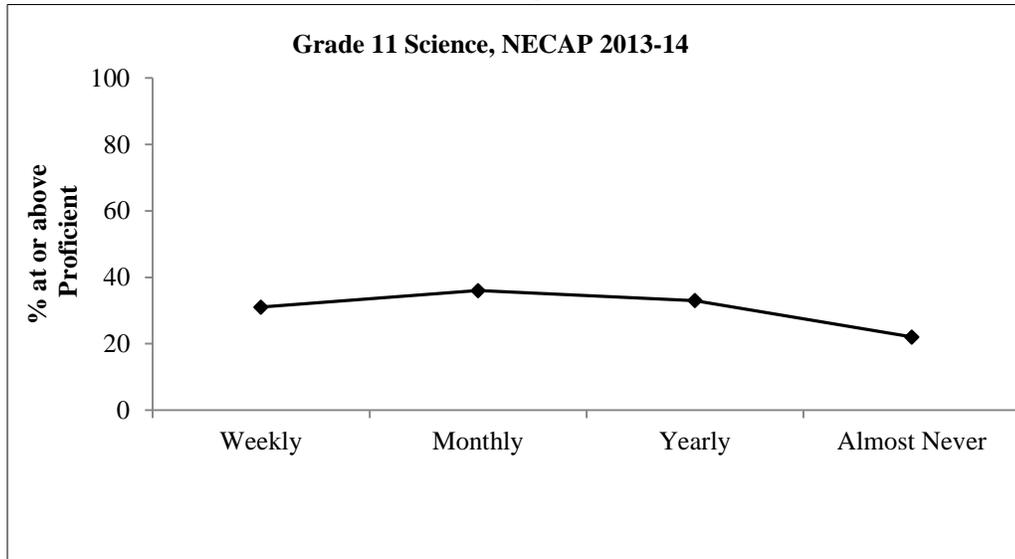
**Figure 10-4. 2013–14 NECAP Science: Questionnaire Responses
Content—Grade 4**



**Figure 10-5. 2013–14 NECAP Science: Questionnaire Responses
Content—Grade 8**



**Figure 10-6. 2013–14 NECAP Science: Questionnaire Responses
Content—High School**

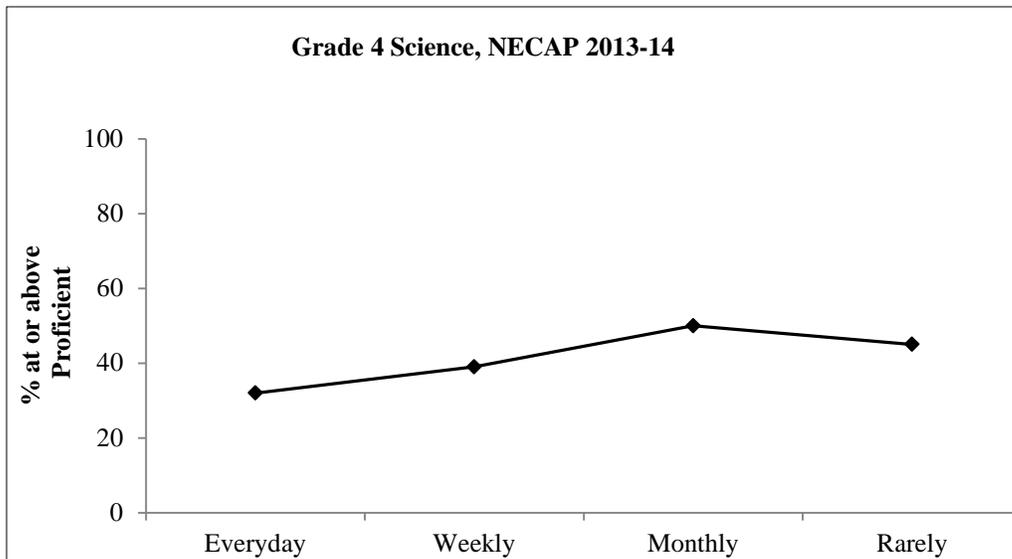


10.1.3 Homework

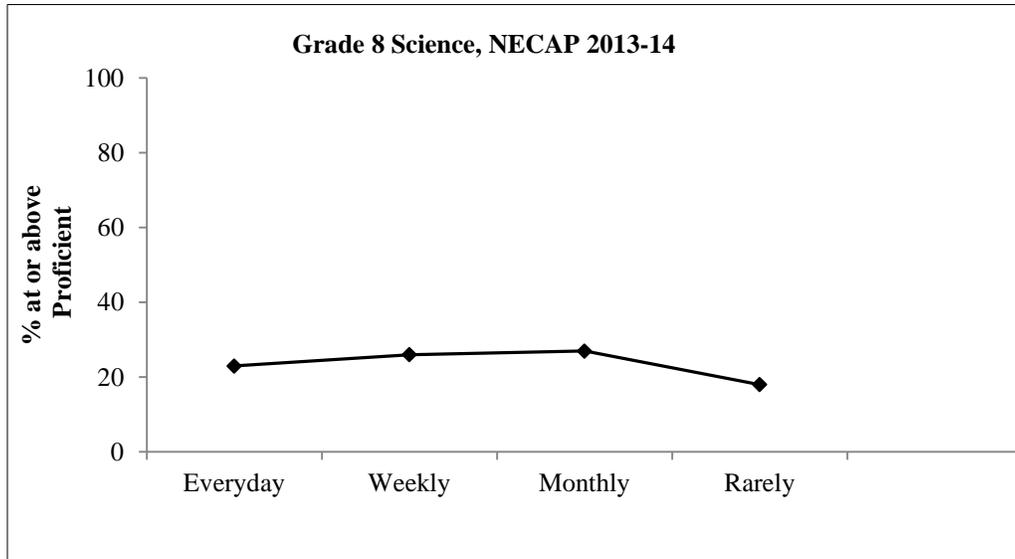
Examinees in all three grades were asked how often they have science homework. Figures 10-7 through 10-9 indicate a strong positive relationship between frequency of homework and NECAP scores for grades 8 and 11. The results for grade 4, on the other hand, indicate a negative relationship between frequency of homework and NECAP scores.

Question: How often do you have science homework?

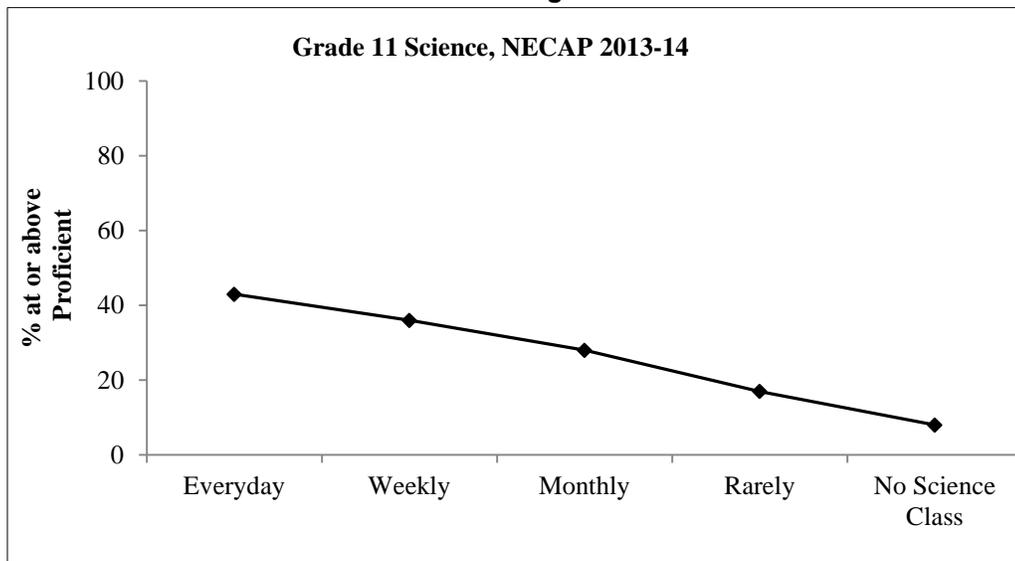
**Figure 10-7. 2013–14 NECAP Science: Questionnaire Responses
Homework—Grade 4**



**Figure 10-8. 2013–14 NECAP Science: Questionnaire Responses
Homework—Grade 8**



**Figure 10-9. 2013–14 NECAP Science: Questionnaire Responses
Homework—High School**

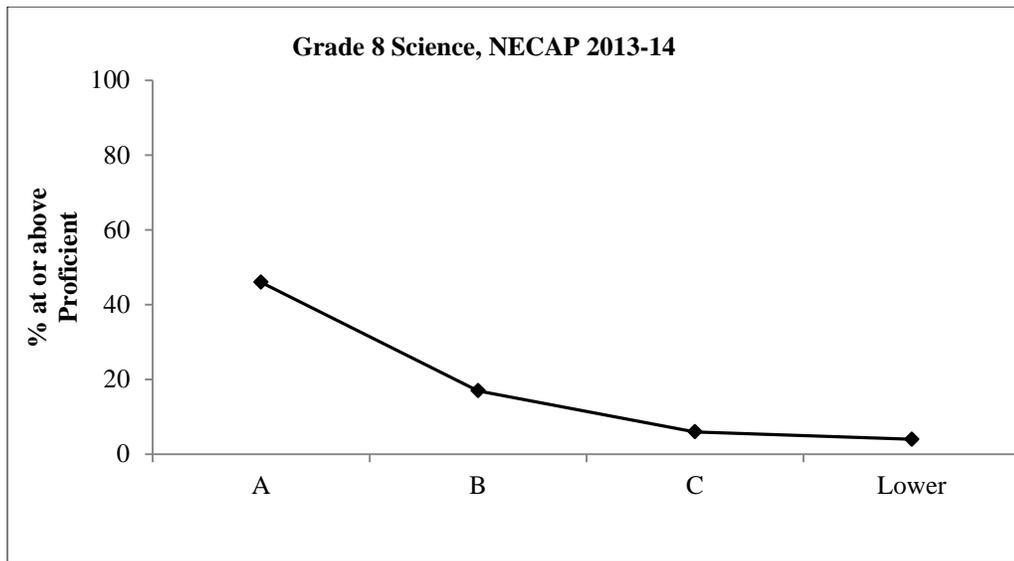


10.1.4 Performance in Science Class

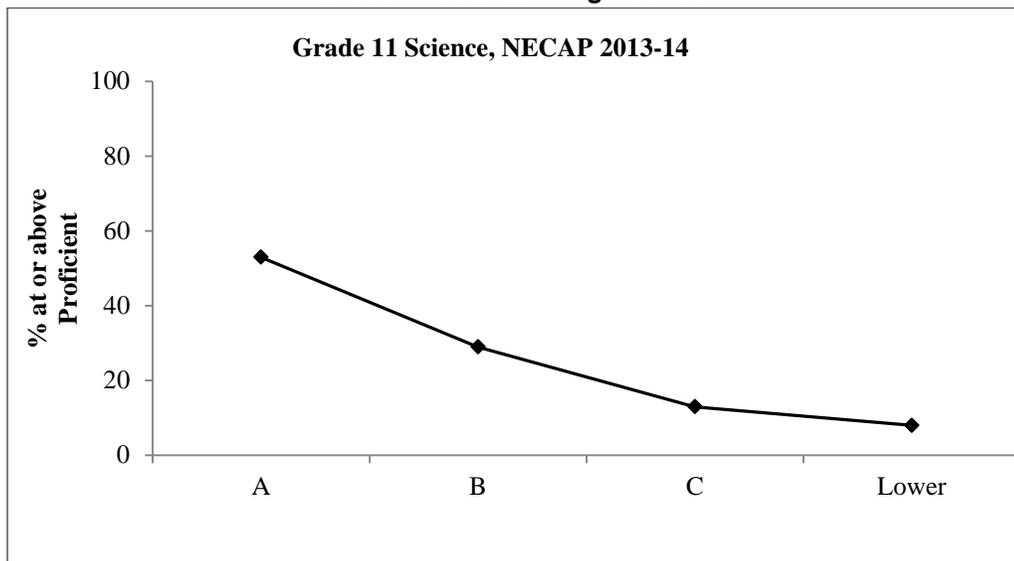
Students in grades 8 and 11 were asked what their most recent science grade was. Figures 10-10 and 10-11 indicate that, for grade 11, there was a strong positive relationship between the most recent science grade and NECAP scores; for grade 8, students with scores of C or higher performed better on NECAP than students with grades lower than C, but the relationship was fairly flat among students who received grades of A, B, or C.

Question: What was your science grade on your most recent report card?

**Figure 10-10. 2013–14 NECAP Science: Questionnaire Responses
Grade in Science—Grade 8**



**Figure 10-11. 2013–14 NECAP Science: Questionnaire Responses
Grade in Science—High School**



The evidence presented in this report supports inferences made about student achievement on the content represented in the NECAP Science assessment targets. As such, the evidence provided also supports the use of NECAP Science results for the purposes of program and instructional improvement and as a component of school accountability.

REFERENCES

- Allen, M. J., & Yen, W. M. (1979). *Introduction to Measurement Theory*. Belmont, CA: Wadsworth, Inc.
- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for educational and psychological testing*. Washington, DC: American Educational Research Association.
- Baker, F. B., & Kim, S-H. (2004). *Item response theory: Parameter estimation techniques* (2nd ed.). New York: Marcel Dekker, Inc.
- Brown, F. G. (1983). *Principles of educational and psychological testing* (3rd ed.). Fort Worth, TX: Holt, Rinehart and Winston.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement, 20*, 37–46.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika, 16*, 297–334.
- Dorans, N. J., & Holland, P. W. (1993). DIF detection and description. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 35–66). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Dorans, N. J., & Kulick, E. (1986). Demonstrating the utility of the standardization approach to assessing unexpected differential item performance on the Scholastic Aptitude Test. *Journal of Educational Measurement, 23*, 355–368.
- Draper, N. R., & Smith, H. (1998). *Applied regression analysis* (3rd ed.). New York: John Wiley & Sons.
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Boston, MA: Kluwer Academic Publishers.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Newbury Park, CA: Sage Publications.
- Hambleton, R. K., & van der Linden, W. J. (1997). *Handbook of modern item response theory*. New York: Springer-Verlag.
- Joint Committee on Testing Practices. (2004). *Code of fair testing practices in education*. Washington, DC: National Council on Measurement in Education.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement, 32*, 179–197.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theories of mental test scores*. Reading, MA: Addison-Wesley.
- Muraki, E., & Bock, R. D. (2003). *PARSCALE 4.1*. Lincolnwood, IL: Scientific Software International.
- Nering, M., & Ostini, R. (2010). *Handbook of polytomous item response theory models*. New York: Routledge.

- Petersen, N. S., Kolen, M. J., & Hoover, H. D. (1989). Scaling, norming, and equating. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 221–262). Washington, DC: American Council on Education/Macmillan.
- Stocking, M. L., & Lord, F. M. (1983). Developing a common metric in item response theory. *Applied Psychological Measurement*, 7, 201–210.
- Stout, W. F. (1987). A nonparametric approach for assessing latent trait dimensionality. *Psychometrika*, 52, 589–617.
- Stout, W. F., Froelich, A. G., & Gao, F. (2001). Using resampling methods to produce an improved DIMTEST procedure. In A. Boomsma, M. A. J. van Duign, & T. A. B. Snijders (Eds.), *Essays on item response theory* (pp. 357–375). New York: Springer-Verlag.
- Zhang, J., & Stout, W. F. (1999). The theoretical DETECT index of dimensionality and its application to approximate simple structure. *Psychometrika*, 64, 213–249.

APPENDICES

APPENDIX A—GUIDELINES FOR THE DEVELOPMENT OF SCIENCE INQUIRY TASKS

NECAP Science Assessment



Guidelines for the Development of Science Inquiry Tasks

Created in Partnership by the New Hampshire,
Rhode Island, and Vermont Departments of
Education



February 2008

Introduction: Inquiry in the NECAP Science Assessment

Defining the NECAP Science Assessment Inquiry Task

Focus – The Science Inquiry Task at every grade level should be rich and engaging. The task may be an experimental question or observational question – it is the quality of the task that is most important. Regardless of the type of task, all Four Broad Areas of Inquiry as defined in the *NECAP Schema for Assessing Scientific Inquiry*, (see column headings in the table on page 4), will be assessed. The task should flow from beginning to end in a purposeful way that allows students to make connections, express their ideas, and provide evidence of scientific thinking.

Design – Inquiry Tasks should be rooted in one or more NECAP Science Assessment Targets (one of which should have INQ code) and over time should address a variety of content domains. For every task at grades four and eight there must be scoreable components from each of the Four Broad Areas of Inquiry. At grade 11, while the focus of the task may be on constructs in the Area of Developing and Evaluating Explanations (column 4), scoreable items from each of the other three Broad Areas of Inquiry should also be included.

Task development will be guided by *Guidelines for the Development of Science Inquiry Tasks (GDIT)*. For each item within a Science Inquiry Task, the developer must identify the Depth of Knowledge (DOK), the Inquiry Construct number, score points, and key elements (scoring notes). Over time, all Inquiry Constructs should be addressed at each grade level. See the Appendix for additional information about the Inquiry Task development process.

Goal – Science Inquiry Tasks will engage students in a range of Depth of Knowledge experiences up to and including strategic thinking (DOK 3). Individual tasks may look different, but each should focus on providing insight into how students engage in scientific thinking. The goal is to encourage the meaningful inclusion of inquiry in classrooms at all levels.

Applying the Guidelines of the Science NECAP Assessment Task in the Classroom

Background – The first version of *Guidelines for Development of Science Inquiry Tasks* was originally created by the Science Specialists from the New Hampshire, Rhode Island, and Vermont Departments of Education to facilitate and refine the development of Inquiry Tasks for the NECAP Science Assessment. It became clear that such a tool would be useful to teachers and local science specialists to guide them in the development of similar tasks for classroom use at all levels. The State Science Specialists have collaborated on this version of *GDIT* to help educators understand and employ the constructs of the Four Broad Areas of Inquiry as they design or evaluate inquiry tasks for classroom instruction and assessment.

Focus - Classroom inquiry tasks should be relevant, engaging and meaningful learning experiences for students. The classroom inquiry tasks included on the state Department of Education website are examples of the kinds of tasks found in the NECAP Science Assessment. In the classroom any inquiry activity should provide regular opportunities for students to experience the science process as defined in the *NECAP Schema for Assessing Scientific Inquiry* (see page 4). Analysis of student performance on classroom inquiry tasks can inform instruction by providing data on student proficiencies within the constructs across the Four Broad Areas of Inquiry. Classroom inquiry tasks might be used as a component of local assessment or as a classroom summative assessment for a specific unit.

Design - While there are many ways to design inquiry experiences and an assessment for the classroom, *GDIT* provides a framework for the development of rich performance assessments that are aligned with this component of the NECAP Science Assessment. *GDIT* offers the necessary details for teachers to develop classroom inquiry tasks that are similar in structure to the NECAP Science Inquiry Tasks. Each classroom inquiry task will include elements from each of the Four Broad Areas of Inquiry, and address specific constructs within each Broad Area. Classroom inquiry tasks can span a class period, a few days or the length of a unit. Classroom inquiry tasks related to units of study provide opportunities for students to become familiar with the format of the NECAP Science Inquiry Tasks and will help to prepare them for the state assessment

Goals - The main goals of *Guidelines for Development of Science Inquiry Tasks* are to help educators:

- encourage the inclusion of engaging and relevant inquiry experiences in classrooms that contribute to increasing the science literacy of the citizens of New Hampshire, Rhode Island and Vermont;
- develop, evaluate and implement rich science tasks that allow students to gain skills across the Four Broad Areas of Inquiry;
- understand the process and parameters used in the development of Inquiry Tasks for the NECAP Science Assessment;
- provide opportunities for students to become familiar with the format and requirements of the NECAP Science Inquiry Tasks.

NECAP Science Inquiry Constructs for all Grade Levels

NECAP Science Schema for Assessing Scientific Inquiry (with DOK levels for constructs)				
Broad Areas of Inquiry to be Assessed	Formulating Questions & Hypothesizing	Planning and Critiquing of Investigations	Conducting Investigations	Developing and Evaluating Explanations
<p>Constructs for each Broad Area of Inquiry</p> <p>(including intended DOK Ceiling Levels, based on Webb Depth of Knowledge Levels for Science – see also Section II)</p> <p>Inquiry Constructs answer the question: What is it about the broad area of Inquiry that we want students to know and be able to do?</p>	<p>1. Analyze information from observations, research, or experimental data for the purpose of formulating a question, hypothesis, or prediction: (DOK 3)</p> <p>1a. Appropriate for answering with scientific investigation</p> <p>1b. For answering using scientific knowledge</p> <p>2. Construct coherent argument in support of a question, hypothesis, prediction (DOK 2 or 3 depending on complexity of argument)</p> <p>3. Make and describe observations in order to ask questions, hypothesize, make predictions related to topic (DOK 2)</p>	<p>4. Identify information/evidence that needs to be collected in order to answer the question, hypothesis, prediction (DOK 2 – routine; DOK 3 non-routine/ more than one dependant variable)</p> <p>5. Develop an organized and logical approach to investigating the question, including controlling variables (DOK 2 – routine; DOK 3 non-routine)</p> <p>6. Provide reasoning for appropriateness of materials, tools, procedures, and scale used in the investigation (DOK 2)</p>	<p>7. Follow procedures for collecting and recording qualitative or quantitative data, using equipment or measurement devices accurately (DOK 1 – use tools; routine procedure; DOK 2 – follow multi-step procedures; make observations)</p> <p>8. Use accepted methods for organizing, representing, and manipulating data (DOK 2 – compare data; display data)</p> <p>9. Collect sufficient data to study question, hypothesis, or relationships (DOK 2 – part of following procedures)</p> <p>10. Summarize results based on data (DOK 2)</p>	<p>11. Analyze data, including determining if data are relevant, artifact, irrelevant, or anomalous (DOK 2 – specify relationships between facts; ordering, classifying data)</p> <p>12. Use evidence to support and justify interpretations and conclusions or explain how the evidence refutes the hypothesis (DOK 3)</p> <p>13. Communicate how scientific knowledge applies to explain results, propose further investigations, or construct and analyze alternative explanations (DOK 3)</p>

NECAP Science Assessment Inquiry Task Flow

Administration of each Science Inquiry Performance Task (Grades 4 and 8) should follow the sequence below:

Prior to start of Session 3:

- Set up materials
- Group students

Standard Flow of NECAP Science Inquiry Performance Tasks: (Grades 4 and 8)

1. Directions read aloud by Test Administrator (basic info)
2. Scenario read aloud by Test Administrator (context)
3. Description of the materials and/or model explained by Test Administrator. Students make a prediction individually
4. Students conduct investigation with partner
5. Students clean up kits/experiment with partner
6. Students return to desks with their own Task Booklet to work individually
7. Test Administrator distributes Student Answer Booklets to students
8. Students copy data from Task Booklet to Student Answer Booklet (non-scored)
9. Students answer eight (8) scored questions in Student Answer Booklet
 - A. For analyzing the prediction, there will be Yes/No check boxes with space for the narrative below.
 - B. At grades 4 and 8, the question where students must graph data will have a hard-coded grid (1/2- inch squares) in the answer box with lines for x and y axis labels as well as a title. At grade 11, use 1/4- inch squares.

Standard Flow of NECAP Science Inquiry Data Analysis Tasks: (Grades 8 and 11)

1. Test Administrator distributes Student Answer Booklets to students
2. Directions read aloud by Test Administrator (basic info)
3. Scenario read aloud by Test Administrator (task context)
4. Students answer questions related to the scenario and complete data analysis in the Student Answer Booklet.
5. Items will require high school students to consider the Inquiry Constructs in relation to a selected data set.
6. Upon completion of the task students sit quietly and read until dismissal.

Broad Area I: Formulating Questions and Hypothesizing

Grade 4

Standard: Task must provide students a scenario that describes objects, organisms, or events within the environment. The scenario must include information relevant to grade 4 students and sufficient for them to construct questions and/or predictions based upon observations, past experiences, and scientific knowledge.

Note: bullets addressing constructs are not all inclusive.

Inquiry Construct:	Items addressing this construct require students to:
<p>1. Analyze information from observations, research, or experimental data for the purpose of formulating a question, hypothesis, or prediction:</p> <p>1a. Appropriate for answering with scientific investigation</p> <p>1b. For answering using scientific knowledge</p> <p>DOK 3</p>	<ul style="list-style-type: none"> • analyze scientific data and use that information to generate a testable question or a prediction that includes a cause and effect relationship; • generate a question or prediction which is reasonable in terms of available evidence; • support a question or prediction with an explanation. <p>Note: Addressing this construct may appear at the beginning of the task, the end, or both.</p>
<p>2. Construct coherent argument in support of a question, hypothesis, prediction</p> <p>DOK 2 or 3 depending on complexity of argument</p>	<ul style="list-style-type: none"> • identify evidence that supports or does not support a question or prediction.
<p>3. Make and describe observations in order to ask questions, hypothesize, make predictions related to topic</p> <p>DOK 2</p>	<ul style="list-style-type: none"> • connect observations to a question or prediction. <p>Note: Items may refer to an existing, new, or student-generated question or prediction.</p>

Broad Area 2: Planning and Critiquing of Investigations

Grade 4

Standard: Task requires students to plan or analyze a simple experiment based upon questions or predictions derived from the scenario. The experiment and related items should emphasize fairness in its design.

Note: The words "procedure" and "plan" are synonymous.

Inquiry Construct:	Items addressing this construct require students to:
<p>4. Identify information and/or evidence that needs to be collected in order to answer the question, hypothesis, prediction</p> <p>DOK 2 (routine) DOK 3 (non-routine or more than one dependant variable)</p>	<ul style="list-style-type: none"> • identify the types of evidence that should be gathered to answer the question; • design an appropriate format, such as data tables or charts, for recording data. <p>Note: These items could appear at the end of the task.</p>
<p>5. Develop an organized and logical approach to investigating the question, including controlling variables</p> <p>DOK 2 (routine) DOK 3 (non-routine)</p>	<ul style="list-style-type: none"> • develop a procedure to gather sufficient evidence (including multiple trials) to answer the question or test the prediction; • develop a procedure that lists steps logically and sequentially; • develop a procedure that changes one variable at a time. <p>Note: These items could appear at the end of the task. Use of the term "variable" should not appear in the item stem.</p>
<p>6. Provide reasoning for appropriateness of materials, tools, procedures, and scale used in the investigation</p> <p>DOK 2</p>	<ul style="list-style-type: none"> • explain why the materials, tools, or procedure for the task are or are not appropriate for the investigation.

Broad Area 3: Conducting Investigations

Grade 4

Standard: The procedure requires the student to demonstrate simple skills (observing, measuring, basic skills involving fine motor movement). The investigation requires the student to use simple scientific equipment (rulers, scales, thermometers) to extend their senses. The procedure provides the student with an opportunity to collect sufficient data to investigate the question, prediction, or relationships. Student is required to organize and represent qualitative or quantitative data using blank graph/chart templates. Student is required to summarize data.

Note: Metric measurements are used for Grade 4, except for those pertaining to weather.

Note: Multiple trials mean repeating the experiment to collect multiple sets of data.

Inquiry Construct	Items addressing this construct require students to:
<p>7. Follow procedures for collecting and recording qualitative or quantitative data, using equipment or measurement devices accurately</p> <p>DOK 1: use tools; routine procedure;</p> <p>DOK 2: follow multi-step procedures; make observations</p>	<ul style="list-style-type: none"> • record precise data and observations that are consistent with the procedure of the investigation; • include appropriate units of all measurements; • use appropriate measurement tools correctly to collect data; • record and label relevant details within a scientific drawing or diagram.
<p>8. Use accepted methods for organizing, representing, and manipulating data</p> <p>DOK 2: compare data; display data</p>	<ul style="list-style-type: none"> • represent data accurately in a graph/ table/ chart; • include titles , labels, keys or symbols as needed; • select a scale appropriate for the range of data to be plotted; • use common terminology to label representations; • identify relationships among variables based upon evidence.
<p>9. Collect sufficient data to study question, hypothesis, or relationships</p> <p>DOK 2 part of following procedures</p>	<ul style="list-style-type: none"> • show understanding of the value of multiple trials; • relate data to original question and prediction; • determine if the quantity of data is sufficient to answer the question or support or refute the prediction.
<p>10. Summarize results based on data</p> <p>DOK 2</p>	<ul style="list-style-type: none"> • consider all data when developing an explanation and/or conclusion; • identify patterns and trends in data.

Broad Area 4: Developing and Evaluating Explanations

Grade 4

Standard: Task must provide the opportunity for students to use data to construct an explanation based on their science knowledge and evidence from experimentation or investigation.

Inquiry Construct	Items addressing this construct require students to:
<p>11. Analyze data, including determining if data are relevant, artifact, irrelevant, or anomalous</p> <p>DOK 2 - specify relationships between facts; ordering, classifying data</p>	<ul style="list-style-type: none"> • identify data relevant to the task or question ; • identify factors that may affect experimental results (e.g. variables, experimental error, environmental conditions); • classify data into meaningful categories.
<p>12. Use evidence to support and justify interpretations and conclusions or explain how the evidence refutes the hypothesis</p> <p>DOK 3</p>	<ul style="list-style-type: none"> • identify data that seem inconsistent ; • use evidence to support or refute a prediction; • use evidence to justify an interpretation of data or trends; • identify and explain differences or similarities between prediction and experimental data; • provide a reasonable explanation that accurately reflects data; • use mathematical reasoning to determine or support conclusions.
<p>13. Communicate how scientific knowledge applies to explain results, propose further investigations, or construct and analyze alternative explanations</p> <p>DOK 3</p>	<ul style="list-style-type: none"> • explain how experimental results compare to accepted scientific understanding; • suggest ways to modify the procedure in order to collect sufficient data; • identify additional data that would strengthen an investigation; • connect the investigation or model to a real world example; • propose new questions, predictions, next steps or technology for further investigations; • design an investigation to further test a prediction.

Broad Area I: Formulating Questions and Hypothesizing

Grade 8

Standard: Task must provide students a scenario that describes objects, organisms, or events to which the student will respond. The task will provide the student with the opportunity to develop their own testable questions or predictions based upon their experimental data, observations, and scientific knowledge. The task could include opportunities for the student to refine and refocus questions or hypotheses related to the scenario using their scientific knowledge and information

Inquiry Construct	Items addressing this construct require students to:
<p>1. Analyze information from observations, research, or experimental data for the purpose of formulating a question, hypothesis, or prediction: (DOK 3)</p> <p>1a. Appropriate for answering with scientific investigation</p> <p>1b. For answering using scientific knowledge</p>	<ul style="list-style-type: none"> • analyze scientific data and use that information to generate a testable question or a prediction that includes a cause and effect relationship; • generate a question or a prediction which is reasonable in terms of available evidence; • support their question or prediction with a <u>scientific explanation</u>; • <u>refine or refocus a question or hypothesis using experimental data, research, or scientific knowledge.</u> <p>Note: Addressing this construct may appear at the beginning of the task, the end, or both.</p>
<p>2. Construct coherent argument in support of a question, hypothesis, prediction DOK 2 or 3 depending on complexity of argument</p>	<ul style="list-style-type: none"> • identify evidence that supports or does not support a question, <u>hypothesis</u> or prediction; • <u>explain the cause and effect relationship within the hypothesis or prediction</u>; • <u>use a logical argument to explain how the hypothesis or prediction is connected to a scientific concept, or observation.</u>
<p>3. Make and describe observations in order to ask questions, hypothesize, make predictions related to topic DOK 2</p>	<ul style="list-style-type: none"> • connect observations to a question or prediction. <p>Note: Items may refer to an existing, new, or student-generated question or prediction.</p>

Broad Area 2: Planning and Critiquing of Investigations

Grade 8

Standard: The task will require students to plan or analyze an experiment or investigation based upon questions, hypothesis, or predictions derived from the scenario. An experiment must provide students with the opportunity to identify and control variables. The task will provide opportunities for students to think critically about experiments and investigations and may ask students to propose alternatives.

Note: Scale refers to proportionality between the model and what it represents or the frequency with which data are collected.

Inquiry Construct	Items addressing this construct require students to:
<p>4. Identify information/evidence that needs to be collected in order to answer the question, hypothesis, prediction</p> <p>DOK 2: routine;</p> <p>DOK 3: non-routine/ more than one dependant variable</p>	<ul style="list-style-type: none"> • identify the types of evidence that should be gathered to answer the question, or <u>support or refute the prediction</u> ; • <u>identify the variables that may affect the outcome of the experiment or investigation;</u> • design an appropriate format for recording data; • <u>evaluate multiple data sets to determine which data are relevant to the question, hypothesis or prediction.</u> <p>Note: These items could appear at the end of the task</p>
<p>5. Develop an organized and logical approach to investigating the question, including controlling variables</p> <p>DOK 2: routine (replicates existing procedure);</p> <p>DOK 3: non-routine (extends, refines, or improves existing procedure)</p>	<ul style="list-style-type: none"> • develop a procedure to gather sufficient evidence (including multiple trials) to answer the question, or test <u>the hypothesis, or prediction;</u> • develop a procedure that lists steps sequentially and logically; • <u>explain which variable will be manipulated or changed (independent) and which variable will be affected by those changes (dependent);</u> • <u>identify variables that will be kept constant throughout the investigation;</u> • <u>use scientific terminology that supports the identified procedures;</u> • evaluate the organization and logical approach of a given procedure including <u>variables, controls, materials, and tools;</u> • <u>evaluate investigation design, including opportunities to collect appropriate and sufficient data.</u> <p>Note: These items could appear at the beginning or the end of the task.</p>
<p>6. Provide reasoning for appropriateness of materials, tools, procedures, and scale used in the investigation</p> <p>DOK 2</p>	<ul style="list-style-type: none"> • explain why the materials, tools, <u>procedure, or scale</u> for a task are appropriate or are inappropriate for the investigation. • <u>evaluate the investigation for the safe and ethical considerations of the materials, tools, and procedures.</u>

Broad Area 3: Conducting Investigations

Grade 8

Standard: The procedure requires the student to demonstrate skills (observing, measuring, basic skills involving fine motor movement) and mathematical understanding. The materials involved in the investigation are authentic to the task required. The procedure provides the student with an opportunity to collect sufficient data to investigate the question, prediction/hypothesis, or relationships. Student is required to organize and represent qualitative or quantitative data. Student is required to summarize data to form a logical argument.

Note: Metric units are used for all Grade 8 measurements.

Note: Multiple trials means repeating the experiment to collect multiple sets of data.

Inquiry Construct	Items addressing this construct require students to:
<p>7. Follow procedures for collecting and recording qualitative or quantitative data, using equipment or measurement devices accurately</p> <p>DOK 1: use tools; routine procedure;</p> <p>DOK 2: follow multi-step procedures; make observations</p>	<ul style="list-style-type: none"> • record precise data and observations that are consistent with the procedure of the investigation; • include appropriate units of all measurements; • use appropriate measurement tools correctly to collect data; • record and label relevant details within a scientific drawing.
<p>8. Use accepted methods for organizing, representing, and manipulating data</p> <p>DOK 2: compare data; display data</p>	<ul style="list-style-type: none"> • represent data accurately in an <u>appropriate</u> graph/table/chart; • include titles, labels, keys or symbols as needed; • select a scale appropriate for the range of data to be plotted; • use <u>scientific</u> terminology to label representations; • identify relationships among variables based upon evidence. <p>Note: The standard practice of graphing in science is to represent the independent on the x-axis and the dependent variable on the y-axis.</p>
<p>9. Collect sufficient data to study question, hypothesis, or relationships</p> <p>DOK 2: part of following procedures</p>	<ul style="list-style-type: none"> • show understanding of the value of multiple trials; • relate data to original question, <u>hypothesis</u> or prediction; • determine if the quantity of data is sufficient to answer the question or support or refute the <u>hypothesis</u> or prediction.
<p>10. Summarize results based on data</p> <p>DOK 2</p>	<ul style="list-style-type: none"> • consider all data when developing an explanation/conclusion; • identify patterns and trends in data.

Broad Area 4: Developing and Evaluating Explanations

Grade 8

Standard Task must provide the opportunity for students to use data to construct an explanation based on their science knowledge and evidence from experimentation or investigation. The task requires students to use qualitative and quantitative data to communicate conclusions and support/refute prediction/hypothesis.

Inquiry Construct	Items addressing this construct require students to:
<p>11. Analyze data, including determining if data are relevant, artifact, irrelevant, or anomalous</p> <p>DOK 2: specify relationships between facts; ordering, classifying data</p>	<ul style="list-style-type: none"> • identify data relevant to the task or question; • identify factors that may affect experimental results (e.g. variables, experimental error, environmental conditions); • classify data into meaningful categories; • <u>compare experimental data to accepted scientific data provided as part of the task;</u> • <u>use mathematical and statistical techniques to analyze data;</u> • <u>provide a reasonable explanation that accurately reflects data;</u> • <u>use content understanding to question data that might seem inaccurate;</u> • evaluate the significance of experimental data.
<p>12. Use evidence to support and justify interpretations and conclusions or explain how the evidence refutes the hypothesis</p> <p>DOK 3</p>	<ul style="list-style-type: none"> • <u>identify and explain data, interpretations or conclusions that seem inaccurate;</u> • use evidence to support or refute <u>question or hypothesis;</u> • use evidence to justify an interpretation of data or trends; • identify and explain differences or similarities between predictions and experimental data; • provide a reasonable explanation that accurately reflects data; • <u>use mathematical computations to determine or support conclusions.</u>
<p>13. Communicate how scientific knowledge applies to explain results, propose further investigations, or construct and analyze alternative explanations</p> <p>DOK 3</p>	<ul style="list-style-type: none"> • explain how experimental results compare to accepted scientific understanding; • <u>recommend changes to procedures to produce data that would provide sufficient data and more accurate analysis;</u> • <u>identify and justify</u> additional data that would strengthen an investigation; • connect the investigation or model to an <u>authentic situation;</u> • propose <u>and evaluate</u> new questions, predictions, next steps or technology for further investigations or <u>alternative explanations;</u> • <u>account for limitations and/or sources of error within the experimental design;</u> • <u>apply experimental results to a new problem or situation.</u>

Broad Area I: Formulating Questions and Hypothesizing

Grade 11

Standard: Task must provide students a scenario with information and detail sufficient for the student to create a testable prediction or hypothesis. Students will draw upon their science knowledge base to advance a prediction or hypothesis using appropriate procedures and controls; this may include an experimental design.

Inquiry Construct	Items addressing this construct require students to:
<p>1. Analyze information from observations, research, or experimental data for the purpose of formulating a question, hypothesis, or prediction.</p> <p>1a. Appropriate for answering with scientific investigation</p> <p>1b. For answering using scientific knowledge</p> <p>DOK 3</p>	<ul style="list-style-type: none"> • analyze scientific data and use that information to generate a testable question, <u>hypothesis</u>, or prediction that includes a cause and effect relationship; • generate a question, <u>hypothesis</u> or a prediction which is reasonable in terms of available evidence; • <u>show connections between hypothesis or prediction and scientific knowledge, observations, or research</u>; • support their question, <u>hypothesis</u>, or prediction with a scientific explanation; • refine or refocus a question or hypothesis using experimental data, research, or scientific knowledge. <p>Note: Addressing this construct may appear at the beginning of the task, the end, or both.</p>
<p>2. Construct coherent argument in support of a question, hypothesis, prediction.</p> <p>DOK 2 or 3: depends on complexity of argument</p>	<ul style="list-style-type: none"> • identify evidence that supports or does not support a question, hypothesis or prediction • explain the cause and effect relationship within the hypothesis or prediction; • use a logical argument to <u>support</u> the hypothesis or prediction using scientific concepts, principles, or observations.
<p>3. Make and describe observations in order to ask questions, hypothesize, make predictions related to topic.</p> <p>DOK 2</p>	<ul style="list-style-type: none"> • connect observations <u>and data</u> to a question, <u>hypothesis</u>, or prediction. <p>Note: Items may refer to an existing, new, or student-generated question, <u>hypothesis</u>, or prediction.</p>

Broad Area 2: Planning and Critiquing of Investigations

Grade 11

Standard: The task will require students to plan or analyze an experiment or investigation based upon questions, hypothesis, or predictions derived from the scenario. An experiment must provide students with the opportunity to identify and control variables. The task will provide opportunities for students to think critically and construct an argument about experiments and investigations and may ask students to propose alternatives. Task will require the student to identify and justify the appropriate use of tools, equipment, materials, and procedures involved in the experiment.

Note: Scale refers to proportionality between the model and what it represents or the frequency with which data are collected.

Inquiry Construct	Items addressing this construct require students to:
<p>4. Identify information/evidence that needs to be collected in order to answer the question, hypothesis, prediction</p> <p>DOK 2: routine;</p> <p>DOK 3: non-routine; more than one dependent variable</p>	<ul style="list-style-type: none"> • identify the types of evidence that should be gathered to answer the question, or support or refute the <u>hypothesis</u> or prediction; • identify the variables that may affect the outcome of the experiment or investigation; • design an appropriate format for recording data <u>and include relevant technology</u>; • evaluate multiple data sets to determine which data are relevant to the question, hypothesis or prediction. <p>Note: These items could appear at the end of the task.</p>
<p>5. Develop an organized and logical approach to investigating the question, including controlling variables</p> <p>DOK 2: routine (replicates existing procedure);</p> <p>DOK 3: non-routine (extends, refines, or improves existing procedure)</p>	<ul style="list-style-type: none"> • develop a procedure to gather sufficient evidence (including multiple trials) to answer the question, or test the hypothesis, or prediction; • develop a procedure that lists steps sequentially and logically and incorporates the use of <u>appropriate technology</u>; • explain which variable will be manipulated or changed (independent) and which variable will be affected by those changes (dependent); • identify variables that will be kept constant throughout the investigation; • distinguish between the control group and the experimental group in an investigation; • use scientific terminology that supports the identified procedures; • evaluate the organization and logical approach of a given procedure including variables, controls, materials, and tools. • <u>evaluate investigation design, including opportunities to collect appropriate and sufficient data.</u> <p>Note: These items could appear at the beginning or the end of the task.</p>

<p>6. Provide reasoning for appropriateness of materials, tools, procedures, and scale used in the investigation</p> <p>DOK 2</p>	<ul style="list-style-type: none">• explain why the materials, tools, procedure, or scale for a task are appropriate or inappropriate for the investigation.• evaluate the investigation for the safe and ethical considerations of the materials, tools, and procedures.
--	--

Broad Area 3: Conducting Investigations

Grade 11

Standard: The procedure requires the student to collect data through observation, inference, and prior scientific knowledge. Mathematics is required for the student to determine and report data. The task scenario is authentic to the realm of the student. The task requires the student to collect sufficient data to investigate the question, prediction/hypothesis, or relationships. Student is required to organize and represent qualitative or quantitative data. Student is required to summarize data to form a logical argument.

Note: Metric units are used for all Grade 11 measurements

Note: Multiple trials mean repeating the experiment to collect multiple sets of data.

Inquiry Construct	Items addressing this construct require students to:
<p>7. Follow procedures for collecting and recording qualitative or quantitative data, using equipment or measurement devices accurately</p> <p>DOK 1: use tools; routine procedure;</p> <p>DOK 2: follow multi-step procedures; make observations</p>	<ul style="list-style-type: none"> • record precise data and observations that are consistent with the procedure of the investigation; • include appropriate units of all measurements; • use appropriate measurement tools correctly to collect data; record and label relevant details within a scientific drawing.
<p>8. Use accepted methods for organizing, representing, and manipulating data</p> <p>DOK 2 : compare data; display data</p>	<ul style="list-style-type: none"> • represent data accurately in an appropriate graph/ table/ chart; • include titles, labels, keys or symbols as needed; • select a scale appropriate for the range of data to be plotted; • use scientific terminology to label representations; • identify relationships among variables based upon evidence. <p>Note: The standard practice of graphing in science is to represent the independent on the x-axis and the dependent variable on the y- axis.</p>
<p>9. Collect sufficient data to study question, hypothesis, or relationships</p> <p>DOK 2 : part of following procedures</p>	<ul style="list-style-type: none"> • show understanding of the value of multiple trials • relate data to original question, hypothesis or prediction; • determine if the quantity of data is sufficient to answer the question or support or refute the hypothesis or prediction.
<p>10. Summarize results based on data</p> <p>DOK 2</p>	<ul style="list-style-type: none"> • consider all data when developing an explanation/conclusion; • identify patterns and trends in data.

Broad Area 4: Developing and Evaluating Explanations

Grade 11

Standard: Task must provide the opportunity for students to use data to construct an explanation based on their science knowledge and evidence from experiment or investigation. The task requires students to use qualitative and quantitative data to communicate conclusions and support/refute prediction/hypothesis. The task provides students the opportunity to recognize and analyze alternative methods and models to evaluate other plausible explanations.

Note: The complexity of the scenario and associated data sets distinguishes this task from an 8th Grade task.

Inquiry Construct	Items addressing this construct require students to:
<p>11. Analyze data, including determining if data are relevant, artifact, irrelevant, or anomalous</p> <p>DOK 2: specify relationships between facts; ordering, classifying data</p>	<ul style="list-style-type: none"> • identify data relevant to the task or question; • identify factors that may affect experimental results (e.g. variables, experimental error, environmental conditions); • <u>analyze data and sort</u> into meaningful categories; • <u>compare experimental data to accepted scientific data provided as part of the task</u>; • <u>use mathematical and statistical techniques to analyze data</u>; • <u>provide a reasonable explanation that accurately reflects data</u>; • <u>use content understanding to question data that might seem inaccurate</u> • evaluate the significance of experimental data.
<p>12. Use evidence to support and justify interpretations and conclusions or explain how the evidence refutes the hypothesis</p> <p>DOK 3</p>	<ul style="list-style-type: none"> • identify and explain data, interpretations or conclusions that seem inaccurate; • use evidence to support or refute question or hypothesis; • use evidence to justify an interpretation of data or trend; • identify and explain differences or similarities between <u>hypothesis</u> and predictions and experimental data; • <u>use evidence to justify a conclusion or explanation based on experimental data</u>; • use mathematical computations to determine or support conclusions; • <u>evaluate potential bias in the interpretation of evidence.</u>

<p>13. Communicate how scientific knowledge applies to explain results, propose further investigations, or construct and analyze alternative explanations</p> <p>DOK 3</p>	<ul style="list-style-type: none"> • explain how experimental results compare to accepted scientific understanding; • recommend changes to procedures to produce data that would provide sufficient data and more accurate analysis; • identify <u>and justify</u> additional data that would strengthen an investigation; • connect the investigation or model to an <u>authentic situation</u>; • propose <u>and evaluate</u> new questions, predictions, next steps or technology for further investigations or <u>alternative explanations</u>; • <u>account for limitations and/or sources of error within the experimental design</u>; • apply experimental results to a new problem or situation; • <u>consider the impact (safety, ethical, social, civic, economic, environmental) of additional investigations.</u>
--	---

APPENDIX

**NECAP Science
Inquiry Task Development Process**

Initial Steps for the Development of an Inquiry Task

1. Identify the NECAP Assessment **TARGET** to be addressed within the major idea for the task.
2. Refer to the **Guidelines for the Development of Science Inquiry Tasks (GDIT)**. Brainstorm constructs that would be addressed under each broad area within the major idea for the task.

Formulating Questions and Hypothesizing	Planning and Critiquing of Investigations	Conducting Investigations	Developing and Evaluating Explanation
--	--	----------------------------------	--

3. Develop a draft **SCENARIO** aligned to the major idea of the task that could generate testable questions.*
4. Identify an authentic **Data Set** (Grades 8 & 11) that applies to the **TARGET** and relates to the **SCENARIO** *

OR

Provide opportunity for **Collection of Data** (Grade 4 & 8) that applies to the **TARGET** and relates to the **SCENARIO** *

* **Note:** *The previous steps are interdependent. The construction of the draft SCENARIO and the identification of a data set, will inform one another. Either may necessitate modifications for alignment, as the task items are being developed.*

Components of the Final Inquiry Task

Each **Inquiry Task** must include:

- A cohesive series of scoreable items, totaling 16-18 points, that assess student understanding in each of the four broad areas of inquiry, as described in the **GDIT**.
- Scoreable items that have sufficient complexity for students to demonstrate scientific thinking related to inquiry.
- An identified DOK level for each scoreable item.
- A scoring rubric for each scoreable item.

APPENDIX B—NECAP SCIENCE COMMITTEE MEMBERS

**Table B-1. 2013–14 NECAP Science: Item Review Committee Participants
August 1 and 2, 2013**

<i>State</i>	<i>Name</i>	<i>School/Association Affiliation</i>	<i>Position</i>
New Hampshire	Mary Fougere	Gilmanton	Science teacher – Grades 7, 8
	Bonnie Skogsholm	Hollis Brookline	Teacher
	Marisa Bozek	Epping MS	Science teacher – Grade 7
	Cheryl Patty	Westmoreland School	Science teacher – Grades 5-8
	Nancy Morse	Hollis Brookline HS	ESOL teacher
	Sonya Roberts	Laconia MS	Science teacher
	Kelly Marcotte	SAU 43	Teacher – Grade 4
	Dennis Vienneau	Moultonborough Academy	Science teacher
	Mary Kate Hartwell	Timberlane Regional MS	Science teacher
	Pauline Corzilius	Lisbon Regional HS	Science teacher
	Debbie Maloney	Hollis/Brookline HS	Chemistry teacher
	Erik Anderson	Sage School	Teacher
Rhode Island	Gail DeRobbio	Daniel D. Waterman	Teacher – Grade 4
	Amy Gearing	Glen Hills	Teacher – Grade 4
	Kiara Tracey	Community School	Teacher – Grade 4
	Debra Turchetti-Ramm	S.D. Barnes ES	Teacher – Grade 4
	Stephen Cormier	Chariho MS	Science teacher
	Desiree Derix	Westerly MS	Science teacher
	Carolyn Higgins	Winman Jr. HS	Science teacher – Grades 7, 8
	Stephen Martin	Lincoln MS	Science teacher – Grade 7
	Karen Finlan	North Kingstown HS	Chemistry teacher
	Helaine Hager	Mt. Pleasant HS	Chemistry teacher
	Thomas Holstein, Jr.	Portsmouth HS	Chemistry teacher
	Heather Taylor	South Kingston School	STEM coordinator
Vermont	Lynn Murphy	Waits River Valley	MS Science teacher
	Sarah Cousino	Irasburg Village School	Teacher – Grades 4, 5
	Jenn Harper	Cavendish Town ES	Teacher – Grade 4
	Ann Thompson	Cavendish Town ES	Teacher – Grade 5
	Debbie Griggs	Blue Mountain Union	Science teacher – Grades 3, 4
	David White	Vermont Science Initiative	Program Coordinator
	Kaija Percy	Lyndon Ed. Alternative	Science teacher
	Susan Steiner	Bellows Falls Union HS	Science teacher
	Greg Renner	Retired Science Teacher	N/A
	Rory Waterman	UVM	Assoc. Professor
	Gay Craig	Champlain Valley Union HS	Science teacher
	Maureen Maidrand	Bennington-Rutland SU	Math/Sci Partnership

**Table B-2. 2013–14 NECAP Science: Bias and Sensitivity Committee Participants
August 1, 2013**

<i>State</i>	<i>Name</i>	<i>School/Association Affiliation</i>	<i>Position</i>
New Hampshire	Enchi Chen	Farmington HS	ESL
	Kathaleen Cobb	Cutler Elementary	Teacher – Grade 4
	Erika Langlais	Gilmanton ES	Special Education teacher
	Christine Leach	Exeter HS	Guidance Counselor
	Diane Bush	Jaffrey Rindge MS	School Counselor
Rhode Island	Barbara Cesana	RI School for Deaf	Math/Science teacher
	Kathleen Beebe	Portsmouth High School	Chemistry teacher
	Kimberly Bontempo	Bradford Elementary School	Teacher – Grade 4
	Darline Berrios	Paul Cuffee School	Teacher – Grade 5
	Alyssa Wood	Sophia Academy	Science teacher
Vermont	Suzanna Buck	Vergennes HS	Special Education
	Sue Vincent	Brattleboro HS	Science teacher
	Aranka Gyuk	Integrated Arts Academy	ELL teacher
	Brenda Seitz	INSPIRE for Autism Inc.	Director of Administrators & Admissions
	Rebekah Thomas	Barrington school district	ELL teacher

APPENDIX C—PARTICIPATION RATES

Table C-1. 2013–14 NECAP Science: Participation Rates

<i>Description</i>	<i>Tested</i>	
	<i>Number</i>	<i>Percent</i>
All Students	91,125	100.00
Male	46,846	51.41
Female	44,277	48.59
Gender Not Reported	2	0.00
Hispanic or Latino	8,673	9.52
American Indian or Alaskan Native	389	0.43
Asian	2,523	2.77
Black or African American	3,625	3.98
Native Hawaiian or Pacific Islander	124	0.14
White (non-Hispanic)	73,389	80.54
Two or More Races (non-Hispanic)	2,328	2.55
No Primary Race/Ethnicity Reported	74	0.08
Currently receiving LEP services	2,777	3.05
Former LEP student – monitoring year 1	766	0.84
Former LEP student – monitoring year 2	351	0.39
LEP: All Other Students	87,231	95.73
Students with an IEP	13,040	14.31
IEP: All Other Students	78,085	85.69
Economically Disadvantaged Students	31,985	35.10
SES: All Other Students	59,140	64.90
Migrant Students	26	0.03
Migrant: All Other Students	91,099	99.97
Students receiving Title 1 Services	14,807	16.25
Title 1: All Other Students	76,318	83.75
Plan 504	720	0.79
Plan 504: All Other Students	90,405	99.21

APPENDIX D—ACCOMMODATION FREQUENCIES

Table D-1. 2013–14 NECAP Science: Accommodation Frequencies

<i>Accommodation</i>	<i>Grade</i>		
	<i>4</i>	<i>8</i>	<i>11</i>
SCIAccomT1	3,407	2,407	1,252
SCIAccomT2	117	69	39
SCIAccomT3	2,412	1,279	673
SCIAccomT4	87	62	40
SCIAccomS1	3,850	2,182	1,917
SCIAccomS2	9	25	59
SCIAccomP1	769	456	198
SCIAccomP2	3,703	2,879	1,763
SCIAccomP3	3,165	1,241	442
SCIAccomP4	338	346	153
SCIAccomP5	1,262	440	198
SCIAccomP6	11	3	6
SCIAccomP7	1,564	511	362
SCIAccomP8	37	17	11
SCIAccomP9	1	1	4
SCIAccomP10	60	95	148
SCIAccomP11	90	30	3
SCIAccomR1	1,354	316	63
SCIAccomR2	33	20	3
SCIAccomR3	173	26	7
SCIAccomR4	38	155	51
SCIAccomR5	96	28	24
SCIAccomR6	3	2	2
SCIAccomR7	38	71	84
SCIAccomO1	6	0	0
SCIAccomM1	4	24	18
SCIAccomM3	1	0	2

APPENDIX E—NECAP TABLE OF STANDARD ACCOMMODATIONS

NECAP Table of Standard Accommodations

Revised August 2009

Any accommodation(s) used for the assessment of an individual student will be the result of a team decision made at the local level. All decisions regarding the use of accommodations must be made on an individual student basis – not for a large group, entire class, or grade level. Accommodations are available to all students on the basis of individual need regardless of disability status and should be consistent with the student’s normal routine during instruction and assessment. This table is not intended to be used as a stand-alone document and should always be used in conjunction with the *NECAP Accommodations Guide*.

T. Timing		
Code	Tests were administered	Details on Delivery of Accommodations
T1	with time to complete a session extended beyond the scheduled administration time within the same day.	NECAP tests are not designed to be timed or speeded tests. The scheduled administration time already includes additional time and the vast majority of students complete the test session within that time period. Extended time within a single sitting may be needed by students who are unable to meet time constraints. A test session may be extended until the student can no longer sustain the activity.
T2	so that only a portion of the test session was administered on a particular day.	In rare and severe cases, the extended time accommodation (T1) may not be adequate for a student not able to complete a test session within a single day. A test session may be administered to a student as two or more “mini-sessions” if procedures are followed to maintain test security and ensure that the student only has access to the items administered on that day (see the <i>NECAP Accommodations Guide</i> for details).
T3	with short, supervised breaks.	Multiple or frequent breaks may be required by a student whose attention span, distractibility, or physical condition, requires shorter working periods.
T4	at the time of day or day of week that takes into account the student’s medical needs or learning style.	Individual scheduling may be used for a student whose school performance is noticeably affected by the time of day or day of the school week on which it is done. This accommodation may not be used specifically to change the order of administration of test sessions. This accommodation must not result in the administration of a test session to an individual student prior to the regularly scheduled administration time for that session for all students.

S. Setting		
Code	Tests were administered	Details on Delivery of Accommodations
S1	in a separate location within the school by trained school personnel.	A student or students may be tested individually or in small groups in an alternative site within the school to reduce distractions for themselves or others, or to increase physical access to special equipment.
S2	in an out-of-school setting by trained school personnel.	Out-of-school testing may be used for a student who is hospitalized or tutored because they are unable to attend school. The test must be administered by trained school personnel familiar with test administration procedures and guidelines. Relatives/guardians of the student may not be used as the test administrator.

P. Presentation		
Code	Tests were administered	Details on Delivery of Accommodations
P1	individually.	Individual or small group testing may be used to minimize distractions for a student or students whose test is administered out of the classroom or so that others will not be distracted by other accommodations being used (e.g., dictation)
P2	in a small group.	
P3	with test and directions read aloud in English or signed to the student. (NOT allowed for the Reading test.)	A reader may be used for a student whose inability to read would hinder performance on the Mathematics, Science, or Writing test. Words must be read as written. Guidelines for reading mathematical symbols must be followed. No translations (with the exception of signed language) or explanations are allowed. Trained personnel may use sign language to administer the test.
P4	with only test directions read aloud or signed to the student.	A reader may be used for a student whose inability to read or locate directions would hinder performance on the test. Note that most directions on the NECAP test occur at the beginning of the test session and are already read aloud by the test administrator. Guidelines for what are and are not “test directions” must be followed. With the exception of sign language and the case of students enrolled in a program where the test administrator routinely presents information in a foreign language, directions may not be translated.
P5	with administrator verification of student understanding following the reading of test directions.	After <u>test directions</u> have been read, the test administrator may ask the student to explain what he/she has been asked to do. If directions have been misunderstood by the student, the <u>test directions</u> may be paraphrased or demonstrated. Test items MUST NOT be paraphrased or explained.
P6	using alternative or assistive technology that is part of the student’s communication system.	The test may be presented through his/her regular communication system to a student who uses alternative or assistive technology on a daily basis. Technology may not be used to “read” the Reading test to the student.
P7	by trained school personnel known to the student other than the student’s classroom teacher.	A student may be more comfortable with a test administrator who works with the student on a regular basis, but is not the student’s regular teacher for the general curriculum or other staff assigned as test administrator. All test administrators must be trained school personnel familiar with test administration and accommodations procedures and guidelines.
P8	using a large-print version of assessment.	Both large-print and Braille versions of the assessment require special preparation and processing and must be pre-ordered. Directions for ordering these materials are included in communications sent to school principals prior to the test.
P9	using Braille version of assessment.	
P10	using a word-to-word translation dictionary for ELL students. (NOT allowed for the Reading test.)	A student with limited English proficiency may have a word-to-word dictionary available for individual use as needed. A word-to-word dictionary is one that does not include any definitions. Information on acceptable dictionaries is provided on the departments’ websites.
P11	using visual or auditory supports.	The test may be presented using visual aids such as visual magnification devices, reduction of visual print by blocking or other techniques, or acetate shields; or auditory devices such as special acoustics, amplification, noise buffers, whisper phones, or calming music.

R. Response		
Code	Tests were administered	Details on Delivery of Accommodations
R1	with a student <u>dictating</u> responses to school personnel. (NOT allowed for the Writing test. See O2 – using a scribe for the Writing test.)	A student may dictate answers to constructed-response or short-answer questions to locally trained personnel or record oral answers in an individual setting so that other students will not benefit by hearing answers or be otherwise disturbed. Policies regarding recorded answers must be followed prior to returning test materials.
R2	with a student <u>dictating</u> responses using alternative or assistive technology/devices that are part of the student’s communication system. (NOT allowed for the Writing test. See O2 – using a scribe for the Writing test.)	Technology is used to permit a student to respond to the test. When using a computer, word processing device, or other assistive technology, spell and grammar checks must be turned off. Policies regarding recorded answers must be followed prior to returning test materials.
R3	with a student using approved tools or devices to minimize distractions.	Noise buffers, place markers, carrels, etc. may be used to minimize distractions for the student. This accommodation does NOT include assistive devices such as templates, graphic organizers, or other devices intended specifically to help students <u>organize thinking or develop a strategy for a specific question.</u>
R4	with a student <u>writing</u> responses using separate paper, a word processor, computer, braille, or similar device.	A student may use technological or other tools (e.g., large-spaced paper) to write responses to constructed-response, short-answer, and extended response items. A key distinction between this accommodation and R2 is that the student using this accommodation is responding in writing rather than dictating. When using a computer, word processing device, or other assistive technology, spell and grammar checks must be turned off, as well as access to the Web. This accommodation is intended for unique individual needs, not an entire class. Policies regarding recorded answers must be followed prior to returning test materials.
R5	with a student indicating responses to multiple-choice items to school personnel.	A student unable to write or otherwise unable to fill-in answers to multiple-choice questions may indicate a response to trained school personnel. The school personnel records the student’s response in the student answer booklet.
R6	with a student responding with the use of visual aids.	Visual aids include any optical or non-optical devices used to enhance visual capability. Examples include magnifiers, special lighting, markers, filters, large-spaced paper, color overlays, etc. An abacus may also be used for student with severe visual impairment or blindness on the Mathematics and Science tests. Note that the use of this accommodation still requires student responses to be recorded in a student answer booklet.
R7	with a student with limited English proficiency responding with use of a word-to-word dictionary. (NOT allowed for the Reading test.)	A student with limited English proficiency may have a word-to-word dictionary available for individual use as needed when responding. A word-to-word dictionary is one that does not include any definitions. Information on acceptable dictionaries is provided on each Department’s website.

O. Other		
These accommodations require DOE approval.		
Code	Tests were administered	Details on Delivery of Accommodations
O1	using other accommodation(s) not on this list, requested by the accommodations team.	An IEP team or other appropriate accommodation team may request that a student be provided an accommodation not included on this standard list of accommodations. Like all other accommodations, these should be consistent with the student's normal routine during instruction and/or assessment. Requests should be made to the DOE when accommodation plans are being made for a student prior to testing. DOE approval must be received for the requested accommodation to be coded as an O1 accommodation. Non-approved accommodations used during test administration will be coded as an M3 modification.
O2	with a scribe used on the Writing test.	The use of a scribe for students dictating a response to the Writing test may only be used under limited circumstances and must be approved by the DOE. When approved as an accommodation, the scribe must follow established guidelines and procedures.

M. Modifications		
All modifications result in impacted items being scored as incorrect.		
Code	Tests were administered	Details on Delivery of Accommodations
M1	using a calculator and/or manipulatives on Session 1 of the Mathematics test or using a scientific or graphing calculator on Session 3 of the Science test	Inappropriate use of a calculator or other tools will result in impacted items being scored as incorrect.
M2	with the test administrator reading the Reading test.	The read aloud accommodation (P3) is not allowed for the Reading test. If it is used, all reading items in the sessions that are read aloud will be scored as incorrect.
M3	using an accommodation on this list not approved for a particular test or an accommodation not included on this list without prior approval of the DOE.	Inappropriate use of an accommodation included on this list or use of another accommodation without prior approval of the DOE will result in impacted items being scored as incorrect.

N. NimbleTools® 2011		
The NECAP Science test was administered using NimbleTools online accommodations.		
Code		
N01	Allow breaks	
N02	Read Aloud Text	
N03	Magnifier	
N04	Custom Masking or Answer Masking	
N05	Color Overlay	
N06	Reverse Contrast	
N07	Font & Background Color Choice	
N08	Auditory Calming	
N09	Read Aloud Text and Graphics	
N10	Microscope	
N11	Magnifying Glass	

Note: English Language Learners may qualify for any of the accommodations listed as appropriate and determined by a team. Refer to the *NECAP Accommodations Guide* for additional information.

APPENDIX F—ITEM-LEVEL CLASSICAL STATISTICS

**Table F-1. 2013–14 NECAP Science: Item-Level Classical Test Theory Statistics—
Grade 4**

<i>Item</i>		<i>Difficulty</i>	<i>Discrimination</i>	<i>Percent Omitted</i>	<i>Item</i>		<i>Difficulty</i>	<i>Discrimination</i>	<i>Percent Omitted</i>
<i>Number</i>	<i>Type</i>				<i>Number</i>	<i>Type</i>			
46316	MC	0.73	0.40	1	221536	MC	0.92	0.36	0
46426	MC	0.80	0.35	1	241791	MC	0.45	0.29	2
46525	MC	0.85	0.33	1	241812	MC	0.84	0.41	1
47974	MC	0.57	0.50	0	241815	CR	0.42	0.59	1
48062	CR	0.38	0.46	0	241826	MC	0.71	0.28	0
50509	MC	0.81	0.25	0	241848	MC	0.63	0.27	0
59271	MC	0.56	0.32	1	241967	MC	0.70	0.31	0
59906	MC	0.69	0.27	0	242407	MC	0.67	0.27	0
59915	MC	0.80	0.45	1	242411	MC	0.68	0.32	0
60292	MC	0.72	0.42	1	242420	MC	0.81	0.34	0
60342	MC	0.77	0.33	0	242423	MC	0.47	0.25	0
60389	MC	0.66	0.32	1	242427	MC	0.40	0.30	2
76741	MC	0.88	0.38	1	242739	MC	0.86	0.30	0
91984	MC	0.76	0.45	1	258700	CR	0.40	0.41	1
99021	CR	0.41	0.53	1	258702	SA	0.46	0.44	0
140765	MC	0.65	0.35	1	258704	SA	0.15	0.31	1
142211	MC	0.46	0.30	1	258705	SA	0.33	0.37	1
142265	MC	0.59	0.31	0	258708	SA	0.19	0.38	2
144894	MC	0.62	0.33	1	258711	SA	0.34	0.28	3
171051	MC	0.48	0.25	0	258729	CR	0.17	0.45	2
174255	MC	0.53	0.24	0	258731	SA	0.43	0.31	2
174283	MC	0.79	0.38	1					
174659	MC	0.78	0.42	0					

**Table F-2. 2013–14 NECAP Science: Item-Level Classical Test Theory Statistics—
Grade 8**

<i>Item</i>		<i>Difficulty</i>	<i>Discrimination</i>	<i>Percent Omitted</i>	<i>Item</i>		<i>Difficulty</i>	<i>Discrimination</i>	<i>Percent Omitted</i>
<i>Number</i>	<i>Type</i>				<i>Number</i>	<i>Type</i>			
46064	MC	0.65	0.38	0	91780	MC	0.76	0.35	0
46068	MC	0.73	0.42	1	144299	MC	0.64	0.35	1
47786	MC	0.47	0.45	0	144301	MC	0.65	0.26	1
48239	MC	0.83	0.40	1	174708	CR	0.34	0.63	1
48252	MC	0.52	0.35	1	174966	MC	0.57	0.36	0
48266	MC	0.72	0.34	1	174976	MC	0.80	0.32	0
48269	MC	0.49	0.28	1	174984	MC	0.84	0.40	0
48343	MC	0.52	0.35	0	175340	MC	0.70	0.35	1
48554	CR	0.37	0.55	1	175377	MC	0.48	0.41	1
50134	MC	0.61	0.28	0	218636	MC	0.40	0.14	1
50140	MC	0.69	0.35	0	218646	MC	0.59	0.38	1
58316	MC	0.65	0.54	1	219112	CR	0.46	0.62	2
59767	MC	0.35	0.37	0	219238	MC	0.49	0.51	1
60201	MC	0.80	0.27	0	219441	MC	0.82	0.46	1
87085	MC	0.36	0.17	0	219445	MC	0.55	0.32	1
91735	MC	0.46	0.31	1	219464	MC	0.92	0.38	0

continued

<i>Item</i>		<i>Difficulty</i>	<i>Discrimination</i>	<i>Percent Omitted</i>
<i>Number</i>	<i>Type</i>			
219466	MC	0.78	0.39	0
220376	MC	0.85	0.27	1
241626	MC	0.50	0.24	0
242842	MC	0.63	0.34	0
258713	SA	0.59	0.44	1
258714	SA	0.60	0.54	1

<i>Item</i>		<i>Difficulty</i>	<i>Discrimination</i>	<i>Percent Omitted</i>
<i>Number</i>	<i>Type</i>			
258715	SA	0.59	0.52	1
258716	SA	0.74	0.49	3
258717	CR	0.67	0.50	2
258719	SA	0.35	0.37	2
258720	CR	0.14	0.39	2
258721	SA	0.12	0.34	3

**Table F-3. 2013–14 NECAP Science: Item-Level Classical Test Theory Statistics—
Grade 11**

<i>Item</i>		<i>Difficulty</i>	<i>Discrimination</i>	<i>Percent Omitted</i>
<i>Number</i>	<i>Type</i>			
46076	MC	0.61	0.28	1
46140	MC	0.73	0.27	1
46184	MC	0.59	0.34	1
47862	MC	0.40	0.29	1
48081	MC	0.76	0.40	1
48115	MC	0.80	0.38	1
48214	MC	0.57	0.19	1
49913	MC	0.88	0.34	1
49914	MC	0.83	0.43	1
49918	MC	0.77	0.41	1
59605	MC	0.84	0.36	1
60693	MC	0.48	0.28	2
61834	MC	0.46	0.27	1
62084	MC	0.60	0.33	1
89544	MC	0.45	0.31	1
89628	MC	0.49	0.21	1
90282	MC	0.47	0.29	1
90297	MC	0.68	0.51	1
91920	MC	0.66	0.48	2
135344	MC	0.58	0.18	1
146766	MC	0.31	0.27	1
146785	MC	0.44	0.35	1
147035	MC	0.59	0.35	1
176787	MC	0.49	0.35	2

<i>Item</i>		<i>Difficulty</i>	<i>Discrimination</i>	<i>Percent Omitted</i>
<i>Number</i>	<i>Type</i>			
176922	MC	0.47	0.18	1
177093	MC	0.59	0.40	1
219134	MC	0.35	0.30	1
219201	MC	0.38	0.21	1
219212	MC	0.48	0.33	1
219607	MC	0.70	0.38	1
219612	MC	0.81	0.47	1
241745	CR	0.26	0.56	8
241752	MC	0.31	0.16	1
242626	MC	0.30	0.20	1
242632	CR	0.39	0.61	4
242637	CR	0.41	0.60	4
258655	CR	0.37	0.61	2
258656	SA	0.47	0.54	2
258657	SA	0.45	0.54	3
258659	SA	0.30	0.45	6
258661	SA	0.47	0.57	3
258662	SA	0.45	0.60	5
258669	SA	0.43	0.55	5
263524	CR	0.36	0.62	6

APPENDIX G—ITEM-LEVEL SCORE POINT DISTRIBUTIONS

Table G-1. 2013–14 NECAP Science: Item-Level Score Distributions for Constructed Response Items by Subject and Grade

Grade	Item Number	Total Possible Points	Percent of Students at Score Points				
			0	1	2	3	4
4	48062	4	24.85	25.39	25.55	18.02	5.80
	99021	4	16.79	29.70	28.19	17.99	5.93
	241815	4	23.24	19.99	25.19	25.29	5.41
	258700	3	35.65	23.54	23.54	16.50	
	258702	2	31.71	43.44	24.49		
	258704	2	72.08	24.14	2.74		
	258705	2	51.84	29.65	17.79		
	258708	2	68.84	20.83	8.59		
	258711	2	42.35	41.47	13.67		
	258729	3	53.26	38.04	5.78	0.47	
	258731	2	23.05	64.89	10.44		
8	48554	4	22.18	23.05	39.21	12.59	1.73
	174708	4	35.85	23.51	17.32	12.78	9.48
	219112	4	11.39	18.52	43.72	19.90	4.43
	258713	2	18.07	45.42	35.91		
	258714	2	23.44	32.43	43.33		
	258715	2	27.67	25.23	46.25		
	258716	2	13.23	20.13	63.96		
	258717	3	12.20	13.59	30.42	41.77	
	258719	2	45.60	33.44	18.60		
	258720	3	67.06	20.77	7.11	2.67	
	258721	2	77.79	14.31	5.01		
11	241745	4	34.75	22.76	24.61	8.57	1.76
	242632	4	26.54	9.72	36.55	18.69	4.10
	242637	4	17.66	28.01	24.74	15.15	10.00
	258655	3	26.19	39.74	23.34	8.30	
	258656	2	36.04	30.00	31.54		
	258657	2	30.41	43.73	22.74		
	258659	2	52.74	23.32	18.03		
	258661	2	23.73	51.14	21.91		
	258662	2	31.97	35.94	27.07		
	258669	2	31.59	39.62	23.35		
	263524	3	31.69	26.36	27.20	8.93	

APPENDIX H—DIFFERENTIAL ITEM FUNCTIONING RESULTS

**Table H-1. 2013–14 NECAP Science: Number of Items Classified as “Low” or “High” DIF
Overall and by Group Favored**

Grade	Group		Item Type	Number of Items	Number “Low”			Number “High”			
	Reference	Focal			Total	Favoring		Total	Favoring		
						Reference	Focal		Reference	Focal	
4	Male	Female	CR	5	0	0	0	0	0	0	
			MC	33	4	2	2	0	0	0	
			SA	6	0	0	0	0	0	0	
	No Disability	Disability	CR	5	1	1	0	0	0	0	
			MC	33	10	10	0	0	0	0	
			SA	6	1	0	1	0	0	0	
	Non-EconDis	EconDis	CR	5	0	0	0	0	0	0	
			MC	33	2	2	0	0	0	0	
			SA	6	0	0	0	0	0	0	
	Non-LEP	LEP	CR	5	0	0	0	0	0	0	
			MC	33	6	6	0	1	1	0	
			SA	6	0	0	0	0	0	0	
	White	Asian	CR	5	1	0	1	0	0	0	
			MC	33	2	2	0	0	0	0	
			SA	6	0	0	0	0	0	0	
		Black	CR	5	0	0	0	0	0	0	
			MC	33	5	5	0	0	0	0	
			SA	6	0	0	0	0	0	0	
	Hispanic	CR	5	1	0	1	0	0	0		
		MC	33	3	3	0	1	1	0		
		SA	6	0	0	0	0	0	0		
	8	Male	Female	CR	5	1	0	1	0	0	0
				MC	33	7	6	1	0	0	0
				SA	6	2	0	2	0	0	0
No Disability		Disability	CR	5	0	0	0	1	1	0	
			MC	33	2	1	1	0	0	0	
			SA	6	4	4	0	0	0	0	
Non-EconDis	EconDis	CR	5	0	0	0	0	0	0		
		MC	33	0	0	0	0	0	0		
		SA	6	0	0	0	0	0	0		

continued

Grade	Group		Item Type	Number of Items	Number "Low"			Number "High"			
	Reference	Focal			Total	Favoring		Total	Favoring		
						Reference	Focal		Reference	Focal	
8	Non-LEP	LEP	CR	5	2	2	0	0	0	0	
			MC	33	7	5	2	0	0	0	
			SA	6	0	0	0	1	1	0	
	White	Asian	CR	5	0	0	0	0	0	0	
			MC	33	2	1	1	0	0	0	
			SA	6	0	0	0	0	0	0	
		Black	CR	5	0	0	0	0	0	0	
			MC	33	6	6	0	0	0	0	
			SA	6	1	1	0	0	0	0	
	Hispanic	CR	5	0	0	0	0	0	0		
		MC	33	5	5	0	0	0	0		
		SA	6	1	1	0	0	0	0		
11	Male	Female	CR	5	1	0	1	0	0	0	
			MC	33	3	3	0	1	1	0	
			SA	6	1	0	1	0	0	0	
	No Disability	Disability	CR	5	3	3	0	0	0	0	
			MC	33	1	0	1	0	0	0	
			SA	6	4	4	0	0	0	0	
	Non-EconDis	EconDis	CR	5	0	0	0	0	0	0	
			MC	33	0	0	0	0	0	0	
			SA	6	0	0	0	0	0	0	
	White	Non-LEP	LEP	CR	5	1	1	0	0	0	0
				MC	33	10	7	3	1	1	0
				SA	6	2	2	0	0	0	0
Asian		CR	5	0	0	0	0	0	0		
		MC	33	3	1	2	1	0	1		
		SA	6	0	0	0	0	0	0		
		Black	CR	5	0	0	0	0	0	0	
			MC	33	4	4	0	1	1	0	
			SA	6	0	0	0	0	0	0	
Hispanic	CR	5	0	0	0	0	0	0			
	MC	33	2	2	0	0	0	0			
	SA	6	0	0	0	0	0	0			

APPENDIX I—ITEM RESPONSE THEORY CALIBRATION RESULTS

**Table I-1. 2013–14 NECAP Science: IRT Parameters for Dichotomous Items—
Grade 4**

Item Number	Parameters and Measures of Standard Error						Item Number	Parameters and Measures of Standard Error					
	a	SE (a)	b	SE (b)	c	SE (c)		a	SE (a)	b	SE (b)	c	SE (c)
171051	0.51419	0.02358	0.65311	0.05259	0.18514	0.01618	242423	0.45816	0.02183	0.57887	0.06468	0.14288	0.01928
174255	0.31996	0.00750	-0.33797	0.02211	0.00000	0.00000	242407	0.57098	0.02345	-0.28535	0.07816	0.30319	0.02220
221536	0.94798	0.01904	-2.15528	0.03907	0.05636	0.02297	174659	0.84129	0.01895	-1.15585	0.04125	0.12672	0.02045
60342	0.74415	0.02403	-0.72159	0.06007	0.34689	0.02037	46426	0.62552	0.01046	-1.73292	0.02385	0.00000	0.00000
241826	0.43777	0.01259	-1.31354	0.09214	0.07271	0.02929	242739	0.59184	0.01127	-2.28044	0.03448	0.00000	0.00000
242411	0.51961	0.01621	-0.87346	0.07953	0.10354	0.02753	76741	0.84157	0.01578	-1.89428	0.03627	0.04560	0.01883
241848	0.50479	0.02254	-0.16087	0.08885	0.24944	0.02456	142211	0.60672	0.02165	0.52711	0.03597	0.14852	0.01251
47974	1.16994	0.02183	-0.18178	0.01444	0.10848	0.00721	142265	0.46368	0.00825	-0.67523	0.01754	0.00000	0.00000
60292	0.82423	0.01912	-0.80143	0.03683	0.14256	0.01687	50509	0.45915	0.01466	-1.95369	0.13737	0.12582	0.04735
60389	0.66704	0.02202	-0.33447	0.05402	0.25776	0.01822	241812	0.88756	0.01889	-1.50492	0.04169	0.09741	0.02297
144894	0.52568	0.01621	-0.53582	0.06502	0.08926	0.02233	59906	0.40131	0.00823	-1.42236	0.02918	0.00000	0.00000
46525	0.65319	0.01149	-2.05814	0.02808	0.00000	0.00000	174283	0.69931	0.01453	-1.43407	0.04531	0.05882	0.02095
46316	0.72341	0.01662	-1.02702	0.04440	0.08849	0.01998	59915	0.94280	0.01881	-1.25808	0.03234	0.08666	0.01768
91984	0.92077	0.01910	-1.01537	0.03169	0.11064	0.01632	140765	0.55964	0.01376	-0.79567	0.05095	0.05579	0.01890
59271	0.60703	0.01995	-0.03627	0.04763	0.15982	0.01652	241791	0.52902	0.02009	0.50721	0.04413	0.11107	0.01480
242427	0.73765	0.02361	0.76038	0.02339	0.14747	0.00841							
242420	0.60537	0.01045	-1.83907	0.02607	0.00000	0.00000							
241967	0.50422	0.01335	-1.11111	0.07023	0.06672	0.02473							

**Table I-2. 2013–14 NECAP Science: IRT Parameters for Polytomous Items—
Grade 4**

Item Number	Parameters and Measures of Standard Error													
	a	SE (a)	b	SE (b)	D0	SE (D0)	D1	SE (D1)	D2	SE (D2)	D3	SE (D3)	D4	SE (D4)
258702	0.61176	0.00459	0.09278	0.01065	1.08504	0.01278	-1.08504	0.01365	0	0				
258704	0.52234	0.00000	2.70782	0.01954	1.54498	0.01533	-1.54498	0.03959	0	0				
258705	0.52203	0.00460	0.91022	0.01354	0.91442	0.01375	-0.91442	0.01745	0	0				
258708	0.66312	0.00670	1.58850	0.01433	0.76013	0.01210	-0.76013	0.01874	0	0				
258711	0.34740	0.00273	1.31257	0.01903	1.84300	0.02008	-1.84300	0.02837	0	0				
258731	0.41936	0.00274	0.63057	0.01674	2.49456	0.01934	-2.49456	0.02660	0	0				

continued

Item Number	Parameters and Measures of Standard Error													
	a	SE (a)	b	SE (b)	D0	SE (D0)	D1	SE (D1)	D2	SE (D2)	D3	SE (D3)	D4	SE (D4)
258700	0.54286	0.00383	0.47508	0.01159	1.31811	0.01375	0.09442	0.01339	-1.41253	0.01733	0	0		
258729	0.72103	0.00579	2.46546	0.01049	2.35356	0.01051	-0.05166	0.01992	-2.30189	0.06817	0	0		
241815	0.92459	0.00552	0.34422	0.00680	1.49918	0.00984	0.66920	0.00847	-0.27257	0.00898	-1.89580	0.01698	0	0
48062	0.62162	0.00368	0.63921	0.00986	2.03460	0.01349	0.76504	0.01167	-0.53680	0.01340	-2.26285	0.02363	0	0
99021	0.76892	0.00442	0.40545	0.00803	2.02303	0.01253	0.60954	0.00977	-0.59890	0.01112	-2.03367	0.01907	0	0

**Table I-3. 2013–14 NECAP Science: IRT Parameters for Dichotomous Items—
Grade 8**

Item Number	Parameters and Measures of Standard Error						Item Number	Parameters and Measures of Standard Error					
	a	SE (a)	b	SE (b)	c	SE (c)		a	SE (a)	b	SE (b)	c	SE (c)
174984	0.91006	0.02211	-1.11796	0.04594	0.22006	0.02228	87085	0.61607	0.03749	1.96148	0.04606	0.23466	0.00857
174966	1.09820	0.03025	0.49290	0.01955	0.30871	0.00759	60201	0.50401	0.01957	-1.33440	0.14488	0.23017	0.04585
91780	0.62989	0.01811	-0.99761	0.07242	0.15799	0.02829	242842	0.49314	0.00844	-0.62850	0.01788	0.00000	0.00000
241626	0.62915	0.02794	0.98939	0.03816	0.27350	0.01167	91735	0.78595	0.02596	0.86970	0.02431	0.21350	0.00862
50134	0.73019	0.02775	0.52041	0.03895	0.34741	0.01183	48266	0.58238	0.01752	-0.80256	0.07365	0.13969	0.02694
50140	0.61405	0.01865	-0.54505	0.06323	0.16819	0.02278	48269	0.71038	0.02599	0.82668	0.03051	0.24117	0.01026
47786	1.42522	0.03052	0.53901	0.01140	0.18215	0.00527	219441	1.11247	0.02379	-0.93140	0.02943	0.19767	0.01588
59767	1.23204	0.02958	0.94694	0.01223	0.13723	0.00454	48252	0.71062	0.02110	0.40496	0.03081	0.16802	0.01145
48343	0.65449	0.01940	0.33019	0.03413	0.12969	0.01260	58316	1.57174	0.02964	-0.08436	0.01220	0.18486	0.00684
174976	0.53059	0.00963	-1.70424	0.02972	0.00000	0.00000	48239	0.77016	0.01351	-1.43003	0.03138	0.03300	0.01385
219238	1.45734	0.02782	0.37466	0.01044	0.13621	0.00513	219466	0.72138	0.01725	-1.07930	0.05252	0.11132	0.02363
175340	0.55264	0.01391	-0.88537	0.06024	0.06627	0.02235	219445	0.48989	0.01470	-0.02218	0.05333	0.05716	0.01736
218646	0.70690	0.01919	-0.01431	0.03602	0.14901	0.01387	144301	0.74132	0.02982	0.47184	0.04276	0.40995	0.01195
218636	0.48454	0.03981	2.19414	0.06755	0.27140	0.01223	144299	0.56596	0.01492	-0.51458	0.05311	0.06946	0.01941
46068	0.74072	0.01584	-0.83610	0.03847	0.07027	0.01722	220376	0.49285	0.01026	-2.32260	0.04405	0.00000	0.00000
175377	1.13154	0.02652	0.55826	0.01485	0.18491	0.00637							
219464	1.01953	0.02291	-1.84622	0.04676	0.11127	0.02967							
46064	0.65997	0.01807	-0.34208	0.04653	0.13549	0.01779							

**Table I-4. 2013–14 NECAP Science: IRT Parameters for Polytomous Items—
Grade 8**

<i>Item Number</i>	<i>Parameters and Measures of Standard Error</i>													
	<i>a</i>	<i>SE (a)</i>	<i>b</i>	<i>SE (b)</i>	<i>D0</i>	<i>SE (D0)</i>	<i>D1</i>	<i>SE (D1)</i>	<i>D2</i>	<i>SE (D2)</i>	<i>D3</i>	<i>SE (D3)</i>	<i>D4</i>	<i>SE (D4)</i>
258713	0.59715	0.00448	-0.46218	0.01106	1.21309	0.01557	-1.21309	0.01263	0	0				
258714	0.81159	0.00662	-0.35832	0.00850	0.68770	0.01107	-0.68770	0.00949	0	0				
258715	0.75662	0.00663	-0.29960	0.00916	0.54629	0.01108	-0.54629	0.01002	0	0				
258716	0.74269	0.00709	-1.04664	0.01131	0.55846	0.01362	-0.55846	0.01060	0	0				
258719	0.51058	0.00429	0.98272	0.01333	1.01409	0.01407	-1.01409	0.01754	0	0				
258721	0.61796	0.00000	2.42445	0.02050	0.78977	0.01440	-0.78977	0.02485	0	0				
258717	0.69568	0.00502	-0.75969	0.00946	1.06074	0.01499	0.15541	0.01179	-1.21615	0.01082	0	0		
258720	0.65406	0.00000	2.37102	0.01379	1.40323	0.01208	-0.03715	0.01777	-1.36608	0.03206	0	0		
174708	1.16793	0.00754	0.72665	0.00567	1.08046	0.00747	0.30670	0.00710	-0.34009	0.00807	-1.04708	0.01105	0	0
219112	0.98272	0.00565	0.36270	0.00653	1.89127	0.01166	0.91932	0.00866	-0.69410	0.00910	-2.11648	0.01738	0	0
48554	0.81954	0.00486	1.03012	0.00779	2.09636	0.01103	1.08648	0.00934	-0.69403	0.01257	-2.48881	0.03181	0	0

**Table I-5. 2013–14 NECAP Science: IRT Parameters for Dichotomous Items—
Grade 11**

<i>Item Number</i>	<i>Parameters and Measures of Standard Error</i>						<i>Item Number</i>	<i>Parameters and Measures of Standard Error</i>					
	<i>a</i>	<i>SE (a)</i>	<i>b</i>	<i>SE (b)</i>	<i>c</i>	<i>SE (c)</i>		<i>a</i>	<i>SE (a)</i>	<i>b</i>	<i>SE (b)</i>	<i>c</i>	<i>SE (c)</i>
49913	0.76130	0.01307	-1.78645	0.02601	0.00000	0.00000	176922	0.75388	0.04363	1.76186	0.03608	0.36109	0.00779
219612	1.09017	0.02274	-0.82902	0.02937	0.14958	0.01629	135344	0.23922	0.00725	-0.63882	0.03822	0.00000	0.00000
219607	0.63554	0.01582	-0.61142	0.05217	0.08059	0.02103	176787	0.67643	0.02166	0.71490	0.03150	0.15614	0.01138
147035	0.92714	0.02829	0.60771	0.02585	0.31415	0.00918	46184	0.51997	0.01567	-0.06711	0.05819	0.07164	0.01965
49914	1.00093	0.02310	-0.91858	0.03848	0.19834	0.02005	59605	0.69005	0.01441	-1.45121	0.05156	0.06137	0.02427
90297	1.13736	0.02324	-0.12879	0.01986	0.17151	0.00983	46076	0.43151	0.01837	-0.09809	0.10815	0.12813	0.03045
49918	0.75688	0.01660	-0.89697	0.04390	0.07953	0.02060	219201	0.69868	0.03574	1.77874	0.03326	0.24483	0.00803
47862	0.79517	0.02888	1.31709	0.02309	0.20156	0.00753	46140	0.40723	0.00847	-1.40311	0.03445	0.00000	0.00000
89628	0.41042	0.02724	1.24033	0.08548	0.21830	0.02236	219212	0.75454	0.02517	0.91140	0.02699	0.21080	0.00955
89544	0.84090	0.02857	1.09996	0.02291	0.23452	0.00789	48214	0.25184	0.00727	-0.43163	0.03300	0.00000	0.00000
48081	0.75933	0.01906	-0.69453	0.04909	0.15124	0.02152	242626	0.47442	0.02952	2.13049	0.04840	0.12874	0.01177
61834	0.77923	0.03092	1.25762	0.02655	0.26711	0.00839	219134	1.08081	0.03443	1.39030	0.01676	0.19430	0.00513
241752	0.49275	0.03745	2.47433	0.06789	0.18085	0.01084	177093	0.75340	0.02034	0.18028	0.03339	0.16647	0.01296

continued

Item Number	Parameters and Measures of Standard Error					
	a	SE (a)	b	SE (b)	c	SE (c)
90282	0.39678	0.00797	0.44031	0.01845	0.00000	0.00000
62084	0.45719	0.00823	-0.34278	0.01828	0.00000	0.00000
146785	0.83200	0.02510	0.96005	0.02155	0.18140	0.00797
48115	0.66250	0.01059	-1.29000	0.02151	0.00000	0.00000
146766	0.84086	0.03095	1.61629	0.02234	0.15296	0.00586

Item Number	Parameters and Measures of Standard Error					
	a	SE (a)	b	SE (b)	c	SE (c)
91920	1.29765	0.02891	0.12872	0.01764	0.26665	0.00833
60693	0.47760	0.02127	0.79821	0.06058	0.13422	0.01852

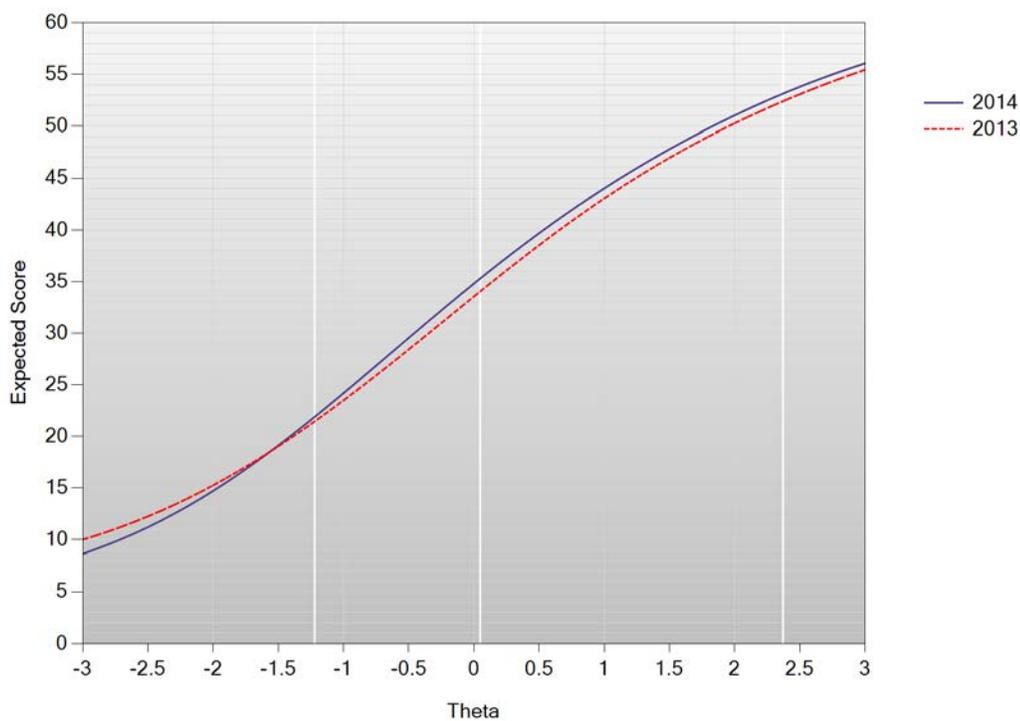
**Table I-6. 2013–14 NECAP Science: IRT Parameters for Polytomous Items—
Grade 11**

Item Number	Parameters and Measures of Standard Error													
	a	SE (a)	b	SE (b)	D0	SE (D0)	D1	SE (D1)	D2	SE (D2)	D3	SE (D3)	D4	SE (D4)
258656	0.85094	0.00714	0.36443	0.00815	0.59059	0.00967	-0.59059	0.00982	0	0				
258657	0.82285	0.00629	0.47347	0.00838	0.91374	0.01025	-0.91374	0.01102	0	0				
258659	0.71324	0.00672	1.18120	0.01098	0.58801	0.01089	-0.58801	0.01336	0	0				
258661	0.91672	0.00681	0.36046	0.00773	1.02556	0.00997	-1.02556	0.01030	0	0				
258662	1.01073	0.00806	0.41820	0.00701	0.65604	0.00858	-0.65604	0.00889	0	0				
258669	0.87013	0.00681	0.51973	0.00799	0.79638	0.00960	-0.79638	0.01042	0	0				
258655	0.99434	0.00650	0.83130	0.00681	1.40702	0.00924	-0.06668	0.00863	-1.34035	0.01351	0	0		
263524	1.06022	0.00727	0.86970	0.00651	1.06020	0.00820	0.12166	0.00803	-1.18187	0.01243	0	0		
241745	0.87729	0.00587	1.58424	0.00781	1.66041	0.00928	0.79919	0.00929	-0.53007	0.01365	-1.92954	0.03037	0	0
242632	0.97538	0.00623	0.85421	0.00687	1.33081	0.00912	0.95924	0.00855	-0.42096	0.00949	-1.86909	0.01866	0	0
242637	0.93027	0.00566	0.65866	0.00687	1.57103	0.01043	0.41861	0.00861	-0.54836	0.00957	-1.44128	0.01329	0	0

APPENDIX J—TEST CHARACTERISTIC CURVE AND TEST INFORMATION FUNCTION CHARTS

Figure J-1. 2013–14 NECAP Science: Grade 4 Charts
Top: Test Characteristic Curve Bottom: Test Information Function

Test Characteristic Curve: Science Grade 4



Test Information Function: Science Grade 4

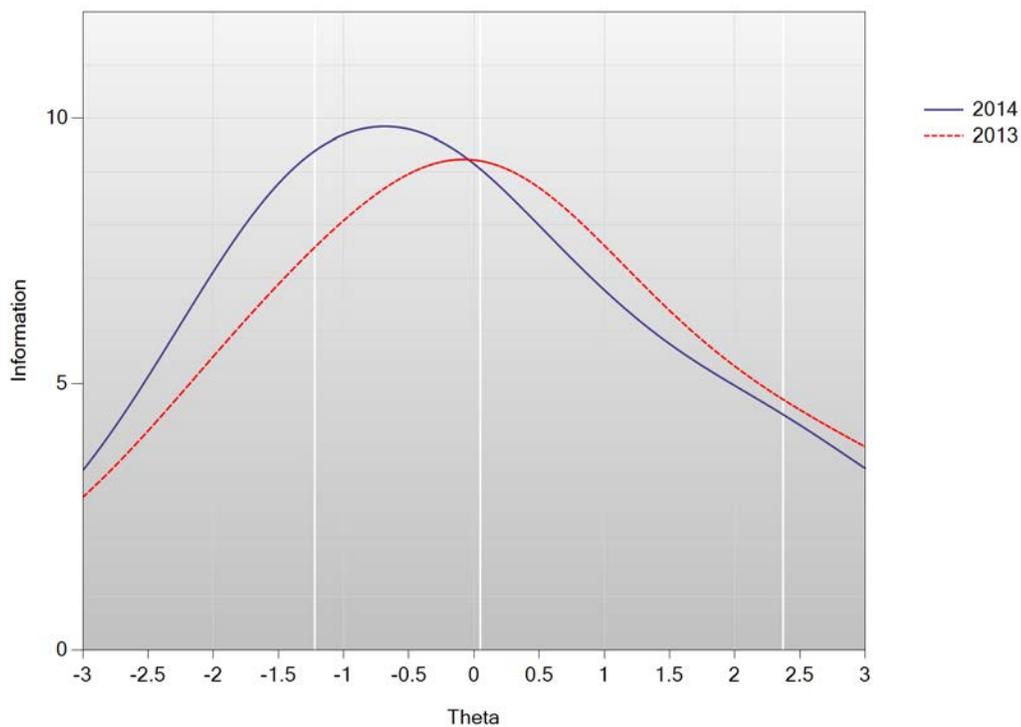
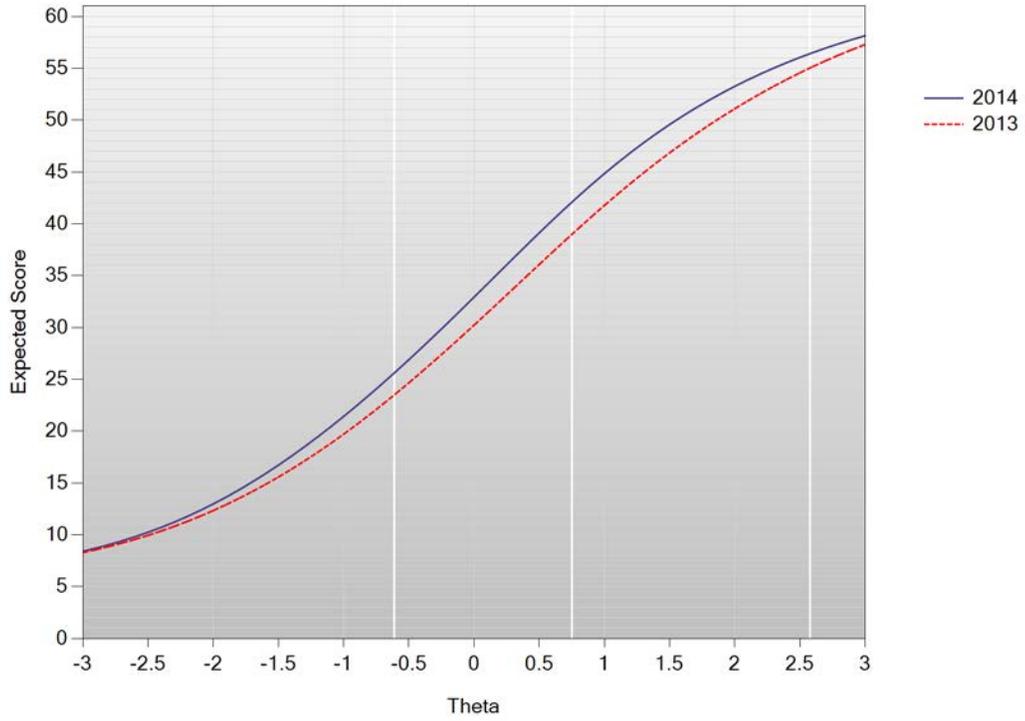


Figure J-2. 2013–14 NECAP Science: Grade 8 Charts
Top: Test Characteristic Curve Bottom: Test Information Function

Test Characteristic Curve: Science Grade 8



Test Information Function: Science Grade 8

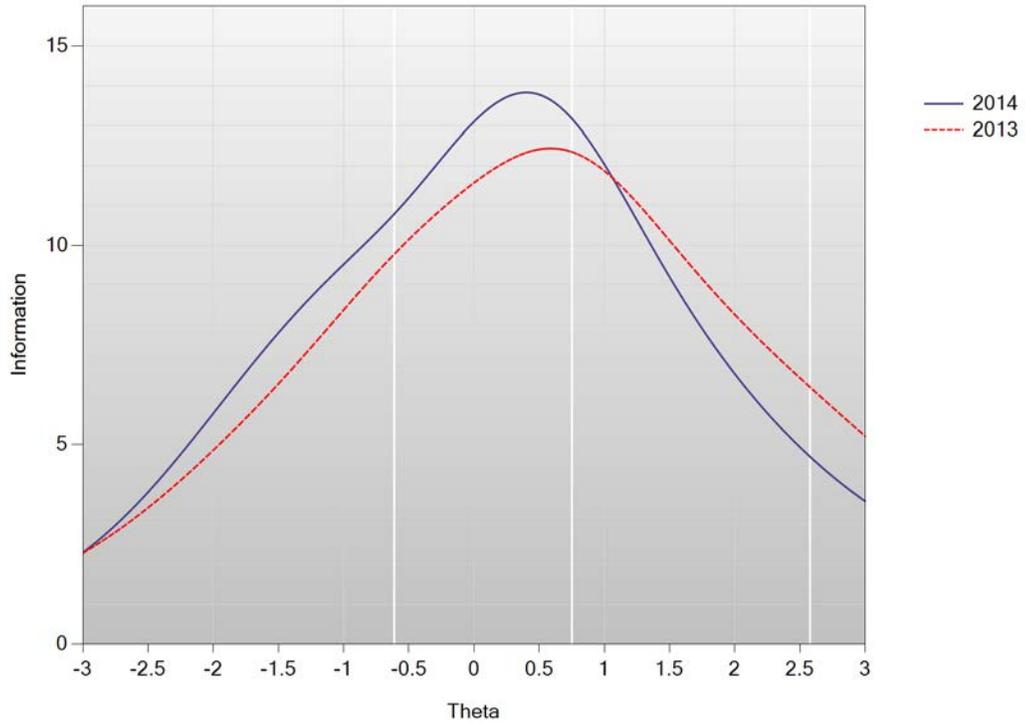
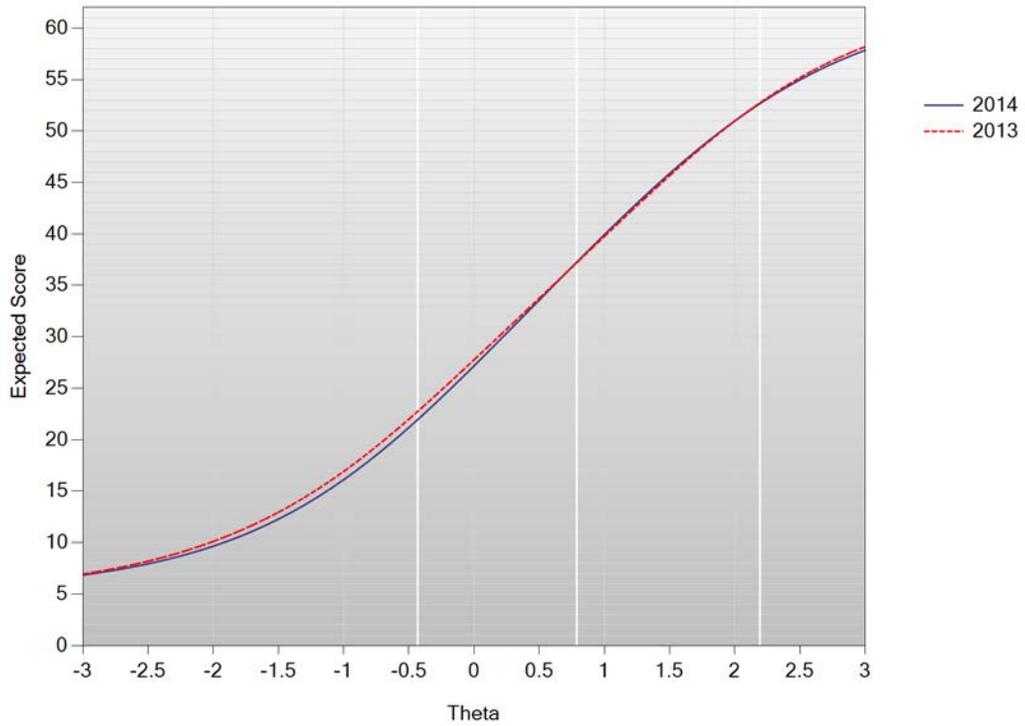
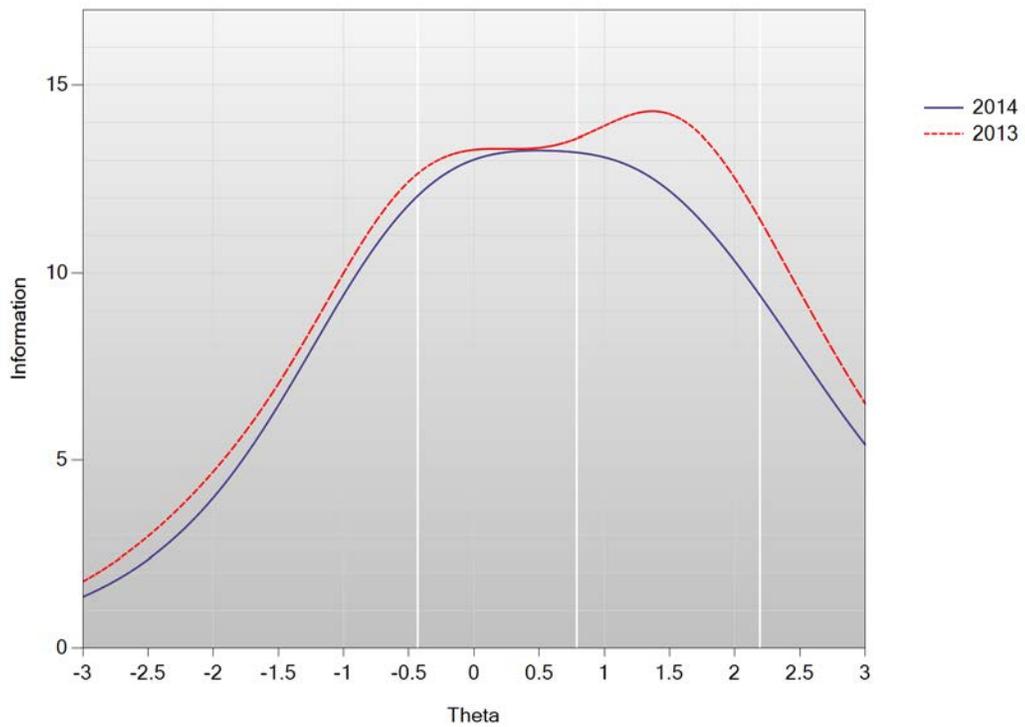


Figure J-3. 2013–14 NECAP Science: Grade 11 Charts
Top: Test Characteristic Curve Bottom: Test Information Function

Test Characteristic Curve: Science Grade 11



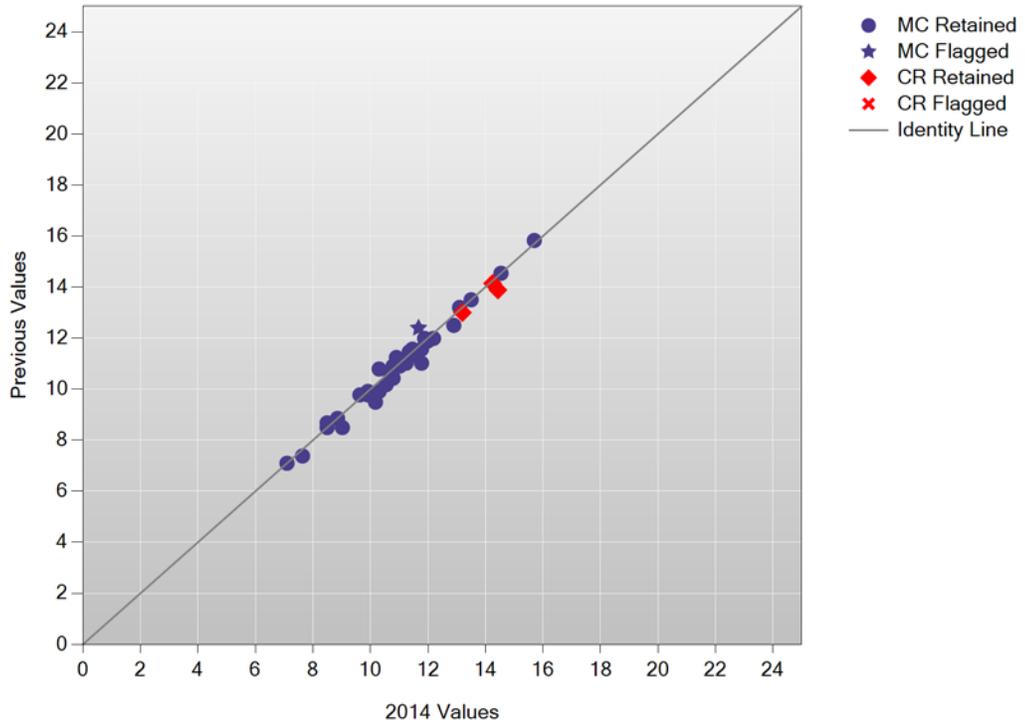
Test Information Function: Science Grade 11



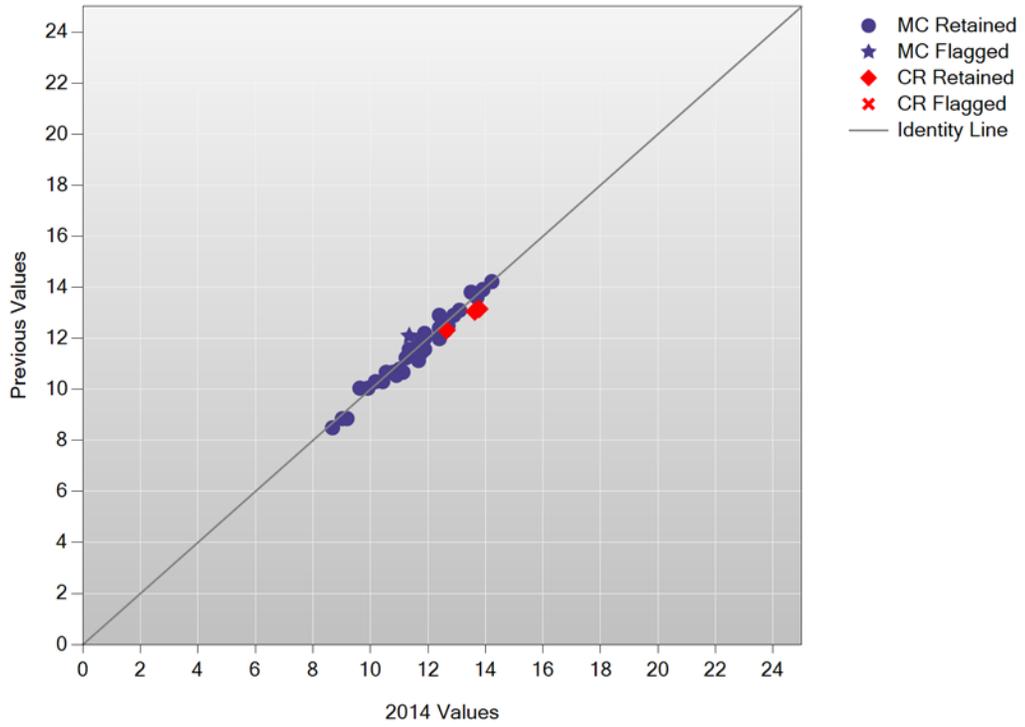
APPENDIX K—DELTA ANALYSES AND RESCORE ANALYSES

Figure K-1. 2013–14 NECAP Science: Delta Plots
Top: Grade 4 Bottom: Grade 8

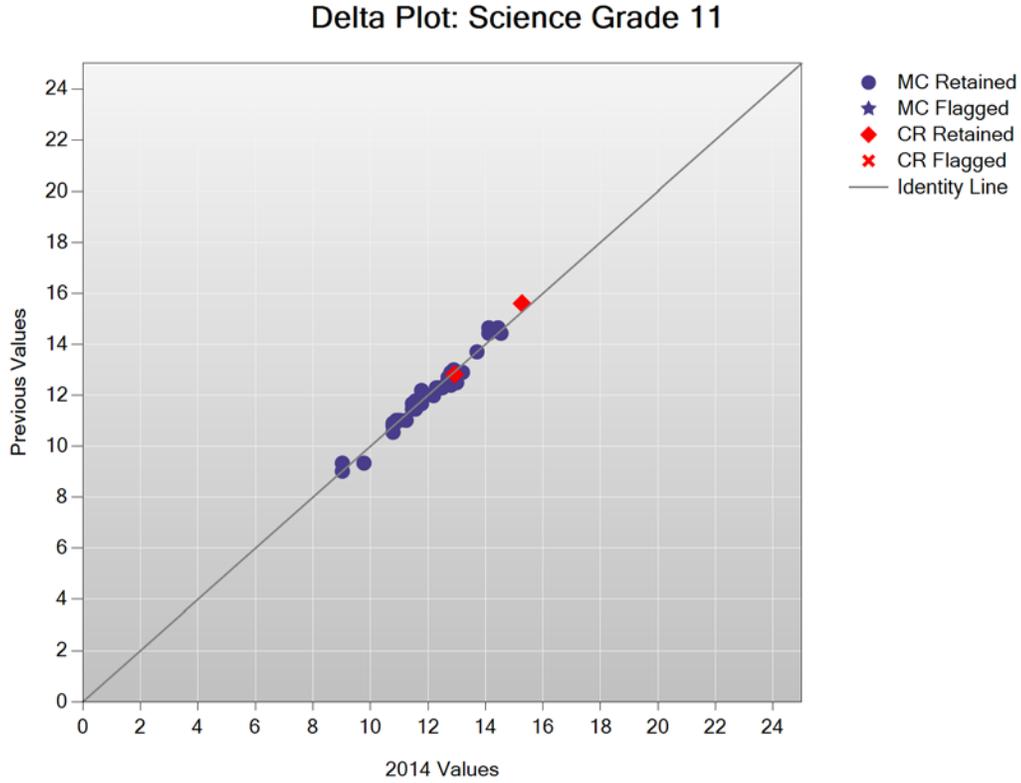
Delta Plot: Science Grade 4



Delta Plot: Science Grade 8



**Figure K-2. 2013–14 NECAP Science: Delta Plot—
Grade 11**



**Table K-1. 2013–14 NECAP Science: Delta Analyses Results—
Grade 4**

Item Number	Difficulty		Delta		Discard	Standard Deviation
	Old	New	Old	New		
141192	0.55000	0.51000	12.49735	12.89972	False	0.55979
142220	0.24000	0.25000	15.82521	15.69796	False	-0.24323
142224	0.69000	0.67000	11.01660	11.24035	False	-0.45411
144866	0.69000	0.62000	11.01660	11.77808	False	2.36870
144907	0.50000	0.48000	13.00000	13.20061	False	-0.47345
174214	0.79000	0.78000	9.77432	9.91123	False	-0.97392
46310	0.87000	0.84000	8.49444	9.02217	False	1.01181
46402	0.71000	0.71000	10.78646	10.78646	False	-0.65184
46416	0.56000	0.63000	12.39612	11.67259	True	3.06350
47361	0.79000	0.77000	9.77432	10.04461	False	-0.27371
47386	0.86000	0.87000	8.67872	8.49444	False	0.42408
47418	0.61000	0.60000	11.88272	11.98661	False	-1.03874
47448	0.71000	0.75000	10.78646	10.30204	False	1.89113
47493	0.70000	0.71000	10.90240	10.78646	False	-0.04920
47531	0.38750	0.37750	14.14336	14.24821	False	-0.91731
47551	0.74000	0.71000	10.42662	10.78646	False	0.22994
47624	0.60000	0.58000	11.98661	12.19243	False	-0.49832
47724	0.93000	0.93000	7.09684	7.09684	False	-0.46190
47760	0.64000	0.62000	11.56616	11.77808	False	-0.48796

continued

<i>Item Number</i>	<i>Difficulty</i>		<i>Delta</i>		<i>Discard</i>	<i>Standard Deviation</i>
	<i>Old</i>	<i>New</i>	<i>Old</i>	<i>New</i>		
47992	0.81000	0.76000	9.48841	10.17479	False	1.89577
48019	0.92000	0.91000	7.37971	7.63698	False	-0.46539
49850	0.76000	0.73000	10.17479	10.54875	False	0.29107
49866	0.85000	0.85000	8.85427	8.85427	False	-0.55237
49894	0.66000	0.65000	11.35015	11.45872	False	-1.04157
49897	0.64000	0.65000	11.56616	11.45872	False	-0.12794
59267	0.68000	0.67000	11.12920	11.24035	False	-1.03944
59434	0.78000	0.75000	9.91123	10.30204	False	0.36599
59640	0.41250	0.36000	13.88447	14.43384	False	1.40283
60352	0.87000	0.87000	8.49444	8.49444	False	-0.53385
61931	0.78000	0.78000	9.91123	9.91123	False	-0.60678
86940	0.45000	0.45000	13.50265	13.50265	False	-0.79167
87135	0.70000	0.69000	10.90240	11.01660	False	-1.03507
87181	0.67000	0.70000	11.24035	10.90240	False	1.09886
88061	0.65000	0.66000	11.45872	11.35015	False	-0.11651
88090	0.79000	0.80000	9.77432	9.63352	False	0.13939
91504	0.60000	0.61000	11.98661	11.88272	False	-0.16827
91509	0.72000	0.72000	10.66863	10.66863	False	-0.64578
91513	0.35000	0.35000	14.54128	14.54128	False	-0.84514
91971	0.48000	0.49000	13.20061	13.10028	False	-0.24940

**Table K-2. 2013–14 NECAP Science: Rescore Analyses Results—
Grade 4**

<i>Item Number</i>	<i>Number of Score Categories</i>	<i>Average Score</i>		<i>Standard Deviation</i>		<i>Effect Size</i>	<i>Discard</i>
		<i>Old</i>	<i>New</i>	<i>Old</i>	<i>New</i>		
144907	4	2.22660	2.06897	1.27345	1.28786	-0.12379	False
59640	4	1.69268	1.51707	1.00401	0.98814	-0.17491	False
47531	4	1.66832	1.76733	1.30208	1.04635	0.07604	False

**Table K-3. 2013–14 NECAP Science: Delta Analyses Results—
Grade 8**

<i>Item Number</i>	<i>Difficulty</i>		<i>Delta</i>		<i>Discard</i>	<i>Standard Deviation</i>
	<i>Old</i>	<i>New</i>	<i>Old</i>	<i>New</i>		
144589	0.59000	0.66000	12.08982	11.35015	True	3.15787
174702	0.49500	0.43750	13.05013	13.62924	False	1.46165
46016	0.72000	0.71000	10.66863	10.78646	False	-1.05786
46026	0.73000	0.70000	10.54875	10.90240	False	0.24407
46039	0.48500	0.42250	13.15043	13.78201	False	1.74999
46041	0.51000	0.56000	12.89972	12.39612	False	1.86421
46070	0.42000	0.45000	13.80757	13.50265	False	0.77784
46072	0.71000	0.69000	10.78646	11.01660	False	-0.43963
46074	0.60000	0.62000	11.98661	11.77808	False	0.22730
46082	0.58000	0.61000	12.19243	11.88272	False	0.78742
46089	0.72000	0.73000	10.66863	10.54875	False	-0.27531
46109	0.41000	0.41000	13.91018	13.91018	False	-0.90293

continued

Item Number	Difficulty		Delta		Discard	Standard Deviation
	Old	New	Old	New		
48184	0.75000	0.76000	10.30204	10.17479	False	-0.23849
48228	0.67000	0.67000	11.24035	11.24035	False	-0.93062
48245	0.77000	0.78000	10.04461	9.91123	False	-0.20733
48267	0.65000	0.62000	11.45872	11.77808	False	0.04550
48268	0.60000	0.61000	11.98661	11.88272	False	-0.34989
48297	0.64000	0.61000	11.56616	11.88272	False	0.02894
48355	0.64000	0.66000	11.56616	11.35015	False	0.26421
48421	0.56000	0.54000	12.39612	12.59827	False	-0.61073
48445	0.55000	0.53000	12.49735	12.69892	False	-0.61496
48449	0.56750	0.53500	12.31995	12.64862	False	0.08794
48456	0.44000	0.43000	13.60388	13.70550	False	-1.17769
48472	0.87000	0.86000	8.49444	8.67872	False	-0.66875
48563	0.77000	0.80000	10.04461	9.63352	False	1.32440
49990	0.56000	0.56000	12.39612	12.39612	False	-0.91863
50012	0.72000	0.68000	10.66863	11.12920	False	0.83255
50016	0.85000	0.84000	8.85427	9.02217	False	-0.76286
50026	0.64000	0.64000	11.56616	11.56616	False	-0.92724
50120	0.49000	0.49000	13.10028	13.10028	False	-0.91133
50511	0.85000	0.83000	8.85427	9.18334	False	0.12608
58352	0.75000	0.74000	10.30204	10.42662	False	-1.01683
60084	0.38000	0.38000	14.22192	14.22192	False	-0.89970
76623	0.51000	0.51000	12.89972	12.89972	False	-0.91341
76626	0.68000	0.63000	11.12920	11.67259	False	1.28452
90304	0.57000	0.54000	12.29450	12.59827	False	-0.04919
91663	0.60000	0.56000	11.98661	12.39612	False	0.53726
91775	0.65000	0.65000	11.45872	11.45872	False	-0.92835

**Table K-4. 2013–14 NECAP Science: Rescore Analyses Results—
Grade 8**

Item Number	Number of Score Categories	Average Score		Standard Deviation		Effect Size	Discard
		Old	New	Old	New		
174702	4	1.82843	1.73039	0.96489	0.91018	-0.10161	False
46039	4	1.91176	1.67647	1.15416	1.01866	-0.20387	False
48449	4	2.34171	2.04020	1.10728	0.89239	-0.27229	False

**Table K-5. 2013–14 NECAP Science: Delta Analyses Results—
Grade 11**

Item Number	Difficulty		Delta		Discard	Standard Deviation
	Old	New	Old	New		
146889	0.25750	0.28500	15.60429	15.27221	False	0.17854
146893	0.53000	0.53000	12.69892	12.69892	False	-1.29965
146924	0.58000	0.62000	12.19243	11.77808	False	1.77766
147033	0.57000	0.56000	12.29450	12.39612	False	-0.64780
169678	0.34000	0.36000	14.64985	14.43384	False	-0.42484
176916	0.84000	0.84000	9.02217	9.02217	False	-0.44431

continued

<i>Item Number</i>	<i>Difficulty</i>		<i>Delta</i>		<i>Discard</i>	<i>Standard Deviation</i>
	<i>Old</i>	<i>New</i>	<i>Old</i>	<i>New</i>		
177100	0.51750	0.50750	12.82448	12.92480	False	-0.50579
242620	0.43000	0.43000	13.70550	13.70550	False	-1.01123
46019	0.50000	0.51000	13.00000	12.89972	False	-0.82651
46094	0.55000	0.50000	12.49735	13.00000	False	2.44001
46096	0.60000	0.58000	11.98661	12.19243	False	0.05115
46121	0.82000	0.84000	9.33854	9.02217	False	1.85518
46148	0.65000	0.65000	11.45872	11.45872	False	-1.14246
46166	0.57000	0.57000	12.29450	12.29450	False	-1.38194
46179	0.69000	0.70000	11.01660	10.90240	False	-0.15301
47884	0.71000	0.71000	10.78646	10.78646	False	-0.94983
47917	0.54000	0.53000	12.59827	12.69892	False	-0.56805
48005	0.57000	0.55000	12.29450	12.49735	False	0.11699
48156	0.82000	0.79000	9.33854	9.77432	False	1.02972
48216	0.52000	0.50000	12.79939	13.00000	False	0.24475
48357	0.36000	0.39000	14.43384	14.11728	False	0.39664
48372	0.36000	0.35000	14.43384	14.54128	False	0.00920
48543	0.51000	0.52000	12.89972	12.79939	False	-0.79731
48908	0.63000	0.65000	11.67259	11.45872	False	0.41201
49902	0.51000	0.48000	12.89972	13.20061	False	1.03107
49903	0.65000	0.64000	11.45872	11.56616	False	-0.84326
49922	0.63000	0.62000	11.67259	11.77808	False	-0.79676
49925	0.62000	0.64000	11.77808	11.56616	False	0.36700
49930	0.35000	0.37000	14.54128	14.32741	False	-0.40996
49935	0.69000	0.69000	11.01660	11.01660	False	-1.01578
61142	0.73000	0.71000	10.54875	10.78646	False	-0.11985
62083	0.70000	0.71000	10.90240	10.78646	False	-0.10717
89407	0.69000	0.67000	11.01660	11.24035	False	-0.09130
89632	0.34000	0.39000	14.64985	14.11728	False	1.96672
91901	0.56000	0.52000	12.39612	12.79939	False	1.66018

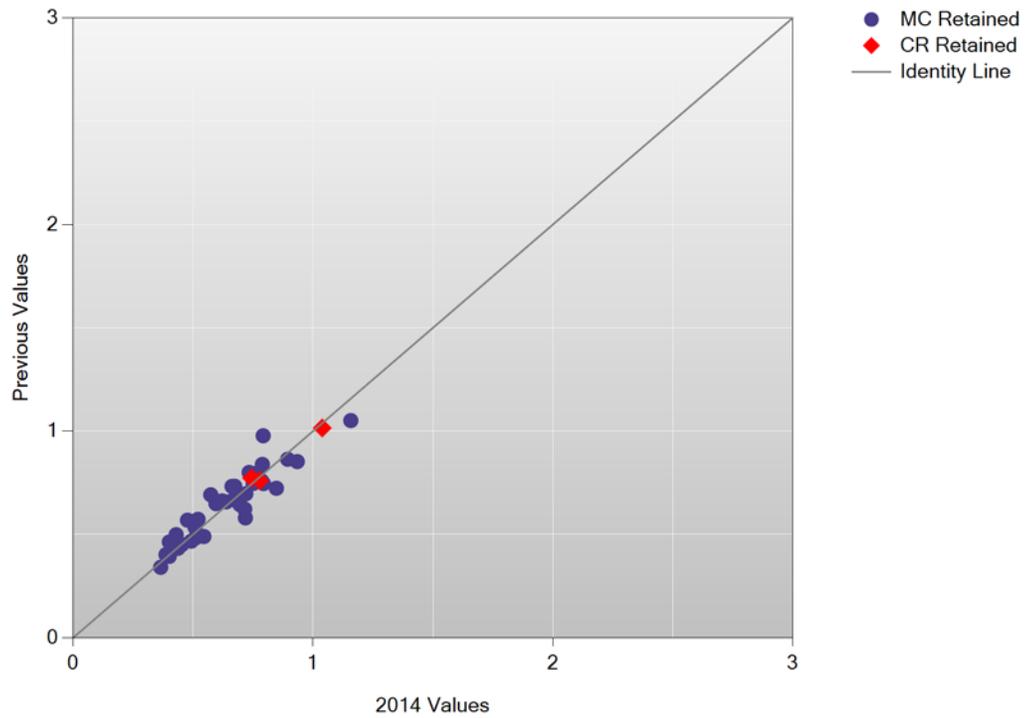
**Table K-6. 2013–14 NECAP Science: Rescore Analyses Results—
Grade 11**

<i>Item Number</i>	<i>Number of Score Categories</i>	<i>Average Score</i>		<i>Standard Deviation</i>		<i>Effect Size</i>	<i>Discard</i>
		<i>Old</i>	<i>New</i>	<i>Old</i>	<i>New</i>		
146889	4	1.02970	1.10396	1.27319	1.28291	0.05832	False
146948	4	1.27941	1.41176	0.93922	1.07229	0.14092	False
177100	4	2.20896	2.33831	0.83434	0.88033	0.15504	False

APPENDIX L—*a*-PLOTS AND *b*-PLOTS

Figure L-1. 2013–14 NECAP Science: Grade 4 Plots
Top: *a*-Plot Bottom: *b*-Plot

A/A Plot: Science Grade 4



B/B Plot: Science Grade 4

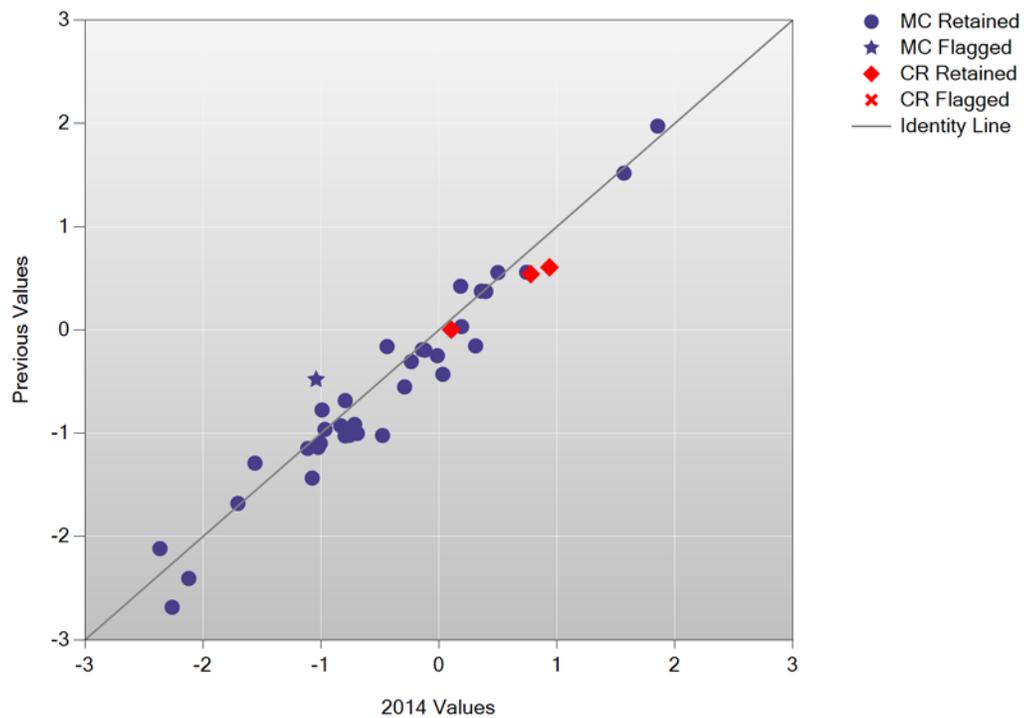
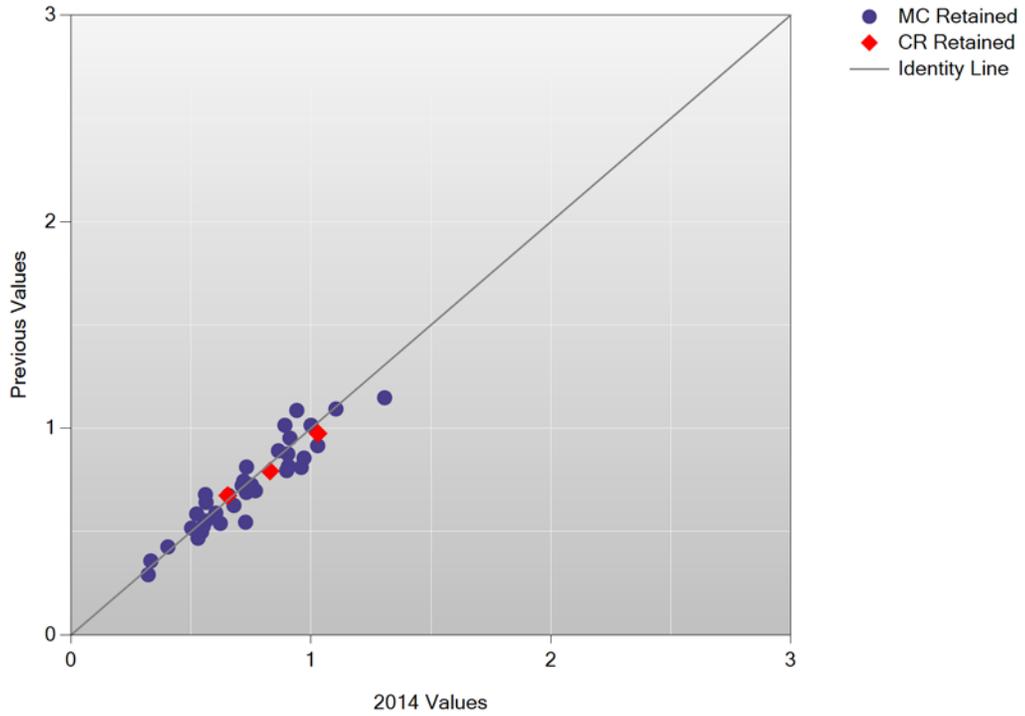


Figure L-2. 2013–14 NECAP Science: Grade 8 Plots

Top: *a*-Plot Bottom: *b*-Plot

A/A Plot: Science Grade 8



B/B Plot: Science Grade 8

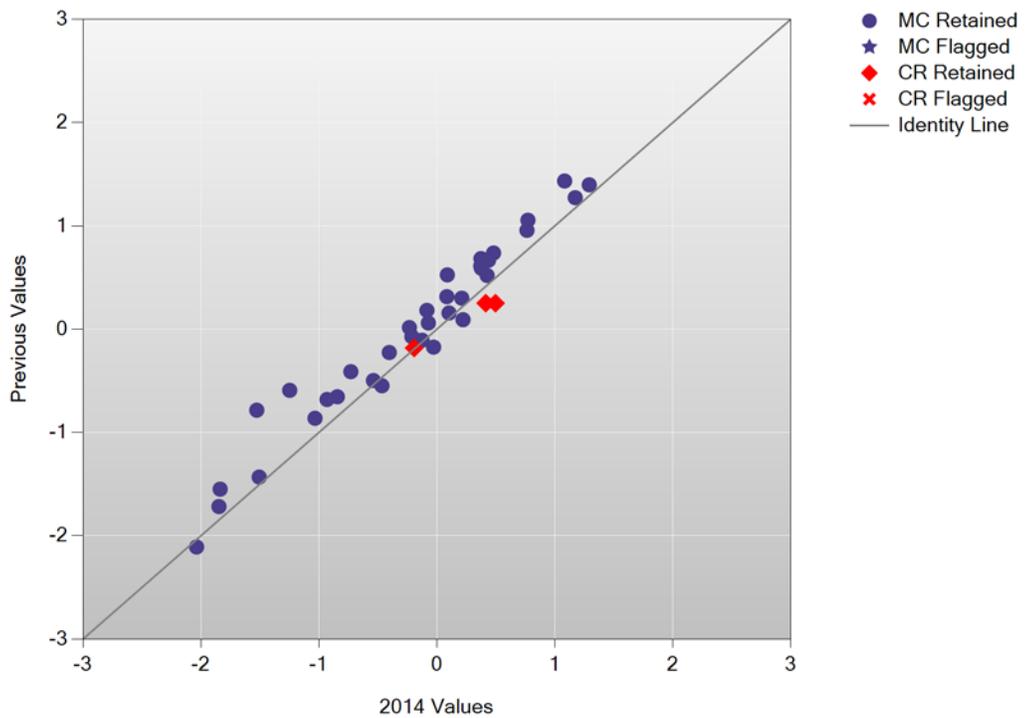
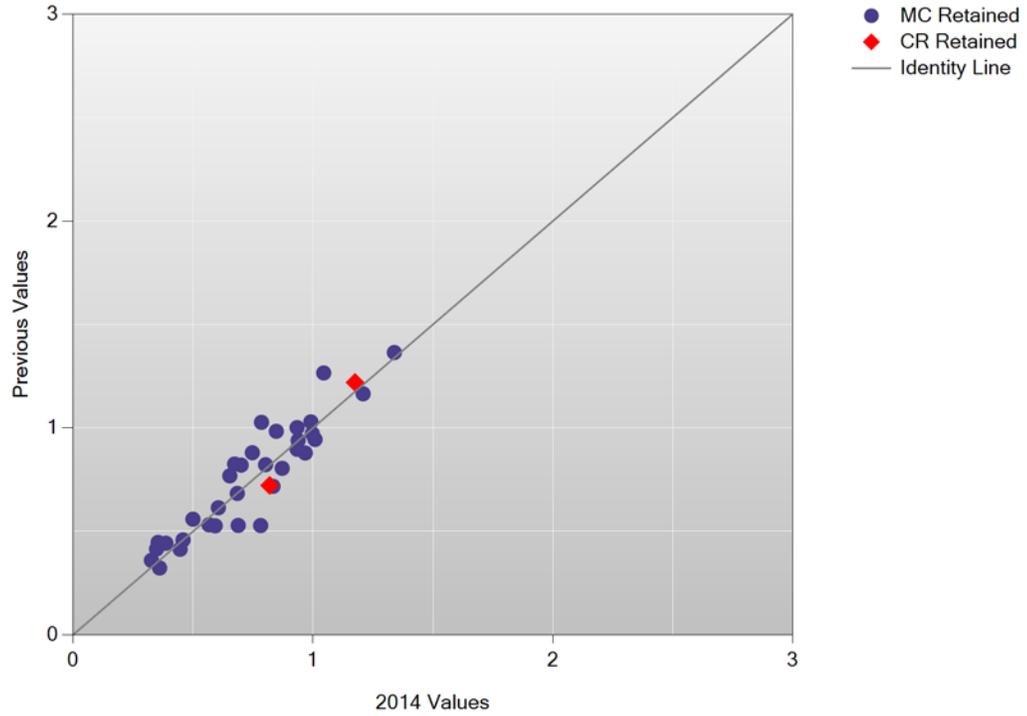


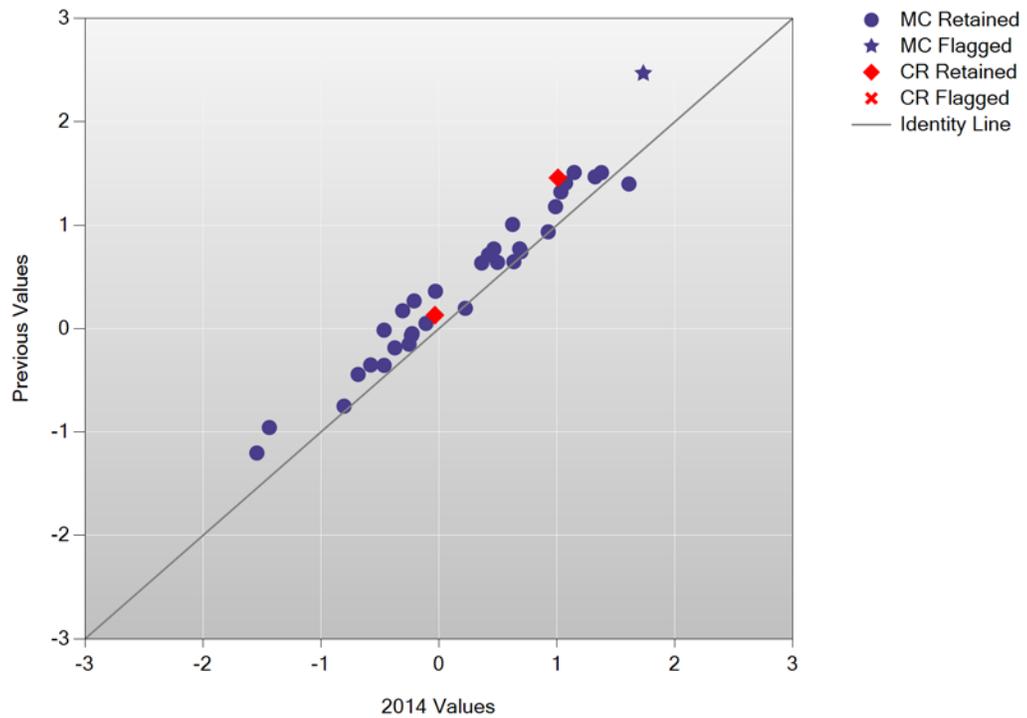
Figure L-3. 2013–14 NECAP Science: Grade 11 Plots

Top: *a*-Plot Bottom: *b*-Plot

A/A Plot: Science Grade 11



B/B Plot: Science Grade 11



APPENDIX M—ACHIEVEMENT LEVEL DISTRIBUTIONS

**Table M-1. 2013–14 NECAP Science: Achievement Level Distributions
by Grade**

<i>Grade</i>	<i>Performance Level</i>	<i>Percent in Level</i>		
		<i>2013–14</i>	<i>2012–13</i>	<i>2011–12</i>
4	4	0.65	0.79	1.29
	3	43.17	45.70	48.90
	2	42.06	39.72	36.47
	1	14.11	13.79	13.33
8	4	0.59	0.60	1.17
	3	23.67	30.27	28.61
	2	52.27	47.45	47.38
	1	23.46	21.68	22.83
11	4	1.70	1.29	2.10
	3	27.99	28.90	30.44
	2	45.97	45.48	44.37
	1	24.34	24.33	23.08

APPENDIX N—RAW TO SCALED SCORE LOOKUP TABLES

**Table N-1. 2013–14 NECAP Science: Raw Score to Scaled Score Correspondence
Grade 4**

Raw Score	2014			2013		
	Scaled Score	Standard Error	Performance Level	Scaled Score	Standard Error	Performance Level
0	400	10.0	1	400	10.0	1
1	400	10.0	1	400	10.0	1
2	400	10.0	1	400	10.0	1
3	400	10.0	1	400	10.0	1
4	400	10.0	1	400	10.0	1
5	400	9.9	1	400	10.0	1
6	401	8.2	1	400	10.0	1
7	405	6.7	1	400	9.5	1
8	408	5.8	1	403	7.8	1
9	411	5.2	1	407	6.6	1
10	413	4.7	1	410	5.8	1
11	414	4.4	1	412	5.3	1
12	416	4.2	1	414	4.9	1
13	418	4.0	1	416	4.7	1
14	419	3.8	1	418	4.4	1
15	420	3.7	1	419	4.2	1
16	421	3.6	1	421	4.1	1
17	422	3.5	1	422	4.0	1
18	424	3.4	1	423	3.9	1
19	425	3.3	1	425	3.8	1
20	426	3.3	1	426	3.7	1
21	426	3.3	1	426	3.6	1
22	428	3.2	2	428	3.6	2
23	429	3.2	2	429	3.5	2
24	429	3.2	2	430	3.5	2
25	430	3.2	2	431	3.4	2
26	431	3.2	2	432	3.4	2
27	432	3.1	2	433	3.3	2
28	433	3.1	2	434	3.3	2
29	434	3.2	2	435	3.3	2
30	435	3.2	2	436	3.3	2
31	436	3.2	2	437	3.3	2
32	437	3.2	2	438	3.3	2
33	438	3.2	2	439	3.3	2
34	439	3.2	2	439	3.3	2
35	439	3.3	2	441	3.3	3
36	441	3.3	3	442	3.3	3
37	442	3.4	3	443	3.3	3
38	443	3.4	3	444	3.3	3
39	444	3.5	3	445	3.4	3
40	445	3.5	3	446	3.4	3
41	446	3.6	3	447	3.5	3
42	447	3.6	3	448	3.5	3
43	448	3.7	3	449	3.6	3
44	449	3.8	3	451	3.7	3

continued

Raw Score	2014			2013		
	Scaled Score	Standard Error	Performance Level	Scaled Score	Standard Error	Performance Level
45	451	3.9	3	452	3.7	3
46	452	4.0	3	453	3.8	3
47	453	4.0	3	454	3.9	3
48	455	4.1	3	456	4.0	3
49	456	4.2	3	457	4.1	3
50	458	4.3	3	459	4.2	3
51	459	4.4	3	460	4.4	3
52	461	4.5	3	462	4.5	3
53	462	4.7	3	464	4.6	4
54	465	4.8	4	466	4.8	4
55	467	5.1	4	468	5.0	4
56	469	5.3	4	470	5.2	4
57	472	5.7	4	473	5.4	4
58	475	6.2	4	476	5.8	4
59	478	6.8	4	479	6.2	4
60	479	7.0	4	479	6.2	4
61	479	7.0	4	479	6.2	4
62	479	7.0	4	479	6.2	4
63	480	7.0	4	480	6.2	4

**Table N-2. 2013–14 NECAP Science: Raw Score to Scaled Score Correspondence
Grade 8**

Raw Score	2014			2013		
	Scaled Score	Standard Error	Performance Level	Scaled Score	Standard Error	Performance Level
0	800	10.0	1	800	10.0	1
1	800	10.0	1	800	10.0	1
2	800	10.0	1	800	10.0	1
3	800	10.0	1	800	10.0	1
4	800	10.0	1	800	10.0	1
5	800	10.0	1	800	10.0	1
6	800	10.0	1	800	9.6	1
7	803	7.9	1	803	7.2	1
8	807	6.0	1	808	5.9	1
9	810	5.1	1	810	5.1	1
10	812	4.4	1	813	4.5	1
11	814	4.0	1	815	4.2	1
12	815	3.7	1	816	3.9	1
13	817	3.5	1	818	3.7	1
14	818	3.3	1	819	3.5	1
15	819	3.2	1	820	3.4	1
16	820	3.1	1	822	3.2	1
17	821	3.0	1	823	3.1	1
18	822	2.9	1	824	3.0	1
19	823	2.9	1	825	3.0	1
20	824	2.8	1	826	2.9	1
21	825	2.7	1	826	2.8	1

continued

Raw Score	2014			2013		
	Scaled Score	Standard Error	Performance Level	Scaled Score	Standard Error	Performance Level
22	826	2.7	1	827	2.8	1
23	827	2.7	1	828	2.7	1
24	827	2.6	1	829	2.7	2
25	828	2.6	1	830	2.6	2
26	829	2.5	2	831	2.6	2
27	830	2.5	2	831	2.6	2
28	830	2.5	2	832	2.5	2
29	831	2.4	2	833	2.5	2
30	832	2.4	2	834	2.5	2
31	832	2.4	2	834	2.5	2
32	833	2.3	2	835	2.4	2
33	834	2.3	2	836	2.4	2
34	834	2.3	2	836	2.4	2
35	835	2.3	2	837	2.4	2
36	836	2.3	2	838	2.4	2
37	836	2.3	2	839	2.4	2
38	837	2.3	2	839	2.4	2
39	838	2.3	2	840	2.4	3
40	838	2.3	2	841	2.4	3
41	839	2.3	2	842	2.4	3
42	839	2.3	2	842	2.4	3
43	841	2.3	3	843	2.5	3
44	841	2.4	3	844	2.5	3
45	842	2.4	3	845	2.6	3
46	843	2.5	3	846	2.6	3
47	844	2.6	3	846	2.7	3
48	845	2.6	3	847	2.7	3
49	846	2.7	3	848	2.8	3
50	847	2.8	3	849	2.8	3
51	848	2.9	3	850	2.9	3
52	849	3.1	3	852	3.0	3
53	850	3.2	3	853	3.1	3
54	852	3.4	3	854	3.2	3
55	853	3.5	3	854	3.3	3
56	854	3.8	3	857	3.5	4
57	856	4.1	4	858	3.6	4
58	859	4.4	4	860	3.9	4
59	861	4.8	4	862	4.2	4
60	864	5.4	4	865	4.7	4
61	867	6.2	4	867	5.3	4
62	867	6.2	4	867	5.3	4
63	880	6.2	4	880	5.3	4

**Table N-3. 2013–14 NECAP Science: Raw Score to Scaled Score Correspondence
Grade 11**

Raw Score	2014			2013		
	Scaled Score	Standard Error	Performance Level	Scaled Score	Standard Error	Performance Level
0	1100	10.0	1	1100	10.0	1
1	1100	10.0	1	1100	10.0	1
2	1100	10.0	1	1100	10.0	1
3	1100	10.0	1	1100	10.0	1
4	1100	10.0	1	1100	10.0	1
5	1100	10.0	1	1100	10.0	1
6	1103	10.0	1	1103	8.9	1
7	1109	6.8	1	1109	6.2	1
8	1113	5.4	1	1112	5.0	1
9	1115	4.6	1	1114	4.3	1
10	1117	4.0	1	1117	3.9	1
11	1119	3.6	1	1118	3.6	1
12	1121	3.4	1	1120	3.3	1
13	1122	3.1	1	1121	3.1	1
14	1123	3.0	1	1122	3.0	1
15	1124	2.8	1	1123	2.8	1
16	1125	2.7	1	1124	2.7	1
17	1126	2.6	1	1125	2.6	1
18	1127	2.6	1	1126	2.6	1
19	1128	2.5	1	1127	2.5	1
20	1128	2.5	1	1128	2.4	1
21	1129	2.4	1	1128	2.4	1
22	1130	2.4	2	1129	2.4	1
23	1131	2.4	2	1130	2.3	2
24	1131	2.4	2	1131	2.3	2
25	1132	2.3	2	1131	2.3	2
26	1133	2.3	2	1132	2.3	2
27	1133	2.3	2	1133	2.3	2
28	1134	2.3	2	1134	2.3	2
29	1135	2.3	2	1134	2.3	2
30	1135	2.3	2	1135	2.3	2
31	1136	2.3	2	1136	2.3	2
32	1137	2.3	2	1136	2.3	2
33	1137	2.3	2	1137	2.3	2
34	1138	2.3	2	1138	2.3	2
35	1139	2.3	2	1138	2.3	2
36	1139	2.3	2	1139	2.3	2
37	1139	2.3	2	1139	2.3	2
38	1140	2.3	3	1141	2.3	3
39	1141	2.3	3	1141	2.2	3
40	1142	2.3	3	1142	2.2	3
41	1142	2.3	3	1143	2.2	3
42	1143	2.3	3	1143	2.2	3
43	1144	2.3	3	1144	2.2	3
44	1145	2.4	3	1145	2.2	3
45	1145	2.4	3	1145	2.2	3

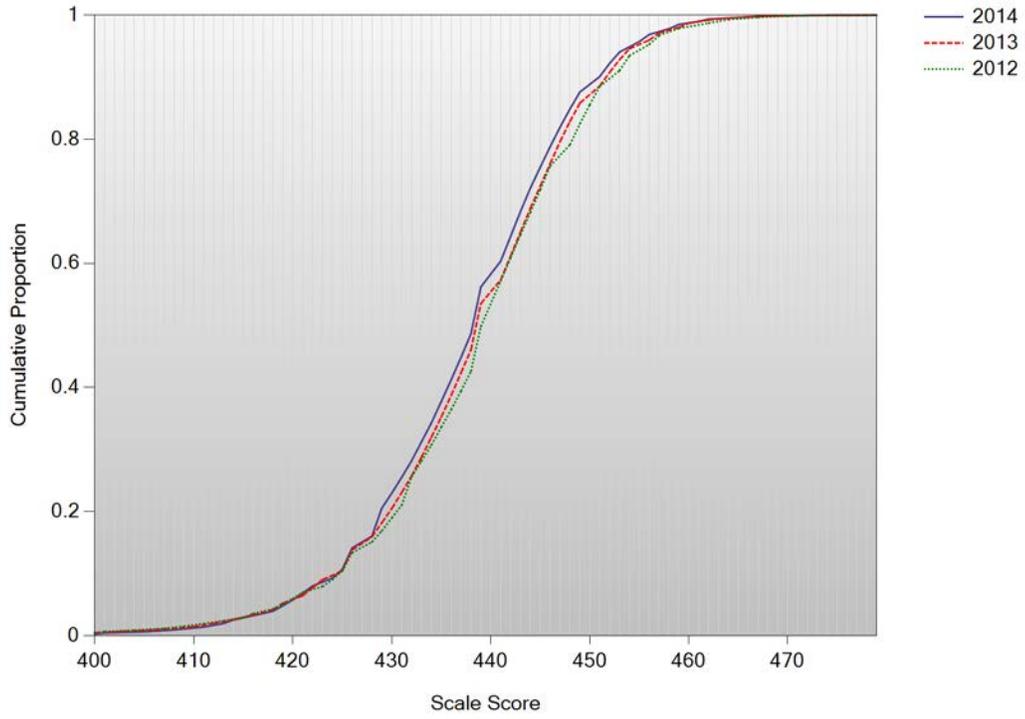
continued

<i>Raw Score</i>	<i>2014</i>			<i>2013</i>		
	<i>Scaled Score</i>	<i>Standard Error</i>	<i>Performance Level</i>	<i>Scaled Score</i>	<i>Standard Error</i>	<i>Performance Level</i>
46	1146	2.4	3	1146	2.2	3
47	1147	2.4	3	1147	2.2	3
48	1148	2.5	3	1148	2.3	3
49	1148	2.5	3	1148	2.3	3
50	1149	2.5	3	1149	2.3	3
51	1150	2.6	3	1150	2.4	3
52	1151	2.7	3	1151	2.4	3
53	1152	2.8	4	1152	2.5	4
54	1153	2.9	4	1153	2.6	4
55	1154	3.0	4	1154	2.7	4
56	1156	3.1	4	1155	2.8	4
57	1157	3.4	4	1157	3.0	4
58	1159	3.6	4	1158	3.2	4
59	1161	4.0	4	1160	3.5	4
60	1163	4.7	4	1162	4.0	4
61	1167	5.7	4	1165	5.0	4
62	1167	5.7	4	1167	5.6	4
63	1180	5.7	4	1180	5.6	4

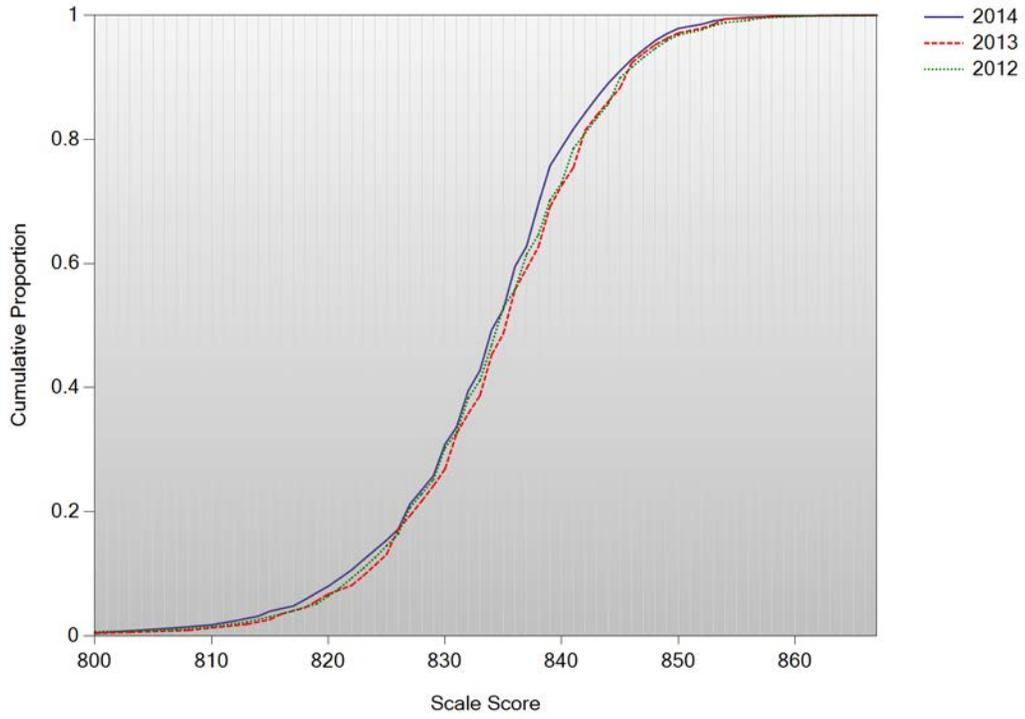
APPENDIX O—SCALED SCORE DISTRIBUTIONS

Figure O-1. 2013–14 NECAP Science: Cumulative Distribution Plots
Top: Grade 4 Bottom: Grade 8

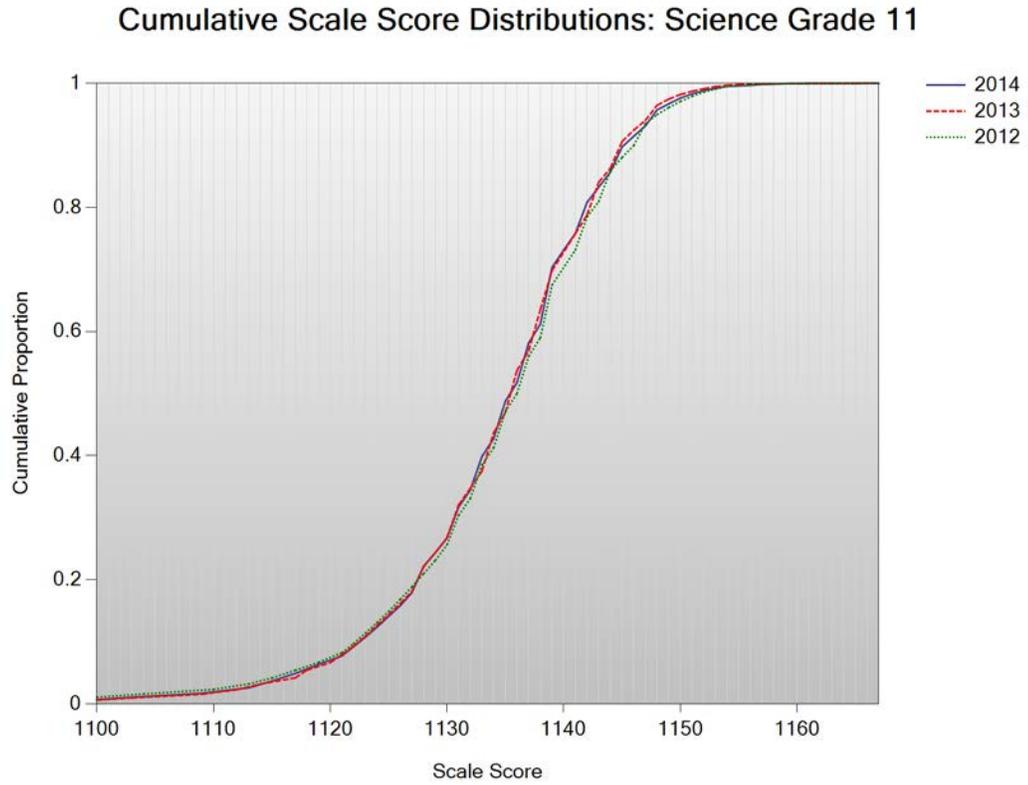
Cumulative Scale Score Distributions: Science Grade 4



Cumulative Scale Score Distributions: Science Grade 8



**Figure O-2. 2013–14 NECAP Science: Cumulative Distribution Plot—
Grade 11**



APPENDIX P—CLASSICAL RELIABILITIES

**Table P-1. 2013–14 NECAP Science: Subgroup Reliabilities—
Grade 4**

<i>Description</i>	<i>Number of Students</i>	<i>Raw Score</i>			<i>Alpha</i>	<i>Standard Error</i>
		<i>Maximum</i>	<i>Mean</i>	<i>Standard Deviation</i>		
All Students	30,843	63	32.99	10.05	0.87	3.61
Male	15,717	63	32.39	10.04	0.87	3.58
Female	15,125	63	33.61	10.02	0.87	3.63
Gender Not Reported	1	63				
Hispanic or Latino	3,351	63	26.56	10.40	0.88	3.57
American Indian or Alaskan Native	137	63	27.55	9.56	0.86	3.60
Asian	917	63	33.71	11.06	0.89	3.68
Black or African American	1,244	63	26.12	10.03	0.87	3.57
Native Hawaiian or Pacific Islander	41	63	31.88	9.33	0.84	3.73
White (non-Hispanic)	24,159	63	34.33	9.44	0.85	3.60
Two or More Races (non-Hispanic)	967	63	30.78	9.86	0.87	3.59
No Primary race/Ethnicity Reported	27	63	29.70	9.73	0.88	3.42
Currently receiving LEP services	1,262	63	20.69	9.57	0.87	3.43
Former LEP student – monitoring year 1	565	63	29.24	8.83	0.83	3.65
Former LEP student – monitoring year 2	140	63	34.21	8.44	0.81	3.64
LEP: All Other Students	28,876	63	33.59	9.73	0.86	3.61
Students with an IEP	4,426	63	25.07	9.95	0.87	3.55
IEP: All Other Students	26,417	63	34.31	9.44	0.85	3.61
Economically Disadvantaged Students	12,164	63	28.61	9.83	0.87	3.60
SES: All Other Students	18,679	63	35.84	9.12	0.84	3.59
Migrant Students	16	63	23.31	8.12	0.82	3.48
Migrant: All Other Students	30,827	63	32.99	10.05	0.87	3.61
Students receiving Title 1 Services	8,782	63	27.90	9.97	0.87	3.58
Title 1: All Other Students	22,061	63	35.01	9.34	0.85	3.60
Plan 504	187	63	33.42	9.43	0.85	3.66
Plan 504: All Other Students	30,656	63	32.99	10.05	0.87	3.61

**Table P-2. 2013–14 NECAP Science: Subgroup Reliabilities—
Grade 8**

<i>Description</i>	<i>Number of Students</i>	<i>Raw Score</i>			<i>Alpha</i>	<i>Standard Error</i>
		<i>Maximum</i>	<i>Mean</i>	<i>Standard Deviation</i>		
All Students	30,769	63	33.85	11.2	0.89	3.65
Male	16,119	63	33.56	11.38	0.9	3.61
Female	14,649	63	34.16	11.00	0.89	3.67
Gender Not Reported	1	63				
Hispanic or Latino	2,875	63	25.36	10.75	0.89	3.62
American Indian or Alaskan Native	133	63	27.49	11.62	0.90	3.69
Asian	787	63	35.91	11.37	0.90	3.66
Black or African American	1,229	63	25.47	10.60	0.88	3.61
Native Hawaiian or Pacific Islander	23	63	26.74	11.73	0.91	3.61
White (non-Hispanic)	24,954	63	35.26	10.61	0.88	3.63
Two or More Races (non-Hispanic)	747	63	32.81	11.54	0.90	3.67
No Primary race/Ethnicity Reported	21	63	18.86	10.91	0.90	3.40

continued

Description	Number of Students	Raw Score			Alpha	Standard Error
		Maximum	Mean	Standard Deviation		
Currently receiving LEP services	883	63	18.02	9.12	0.86	3.40
Former LEP student – monitoring year 1	90	63	28.37	8.53	0.82	3.58
Former LEP student – monitoring year 2	75	63	28.71	8.85	0.83	3.63
LEP: All Other Students	29,721	63	34.35	10.92	0.89	3.64
Students with an IEP	4,622	63	23.02	9.64	0.87	3.51
IEP: All Other Students	26,147	63	35.76	10.34	0.88	3.62
Economically Disadvantaged Students	11,119	63	28.34	10.69	0.88	3.64
SES: All Other Students	19,650	63	36.96	10.24	0.88	3.60
Migrant Students	9	63				
Migrant: All Other Students	30,760	63	33.85	11.20	0.89	3.65
Students receiving Title 1 Services	3,431	63	24.99	10.64	0.89	3.60
Title 1: All Other Students	27,338	63	34.96	10.77	0.89	3.64
Plan 504	243	63	34.19	10.45	0.88	3.63
Plan 504: All Other Students	30,526	63	33.84	11.21	0.89	3.65

**Table P-3. 2013–14 NECAP Science: Subgroup Reliabilities—
Grade 11**

Description	Number of Students	Raw Score			Alpha	Standard Error
		Maximum	Mean	Standard Deviation		
All Students	29,513	63	30.42	11.67	0.90	3.74
Male	15,010	63	29.89	12.07	0.91	3.68
Female	14,503	63	30.97	11.22	0.89	3.79
Gender Not Reported	0	63				
Hispanic or Latino	2,447	63	23.72	10.63	0.88	3.64
American Indian or Alaskan Native	119	63	26.13	12.16	0.91	3.68
Asian	819	63	32.31	12.09	0.90	3.79
Black or African American	1,152	63	23.65	10.66	0.88	3.63
Native Hawaiian or Pacific Islander	60	63	29.32	10.66	0.87	3.77
White (non-Hispanic)	24,276	63	31.40	11.46	0.89	3.75
Two or More Races (non-Hispanic)	614	63	29.78	11.79	0.90	3.73
No Primary race/Ethnicity Reported	26	63	20.50	11.39	0.91	3.41
Currently receiving LEP services	632	63	16.27	8.29	0.84	3.29
Former LEP student – monitoring year 1	111	63	24.26	8.10	0.80	3.66
Former LEP student – monitoring year 2	136	63	24.54	8.32	0.81	3.64
LEP: All Other Students	28,634	63	30.78	11.55	0.89	3.75
Students with an IEP	3992	63	19.13	9.21	0.87	3.37
IEP: All Other Students	25,521	63	32.18	11.01	0.88	3.76
Economically Disadvantaged Students	8,702	63	25.36	10.89	0.89	3.66
SES: All Other Students	20,811	63	32.53	11.33	0.89	3.75
Migrant Students	1	63				
Migrant: All Other Students	29,512	63	30.42	11.67	0.90	3.74
Students receiving Title 1 Services	2,594	63	23.91	10.79	0.89	3.64
Title 1: All Other Students	26,919	63	31.04	11.56	0.89	3.75
Plan 504	290	63	32.22	10.78	0.88	3.77
Plan 504: All Other Students	29,223	63	30.40	11.68	0.90	3.74

**Table P-4. 2013–14 NECAP Science: Reliabilities
by Reporting Category**

Grade	Reporting Category	Number of Items	Raw Score			Alpha	Standard Error
			Maximum	Mean	Standard Deviation		
4	ESS	12	15	9.18	2.95	0.66	1.73
	INQ	8	18	5.51	3.13	0.65	1.86
	LS	12	15	9.49	3.00	0.66	1.75
	PS	12	15	8.81	2.99	0.63	1.82
8	ESS	12	15	8.54	2.88	0.65	1.71
	INQ	8	18	8.40	3.69	0.73	1.93
	LS	12	15	8.96	3.01	0.69	1.67
	PS	12	15	7.94	3.44	0.70	1.87
11	ESS	12	15	6.99	2.82	0.59	1.81
	INQ	8	18	7.32	4.47	0.84	1.82
	LS	12	15	7.20	3.15	0.64	1.89
	PS	12	15	8.91	3.24	0.71	1.74

APPENDIX Q—INTERRATER CONSISTENCY

Table Q-1. 2013–14 NECAP Science: Item-Level Interrater Consistency Statistics by Grade

Grade	Item Number	Number of		Percent		Correlation	Percent of Third Scores
		Score Categories	Examinee Scores	Exact	Adjacent		
4	241815	5	627	75.76	22.49	0.90	1.75
	258700	4	637	74.73	22.61	0.86	2.67
	258702	3	633	73.93	24.64	0.73	1.42
	258704	3	625	79.04	20.64	0.56	0.32
	258705	3	628	77.39	21.97	0.79	0.64
	258708	3	620	86.61	12.10	0.79	1.29
	258711	3	615	78.05	21.14	0.75	0.81
	258729	4	613	81.73	17.46	0.74	0.82
	258731	3	613	74.06	24.96	0.55	0.98
	48062	5	628	53.34	40.13	0.77	5.73
	99021	5	607	66.72	28.67	0.82	4.45
8	174708	5	646	62.38	33.44	0.86	4.02
	219112	5	614	61.56	35.67	0.76	2.61
	258713	3	633	70.14	28.75	0.69	0.95
	258714	3	601	81.20	18.64	0.85	0.17
	258715	3	614	93.16	6.03	0.94	0.81
	258716	3	597	84.76	14.57	0.83	0.67
	258717	4	626	79.23	20.13	0.89	0.64
	258719	3	620	67.58	28.55	0.64	3.71
	258720	4	608	78.62	19.74	0.75	1.64
	258721	3	612	86.27	11.44	0.60	2.29
48554	5	650	56.62	37.85	0.70	5.23	
11	241745	5	558	70.43	26.52	0.82	2.69
	242632	5	575	82.78	16.35	0.93	0.87
	242637	5	581	84.17	15.15	0.94	0.69
	258655	4	597	67.34	30.49	0.75	2.18
	258656	3	582	76.98	21.82	0.80	1.03
	258657	3	582	58.59	38.49	0.52	2.92
	258659	3	561	78.25	20.32	0.79	1.43
	258661	3	583	70.84	28.30	0.67	0.86
	258662	3	552	68.84	27.54	0.66	3.44
	258669	3	568	82.75	15.85	0.82	1.41
263524	4	551	61.89	35.75	0.77	2.36	

APPENDIX R—DECISION ACCURACY AND CONSISTENCY RESULTS

Table R-1. 2013–14 NECAP Science: Summary of Decision Accuracy (and Consistency) Results by Subject and Grade—Overall and Conditional on Performance Level

Grade	Overall	Kappa	Conditional on Level			
			Substantially Below Proficient	Partially Proficient	Proficient	Proficient with Distinction
4	0.81 (0.74)	0.58	0.81 (0.68)	0.79 (0.72)	0.85 (0.79)	0.72 (0.41)
8	0.83 (0.76)	0.61	0.85 (0.76)	0.82 (0.78)	0.81 (0.72)	0.73 (0.43)
11	0.82 (0.75)	0.62	0.87 (0.79)	0.81 (0.75)	0.82 (0.75)	0.64 (0.34)

Table R-2. 2013–14 NECAP Science: Summary of Decision Accuracy (and Consistency) Results by Subject and Grade—Conditional on Cutpoint

Grade	Substantially Below Proficient / Partially Proficient			Partially Proficient / Proficient			Proficient / Proficient with Distinction		
	Accuracy (Consistency)	False		Accuracy (Consistency)	False		Accuracy (Consistency)	False	
		Positive	Negative		Positive	Negative		Positive	Negative
4	0.94 (0.91)	0.02	0.04	0.89 (0.84)	0.06	0.05	0.99 (0.99)	0.01	0
8	0.92 (0.88)	0.03	0.05	0.92 (0.88)	0.05	0.03	0.99 (0.99)	0.01	0
11	0.93 (0.90)	0.03	0.04	0.91 (0.88)	0.05	0.04	0.99 (0.98)	0.01	0

APPENDIX S—SAMPLE REPORTS

Student Michael Laststar	Grade 4	School Demonstration School 1	District Demonstration District A	State NH
------------------------------------	-------------------	---	---	--------------------

Spring 2014 - Grade 4 NECAP Science Test Results

Achievement Level	Scaled Score	This Student's Achievement Level and Scaled Score				
		Below	Partial	Proficient	Distinction	
Proficient	447	400	427	440	463	480

Interpretation of Graphic Display

The line (|) represents the student's score. The bar (▬) surrounding the score represents the probable range of scores for the student if he or she were to be tested many times. This statistic is called the standard error of measurement. See the reverse side for the achievement level descriptions.

This Student's Achievement Level Compared to Other End of Grade 4 Students by School, District, and State				
	Student	School	District	State
Proficient with Distinction		0%	0%	1%
Proficient	✓	59%	48%	45%
Partially Proficient		26%	36%	44%
Substantially Below Proficient		15%	16%	10%

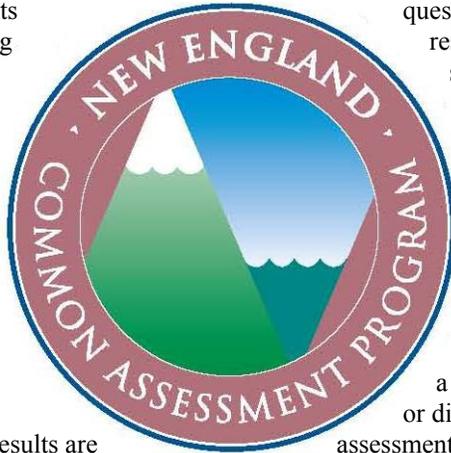
This Student's Performance in Science Domains						
	Possible Points	Student	Average Points Earned			Students at Beginning of Proficient
			School	District	State	
Physical Science	15	12	9.3	8.8	9.0	7.8-11.4
Earth Space Science	15	11	9.2	8.9	9.5	8.4-11.9
Life Science	15	11	9.4	9.0	9.8	8.7-12.2
Inquiry	18	8	5.9	5.5	5.5	3.9-7.6

Description of the Inquiry Task

There are many interesting and essential facts, formulas, and processes that students should know across the three content domains of science. But science is more than content. Inquiry skills are skills that all students should have in addition to the content. Inquiry skills are the ability to formulate questions and hypothesize, plan investigations and experiments, conduct investigations and experiments, and evaluate results. These are the broad areas that constitute scientific inquiry. Content from Physical Science, Earth Space Science, and Life Science forms the basis of each NECAP Science Inquiry Task. Instead of measuring student knowledge of content, inquiry tasks measure the student's ability to make connections, express ideas, and provide evidence of scientific thinking.

The grade 4 inquiry task, *Toy Skateboard Roll*, required students to read a story about two students riding skateboards down a ramp. In the story, the students discuss how to change the ramp to help them roll farther. They develop a research question based on their experience riding down the ramp on their skateboards. In the investigation, the students use a model to investigate if raising the height of a ramp would cause the skateboard to roll farther from the end of the ramp. Students worked with partners to complete the task and then answered questions on their own.

About The New England Common Assessment Program



This report highlights results from the Spring 2014 New England Common Assessment Program (NECAP) science tests. The NECAP tests are administered to students in New Hampshire, Rhode Island, and Vermont as part of each state's statewide assessment program.

NECAP science test results are used primarily for program evaluation, school improvement and public reporting. Achievement level results are used in the state accountability system required under No Child Left Behind (NCLB). More detailed school and district results are used by schools to help improve curriculum and instruction. Individual student results are used to support information gathered through classroom instruction and assessments.

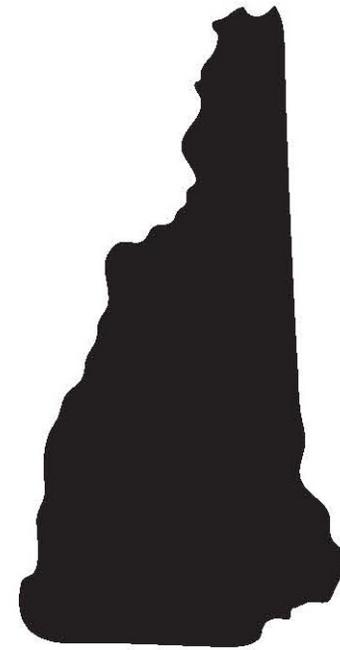
The NECAP science tests are administered to students in grades 4, 8, and 11. The tests are designed to measure student performance on standards developed and adopted by the three states. Specifically, the tests are designed to measure the content and skills that students are expected to have as they complete the K-4, 5-8, and 9-11 grade spans—in other words, the content and skills that students have learned through the end of the tested grade.

Each test contains a mix of multiple-choice and constructed-response

questions. Constructed-response questions require students to develop their own answers to questions. The science test also includes an inquiry session that requires students to answer questions based on results of an actual scientific investigation. This report contains a variety of school- and/or district-, and state-level assessment results for the NECAP

science tests administered at a grade level. Achievement level distributions and mean scaled scores are provided for all students tested as well as for subgroups of students classified by demographics or program participation. The report also contains comparative information on school and district performance on four specific science domains.

In addition to this report of grade level results, schools and districts will also receive Item Analysis Reports, released item support materials, and student-level data files containing NECAP results. Districts will also receive a Summary Report that will show results for all district schools. Together, these reports and data constitute a rich source of information to support local decisions in curriculum, instruction, assessment, and professional development. Over time, this information can also strengthen the school's and district's evaluation of their ongoing improvement efforts.



Spring 2014 Grade 4 NECAP Science Test

District Results

District: Demonstration District A

Code: DEM-DEA



Spring 2014 - Grade 4 NECAP Science Test

Grade Level Summary Report

District: Demonstration District A
State: New Hampshire
Code: DEM-DEA

Schools and districts administered all NECAP tests to every enrolled student with the following exceptions: students who participated in the alternate assessment for the 2013-14 school year, students who withdrew from the school after May 5, 2014, students who enrolled in the school after

May 5, 2014, students for whom a special consideration was granted through the state Department of Education, and other students for reasons not approved. On this page, and throughout this report, results are only reported for groups of students that are larger than nine (9).

PARTICIPATION in NECAP	Number			Percentage				
	School	District	State	School	District	State		
Students enrolled on or after May 5		80	13,948		100	100		
		Science				Science		
Students tested		77	13,799		96	99		
With an approved accommodation		19	2,978		25	22		
Current LEP Students		3	282		4	2		
With an approved accommodation		2	150		67	53		
IEP Students		14	2,078		18	15		
With an approved accommodation		11	1,584		79	76		
Students not tested in NECAP		3	149		4	1		
State Approved		2	102		67	68		
Alternate Assessment		1	97		50	95		
Withdrew After May 5		0	0		0	0		
Enrolled After May 5		0	0		0	0		
Special Consideration		1	5		50	5		
Other		1	47		33	32		

NECAP RESULTS

	District											State													
	Enrolled	NT Approved	NT Other	Tested	Level 4		Level 3		Level 2		Level 1		Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score
	N	N	N	N	N	%	N	%	N	%	N	%		N	%	%	%	%		N	%	%	%	%	
SCIENCE	80	2	1	77	0	0	37	48	28	36	12	16	437	13,799	1	45	44	10	439						

Level 4 = Proficient with Distinction; Level 3 = Proficient; Level 2 = Partially Proficient; Level 1 = Substantially Below Proficient

Note: Throughout this report, percentages may not total 100 since each percentage is rounded to the nearest whole number.

Note: Some numbers may have been left blank because fewer than ten (10) students were tested.



Spring 2014 - Grade 4 NECAP Science Test

Science Results

District: Demonstration District A
 State: New Hampshire
 Code: DEM-DEA

Proficient with Distinction (Level 4)

Students performing at this level demonstrate the knowledge and skills as described in the content standards for this grade span. Errors made by these students are few and minor and do not reflect gaps in knowledge and skills.

(Scaled Score 463–480)

Proficient (Level 3)

Students performing at this level demonstrate the knowledge and skills as described in the content standards for this grade span with only minor gaps. It is likely that any gaps in knowledge and skills demonstrated by these students can be addressed by the classroom teacher during the course of classroom instruction.

(Scaled Score 440–462)

Partially Proficient (Level 2)

Students performing at this level demonstrate gaps in knowledge and skills as described in the content standards for this grade span. Additional instructional support may be necessary for these students to achieve proficiency on the content standards.

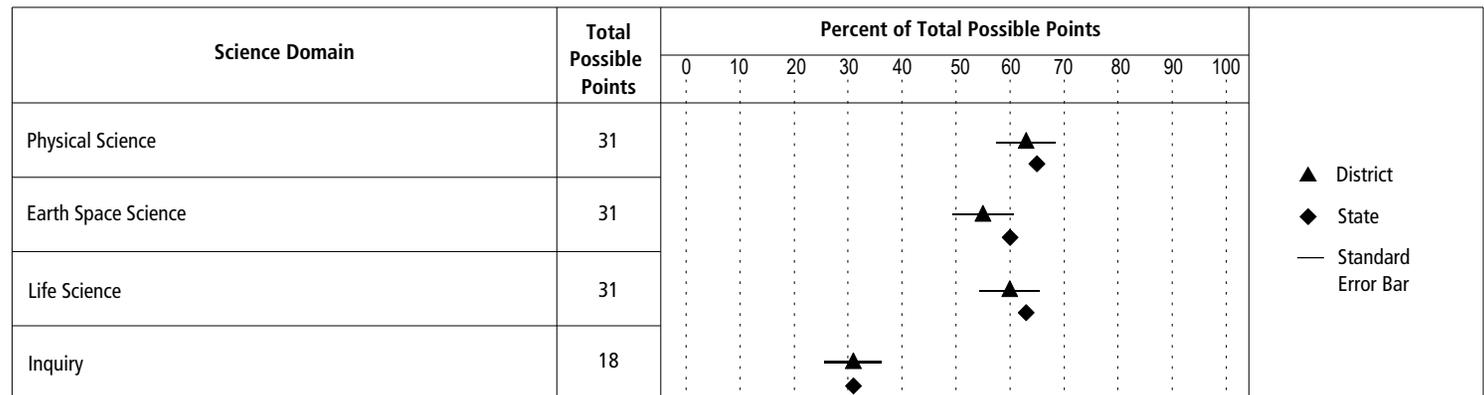
(Scaled Score 427–439)

Substantially Below Proficient (Level 1)

Students performing at this level demonstrate extensive and significant gaps in knowledge and skills as described in the content standards for this grade span. Additional instructional support is necessary for these students to achieve proficiency on the content standards.

(Scaled Score 400–426)

	Enrolled	NT Approved	NT Other	Tested	Level 4		Level 3		Level 2		Level 1		Mean Scaled Score
	N	N	N	N	N	%	N	%	N	%	N	%	
School													
2011-12													
2012-13													
2013-14													
Cumulative Total													
District													
2011-12	86	2	1	83	1	1	39	47	32	39	11	13	439
2012-13	80	3	1	76	2	3	35	46	29	38	10	13	439
2013-14	80	2	1	77	0	0	37	48	28	36	12	16	437
Cumulative Total	246	7	3	236	3	1	111	47	89	38	33	14	438
State													
2011-12	14,234	158	23	14,053	147	1	7,261	52	5,235	37	1,410	10	440
2012-13	13,919	165	26	13,728	93	1	6,855	50	5,416	39	1,364	10	440
2013-14	13,948	102	47	13,799	84	1	6,238	45	6,055	44	1,422	10	439
Cumulative Total	42,101	425	96	41,580	324	1	20,354	49	16,706	40	4,196	10	440





Spring 2014 - Grade 4 NECAP Science Test

Disaggregated Science Results

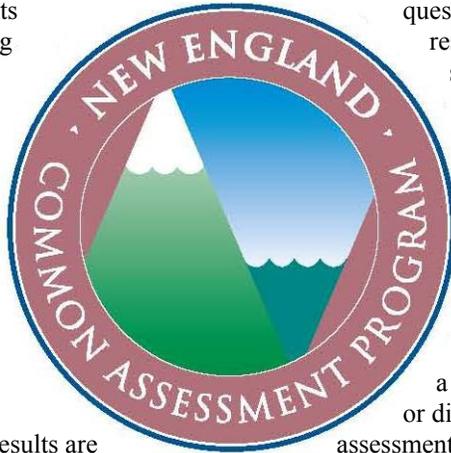
District: Demonstration District A
 State: New Hampshire
 Code: DEM-DEA

REPORTING CATEGORIES	District												State												
	Enrolled	NT Approved	NT Other	Tested	Level 4		Level 3		Level 2		Level 1		Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score
	N	N	N	N	N	%	N	%	N	%	N	%		N	%	%	%	%		N	%	%	%	%	
All Students	80	2	1	77	0	0	37	48	28	36	12	16	437	13,799	1	45	44	10	439						
Gender																									
Male	42	1	0	41	0	0	16	39	17	41	8	20	436	6,931	<1	44	45	11	438						
Female	37	1	1	35	0	0	21	60	10	29	4	11	439	6,867	1	47	43	9	440						
Not Reported	1	0	0	1										1											
Race/Ethnicity																									
Hispanic or Latino	5	1	0	4										732	<1	27	50	23	434						
Not Hispanic or Latino																									
American Indian or Alaskan Native	1	0	0	1										45	0	36	49	16	436						
Asian	2	0	0	2										456	3	51	35	11	441						
Black or African American	3	0	0	3										244	<1	24	45	30	433						
Native Hawaiian or Pacific Islander	2	0	1	1										14	0	57	36	7	441						
White	64	1	0	63	0	0	29	46	25	40	9	14	437	12,018	1	47	44	9	439						
Two or more races	2	0	0	2										289	1	34	51	15	437						
No Race/Ethnicity Reported	1	0	0	1										1											
LEP Status																									
Current LEP student	3	0	0	3										282	0	12	41	47	428						
Former LEP student - monitoring year 1	1	0	0	1										233	0	23	58	19	434						
Former LEP student - monitoring year 2	1	0	0	1										37	0	51	43	5	441						
All Other Students	75	2	1	72	0	0	36	50	27	38	9	13	438	13,247	1	46	44	9	439						
IEP																									
Students with an IEP	16	2	0	14	0	0	3	21	3	21	8	57	424	2,078	<1	21	51	28	432						
All Other Students	64	0	1	63	0	0	34	54	25	40	4	6	440	11,721	1	50	43	7	440						
SES																									
Economically Disadvantaged Students	24	2	0	22	0	0	7	32	9	41	6	27	433	4,001	<1	26	54	19	435						
All Other Students	56	0	1	55	0	0	30	55	19	35	6	11	439	9,798	1	53	40	7	441						
Migrant																									
Migrant Students	1	0	0	1										4											
All Other Students	79	2	1	76	0	0	37	49	28	37	11	14	437	13,795	1	45	44	10	439						

Level 4 = Proficient with Distinction; Level 3 = Proficient; Level 2 = Partially Proficient; Level 1 = Substantially Below Proficient

Note: Some numbers may have been left blank because fewer than ten (10) students were tested.

About The New England Common Assessment Program



This report highlights results from the Spring 2014 New England Common Assessment Program (NECAP) science tests. The NECAP tests are administered to students in New Hampshire, Rhode Island, and Vermont as part of each state's statewide assessment program.

NECAP science test results are used primarily for program evaluation, school improvement and public reporting. Achievement level results are used in the state accountability system required under No Child Left Behind (NCLB). More detailed school and district results are used by schools to help improve curriculum and instruction. Individual student results are used to support information gathered through classroom instruction and assessments.

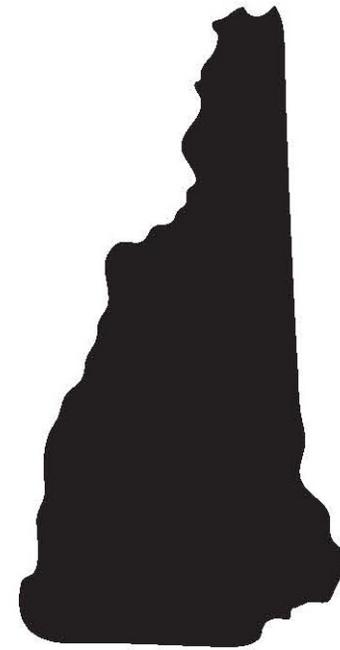
The NECAP science tests are administered to students in grades 4, 8, and 11. The tests are designed to measure student performance on standards developed and adopted by the three states. Specifically, the tests are designed to measure the content and skills that students are expected to have as they complete the K-4, 5-8, and 9-11 grade spans—in other words, the content and skills that students have learned through the end of the tested grade.

Each test contains a mix of multiple-choice and constructed-response

questions. Constructed-response questions require students to develop their own answers to questions. The science test also includes an inquiry session that requires students to answer questions based on results of an actual scientific investigation. This report contains a variety of school- and/or district-, and state-level assessment results for the NECAP

science tests administered at a grade level. Achievement level distributions and mean scaled scores are provided for all students tested as well as for subgroups of students classified by demographics or program participation. The report also contains comparative information on school and district performance on four specific science domains.

In addition to this report of grade level results, schools and districts will also receive Item Analysis Reports, released item support materials, and student-level data files containing NECAP results. Districts will also receive a Summary Report that will show results for all district schools. Together, these reports and data constitute a rich source of information to support local decisions in curriculum, instruction, assessment, and professional development. Over time, this information can also strengthen the school's and district's evaluation of their ongoing improvement efforts.



Spring 2014 Grade 8 NECAP Science Test

District Results

District: Demonstration District A

Code: DEM-DEA



Spring 2014 - Grade 8 NECAP Science Test

Grade Level Summary Report

District: Demonstration District A
State: New Hampshire
Code: DEM-DEA

Schools and districts administered all NECAP tests to every enrolled student with the following exceptions: students who participated in the alternate assessment for the 2013-14 school year, students who withdrew from the school after May 5, 2014, students who enrolled in the school after

May 5, 2014, students for whom a special consideration was granted through the state Department of Education, and other students for reasons not approved. On this page, and throughout this report, results are only reported for groups of students that are larger than nine (9).

PARTICIPATION in NECAP	Number			Percentage			
	School	District	State	School	District	State	
Students enrolled on or after May 5		83	14,662		100	100	
		Science				Science	
Students tested		78	14,414		94	98	
With an approved accommodation		9	2,016		12	14	
Current LEP Students		2	255		3	2	
With an approved accommodation		2	123		100	48	
IEP Students		9	2,239		12	16	
With an approved accommodation		5	1,491		56	67	
Students not tested in NECAP		5	248		6	2	
State Approved		2	152		40	61	
Alternate Assessment		1	138		50	91	
Withdrew After May 5		0	0		0	0	
Enrolled After May 5		0	0		0	0	
Special Consideration		1	14		50	9	
Other		3	96		60	39	

NECAP RESULTS

	District											State														
	Enrolled	NT Approved	NT Other	Tested		Level 4		Level 3		Level 2		Level 1		Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score
	N	N	N	N	N	%	N	%	N	%	N	%	N	%	N	%	%	%	%	%	N	%	%	%	%	%
SCIENCE	83	2	3	78	0	0	19	24	46	59	13	17	835	14,414	<1	24	55	20	835							

Level 4 = Proficient with Distinction; Level 3 = Proficient; Level 2 = Partially Proficient; Level 1 = Substantially Below Proficient

Note: Throughout this report, percentages may not total 100 since each percentage is rounded to the nearest whole number.

Note: Some numbers may have been left blank because fewer than ten (10) students were tested.



Spring 2014 - Grade 8 NECAP Science Test

Disaggregated Science Results

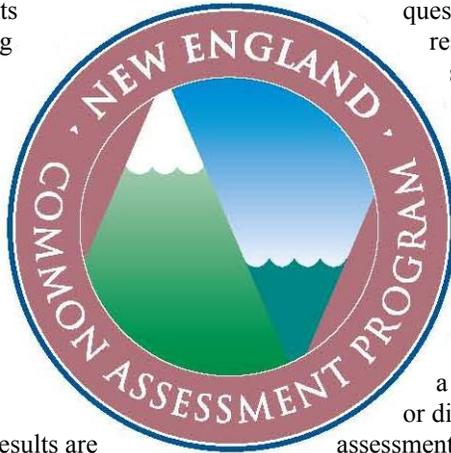
District: Demonstration District A
 State: New Hampshire
 Code: DEM-DEA

REPORTING CATEGORIES	District												State												
	Enrolled	NT Approved	NT Other	Tested	Level 4		Level 3		Level 2		Level 1		Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score
	N	N	N	N	N	%	N	%	N	%	N	%		N	%	%	%	%		N	%	%	%	%	
All Students	83	2	3	78	0	0	19	24	46	59	13	17	835	14,414	<1	24	55	20	835						
Gender																									
Male	40	0	2	38	0	0	10	26	22	58	6	16	835	7,581	<1	24	55	21	834						
Female	43	2	1	40	0	0	9	23	24	60	7	18	835	6,833	1	25	56	18	835						
Not Reported	0	0	0	0									0												
Race/Ethnicity																									
Hispanic or Latino	5	0	1	4										648	<1	8	50	42	829						
Not Hispanic or Latino																									
American Indian or Alaskan Native	1	0	0	1										41	0	22	49	29	833						
Asian	2	0	0	2										402	1	35	50	14	837						
Black or African American	1	0	0	1										279	<1	8	47	44	828						
Native Hawaiian or Pacific Islander	1	0	0	1										10	0	10	70	20	833						
White	71	2	2	67	0	0	16	24	42	63	9	13	836	12,839	<1	25	56	18	835						
Two or more races	2	0	0	2										195	1	23	55	21	834						
No Race/Ethnicity Reported	0	0	0	0										0											
LEP Status																									
Current LEP student	3	0	1	2										255	0	2	20	78	822						
Former LEP student - monitoring year 1	1	0	0	1										29	0	3	59	38	830						
Former LEP student - monitoring year 2	1	0	0	1										28	0	4	71	25	831						
All Other Students	78	2	2	74	0	0	19	26	44	59	11	15	835	14,102	<1	25	56	19	835						
IEP																									
Students with an IEP	12	1	2	9										2,239	0	4	40	57	827						
All Other Students	71	1	1	69	0	0	18	26	45	65	6	9	837	12,175	1	28	58	13	836						
SES																									
Economically Disadvantaged Students	28	0	2	26	0	0	4	15	16	62	6	23	832	3,917	<1	10	55	35	831						
All Other Students	55	2	1	52	0	0	15	29	30	58	7	13	836	10,497	1	30	56	14	836						
Migrant																									
Migrant Students	1	0	1	0										4											
All Other Students	82	2	2	78	0	0	19	24	46	59	13	17	835	14,410	<1	24	55	20	835						

Level 4 = Proficient with Distinction; Level 3 = Proficient; Level 2 = Partially Proficient; Level 1 = Substantially Below Proficient

Note: Some numbers may have been left blank because fewer than ten (10) students were tested.

About The New England Common Assessment Program



This report highlights results from the Spring 2014 New England Common Assessment Program (NECAP) science tests. The NECAP tests are administered to students in New Hampshire, Rhode Island, and Vermont as part of each state's statewide assessment program.

NECAP science test results are used primarily for program evaluation, school improvement and public reporting. Achievement level results are used in the state accountability system required under No Child Left Behind (NCLB). More detailed school and district results are used by schools to help improve curriculum and instruction. Individual student results are used to support information gathered through classroom instruction and assessments.

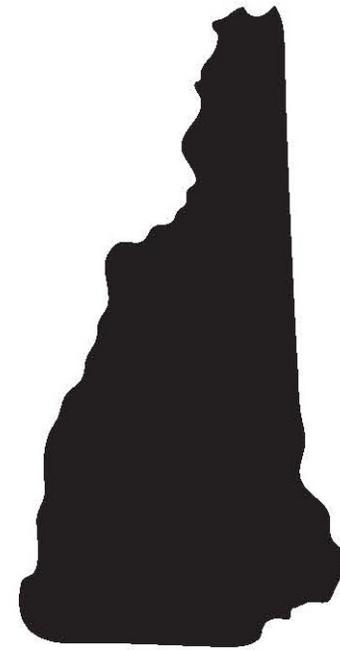
The NECAP science tests are administered to students in grades 4, 8, and 11. The tests are designed to measure student performance on standards developed and adopted by the three states. Specifically, the tests are designed to measure the content and skills that students are expected to have as they complete the K-4, 5-8, and 9-11 grade spans—in other words, the content and skills that students have learned through the end of the tested grade.

Each test contains a mix of multiple-choice and constructed-response

questions. Constructed-response questions require students to develop their own answers to questions. The science test also includes an inquiry session that requires students to answer questions based on results of an actual scientific investigation. This report contains a variety of school- and/or district-, and state-level assessment results for the NECAP

science tests administered at a grade level. Achievement level distributions and mean scaled scores are provided for all students tested as well as for subgroups of students classified by demographics or program participation. The report also contains comparative information on school and district performance on four specific science domains.

In addition to this report of grade level results, schools and districts will also receive Item Analysis Reports, released item support materials, and student-level data files containing NECAP results. Districts will also receive a Summary Report that will show results for all district schools. Together, these reports and data constitute a rich source of information to support local decisions in curriculum, instruction, assessment, and professional development. Over time, this information can also strengthen the school's and district's evaluation of their ongoing improvement efforts.



Spring 2014 Grade 11 NECAP Science Test

District Results

District: Demonstration District A

Code: DEM-DEA



Spring 2014 - Grade 11 NECAP Science Test

Grade Level Summary Report

District: Demonstration District A
State: New Hampshire
Code: DEM-DEA

Schools and districts administered all NECAP tests to every enrolled student with the following exceptions: students who participated in the alternate assessment for the 2013-14 school year, students who withdrew from the school after May 5, 2014, students who enrolled in the school after

May 5, 2014, students for whom a special consideration was granted through the state Department of Education, and other students for reasons not approved. On this page, and throughout this report, results are only reported for groups of students that are larger than nine (9).

PARTICIPATION in NECAP	Number			Percentage			
	School	District	State	School	District	State	
Students enrolled on or after May 5		83	14,167		100	100	
		Science				Science	
Students tested		78	13,658		94	96	
With an approved accommodation		6	1,459		8	11	
Current LEP Students		2	169		3	1	
With an approved accommodation		0	72		0	43	
IEP Students		9	1,950		12	14	
With an approved accommodation		5	1,225		56	63	
Students not tested in NECAP		5	509		6	4	
State Approved		2	94		40	18	
Alternate Assessment		1	82		50	87	
Withdrew After May 5		0	0		0	0	
Enrolled After May 5		0	0		0	0	
Special Consideration		1	12		50	13	
Other		3	415		60	82	

NECAP RESULTS

	District											State													
	Enrolled	NT Approved	NT Other	Tested	Level 4		Level 3		Level 2		Level 1		Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score
	N	N	N	N	N	%	N	%	N	%	N	%		N	%	%	%	%		N	%	%	%	%	
SCIENCE	83	2	3	78	1	1	27	35	33	42	17	22	1136	13,658	1	28	47	24	1135						

Level 4 = Proficient with Distinction; Level 3 = Proficient; Level 2 = Partially Proficient; Level 1 = Substantially Below Proficient

Note: Throughout this report, percentages may not total 100 since each percentage is rounded to the nearest whole number.

Note: Some numbers may have been left blank because fewer than ten (10) students were tested.



Spring 2014 - Grade 11 NECAP Science Test

Disaggregated Science Results

District: Demonstration District A
 State: New Hampshire
 Code: DEM-DEA

REPORTING CATEGORIES	District												State												
	Enrolled	NT Approved	NT Other	Tested	Level 4		Level 3		Level 2		Level 1		Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score	Tested	Level 4	Level 3	Level 2	Level 1	Mean Scaled Score
	N	N	N	N	N	%	N	%	N	%	N	%		N	%	%	%	%		N	%	%	%	%	
All Students	83	2	3	78	1	1	27	35	33	42	17	22	1136	13,658	1	28	47	24	1135						
Gender																									
Male	41	1	1	39	1	3	12	31	17	44	9	23	1135	7,027	1	27	45	27	1134						
Female	41	1	1	39	0	0	15	38	16	41	8	21	1137	6,631	1	28	50	20	1135						
Not Reported	1	0	1	0										0											
Race/Ethnicity																									
Hispanic or Latino	6	0	1	5										497	0	8	45	47	1129						
Not Hispanic or Latino																									
American Indian or Alaskan Native	1	0	0	1										34	0	18	32	50	1131						
Asian	3	0	0	3										387	3	35	41	20	1136						
Black or African American	2	0	0	2										247	0	10	48	43	1129						
Native Hawaiian or Pacific Islander	1	0	0	1										14	0	14	43	43	1130						
White	68	2	1	65	1	2	25	38	27	42	12	18	1137	12,324	1	29	48	22	1135						
Two or more races	1	0	0	1										155	1	28	44	27	1134						
No Race/Ethnicity Reported	1	0	1	0										0											
LEP Status																									
Current LEP student	2	0	0	2										169	0	2	13	85	1122						
Former LEP student - monitoring year 1	1	0	0	1										42	0	2	62	36	1130						
Former LEP student - monitoring year 2	1	0	0	1										76	0	3	67	30	1131						
All Other Students	79	2	3	74	1	1	27	36	31	42	15	20	1136	13,371	1	28	47	23	1135						
IEP																									
Students with an IEP	10	1	0	9										1,950	<1	5	32	63	1126						
All Other Students	73	1	3	69	1	1	26	38	31	45	11	16	1137	11,708	2	32	50	17	1136						
SES																									
Economically Disadvantaged Students	19	0	0	19	0	0	2	11	10	53	7	37	1130	2,869	1	13	47	39	1131						
All Other Students	64	2	3	59	1	2	25	42	23	39	10	17	1138	10,789	2	32	47	19	1136						
Migrant																									
Migrant Students	1	0	0	1										1											
All Other Students	82	2	3	77	1	1	27	35	33	43	16	21	1136	13,657	1	28	47	24	1135						

Level 4 = Proficient with Distinction; Level 3 = Proficient; Level 2 = Partially Proficient; Level 1 = Substantially Below Proficient

Note: Some numbers may have been left blank because fewer than ten (10) students were tested.



District Summary

2013-2014 Students

District: Demonstration District A
State: New Hampshire
Code: DEM-DEA

Science	Enrolled	NT Approved	NT Other	Tested	Achievement Level								Mean Score
	N	N	N	N	Level 4		Level 3		Level 2		Level 1		
					N	%	N	%	N	%	N	%	
Demonstration District A	246	6	7	233	1	<1	83	36	107	46	42	18	
Grade 4	80	2	1	77	0	0	37	48	28	36	12	16	437
Grade 8	83	2	3	78	0	0	19	24	46	59	13	17	835
Grade 11	83	2	3	78	1	1	27	35	33	42	17	22	1136

Level 4 = Proficient with Distinction; Level 3 = Proficient; Level 2 = Partially Proficient; Level 1 = Substantially Below Proficient

APPENDIX T—INTERACTIVE REPORTS



CONFIDENTIAL

Spring 2014 - Grade 04 NECAP Tests

Item Analysis Report - Science

School: Demonstration School 1
 District: Demonstration District A
 State: Vermont
 Code: DEMOA-DEMO1

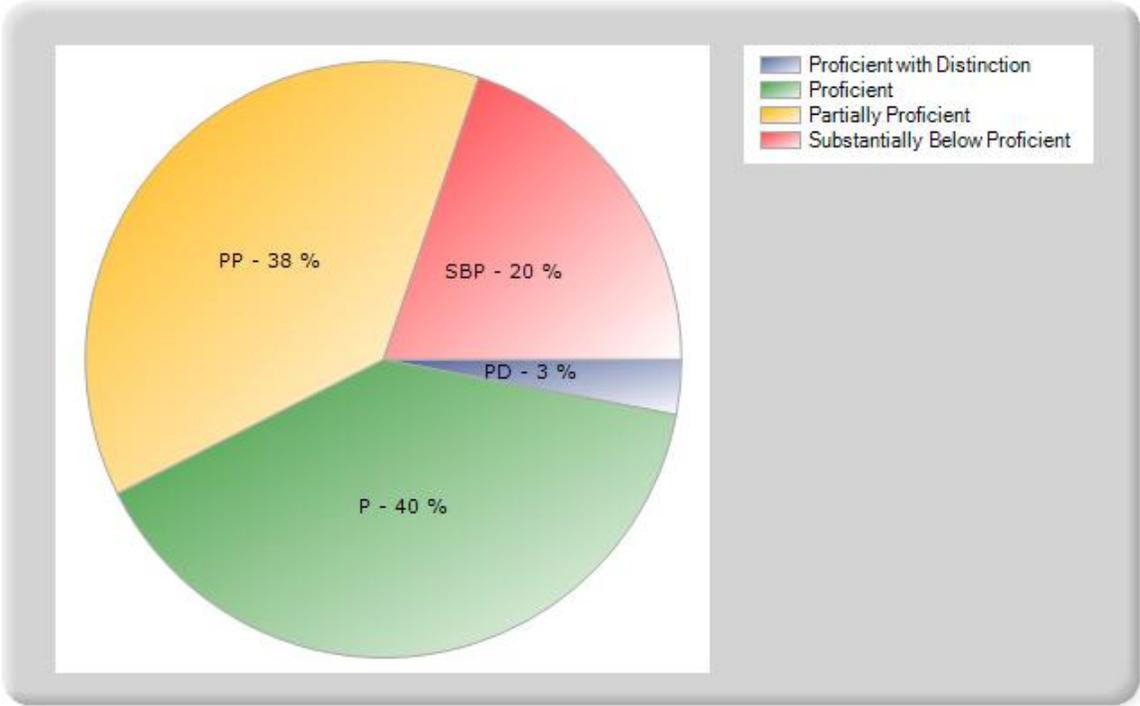
Name/Student ID	Item Number	Released Items										Total Test Results						
		1	2	3	4	5	6	7	8	9	10	Domain Points Earned				Total Points Earned	Scaled Score	Achievement Level
		Science Domain	PS	PS	PS	ESS	ESS	ESS	ESS	LS	LS	LS	Physical Science	Earth Space Science	Life Science			
		Assessment Target	1-1	2-6	3-8	1-3	1-4	1-6	1-2	1-1	2-5	3-7						
		Depth of Knowledge Code	2	2	2	2	2	2	2	2	2	2						
		Item Type	MC	MC	MC	MC	MC	MC	CR	MC	MC	MC						
Correct MC Response	B	A	C	A	D	C		B	B	B								
Total Possible Points	1	1	1	1	1	1	4	1	1	1	15	15	15	18	63			
Austin, Samantha	D043038	+	C	+	+	+	+	2	+	+	C	11	12	11	10	44	449	3
Bagge, Connor A	D043047	A	C	B	+	+	+	3	+	D	+	6	10	11	6	33	438	2
Caldwell, Terrence	D043024	+	B	+	+	+	+	1	+	+	C	10	11	13	5	39	444	3
Cardoza, Devyn	D043045	+	B	B	+	+	+	1	C	D	+	5	7	5	4	21	426	1
Carson, Essence S	D043018	+	D	+	C	C	+	1	D	A	A	12	7	10	11	40	445	3
Liberty, Sydney	D043002	+	D	A	+	C	+	1	A	D	C	7	8	6	11	32	437	2
Lord, Hannah	D043037	A	+	+	+	+	+	0	+	D	+	10	7	7	2	26	431	2
Lutskiy, Dennis L	D043012	A	C	A	+	A	+	0	A	A	+	2	5	7	1	15	420	1
Mansfield, Stephanie	D043005	+	+	+	+	+	+	1	+	+	+	12	9	10	3	34	439	2
Marshall, Shalane	D043048	+	C	+	+	C	+	3	+	C	+	10	12	12	6	40	445	3
Martin, Hugh L	D043025	D	B	+	+	+	+	2	+	A	C	8	9	7	2	26	431	2
Masten, Dakota J	D043006	+	C	+	+	A	+	0	+	+	C	7	7	5	2	21	426	1
Matson, Miranda	D043023	D	D	D	D	+	+	0	+	A	D	5	5	4	1	15	420	1
Mccarthy, Dylan	D043050	C	D	+	+	C	+	3	+	A	+	10	11	7	9	37	442	3
Patterson, Justin H	D043031	+	D	B	+	+	+	2	A	+	C	12	10	11	5	38	443	3
Perez, Eric L	D043011	+	B	+	+	+	+	3	+	+	+	9	10	12	4	35	439	2
Rileymcnary, Maille	D043004											0	0	0	0	0		W
Roberts, Brandon S	D043051	D	+	B	+	+	+	3	D	+	+	9	10	12	8	39	444	3
Saylor, Brenda	D043022	A	+	A	B	+	B	0	C	C	D	4	3	3	1	11	414	1
Smith, Dezmond D	D043014											2	0	0	2	4	400	1
Snow, Mackayla M	D043029	+	B	+	+	+	+	1	+	+	C	12	11	13	8	44	449	3
Vettese, Mariah O	D043021	+	B	+	+	A	D	2	+	C	+	9	9	9	8	35	439	2
Vogel, Dean A	D043007	C	+	+	+	+	D	1	A	+	+	10	6	9	6	31	436	2
Item Number	1	2	3	4	5	6	7	8	9	10								
Percent Correct/Average Score: Group	55	23	59	82	64	82	1.4	59	41	50	8.3	8.1	8.4	5.2				
Percent Correct/Average Score: School	55	23	59	82	64	82	1.4	59	41	50	8.3	8.1	8.4	5.2				
Percent Correct/Average Score: District	56	42	50	79	67	77	1.5	67	50	46	8.5	8.5	9.0	5.4				
Percent Correct/Average Score: State	58	47	62	86	60	81	1.7	69	60	65	8.8	9.3	9.4	5.4				



Achievement Level Summary

District: Demonstration District A
School: Demonstration School 1
Grade: 11
Date: 1/6/2015 2:18:24 PM

Science



Achievement Level	Count	Percentage %*
Proficient with Distinction	1	3
Proficient	16	40
Partially Proficient	15	38
Substantially Below Proficient	8	20

*Percentages may not total exactly 100% due to applied rounding.



Science Released Items Summary Data

District: Demonstration District A

School: Demonstration School 1

Grade: 11

Date: 1/6/2015 2:19:35 PM

Multiple Choice

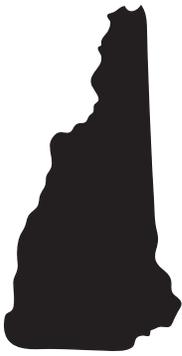
Released Item	Domain	Assessment Target	Correct (#)	A (#)	B (#)	C (#)	D (#)	IR (#)	Correct Response
1	PS	1-1	31	2	7	31	0	0	C
2	PS	2-6	32	0	3	32	5	0	C
3	PS	3-10	24	3	24	8	4	1	B
5	ESS	1-2	23	23	6	1	10	0	A
6	ESS	3-6	24	4	9	24	2	1	C
7	ESS	3-8	25	1	4	10	25	0	D
8	LS	2-3	34	3	34	0	1	2	B
9	LS	4-9	22	3	4	10	22	1	D
10	LS	4-10	23	2	23	7	7	1	B

Constructed Response

Released Item	Domain	Assessment Target	Point Value	Average Score
4	PS	1-3	4	1.8

Inquiry Task

Released Item	Domain	Inquiry Construct	Point Value	Average Score
1	INQ	1	3	1.3
2	INQ	3	2	1.1
3	INQ	9	2	1.1
4	INQ	10	2	1.0
5	INQ	5	2	0.5
6	INQ	12	2	1.1
7	INQ	11	3	1.2
8	INQ	13	2	1.0



CONFIDENTIAL

Student Name
Samantha Barnhart

Longitudinal Data Report

Year	Enrolled Grade	School Name	Administration	Test Name	Content Area	Score	Achievement Level
0910	11	Demonstration School 2	NECAP Spring 2010 - Science	Grade 11 Science	sci	1125	Substantially Below Proficient
1011	11	Demonstration School 2	NECAP Fall 2010	Grade 11 Mathematics	mat	1137	Partially Proficient
1011	11	Demonstration School 2	NECAP Fall 2010	Grade 11 Reading	rea	1144	Proficient
1011	11	Demonstration School 1	NECAP Spring 2011 - Science	Grade 11 Science	sci	1138	Partially Proficient
1011	11	Demonstration School 2	NECAP Fall 2010	Grade 11 Writing	wri	8	Proficient
1112	11	Demonstration School 1	NECAP Fall 2011	Grade 11 Mathematics	mat	1108	Substantially Below Proficient
1112	11	Demonstration School 1	NECAP Fall 2011	Grade 11 Reading	rea	1120	Substantially Below Proficient
1112	11	Demonstration School 1	NECAP Spring 2012 - Science	Grade 11 Science	sci	1147	Proficient
1112	11	Demonstration School 1	NECAP Fall 2011	Grade 11 Writing	wri	2	Substantially Below Proficient
1213	11	Demonstration School 2	NECAP Fall 2012	Grade 11 Mathematics	mat	1147	Proficient
1213	11	Demonstration School 2	NECAP Fall 2012	Grade 11 Reading	rea	1162	Proficient with Distinction
1213	11	Demonstration School 2	NECAP Spring 2013 - Science	Grade 11 Science	sci	1143	Proficient
1213	11	Demonstration School 2	NECAP Fall 2012	Grade 11 Writing	wri	8	Proficient
1314	11	Demonstration School 1	NECAP Fall 2013	Grade 11 Mathematics	mat	1143	Proficient
1314	11	Demonstration School 1	NECAP Fall 2013	Grade 11 Reading	rea	1151	Proficient
1314	11	Demonstration School 1	NECAP Spring 2014 - Science	Grade 11 Science	sci	1117	Substantially Below Proficient
1314	11	Demonstration School 1	NECAP Fall 2013	Grade 11 Writing	wri	7	Proficient

Note: This report returns as many years of NECAP data as are available for this student beginning with 08-09.

APPENDIX U—ANALYSIS AND REPORTING DECISION RULES

Data Analysis and Static Reporting Decision Rules
NECAP
Spring 13-14 Administration

This document specifies rules for data analysis and static reporting requirements. The final student level data set used for analysis and reporting is described in the “Data Processing Specifications.” This document is considered a draft until the NECAP State Departments of Education (DOE) signs off. If there are rules that need to be added or modified after said sign-off, DOE sign off will be obtained for each rule. Details of these additions and modifications will be in the Addendum section.

I. General Information

NECAP is administered in the fall and spring. This document incorporates fall and spring rules so that changes are carried to future administrations. In the spring, students are reported based on the spring school/district (referred to as testing school/district). In the spring, students are not reported based on the teaching school. Rules pertaining to the teaching school/district can be ignored for spring administrations. For more information regarding discode, schcode, sprdiscode, sprschcode, senddiscode, and sprsenddiscode, please refer to the data processing specifications and demographic data specification.

This document is the official rules for the current reporting administration.

A. Spring Tests Administered

Grade	Subject	Test items used for Scaling	Item Reporting Categories (Subtopic and Subcategory Source)
04	Science	Common	Cat3
08	Science	Common	Cat3
11	Science	Common	Cat3

B. Reports Produced:

1. Student Report
 - a. Testing School District
 - I. Parent Copy
 - II. School Copy
2. Interactive Reporting (Only the data analysis requirements are outlined in this document)
 - a. Item Analysis
 - b. Achievement Level Summary
 - c. Item Information
 - d. Student Longitudinal
3. Grade Level School/District/State Results
 - a. Testing School District
4. District/State Summary
 - a. Testing School District
5. Writing Prompt CDs

C. Files Produced:

1. Preliminary State Results
2. State Student Released Item Data
3. State Student Raw Data
4. State Student Scored Data

5. District Student Data
6. School Student Data
7. Common Item Information
8. Grade Level Results Report Disaggregated and Historical Data
9. State Standard Deviations and Average Scaled Scores
10. Grade Level Results Report Participation Category Data
11. Grade Level Results Report Subtopic Data
12. Summary Results Data
13. Released Item Percent Responses Data
14. Invalidated Students Original Score
15. Student Questionnaire Summary
16. TCTA Questionnaire Raw Data
17. TCTA Questionnaire Frequency Distribution
18. Scaled Score Lookup
19. Subtopic Average Points Earned (For Program Management)
20. Item Stats for Inquiry Task Items (For Program Management)
21. Memo Shipping files (For Program Management)
22. Released Item Info File (For Program Management)
23. CD print file.

D. School Type:

Testing School Type: SchType	Source: ICORE SubTypeID	Description	States
PUB	1,12,13	Public School	ME, NH, RI, VT
CHA	11	Charter School	NH, RI, ME
PSP	19	Public Special Purpose	ME
PSE	15	Public Special Education	ME
INS	7	Institution	VT
OTH	9	Other	VT
OOD	4	Out-of-District Private Providers	NH
OUT	8	Out Placement	RI
PSN	23	Private Special Purpose	ME
BIG	6	Private with >60% Publicly Funded	ME
PRI	3	Private School	RI, VT

School Type Impact on Data Analysis and Reporting				
Level	Testing		Teaching (Fall Only)	
	Impact on Analysis	Impact on Reporting	Impact on Analysis	Impact on Reporting
Student	n/a	Report students based on testing discode and schcode. District data will be blank for students tested at BIG, PSN, PRI, OOD, OUT, INS, or OTH schools. Always print tested year state data.	n/a	n/a
School	Do not exclude any students based on school type using testing school code for aggregations	Generate a report for each school with at least one student enrolled using the tested school aggregate denominator. District data will be blank for BIG, PSN, PRI, OOD, OUT, INS, or OTH schools. Always print tested year state data.	Exclude students who do not have a teaching school code.	Generate a report for each school with at least one student enrolled using the teaching school aggregate denominator. District data will be blank for BIG, PSN, PRI, OOD, OUT, INS, or OTH schools. Always print tested year state data.
District	For OUT, OOD, BIG, and PSN schools, aggregate using the sending district. If OUT, OOD, BIG, or PSN student does not have a sending district, do not include in aggregations. Do not include students tested at PRI, INS, or OTH schools	Generate a report for each district with at least one student enrolled using the tested district aggregate denominator. Always report tested year state data.	For OUT, OOD, BIG, and PSN teaching schools, aggregate using the spring sending district. If OUT, OOD, BIG, or PSN teaching school student does not have a teaching sending district, do not include in aggregations. Do not include students taught at PRI, INS, or OTH schools	Generate a report for each district with at least one student enrolled using the teaching district aggregate denominator. Always report tested year state data.
State	Do not include students tested at PRI schools for NH and RI. Include all students for VT and ME.	Always report testing year state data.	n/a	n/a

E. Student Status

StuStatus	Description
1	Homeschooled
2	Privately Funded
3	Exchange Student
4	Excluded State
0	Publicly Funded

StuStatus impact on Data Analysis and Reporting		
Level	Impact on Analysis	Impact on Reporting
Student	n/a	School and District data will be blank for students with a StuStatus value of 1,2 or 3. Always print tested year state data. For StuStatus values of 1, 2, and 3 print the description from the table above for the school and district names.
School	Exclude all students with a StuStatus value of 1, 2 or 3.	Students with a StuStatus value of 1, 2 or 3 are excluded from Interactive Reporting.
District	Exclude all students with a StuStatus value of 1, 2 or 3.	n/a
State	Exclude all students with a StuStatus value of 1, 2, 3 or 4.	n/a.

F. Requirements To Report Aggregate Data(Minimum N)

Calculation Description	Rule
Number and Percent at each achievement level, mean score by disaggregated category and aggregate level	If the number of tested students included in the denominator is less than 10, then do not report.
Content Area Subcategories Average Points Earned based on common items only by aggregate level	If the number of tested students included in the denominator is less than 10, then do not report.
Aggregate data on Item Analysis report	No required minimum number of students
Number and Percent of students in a participation category by aggregate level	No required minimum number of students
Content Area Subtopic Percent of Total Possible Points and Standard Error Bar and Grade 11 Writing Distribution of Score Points Across Prompts	If any item was not administered to at least one tested student included in the denominator or the number of tested students included in the denominator is less than 10, then do not report
Content Area Cumulative Total Enrollment, Not tested, Tested, Number and Percent at each achievement level, mean score	Suppress all cumulative total data if at least one reported year has fewer than 10 tested students. Spring: The reported years are 1112, 1213, and 1314.

G. Special Forms:

1. Form 00 is created for students whose matrix scores will be ignored for analysis. Such students include Braille or administration issues resolved by program management.

H. Other Information

1. NH, RI, and VT participate in NECAP testing for Grades 03-08 and 11. ME only participates in NECAP testing for Grades 03-08.
2. Grade 12 students are allowed to participate in the NECAP Grade 11 test under the following circumstances: RI students trying to improve prior NECAP score, and NH, RI, and VT students taking the NECAP Grade 11 test for the first time.
 - a. RI students trying to improve are identified as StuGrade=12 and Grade=11. They only receive a student report. They are not listed on a roster or included in any aggregations. Do not print tested school and district aggregate data on the student report.
 - b. For students taking NECAP for the first time the StuGrade in the student demographics file will be 11 and the remaining decision rules apply.
3. Plan504 data not available for NH and VT; therefore 504 Plan section will be suppressed for NH and VT.
4. To calculate Title1 data for writing using Title1rea variable.
5. Title 1 data are not available for VT; therefore Title 1 section will be suppressed for VT.
6. Title 1 Science data are not available for NH; therefore, Title 1 section will be suppressed for NH on Science specific reports. Title 1 Reading and Math data are available for NH and should not be suppressed.
7. Testing level is defined by the variables discode and schcode. Teaching level is defined by the variables sprdiscode and sprschcode. Every student will have testing district and school codes. In the fall, some students will have a teaching school code and some students will have a teaching district code. In the spring, no students will have a teaching school/district.
8. A non-public district code is a district code associated with a school that is type BIG, PSN, PRI, OOD, OUT, INS, or OTH. Non-public testing sending district codes will be ignored. . For example: For RI, senddiscode of 88 is ignored. For NH, senddiscode of 000 is ignored.
9. Only students with a testing school type of OUT, OOD, BIG, or PSN are allowed to have a testing sending district code. Testing sending district codes will be blanked for students at any other testing school types.
10. Only students with a teaching school type of OUT, OOD, BIG, or PSN are allowed to have a spring sending district code. Spring sending district codes will be blanked for students at any other teaching school types.
11. If students have a teaching district code and no teaching school, then ignore teaching district codes that are associated with schools that are BIG, PSN, PRI, OOD, OUT, INS, or OTH.

II. Student Participation / Exclusions

A. Test Attempt Rules by content area

1. A multiple choice item has been answered by a student if the response is A, B, C, D, or * (*=multiple responses)
2. An open response item has been answered if it is not scored blank 'B'

B. Session Attempt Rules by content area

1. A session was attempted if any multiple choice item or non-field test open response item has been answered in the session. (Use original item responses – see special circumstances section II.F)

C. Not Tested Reasons by content area

1. Not Tested State Approved Alternate Assessment
 - a. If a student links to the demographic file has content area not tested status of “Not Tested State Approved Alternate Assessment” is identified as “Not Tested State Approved Alternate Assessment” for the content area.
 - b. If a student is identified as receiving an alternate assessment achievement level, then the student’s record will be updated as outlined in the `NECAP1314StudentDemographicFileDescription.doc`.
 2. Not Tested State Approved First Year LEP (reading and writing only)
 - a. If a student links to the demographic file has content area not tested status of “Not Tested State Approved First Year LEP” or does not link to the demographic file has content area “First Year LEP blank or partially blank reason” marked, then the student is identified as “Not Tested State Approved First Year LEP”.
 3. Not Tested State Approved Special Consideration
 - a. If a student links to the demographic data file has content area “Not Tested State Approved Special Consideration” indicated or does not link to the demographic data file and has content area “Special Consideration blank or partially blank reason” marked, then the student is identified as “Not Tested State Approved Special Consideration”.
 4. Not Tested State Approved Withdrew After
 - a. If a student links to the demographic data file has content area not tested status of “Not Tested Withdrew After” and at least one content area session was not attempted or does not link to the demographic file has content area “Withdrew After blank or partially blank reason” marked and at least one content area session was not attempted, then the student is identified as “Not Tested State Approved Withdrew After”.
 5. Not Tested State Approved Enrolled After
 - a. If a student links to the demographic data file has content area not tested status of “Not Tested Enrolled After” and at least one content area session was not attempted or does not link to the demographic file has content area “Enrolled After blank or partially blank reason” marked and at least one content area session was not attempted, then the student is identified as “Not Tested State Approved Enrolled After”.
 6. Not Tested Other
 - a. If content area test was not attempted, the student is identified as “Not Tested Other”.
- D. Not Tested Reasons Hierarchy by content area: if more than one reason for not testing at a content area is identified then select the first category indicated in the order of the list below.
1. Not Tested State Approved Alternate Assessment
 2. Not Tested State Approved Special Consideration
 3. Not Tested State Approved Enrolled After
 4. Not Tested State Approved Withdrew After
 5. Not Tested Other
- E. Special Circumstances by content area
1. Item invalidation flags are provided to the DOE during data processing test clean up. The item invalidation flag variables are initially set using the rules below. The final values used for reporting are provided back to Measured Progress by the DOE and used in reporting..
 - a. If `sciaccomM1` is marked, then mark `sciInvSes3`.
 - b. If `sciaccomM3` is marked, then mark `sciInvSes1`, `sciInvSes2`, and `sciInvSes3`.

2. A student is identified as content area tested if the student does not have any content area not tested reasons identified. Tested students are categorized in one of the four tested participation statuses: “Tested Damaged SRB”, “Tested with Non-Standard Accommodations”, “Tested Incomplete”, and “Tested”.
 - a. Students with a common item response of ‘X’ are identified as “Tested Damaged SRB”.
 - b. Students identified as content area tested, are not identified as “Tested Damaged SRB”, and have at least one of the content area invalidation session flags marked will be identified as “Tested with Non-Standard Accommodations”.
 - c. Students identified as content area tested, are not identified as “Tested Damaged SRB”, and not identified as “Tested with Non-Standard Accommodations” and did not attempt all sessions in the test are considered to be “Tested Incomplete.”
 - d. All other tested students are identified as “Tested”.
 3. For students identified as “Tested Damaged SRB”, the content area subcategories with at least one damaged item will not be reported. The school and district averages will be suppressed for the impacted subcategories on the student report. These students are excluded from all raw score aggregations (item, subcategory, and total raw score). They are included in participation, achievement level, and scaled score aggregations.
 4. For students identified as “Tested with Non-Standard Accommodations” the content area sessions item responses which are marked for invalidation will be treated as a non-response
 5. Students identified as tested in a content area will receive released item scores, scaled score, scale score bounds, achievement level, raw total score and subcategory scores.
 6. Students identified as not tested in a content area will not receive a scaled score, scaled score bounds, achievement level, writing annotations (where applicable). They will receive released item scores, raw total score, and subcategory scores.
 7. Item scores for students with an invalidation flag marked and have a not tested status will be blanked out based on the invalidation flag. For example, if the student is identified as “Not Tested: State Approved Alternate Assessment” and has sciInvSes3 marked, then all science session 1 item responses will be reported as a blank.
- F. Student Participation Status Hierarchy by content area
1. Not Tested: State Approved Alternate Assessment
 2. Not Tested: State Approved Special Consideration
 3. Not Tested: State Approved Enrolled After
 4. Not Tested: State Approved Withdrew After
 5. Not Tested: Other
 6. Tested Damaged SRB
 7. Tested with Non-Standard Accommodations
 8. Tested Incomplete
 9. Tested

G. Student Participation Summary

Participation Status	Description	Raw Score (*)	Scaled Score (&)	Ach. Level	Student Report Text	Ach. Level	Roster Ach. Level Text
Z	Tested Damaged SRB(**)	✓	✓	✓	Substantially Below Proficient, Partially Proficient, Proficient, or Proficient with Distinction		1,2,3, or 4
A	Tested	✓	✓	✓	Substantially Below Proficient, Partially Proficient, Proficient, or Proficient with Distinction		1,2,3, or 4
B	Tested Incomplete(%)	✓	✓	✓	Substantially Below Proficient, Partially Proficient, Proficient, or Proficient with Distinction		1,2,3, or 4
C	Tested with Non-Standard Accommodations (%%)	✓	✓	✓	Substantially Below Proficient, Partially Proficient, Proficient, or Proficient with Distinction		1,2,3, or 4
D	Not Tested State Approved Alternate Assessment	✓			Alternate Assessment		A
E	Not Tested State Approved First Year LEP (Reading and Writing only)	✓			First Year LEP		L
F	Not Tested State Approved Enrolled After	✓			Spring: Enrolled After May 5, 2014		E
G	Not Tested State Approved Withdrew After	✓			Spring: Withdrew After May 5, 2014		W
H	Not Tested State Approved Special Consideration	✓			Special Consideration		S
I	Not Tested Other	✓			Not Tested		N

If a student has a participation status of Alternate Assessment for all subjects assessed at the grade level, a Parent Letter is not produced.

(*) Raw scores are not printed on student report for students with a not tested status.

(**) Raw scores for Tested damaged SRB students will be reported based on the set of non-damaged items. Subcategory scores will not be reported if it includes a damaged item.

(%) Tested incomplete students will be identified on the student report with a footnote.

(%%) Tested with Non-standard accommodations students will be identified on student report with a footnote. The invalidated items will be stored as a '-' for item analysis.

III. Calculations

A. Rounding

1. All percents are rounded to the nearest whole number
2. All mean scaled scores are rounded to the nearest whole number
3. All mean raw scores are rounded to the nearest tenth.

4. Content Area Subcategories: Average Points Earned (student report): round to the nearest tenth.
 5. Round non-multiple choice average item scores to the nearest tenth.
- B. Students included in calculations based on participation status
1. For number and percent of students enrolled, tested, and not tested categories include all students not excluded by other decision rules.
 2. For number and percent at each achievement level, average scaled score, subtopic percent of total possible points and standard error, subtopic distribution across writing prompts, subcategories average points earned, percent/correct average score for each released item include all tested students not excluded by other decision rules.
 3. Students identified as Tested Damaged SRB are excluded from all raw score aggregations (item, subcategory, and total raw score). They are included in participation, achievement level, and scaled score aggregations.
- C. Raw scores
1. For all analyses, non-response for an item by a tested student is treated as a score of 0. Items identified as damaged (response of 'X') will be excluded for student identified as "Tested Damaged SRB".
 2. Content Area Total Points: Sum the points earned by the student for the common items.
- D. Item Scores
1. For all analysis, non-response for an item by a tested student is treated as a score of 0.
 2. For multiple choice released item data store a '+' for correct response, or A,B,C,D,* or blank
 3. For open response released items, store the student score. If the score is not numeric ('B'), then store it as blank.
 4. For students identified as content area tested with non-standard accommodations, then store the released item score as '-' for invalidated items.
 5. For all writing prompt scores, the final score of record is the sum of scorer 1 and scorer 2. If both scorers give the student a B, then the final score is B. If both scorers give the student an O or F, then the final score is 0.
- E. Scaling
1. Scale Form creation

Scaling is accomplished by defining the unique set of test forms for the grade/subject. This is accomplished as follows:

 - a. Translate each form and position into the unique item number assigned to the form/position.
 - b. Order the items by
 - I.* Type – multiple-choice, short-answer, constructed- response, extended-response, writing prompt.
 - II.* Form – common, then by ascending form number.
 - III.* Position
 - c. If an item number is on a form, then set the value for that item number to '1', otherwise set to '.'. Set the Exception field to '0' to indicate this is an original test form.
 - d. If an item number contains an 'X' (item is not included in scaling) then set the item number to '.'. Set the Exception field to '1' to indicate this is not an original test form.

- e. Compress all of the item numbers together into one field in the order defined in step II to create the test for the student.
 - f. Select the distinct set of tests from the student data and order them by the exception field and the descending test field.
 - g. Check to see if the test has already been assigned a scale form by looking in the tblScaleForm table. If the test exists then assign the existing scale form. Otherwise assign the next available scale form number. All scale form numbering starts at 01 and increments by 1 up to 99.
2. Scaled Score assignment
- a. Psychometrics provides data analysis with a lookup table for each scale form. The lookup table contains the raw score and the resulting scaled score.
- F. SubTopic Item Scores
1. Identify the Subtopic
- a. Spring: NECAP science item information is stored in IABS, including inquiry items.
 - I. Program management provided Data Analysis with “IABS Export Codes for NECAP SCI Reporting.doc” which contains the crosswalk between IABS item information and reporting.
 - II. Program management provided Data Analysis with “2010 IABS_Released ItemsSCI for Tara.xls” which contains released item order. Inquiry items are listed at the end in the order they are in the test booklet.
2. Student Content Area Subcategories (student report): Subtopic item scores at the student level is the sum of the points earned by the student for the common items in the subtopic.
3. Content Area Subtopic (grade level results report): Subtopic scores are based on all unique common and matrix items.
- a. Percent of Total Possible Points:
 - I. For each unique common and matrix item calculate the average student score as follows: (sum student item score/number of tested students administered the item).
 - II. $100 * (\text{Sum the average score for items in the subtopic}) / (\text{Total Possible Points for the subtopic})$ rounded to the nearest whole number.
 - b. Standard Error Bar: Before multiplying by 100 and rounding the Percent of Total Possible points (ppe) calculate standard error for school, district and state: $100 * (\text{square root } ((\text{ppe}) * (1 - \text{ppe}) / \text{number of tested students}))$ rounded to the nearest tenth. For the lower bound and upper bound round the Percent of Total Possible Points +/- Rounded Standard Error to the nearest hundredth.
- G. Cumulative Total
- 1. Include the yearly results where the number tested is greater than or equal to 10
 - 2. Cumulative total N (Enrolled, Not Tested Approved, Not Tested Other, Tested, at each achievement level) is the sum of the yearly results for each category where the number tested is greater than or equal to 10.
 - 3. Cumulative percent for each achievement level is $100 * (\text{Number of students at the achievement level cumulative total} / \text{number of students tested cumulative total})$ rounded to the nearest whole number.
 - 4. Cumulative mean scaled score is a weighted average. For years where the number tested is greater than or equal to 10, $(\text{sum of (yearly number tested * yearly mean scaled score)}) / (\text{sum of yearly number tested})$ rounded to the nearest whole number.
- H. Participation

1. For participation calculate the number and percent of students in each of the following categories by school, district, and state according to schtype and stustatus decision rules.
2. Note that a student is tested with approved accommodations if one is tested, has a non-M accommodation marked, and does not have the M2 or M3 accommodation marked for that subject.
 - a. For Students Enrolled, Students Tested, and Students Not Tested the denominator will be the number of students enrolled
 - b. For Students Tested with approved Accommodations, Current LEP Students Tested (LEP=1), and IEP Students Tested the denominator will be the number of students tested.
 - c. For Current LEP Students Tested with approved accommodations (LEP=1 the denominator will be the number of current LEP students tested.
 - d. For IEP Students Tested with approved accommodations the denominator will be the number of IEP students tested.
 - e. For Students Not Tested State Approved and Not Tested Other the denominator will be the number of students not tested.
 - f. For Students Not Tested Alternate Assessment, First Year LEP, Withdrew After October 1, Enrolled After October 1, and Special Considerations the denominator will be the number of students not tested state approved.
- I. Average Points Earned Students at Proficient Level (Range)
 1. Select all students across the states with Y40 scaled score, where Y=grade. Average the content area subcategories across the students. Add and subtract one standard error of measurement to get the range and round to the nearest tenth.

IV. Report Specific Rules

A. Student Report

1. Student header Information
 - a. If “FNAME” or “LNAME” is not missing then print “FNAME MI LNAME”. Otherwise, print “No Name Provided”.
 - b. Print the student’s single digit tested grade
 - c. For school and district name do the following.
 - I. For students with a stustatus value of 0 or 4, print the abbreviated tested school and district ICORE name based on school type decision rules.
 - II. Otherwise, for the school and district names print the “Description” in the StuStatus table presented earlier in this document.
 - d. Print “ME”, “NH”, “RI”, or “VT” for state.
2. Test Results by content area
 - a. Always display the cut scores in the graphic display.
 - a. For students identified as “Not Tested”, print the not tested reason in the achievement level, leave scaled score and graphic display blank.
 - b. For students identified as tested for the content area then do the following
 - I. Print the complete achievement level name the student earned
 - II. Print the scaled score the student earned
 - III. Print a vertical black bar for the student scaled score with gray horizontal bounds in the graphic display

- IV. For students identified as “Tested with a non-standard accommodation” for a content area, print ‘**’ after the content area earned achievement level and after student points earned for each subcategory.
 - V. For students identified as “Tested Incomplete” for a content area, place a section symbol after content area earned scaled score.
3. This Student’s Achievement Compared to Other Students by content area
 - a. For tested students, print a check mark in the appropriate achievement level in the content area student column. For not tested students leave student column blank
 - b. For percent of students with achievement level by school, district, and state print aggregate data based on student status, StuGrade, school type and minimum N rules.
 4. This Student’s Performance in Content Area Subcategories by content area
 - a. Always print total possible points and students at proficient average points earned range.
 - b. For students identified as not tested then leave student scores blank
 - c. For students identified as tested do the following
 - I. Print school, district, and state aggregate data for subcategories based on student status, StuGrade, school type and minimum N rules.
 - II. For students identified as “Tested Damaged SRB” do not report student, school, and district aggregate data for subcategories that have at least one damaged item. Print Points Possible and state aggregate data.
 - III. Otherwise, always print student subcategory scores
 - IV. If the student is identified as tested with a non-standard accommodation for the content area then place ‘**’ after the student points earned for each subcategory.
 5. Footer information
 - a. Footnotes
 - I. If the student received a participation status of “Tested with a non-standard accommodation” for any content area then print “** Student received no credit for parts of the test that were administered under non-standard conditions.”
 - II. If the student received a participations status of “Tested Incomplete” for any content area then print “§ This score should be viewed with caution because the student did not complete all parts of the test.”
 - III. If both footnotes should appear, the print I. above II.
 - b. For NH the SAU, district, and school codes should appear at the bottom right of the page separated by ‘-’.
 - c. For ME, RI, and VT district and school codes should appear at the bottom right of the page separated by ‘-’.
- B. Grade Level School/District/State Results
1. Reports are run by testing state, testing district, testing school using the aggregate school and district codes described in the school type table.
 2. Exclude students based on stugrade=12, student status, school type and participation status decision rules for aggregations.
 3. The reports will be collated as follows:
 - a. Page 1 is the Title page.
 - b. Page 2 is the Participation Results
 - c. Page 3 is the Historical and Subtopic Results

- d. Page 4 is the Disaggregated Results
4. Report Header Information
 - a. “Spring YYYY Grade XX NECAP Science Tests” where XX is the single digit grade level and YYYY is the year, will print as the title.
 - b. Teaching level reports will have the following subtitle: “Grade XX-1 Students in (YYYY-1)-(YYYY)”.
 - c. Testing level reports will have the following subtitle: “Grade XX Students in (YYYY)-(YYYY+1)”.
 - d. Use abbreviated school and district name from ICORE based on school type decision rules.
 - e. Print “Maine”, “New Hampshire”, “Rhode Island”, or “Vermont” to reference the state. The state graphic is printed on the first page.
 - f. For NH print SAU, district, and school codes separated by ‘-’ for Code on first page for school level. Print SAU and district codes separated by ‘-’ for the district level. Print the full state name for the state level.
 - g. For ME, RI, and VT print district and school codes separated by ‘-’ for Code on first page for the school level. Print the district code for the district level. Print the full state name for the state level.
 5. For achievement level and participation category data if the number of students in an achievement level or participation category does not equal 0, and the percent of students is 0 then format the percent as <1.
 6. Report Section: Participation in NECAP
 - a. For testing level reports always print number and percent based on school type decision rules.
 - b. For the teaching level reports leave the section blank.
 7. Report Section: NECAP Results by content area
 - a. For the testing level report always print based on minimum N-size and school type decision rules.
 - b. For the teaching level report leave Enrolled, NT Approved, and NT Other blank. Print Tested, number and percent at each achievement level, mean scaled score based on minimum N-size and school type decision rules.
 8. Report Section: Historical NECAP Results by content area
 - a. For tested level report always print current year, prior years, and cumulative total results based on minimum N-size and school type decision rules.
 - b. For teaching level report leave Enrolled, NT Approved, and NT Other blank. Print Tested, number and percent at each achievement level, mean scores based on minimum N-size and school type decision rules.
 - c. Bold current year data.
 9. Report Section: Subtopic Results by content area
 - a. For testing and teaching level reports always print based on minimum N-size and school type decision rules
 10. Report Section: Disaggregated Results by content area
 - a. For testing level report always print based on minimum N-size and school type decision rules.

- b. For teaching level report leave Enrolled, NT Approved, and NT Other blank. Print Tested, number and percent at each achievement level, mean scores based on minimum N-size and school type decision rules.
- C. School/District/State Summary(School Level is run in the Fall Only)
1. Report Header Information
 - a. Use abbreviated school and district name from ICORE based on school type decision rules.
 - b. Print “Maine”, “New Hampshire”, “Rhode Island”, or “Vermont” to reference the state.
 - c. For NH print SAU, district, and school codes separated by ‘-‘ for Code on first page for school level. Print SAU and district codes separated by ‘-‘for the district level. Print the full state name for the state level.
 - d. For ME, RI, and VT print district and school codes separated by ‘-‘ for Code on first page for the school level. Print the district code for the district level. Print the full state name for the state level.
 2. Reports are run by testing state, testing district, testing school (Fall Only) using the aggregate school and district codes described in the school type table
 3. Exclude students based on StuGrade=12, student status, school type and participation status decision rules for aggregations.
 4. For achievement level and participation category data if the number of students in an achievement level or participation category does not equal 0, and the percent of students is 0 then format the percent as <1.
 5. For testing level report print entire aggregate group across grades tested and list grades tested results based on minimum N-size and school type decision rules. Mean scores across the grades is not calculated.
 6. For the teaching level report leave Enrolled, NT Approved, and NT Other blank. Print Tested, number and percent at each achievement level, mean scaled score based on minimum N-size and school type decision rules. Mean scores across the grades is not calculated.
 7. Printed Grade Column
 - a. For the all grades row, display the school, district, or state name.
 - b. For grades 3-8 and 11 rows print Grade X (no leading zero).

V. Data Requirements Interactive Reporting

- A. Student Level
1. Refer to Sections II and III. D for decision rules on how student test data will be stored.
 2. Students will be loaded into the Interactive System based off of the Interactive flag in tblStuDemo. Students with Interactive flag set to 0 will not be loaded into the system. Students with Interactive set to 1 will be loaded.
 - a. Students with StuStatus value of 1, 2 or 3 or RI StuGrade=12 will have the Interactive flag set to 0.
 - b. All others will have Interactive=1.
 3. The Included flag will determine which students are included in school level aggregations. Students with Included=0 are excluded from all aggregations. Students with Included=2 will be included in Performance Level aggregations and excluded from raw score aggregations (item, subcategory, and total raw score). Students with Included=1 will be included in all school level aggregations.

- a. Students with a Not Tested Participation Status, StuStatus=1, 2, or 3, or RI StuGrade=12 will have their Included flag set to 0.
 - b. Students who do fall into the above group and have Participation Status of Tested Damaged SRB will have their Included flag set to 2.
 - c. All other students will have their Included flag set to 1.
4. Longitudinal Data
- a. Only students with a valid StudentID and Interactive flag=1 will be loaded.
 - b. The complete achievement level name or not tested reason will be stored .
- B. Aggregate Level
- 1. Data Analysis will compute Item Averages for the whole group only at the testing and teaching (Fall only) School and District Levels.
 - 2. Data Analysis will compute Item Averages for all of the filter combinations that exist at the State Level.
 - 3. Data Analysis will create a lookup table with all of the possible filter combinations. It will contain the variable Filter with length 5. Each position represents one of the filter variables. It will contain all the possible combinations of the values plus nulls for when variables are not selected. The first position will be Gender, second Ethnic, third IEP, fourth LEP, and fifth EconDis.
 - 4. Data Analysis will compute Item Averages, Achievement Level Summary, and Item Summary data for the filter combinations for a sample of schools for quality assurance review.
 - a. For this sample, percents will be rounded to the nearest whole number and open response average scores will be rounded to the nearest tenth.
 - b. For the Item Summary data, item responses other than A, B, C, and D will be counted in the IR column.

VI. Data File Rules

In the file names GR refers to the two digit grade (03-08, 11), YYYY refers to the year, DDDDD refers to the district code, and SS refers to two letter state code. Refer to the tables at the end of this section for filenames and layouts. Teaching level data files will be produced in the Fall Only.

- A. Preliminary State Results
- 1. A PDF file will be created for each state containing preliminary state results for each grade and subject and will list historical state data for comparison.
 - 2. The file name will be SSPreliminaryResultsDATE.pdf
- B. State Student Released Item Data
- 1. A CSV file will be created for each state for grades 3-8 and one for grade 11.
 - 2. One CSV file will be created for each state in the Spring.
 - 3. Accommodation Flags
 - a. If the student has at least 1 standard accommodation marked (excluding M) for a given subject then set [sub]STDaccom flag to '1'. Otherwise set it to '0'.
 - b. For each group of accommodations (S, T, P, R, and O) if a student has any accommodation in that group marked set [sub]Accom[group]='1'. Otherwise set it to '0'.
 - c. If a student has the M2 accommodation marked, then set [sub]AccomM2='1'. Otherwise set it to '0'.

- d. If a student has the M3 accommodation marked, then set [sub]AccomM3='1'. Otherwise set it to '0'.
4. Exclusion Rules
- a. NH: If the student has a testing school type of 'PRI' or StuStatus is 1,2 ,3,or 4 then exclude the student
 - b. RI: If testing school type is PRI and teaching school type is PRI or blank, then exclude the student.
 - c. VT: Do not exclude any students
 - d. ME: If the student has a StuStatus is 1,2 ,3,or 4 then exclude the student
- C. State Student Raw Data
- 1. A CSV file will be created for each state by grade span. The grade spans are 3-4, 5-8, and 11. In the spring, all grades will be combined.
 - 2. If the student has at least 1 standard accommodation marked (excluding M) for a given subject then set [sub]STDaccom flag to '1'. Otherwise set it to '0'.
 - 3. Exclusion Rules
 - a. NH: If the student has a testing school type of 'PRI' or StuStatus is 1,2 ,3,or 4 then exclude the student
 - b. RI: If testing school type is PRI and teaching school type is PRI or blank, then exclude the student.
 - c. VT: Do not exclude any students
 - d. ME: If the student has a StuStatus is 1,2 ,3,or 4 then exclude the student.
- D. State Student Scored Data
- 1. A CSV file will be created for each state including all grades.
 - 2. Exclusion Rules
 - a. NH: If the student has a testing school type of 'PRI' or StuStatus is 1,2 ,3,or 4 then exclude the student
 - b. RI: If testing school type is PRI and teaching school type is PRI or blank, then exclude the student.
 - c. VT: Do not exclude any students
 - d. ME: If the student has a StuStatus is 1,2 ,3,or 4 then exclude the student.
- E. District Student Data
- 1. Testing and teaching CSV files will be created for each state and grade and district.
 - 2. Students with the Discode or SendDiscode will be in the district grade specific CSV file for the testing year.
 - 3. For ME, NH, and RI only public school districts will receive district data files. (Districts with at least one school with schoolsubtypeID=1, 11, 19, or 15 in ICORE)
 - 4. Accommodation Flags
 - a. If the student has at least 1 standard accommodation marked (excluding M) for a given subject then set [sub]STDaccom flag to '1'. Otherwise set it to '0'.
 - b. For each group of accommodations (S, T, P, R, and O) if a student has any accommodation in that group marked set [sub]Accom[group]='1'. Otherwise set it to '0'.

- c. If a student has the M2 accommodation marked, then set [sub]AccomM2='1'. Otherwise set it to '0'.
 - d. If a student has the M3 accommodation marked, then set [sub]AccomM3='1'. Otherwise set it to '0'.
 - 5. Exclusion Rules
 - a. NH & RI: If the student has a StuStatus value of 1,2, or 3 then exclude the student
 - b. VT: If the student has a StuStatus value of 1, then exclude the student.
 - c. ME: If the student has a StuStatus is 1, 2, or 3 then exclude the student.
- F. School Student Data
 - 1. Testing and teaching CSV files will be created for each state and grade and school.
 - 2. Students with the SchCode will be in the school grade specific CSV file for the testing year.
 - 3. Accommodation Flags
 - a. If the student has at least 1 standard accommodation marked (excluding M) for a given subject then set [sub]STDaccom flag to '1'. Otherwise set it to '0'.
 - b. For each group of accommodations (S, T, P, R, and O) if a student has any accommodation in that group marked set [sub]Accom[group]='1'. Otherwise set it to '0'.
 - c. If a student has the M2 accommodation marked, then set [sub]AccomM2='1'. Otherwise set it to '0'.
 - d. If a student has the M3 accommodation marked, then set [sub]AccomM3='1'. Otherwise set it to '0'.
 - 4. Exclusion Rules
 - a. NH & RI: If the student has a StuStatus value of 1,2 or 3, then exclude the student
 - b. VT: If the student has a StuStatus value of 1, then exclude the student.
 - c. ME: If the student has a StuStatus is 1, 2, or 3 then exclude the student.
- G. Common Item Information
 - 1. An excel file will be created containing item information for common items: grade, subject, released item number, item analysis heading data, raw data item name, item type, key, and point value.
- H. State Standard Deviations and Averages Scaled Scores
 - 1. A csv file will be created for each state containing the standard deviations and average scale scores for disaggregated subgroups by subject.
 - 2. Exclude students based on state aggregation StuGrade, StuStatus, and SchType decision rules.
 - 3. Data will be suppressed based on minimum N-size and report type decision rules.
 - 4. Average scaled score will be rounded to the nearest whole number. Standard deviations will be rounded to the nearest tenth.
- I. Grade Level Results Report Disaggregated and Historical Data
 - 1. Teaching and testing CSV files will be created for each state containing the grade level results disaggregated and historical data.
 - 2. Data will be suppressed based on minimum N-size and report type decision rules.
 - 3. Private schools are excluded from NH & RI files.
- J. Grade Level Results Report Participation Category Data

1. Testing CSV file will be created for each state containing the grade level results participation data.
 2. Private schools are excluded from NH & RI files.
- K. Grade Level Results Report Subtopic Data
1. Teaching and testing CSV files will be created for each state containing the grade level results subtopic.
 2. Data will be suppressed based on minimum N-size and report type decision rules.
 3. Private schools are excluded from NH & RI files.
- L. Summary Results Data
1. Teaching and testing CSV files will be created for each state containing the school, district and state summary data.
 2. Data will be suppressed based on minimum N-size and report type decision rules.
 3. Private schools are excluded from NH & RI files.
- M. Released Item Percent Responses Data
1. The CSV files will only contain state level aggregation for released items.
 2. CSV files will be created for each state and grade containing the released item analysis report state data.
- N. Invalidated Students Original Score
1. A CSV file will be created for each state including all grades.
 2. Original raw scores for students whose responses were invalidated for reporting will be provided.
 3. Exclusion Rules
 - a. NH: If the student has a testing school type of 'PRI' or StuStatus is 1,2, 3, or 4 then exclude the student
 - b. RI: If testing school type is PRI and teaching school type is PRI or blank, then exclude the student.
 - c. VT: Do not exclude any students
 - d. ME: If the student has a StuStatus is 1,2 ,3,or 4 then exclude the student.
- O. Student Questionnaire Summary
1. One CSV file will be created for each state containing percent of students at each response, percent of students at each achievement level, and average scaled score, by student questionnaire response.
 2. Only include students who are included in state level aggregations.
 3. Data will be suppressed based on minimum N-size and report type decision rules.
- P. TCTA Questionnaire Raw Data
1. One CSV file will be created for each state containing raw TC Questionnaire data.
 2. One CSV file will be created for each state containing raw TA Questionnaire data.
- Q. TCTA Questionnaire Frequency Distribution
1. One CSV file will be created for each state containing the distribution of responses of TC Questionnaire raw data.

2. One CSV file will be created for each state containing the distribution of responses of TA Questionnaire raw data.
- R. Scaled Score Lookup
1. One CSV file and one excel file will be created containing the scaled score lookup data.
- S. Subtopic Average Points Earned (For Program Management)
1. One excel file will be created containing four worksheets. The first worksheet contains the total possible points for each subtopic as reported on the item analysis report and the range for students who are just proficient. The remaining three worksheets contain state average subtopic scores as reported on the item analysis report.
 2. Program management uses this file to create a document which is provided to the schools.
- T. Item Stats for Inquiry Task Items (For Program Management)
1. Since Inquiry Task Items are not stored in IABS, one CSV file will be created containing item stats for Inquiry Task items.
 2. All three states are included in the calculations.
- U. Memo Shipping Files (For Program Management)
1. Provide PM in excel list of schools and districts that tested regardless of grade.
- V. CD Print File
- A. Spring Table Data File Deliverables

Data File	Layout	File Name
Preliminary State Results	N/A	Included in Equating Report
State Student Released Item Data	NECAP1314SpringStudentReleasedItemLayout.xls	NECAP1314SpringStateStudentReleasedItem.csv
State Student Raw Data	NECAP1314SpringStateStudentRawLayout.xls	NECAP1314SpringStateStudentRaw.csv
State Student Scored Data	NECAP1314SpringStateStudentScoredLayout.xls	NECAP1314SpringStateStudentScored.csv
District Student Data	NECAP1314SpringStudentReleasedItemLayout.xls	NECAP1314SpringDistrictSlice[GR]_[District Code].csv
School Student Data	NECAP1314SpringStudentReleasedItemLayout.xls	NECAP1314SpringSchoolSlice[GR]_[District Code][School Code].csv
Common Item Information	NECAP1314SpringCommonItemInformationLayout.xls	NECAP1314SpringCommonItemInformation.csv
State Standard Deviations and Average Scaled Scores	NECAP1314SpringStateStandardDeviationsLayout.xls	NECAP1314SpringStateStandardDeviations.csv
Grade Level Results Report Disaggregated and Historical Data	NECAP1314SpringResultsReport DisaggregatedandHistoricalLayout.xls	NECAP1314SpringResultsReportDisaggregatedandHistorical.csv

Grade Level Results Report Participation Category Data	NECAP1314SpringResultsReport ParticipationLayout.xls	NECAP1314SpringResultsReportParticipation.csv
Grade Level Results Report Subtopic Data	NECAP1314SpringResultsReport SubtopicLayout.xls	NECAP1314SpringResultsReportSubtopic.csv
Summary Results Data	NECAP1314SpringSummaryResultsLayout.xls	NECAP1314SpringSummaryResults.csv
Released Item Percent Responses Data	NECAP1314SpringReleasedItemPercentResponsesLayout.xls	NECAP1314SpringReleasedItemPercentResponses.csv
Invalidated Students Original Score	NECAP1314SpringStateInvalidatedStudent OriginalScoredLayout.xls	NECAP1314SpringStateInvalidatedStudent OriginalScored.csv
Student Questionnaire Summary	NECAP1314SpringStudentQuestionnaireSummaryLayout.xls	NECAP1314SpringStudentQuestionnaireSummary.csv
TCTA Questionnaire Raw Data	NECAP1314SpringTCQuestionnaireRawLayout.xls NECAP1314SpringTAQuestionnaireRawLayout.xls	NECAP1314SpringTCQuestionnaireRaw.csv NECAP1314SpringTAQuestionnaireRaw.csv
TCTA Questionnaire Frequency Distribution	NECAP1314SpringTCTAQuestionnaireFreqLayout.xls	NECAP1314SpringTCTAQuestionnaireFreq.csv
Scaled Score Lookup	NECAP1314SpringScaleScoreLookupLayout.xls	NECAP1314SpringScaleScoreLookup.xls NECAP1314SpringScaleScoreLookup.csv
Subtopic Average Points Earned (For Project Management)	N/A	NECAP1314SpringSubtopicAvgPointsEarned.xls
Item Stats for Inquiry Task Items (For Program Management)	N/A	NECAP1314SpringInquiryItemStats.csv
Memo Shipping Files (For Program Management)	N/A	TBD