



Spring 2015

Nebraska State Accountability (NeSA)

Reading, Mathematics, and Science

Technical Report

September 2015

Prepared by Data Recognition Corporation





2015 NEBRASKA STATE ACCOUNTABILITY (NeSA) TECHNICAL REPORT

TABLE OF CONTENTS

1. BACKGROUND

1.1. Purpose and Organization of This Report	1
1.2. Background of the Nebraska State Accountability (NeSA)	1
• Previous Nebraska Assessment (STARS)	
• Purpose of the NeSA	
• Phase-In Schedule for NeSA	
• Advisory Committees	

2. ITEM AND TEST DEVELOPMENT

2.1. Content Standards.....	3
2.2. Test Blueprints.....	4
2.3. Multiple-Choice Items.....	4
2.4. Passage Selection.....	4
2.5. Item Development and Review	5
• Item Writer Training	
• Item Writing	
• Item Review	
• Editorial Review of Items	
• Review of the Online Items	
• Universally Designed Assessments	
• Depth of Knowledge	
• Item Content Review	
• Sensitivity and Bias Review	
2.6. Item Banking	16
2.7. The Operational Forms Construction Process.....	16
• DIF in Operational Forms Construction	
• Review of the Items and Test Forms	
2.8. Reading Assessment.....	18
• Test Design	
• Psychometric Targets	
• Equating Design	
2.9. Mathematics Assessment.....	19
• Test Design	
• Psychometric Targets	

- Equating Design

2.10 Science Assessment20

- Test Design
- Psychometric Targets
- Equating Design

3. STUDENT DEMOGRAPHICS

3.1. Demographics and Accommodations22

3.2. Students Tested and Mode Summary Data31

3.3. Testing Time.....32

4. CLASSICAL ITEM STATISTICS

4.1. Item Difficulty35

4.2. Item-Total Correlation.....36

4.3. Percent Selecting Each Response Option.....38

4.4. Point-Biserial Correlations of Response Options.....38

4.5. Percent of Students Omitting an Item38

5. RASCH ITEM CALIBRATION

5.1. Description of the Rasch Model39

5.2. Checking Rasch Assumptions39

- Unidimensionality
- Local Independence
- Item Fit

5.3. Rasch Item Statistics.....48

6. EQUATING AND SCALING

6.1. Equating.....50

6.2. Scaling52

7. FIELD TEST ITEM DATA SUMMARY

7.1. Classical Item Statistics57

7.2. Differential Item Functioning.....59

8. RELIABILITY

8.1. Coefficient Alpha66

8.2. Standard Error of Measurement67

8.3. Conditional Standard Error of Measurement (CSEM).....68

8.4. Decision Consistency and Accuracy69

9. VALIDITY

9.1. Evidence Based on Test Content.....72

9.2. Evidence Based on Internal Structure73

- Item-Test Correlation
- Item Response Theory Dimensionality
- Strand Correlations

9.3. Evidence Related to the Use of the Rasch Model78

10. REFERENCES 79

11. APPENDICES

A. NeSA-R Test Blueprint.....83

B. NeSA-M Test Blueprint.....97

C. NeSA-S Test Blueprint.....127

D. Confidentiality Agreement144

E. Fairness in Testing Manual.....145

F. Reading Key Verification and Foil Analysis.....162

G. Mathematics Key Verification and Foil Analysis180

H. Science Key Verification and Foil Analysis.....200

I. Overview of Rasch Measurement.....209

J. Reading, Mathematics, and Science Operational Form Calibration Summaries.....213

K. Reading Item Bank Difficulties220

L. Mathematics Item Bank Difficulties.....233

M. Science Item Bank Difficulties247

N. Reading Pre- and Post-Equating Summary253

O. Mathematics Pre- and Post-Equating Summary259

P. Science Pre- and Post-Equating Summary265

Q. Reading Raw-to-Scale Conversion Tables and Distributions of Ability.....269

R. Mathematics Raw-to-Scale Conversion Tables and Distributions of Ability282

S. Science Raw-to-Scale Conversion Tables and Distributions of Ability.....297

T. Reading Field Test Differential Item Functioning.....302

U. Mathematics Field Test Differential Item Functioning316

V. Science Field Test Differential Item Functioning330

W. Reading, Mathematics, and Science Analysis and Demographic Summary Sheets.....336

X. Reading, Mathematics, and Science Strand Reliability and SEM.....353



1. BACKGROUND

1.1 PURPOSE AND ORGANIZATION OF THIS REPORT

This report documents the technical aspects of the 2015 Nebraska State Accountability Reading (NeSA-R), Mathematics (NeSA-M), and Nebraska Science (NeSA-S) operational tests, along with the NeSA-R, NeSA-M and NeSA-S embedded field tests, covering details of item and test development processes, administration procedures, and psychometric methods and summaries.

1.2 BACKGROUND OF THE NEBRASKA STATE ACCOUNTABILITY (NE SA)

Previous Nebraska Assessments: In previous years, Nebraska administered a blend of local and state-generated assessments to meet No Child Left Behind (NCLB) requirements called STARS (School-based Teacher-led Assessment and Reporting System). STARS was a decentralized local assessment system that measured academic content standards in reading, mathematics, and science. The state reviewed every local assessment system for compliance and technical quality. The Nebraska Department of Education (NDE) provided guidance and support for Nebraska educators by training them to develop and use classroom-based assessments. For accreditation, districts were also required to administer national norm-referenced tests (NRT).

As a component of STARS, the NDE administered one writing assessment annually in grades 4, 8, and 11. In addition, the NDE provided an alternate assessment for students severely challenged by cognitive disabilities.

Purpose of the NeSA: Legislative Bill 1157 passed by the 2008 Nebraska Legislature (<http://www.legislature.ne.gov/laws/statutes.php?statute=79-760.03>) required a single statewide assessment of the Nebraska academic content standards for reading, mathematics, science, and writing in Nebraska's K-12 public schools. The new assessment system was named NeSA (Nebraska State Accountability), with NeSA-R for reading assessments, NeSA-M for mathematics, NeSA-S for science, and NeSA-W for writing (Complete documentation of the technical details for NeSA-W are presented in a separate document labeled *NeSA 2015 Writing Test Technical Report*). The assessments in reading and mathematics were administered in grades 3-8 and 11; science was administered in grades 5, 8, and 11.

NeSA replaced previous school-based assessments for purposes of local, state, and federal accountability. The NeSA RMS consists entirely of multiple choice items and will be administered, to the extent practicable, online. In January 2009, the NDE contracted with Data Recognition Corporation (DRC) to support the Department of Education with the administration, record keeping, and reporting of statewide student assessment and accountability.

Phase-In Schedule for NeSA: The NDE prescribed such assessments starting in the 2009-2010 school year to be phased in as shown in Table 1.2.1. The state intends to use the expertise and experience of

in-state educators to participate, to the maximum extent possible, in the design and development of the new statewide assessment system.

Table 1.2.1: NeSA Administration Schedule

Subject	Administration Year		Grades
	Field Test	Operational	
Reading	2009	2010	3 through 8 plus high school
Mathematics	2010	2011	3 through 8 plus high school
Science	2011	2012	5, 8 and 11

Advisory Committees: Legislative Bill 1157 added a governor-appointed Technical Advisory Committee (TAC) with three nationally recognized experts in educational assessment, one Nebraska administrator, and one Nebraska teacher. The TAC reviewed the development plan for the NeSA, and provided technical advice, guidance, and research to help the NDE make informed decisions regarding standards, assessment, and accountability.

2. ITEM AND TEST DEVELOPMENT

2.1 CONTENT STANDARDS

In April of 2008, the Nebraska Legislature passed into state law Legislative Bill 1157. This action changed previous provisions related to standards, assessment, and reporting. Specific to standards, the legislation stated:

- The State Board of Education shall adopt measurable academic content standards for at least the grade levels required for statewide assessment. The standards shall cover the content areas of reading, writing, mathematics, science, and social studies. The standards adopted shall be sufficiently clear and measurable to be used for testing student performance with respect to mastery of the content described in the state standards.
- The State Board of Education shall develop a plan to review and update standards for each content area every five years.
- The State Board of Education shall review and update the standards in reading by July 1, 2009, the standards in mathematics by July 1, 2010, and these standards in all other content areas by July 1, 2013.

The Nebraska Language Arts Standards are the foundation for NeSA-R. This assessment instrument is comprised of items that address standards for grades 3–8 and 12. The standards are assessed at grade-level with the exception of grade 12. The grade 12 standards are assessed on the NeSA tests at grade 11. The reading standards for each grade are represented in items that are distributed between two reporting categories: Vocabulary and Comprehension. The Vocabulary standards include word structure, context clues, and semantic relationships. The Comprehension standards include author's purpose, elements of narrative text, literary devices, main idea, relevant details, text features, genre, and generating questions while reading.

The mathematics component of the NeSA is composed of items that address indicators in grades 3–8 and high school. The standards are assessed at grade level with the exception of high school. The high school standards are assessed on the NeSA-M at grade 11. The assessable standards for each grade level are distributed among the four reporting categories: Number Sense Concepts, Geometric/Masurement Concepts, Algebraic Concepts, and Data Analysis/Probability Concepts. The National Council of Teachers of Mathematics (NCTM) and the National Assessment of Educational Progress (NAEP) standards are the foundation of the Nebraska Mathematics standards.

The science component of the NeSA is composed of items that address indicators in grade-band strands 3–5, 6–8, and 9–12. The NeSA-S assesses the standards for each grade-band strand at a specific grade: 3–5 strand at grade 5, 6–8 strand at grade 8, and 9–12 strand at grade 11. The assessable standards for each grade level are distributed among the four reporting categories: Inquiry, The Nature of Science, and Technology; Physical Science; Life Science; and Earth and Space Sciences.

2.2 TEST BLUEPRINTS

The test blueprints for each assessment include lists of all the standards, organized by reporting categories. The test blueprints also contain the Depth of Knowledge (DOK) level assigned to each standard and the range of test items to be part of the assessment by indicator. The NeSA-R test blueprint was developed and approved in fall 2009 (Appendix A). The NeSA-M test blueprint was developed and approved in fall 2010 (Appendix B). The NeSA-S test blueprint was developed and approved in fall 2011 (Appendix C).

2.3 MULTIPLE-CHOICE ITEMS

Each assessment incorporates multiple-choice (MC) items to assess the content standards. Students are required to select a correct answer from four response choices with a single correct answer. Each MC item is scored as right or wrong and has a value of one raw score point. MC items are used to assess a variety of skill levels in relation to the tested standards.

2.4 PASSAGE SELECTION

All items in the reading assessment were derived from a selection of narrative and informational passages. Passages acquired were “authentic” in that they were purchased from the test vendor that commissioned experienced passage writers to provide quality pieces of text. Passages were approved by a group of reading content specialists that have teaching experience at specific grade levels. These experts were given formal training on the specific requirements of the Nebraska assessment of reading. The group, under the facilitation of the NDE test development team, screened and edited passages for:

- interest and accuracy of information in a passage to a particular grade level;
- grade-level appropriateness of passage topic and vocabulary;
- rich passage content to support the development of high-quality test questions;
- bias, sensitivity, and fairness issues; and
- readability considerations and concerns.

Passages that were approved moved forward for the development of test items.

The readability of a passage was an evaluative process made by Nebraska educators, the NDE’s test development team, DRC’s reading content specialists, and other individuals who understand each particular grade level and children of a particular age group. In addition, formal readability programs were also used by DRC to provide a “snapshot” of a passage’s reading difficulty based on sentence structure, length of words, etc. All of this information, along with the classroom context and content appropriateness of a passage, was taken into consideration when placing a passage at a particular grade.

2.5 ITEM DEVELOPMENT AND REVIEW

The most significant considerations in the item and test development process are: aligning the items to the grade level indicators; determining the grade-level appropriateness; DOK; estimated difficulty level; and determining style, accuracy, and correct terminology. In addition, the *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014) and *Universal Design* (Thompson, Johnstone, & Thurlow, 2002) guided the following steps in the item development process:

- Analyze the grade-level indicators and test blueprints.
- Analyze item specifications and style guides.
- Select qualified item writers.
- Develop item-writing workshop training materials.
- Train Nebraska educators to write items.
- Write items that match the standards, are free of bias, and address fairness and sensitivity concerns.
- Conduct and monitor internal item reviews and quality processes.
- Prepare passages (reading only) and items for review by a committee of Nebraska educators (content and bias/sensitivity).
- Select and assemble items for field testing.
- Field test items, score the items, and analyze the data.
- Review items and associated statistics after field testing, including bias statistics.
- Update item bank.

Item Writer Training: The test items were written by Nebraska educators who were recommended for the process by an administrator. Three criteria were considered in selecting the item writers: educational role, geographic location, and experience with item writing.

Prior to developing items for NeSA, a cadre of item writers was trained with regard to:

- Nebraska content standards and test blueprints;
- cognitive levels, including Depth of Knowledge (DOK);
- principles of Universal Design;
- skill-specific and balanced test items for the grade level;
- developmentally appropriate structure and content;
- item-writing technical quality issues;
- bias, fairness, and sensitivity issues; and
- style considerations and item specifications.

Item Writing: To ensure that all test items met the requirements of the approved target content test blueprint and were adequately distributed across subcategories and levels of difficulty, item writers were asked to document the following specific information as each item was written:

- **Alignment to the Nebraska Standards:** There must be a high degree of match between a particular question and the standard it is intended to measure. Item writers were asked to clearly indicate which standard each item was measuring.
- **Estimated Difficulty Level:** Prior to field testing items, the item difficulties were not known, and writers could only make approximations as to how difficult an item might be. The estimated difficulty level was based upon the writer's own judgment as directly related to his or her classroom teaching and knowledge of the curriculum for a given content area and grade level. The purpose for indicating estimated difficulty levels as items were written was to help ensure that the pool of items would include a range of difficulty (easy, medium, and challenging).
- **Appropriate Grade Level, Item Context, and Assumed Student Knowledge:** Item writers were asked to consider the conceptual and cognitive level of each item. They were asked to review each item to determine whether or not the item was measuring something that was important and could be successfully taught and learned in the classroom.
- **MC Item Options and Distractor Rationale:** Writers were instructed to make sure that each item had only one clearly correct answer. Item writers submitted the answer key with the item. All distractors were plausible choices that represented common errors and misconceptions in student reasoning.
- **Face Validity and Distribution of Items Based upon DOK:** Writers were asked to classify the DOK of each item, using a model based on Norman Webb's work on DOK (Webb, 2002). Items were classified as one of four DOK categories: recall (DOK Level 1), skill/concept (DOK Level 2), strategic thinking (DOK Level 3), and extended thinking (DOK Level 4).
- **Readability:** Writers were instructed to pay careful attention to the readability of each item to ensure that the focus was on the concepts; not on reading comprehension of the item. Resources writers used to verify the vocabulary level were the *EDL Core Vocabularies* (Taylor, Frackenpohl, White, Nieroroda, Browning, & Brisner, 1989) and the *Children's Writer's Word Book* (Mogilner, 1992). In addition, every test item was reviewed by grade-level experts. They reviewed each item from the perspective of the students they teach, and they determined the validity of the vocabulary used.
- **Grammar and Structure for Item Stems and Item Options:** All items were written to meet technical quality, including correct grammar, syntax, and usage in all items, as well as parallel construction and structure of text associated with each MC item.

Item Review: Throughout the item development process, independent panels of reading content experts reviewed the items. The following guidelines for reviewing assessment items were used during each review process.

A quality item should:

- have only one clear correct answer and contain answer choices that are reasonably parallel in length and structure;
- have a correctly assigned content code (item map);
- measure one main idea or problem;
- measure the objective or curriculum content standard it is designed to measure;
- be at the appropriate level of difficulty;
- be simple, direct, and free of ambiguity;
- make use of vocabulary and sentence structure that is appropriate to the grade level of the student being tested;
- be based on content that is accurate and current;
- when appropriate, contain stimulus material that are clear and concise and provide all information that is needed;
- when appropriate, contain graphics that are clearly labeled;
- contain answer choices that are plausible and reasonable in terms of the requirements of the question, as well as the students' level of knowledge;
- contain distractors that relate to the question and can be supported by a rationale;
- reflect current teaching and learning practices in the content area; and
- be free of gender, ethnic, cultural, socioeconomic, and regional stereotyping bias.

Following each review process, the item writer group and the item review panel discussed suggestions for revisions related to each item. Items were revised only when both groups agreed on the proposed change.

Editorial Review of Items: After items were written and reviewed, the NDE test development specialists reviewed each item for item quality, making sure that the test items were in compliance with guidelines for clarity, style, accuracy, and appropriateness for Nebraska students. Additionally, DRC test development content experts worked collaboratively with the NDE to review and revise the items prior to field testing to ensure highest level of quality possible.

Review of the Online Items: All items for online assessment were reviewed by the NDE and DRC. In addition to DRC's standard review process to which all items are subjected, and to ensure comparability with paper and pencil versions, all items were reviewed for formatting and scrolling concerns.

Universally Designed Assessments: Universally designed assessments allow participation of the widest possible range of students and result in valid inferences about performance of all students who participate and are based on the premise that each child in school is a part of the population to be tested, and that testing results should not be affected by disability, gender, race, or English language ability (Thompson, Johnstone, & Thurlow, 2002). The NDE and DRC are committed to the development of items and tests that are fair and valid for all students. At every stage of the item and

test development process, procedures ensure that items and tests are designed and developed using the elements of universally designed assessments that were developed by the National Center on Educational Outcomes (NCEO).

Federal legislation addresses the need for universally designed assessments. The *No Child Left Behind Act* (Elementary and Secondary Education Act) requires that each state must “provide for the participation in [statewide] assessments of all students” [Section 1111(b) (3) (C) (ix) (I)]. Both Title 1 and IDEA regulations call for universally designed assessments that are accessible and valid for all students including students with disabilities and students with limited English proficiency. The NDE and DRC recognize that the benefits of universally designed assessments not only apply to these groups of students, but to all individuals with wide-ranging characteristics.

The NDE test development team and Nebraska item writers have been fully trained in the elements of Universal Design as it relates to developing large-scale statewide assessments. Additionally, the NDE and DRC partner to ensure that all items meet the Universal Design requirements during the item review process.

After a review of research relevant to the assessment development process and the principles of Universal Design (Center for Universal Design, 1997), NCEO has produced seven elements of Universal Design as they apply to assessments (Thompson, Johnstone, & Thurlow, 2002).

Inclusive Assessment Population

When tests are first conceptualized, they need to be thought of in the context of who will be tested. If the test is designed for state, district, or school accountability purposes, the target population must include every student except those who will participate in accountability through an alternate assessment. The NDE and DRC are fully aware of increased demands that statewide assessment systems must include and be accountable for ALL students.

Precisely Defined Constructs

An important function of well-designed assessments is that they actually measure what they are intended to measure. The NDE item writers and DRC carefully examine what is to be tested and design items that offer the greatest opportunity for success within those constructs. Just as universally designed architecture removes physical, sensory, and cognitive barriers to all types of people in public and private structures, universally designed assessments must remove all non-construct-oriented cognitive, sensory, emotional, and physical barriers.

Accessible, Non-biased Items

The NDE conducts both internal and external review of items and test specifications to ensure that they do not create barriers because of lack of sensitivity to disability, cultural, or other subgroups. Items and test specifications are developed by a team of individuals who understand the varied characteristics of items that might create difficulties for any group of students. Accessibility is

incorporated as a primary dimension of test specifications, so that accessibility is woven into the fabric of the test rather than being added after the fact.

Amenable to Accommodations

Even though items on universally designed assessments will be accessible for most students, there will still be some students who continue to need accommodations. Thus, another essential element of any universally designed assessment is that it is compatible with accommodations and a variety of widely used adaptive equipment and assistive technology. The NDE and DRC work to ensure that state guidelines on the use of accommodations are compatible with the assessment being developed.

Simple, Clear, and Intuitive Instructions and Procedures

Assessment instructions should be easy to understand, regardless of a student's experience, knowledge, language skills, or current concentration level. Directions and questions need to be in simple, clear, and understandable language. Knowledge questions that are posed within complex language certainly invalidate the test if students cannot understand how they are expected to respond to a question.

Maximum Readability and Comprehensibility

A variety of guidelines exist to ensure that text is maximally readable and comprehensible. These features go beyond what is measured by readability formulas. Readability and comprehensibility are affected by many characteristics, including student background, sentence difficulty, organization of text, and others. All of these features are considered as the NDE develops the text of assessments.

Plain language is a concept now being highlighted in research on assessments. Plain language has been defined as language that is straightforward and concise. The following strategies for editing text to produce plain language are used during the NDE's editing process:

- Reduce excessive length.
- Use common words.
- Avoid ambiguous words.
- Avoid irregularly spelled words.
- Avoid proper names.
- Avoid inconsistent naming and graphic conventions.
- Avoid unclear signals about how to direct attention.
- Mark all questions.
- Maximum legibility.

Legibility is the physical appearance of text, the way that the shapes of letters and numbers enable people to read text easily. Bias results when tests contain physical features that interfere with a student's focus on or understanding of the constructs that test items are intended to assess. DRC

works closely with the NDE to develop a style guide that includes dimensions of style that are consistent with universal design.

DOK: Interpreting and assigning DOK levels to both objectives within standards and assessment items is an essential requirement of alignment analysis. Four levels of DOK are used for this analysis. The NeSA assessments include items written at levels 1, 2, and 3. Level 4 items are not included due to the test being comprised of only MC items.

Reading Level 1

Level 1 requires students to receive or recite facts or to use simple skills or abilities. Oral reading that does not include analysis of the text as well as basic comprehension of a text is included. Items require only a shallow understanding of text presented and often consist of verbatim recall from text or simple understanding of a single word or phrase. Some examples that represent, but do not constitute all of, Level 1 performance are:

- Support ideas by reference to details in the text.
- Use a dictionary to find the meaning of words.
- Identify figurative language in a reading passage.

Reading Level 2

Level 2 includes the engagement of some mental processing beyond recalling or reproducing a response; it requires both comprehension and subsequent processing of text or portions of text. Intersentence analysis of inference is required. Some important concepts are covered, but not in a complex way. Standards and items at this level may include words such as summarize, interpret, infer, classify, organize, collect, display, compare, and determine whether fact or opinion. Literal main ideas are stressed. A Level 2 assessment item may require students to apply some of the skills and concepts that are covered in Level 1. Some examples that represent, but do not constitute all of, Level 2 performance are:

- Use context cues to identify the meaning of unfamiliar words.
- Predict a logical outcome based on information in a reading selection.
- Identify and summarize the major events in a narrative.

Reading Level 3

Deep knowledge becomes more of a focus at Level 3. Students are encouraged to go beyond the text; however, they are still required to show understanding of the ideas in the text. Students may be encouraged to explain, generalize, or connect ideas. Standards and items at Level 3 involve reasoning and planning. Students must be able to support their thinking. Items may involve abstract theme identification, inference across an entire passage, or students' application of prior knowledge. Items may also involve more superficial connections between texts. Some examples that represent, but do not constitute all of, Level 3 performance are:

- Determine the author’s purpose and describe how it affects the interpretation of a reading selection.
- Summarize information from multiple sources to address a specific topic.
- Analyze and describe the characteristics of various types of literature.

Reading Level 4

Higher-order thinking is central and knowledge is deep at Level 4. The standard or assessment item at this level will probably be an extended activity, with extended time provided. The extended time period is not a distinguishing factor if the required work is only repetitive and does not require applying significant conceptual understanding and higher-order thinking. Students take information from at least one passage and are asked to apply this information to a new task. They may also be asked to develop hypotheses and perform complex analyses of the connections among texts. Some examples that represent, but do not constitute all of, Level 4 performance are:

- Analyze and synthesize information from multiple sources.
- Examine and explain alternative perspectives across a variety of sources.
- Describe and illustrate how common themes are found across texts from different cultures.

Mathematics Level 1

Level 1 includes the recall of information such as a fact, definition, term, or a simple procedure, as well as performing a simple algorithm or applying a formula. That is, in mathematics, a one-step, well-defined, and straight algorithmic procedure should be included at this lowest level. Other key words that signify a Level 1 include “identify,” “recall,” “recognize,” “use,” and “measure.” Verbs such as “describe” and “explain” could be classified at different levels, depending on what is to be described and explained.

Mathematics Level 2

Level 2 includes the engagement of some mental processing beyond a habitual response. A Level 2 assessment item requires students to make some decisions as to how to approach the problem or activity, whereas Level 1 requires students to demonstrate a rote response, perform a well-known algorithm, follow a set procedure (like a recipe), or perform a clearly defined series of steps. Keywords that generally distinguish a Level 2 item include “classify,” “organize,” “estimate,” “make observations,” “collect and display data,” and “compare data.” These actions imply more than one step. For example, to compare data requires first identifying characteristics of the objects or phenomenon and then grouping or ordering the objects. Some action verbs, such as “explain,” “describe,” or “interpret” could be classified at different levels depending on the object of the action. For example, if an item required students to explain how light affects mass by indicating there is a relationship between light and heat, this is considered a Level 2. Interpreting information from a simple graph, requiring reading information from the graph, also is a Level 2. Interpreting information from a complex graph that requires some decisions on what features of the graph need to be considered and how information from the graph can be aggregated is a Level 3. Caution is

warranted in interpreting Level 2 as only skills because some reviewers will interpret skills very narrowly, as primarily numerical skills. Such interpretation excludes from this level other skills, such as visualization skills and probability skills, which may be more complex simply because they are less common. Other Level 2 activities include explaining the purpose and use of experimental procedures; carrying out experimental procedures; making observations and collecting data; classifying, organizing, and comparing data; and organizing and displaying data in tables, graphs, and charts.

Mathematics Level 3

Level 3 requires reasoning, planning, using evidence, and a higher level of thinking than the previous two levels. In most instances, requiring students to explain their thinking is a Level 3. Activities that require students to make conjectures are also at this level. The cognitive demands at Level 3 are complex and abstract. The complexity does not result from the fact that there are multiple answers, a possibility for both Levels 1 and 2, but because the task requires more demanding reasoning. An activity, however, that has more than one possible answer and requires students to justify the response they give would most likely be a Level 3. Other Level 3 activities include drawing conclusions from observations, citing evidence and developing a logical argument for concepts, explaining phenomena in terms of concepts, and using concepts to solve problems.

Mathematics Level 4

Level 4 requires complex reasoning, planning, developing, and thinking most likely over an extended period of time. The extended time period is not a distinguishing factor if the required work is only repetitive and does not require applying significant conceptual understanding and higher-order thinking. For example, if a student has to take the water temperature from a river each day for a month and then construct a graph, this would be classified as a Level 2. However, if the student were to conduct a river study that requires taking into consideration a number of variables, this would be a Level 4. At Level 4, the cognitive demands of the task should be high and the work should be very complex. Students should be required to make several connections—relate ideas *within* the content area or *among* content areas—and have to select one approach among many alternatives on how the situation should be solved, in order to be at this highest level. Level 4 activities include designing and conducting experiments, making connections between a finding and related concepts and phenomena, combining and synthesizing ideas into new concepts, and critiquing experimental designs.

Science Level 1

Level 1 (Recall and Reproduction) requires the recall of information, such as a fact, definition, term, or a simple procedure, as well as performance of a simple science process or procedure. Level 1 only requires students to demonstrate a rote response, use a well-known formula, follow a set procedure (like a recipe), or perform a clearly defined series of steps. A “simple” procedure is well defined and typically involves only one step. Verbs such as “identify,”

“recall,” “recognize,” “use,” “calculate,” and “measure” generally represent cognitive work at the recall and reproduction level. Simple word problems that can be directly translated into and solved by a formula are considered Level 1. Verbs such as “describe” and “explain” could be classified at different DOK levels, depending on the complexity of what is to be described and explained. A student answering a Level 1 item either knows the answer or does not: that is, the item does not need to be “figured out” or “solved.” In other words, if the knowledge necessary to answer an item automatically provides the answer to it, then the item is at Level 1. If the knowledge needed to answer the item is not automatically provided in the stem, the item is at least at Level 2. Some examples that represent, but do not constitute all of, Level 1 performance are:

- Recall or recognize a fact, term, or property.
- Represent in words or diagrams a scientific concept or relationship.
- Provide or recognize a standard scientific representation for simple phenomenon.
- Perform a routine procedure, such as measuring length.

Science Level 2

Level 2 (Skills and Concepts) includes the engagement of some mental processing beyond recalling or reproducing a response. The content knowledge or process involved is **more complex** than in Level 1. Items require students to make some decisions as to how to approach the question or problem. Keywords that generally distinguish a Level 2 item include “classify,” “organize,” “estimate,” “make observations,” “collect and display data,” and “compare data.” These actions imply **more than one step**. For example, to compare data requires first identifying characteristics of the objects or phenomena and then grouping or ordering the objects. Level 2 activities include making observations and collecting data; classifying, organizing, and comparing data; and organizing and displaying data in tables, graphs, and charts. Some action verbs, such as “explain,” “describe,” or “interpret,” could be classified at different DOK levels, depending on the complexity of the action. For example, interpreting information from a simple graph, requiring reading information from the graph, is a Level 2. An item that requires interpretation from a complex graph, such as making decisions regarding features of the graph that need to be considered and how information from the graph can be aggregated, is at Level 3. Some examples that represent, but do not constitute all of, Level 2 performance are:

- Specify and explain the relationship between facts, terms, properties, or variables.
- Describe and explain examples and non-examples of science concepts.
- Select a procedure according to specified criteria and perform it.
- Formulate a routine problem, given data and conditions.
- Organize, represent, and interpret data.

Science Level 3

Level 3 (Strategic Thinking) requires reasoning, planning, using evidence, and a higher level of thinking than the previous two levels. The cognitive demands at Level 3 are complex and abstract. The complexity does not result only from the fact that there could be multiple answers, a possibility for both Levels 1 and 2, but because the multi-step task requires more demanding reasoning. In most instances, requiring students to explain their thinking is at Level 3; requiring a very simple explanation or a word or two should be at Level 2. An activity that has more than one possible answer and requires students to justify the response they give would most likely be a Level 3. Experimental designs in Level 3 typically involve more than one dependent variable. Other Level 3 activities include drawing conclusions from observations; citing evidence and developing a logical argument for concepts; explaining phenomena in terms of concepts; and using concepts to solve non-routine problems. Some examples that represent, but do not constitute all of, Level 3 performance are:

- Identify research questions and design investigations for a scientific problem.
- Solve non-routine problems.
- Develop a scientific model for a complex situation.
- Form conclusions from experimental data.

Science Level 4

Level 4 (Extended Thinking) involves high cognitive demands and complexity. Students are required to make several connections—relate ideas within the content area or among content areas—and have to select or devise one approach among many alternatives to solve the problem. Many on-demand assessment instruments will not include any assessment activities that could be classified as Level 4. However, standards, goals, and objectives can be stated in such a way as to expect students to perform extended thinking. “Develop generalizations of the results obtained and the strategies used and apply them to new problem situations,” is an example of a grade 8 objective that is a Level 4. Many, but not all, performance assessments and open-ended assessment activities requiring significant thought will be Level 4.

Level 4 requires complex reasoning, experimental design and planning, and probably will require an extended period of time either for the science investigation required by an objective, or for carrying out the multiple steps of an assessment item. However, the extended time period is not a distinguishing factor if the required work is only repetitive and does not require applying significant conceptual understanding and higher-order thinking. For example, if a student has to take the water temperature from a river each day for a month and then construct a graph, this would be classified as a Level 2 activity. However, if the student conducts a river study that requires taking into consideration a number of variables, this would be a Level 4. Some examples that represent, but do not constitute all of, Level 4 performance are:

- Based on data provided from a complex experiment that is novel to the student, deduce the fundamental relationship between a controlled variable and an experimental variable.
- Conduct an investigation, from specifying a problem to designing and carrying out an experiment, to analyzing its data and forming conclusions.

Source of Challenge Criterion

Source of Challenge criterion is only used to identify items where the major cognitive demand is inadvertently placed and is other than the targeted skill, concept, or application. Cultural bias or specialized knowledge could be reasons for an item to have a source of challenge problem. Such items' characteristics may cause some students to not answer an assessment item or answer an assessment item incorrectly or at a lower level even though they have the understanding and skills being assessed.

Item Content Review: Prior to field testing, all newly developed test passages/items were submitted to grade-level content committees for review. The content committees consisted of Nebraska educators from school districts throughout the state. The primary responsibility of the content committees was to evaluate items with regard to quality and content classification, including grade-level appropriateness, estimated difficulty, DOK, and source of challenge. They also suggested revisions, if appropriate. The committees also reviewed the items for adherence to the principles of universal design, including language demand and issues of bias, fairness, and sensitivity.

Item review committee members were selected by the NDE. The NDE test development team members facilitated the process. Training was provided by the NDE and included how to review items for technical quality and content quality, including DOK and adherence to principles of universal design. In addition, training included providing committee members with the procedures for item review.

Committee members reviewed the items for quality and content, as well as for the following categories:

- Indicator (standard) Alignment
- Difficulty Level (classified as Low, Medium, or High)
- DOK (classified as Recall, Application, or Strategic Thinking)
- Correct Answer
- Quality of Graphics
- Appropriate Language Demand
- Freedom from Bias (classified as Yes or No)

Committee members were asked to flag items that needed revision and to denote suggested revisions on the flagged item cards.

Security was addressed by adhering to a strict set of procedures. Items in binders did not leave the meeting rooms and were accounted for at the end of each day before attendees were dismissed. All

attendees, with the exception of the NDE staff, were required to sign a Confidentiality Agreement (Appendix D).

Sensitivity and Bias Review: Prior to field testing items, all newly developed test items were submitted to a Bias and Sensitivity Committee for review. The committee's primary responsibility was to evaluate passages and items as to acceptability with regard to bias and sensitivity issues. They also made recommendations for changes or deletion of items in order to remove the area of concern. The bias/sensitivity committee was composed of Nebraska educators who represented the diversity of students. All committee members were trained by an NDE test development lead to review items for bias and sensitivity issues using *Fairness in Testing* training manual developed by DRC (Appendix E).

All passages/items were read by all of the respective committee members. Each member noted bias and/or sensitivity comments on a review form. All comments were then compiled and the actions taken on these items were recorded by the NDE. Committee members were required to sign a Confidentiality Agreement and strict security measures were in place to ensure that secure materials remained guarded (Appendix D).

2.6 ITEM BANKING

DRC maintains an item bank (IDEAS) that provides a repository of item image, history, statistics, and usage. IDEAS includes a record of all newly created items together with item data from each item field test. It also includes all data from the operational administration of the items. Within IDEAS, DRC:

- updates the Nebraska item bank after each administration;
- updates the Nebraska item bank with newly developed items;
- monitors the Nebraska item bank to ensure an appropriate balance of items aligned with content standards, goals, and objectives;
- monitors item history statistics; and
- monitors the Nebraska item bank for an appropriate balance of DOK levels.

2.7 THE OPERATIONAL FORM CONSTRUCTION PROCESS

The Spring 2015 operational forms were constructed in Lincoln, Nebraska in August 2014 (Reading and Mathematics) and early September 2014 (Science). The forms were constructed by NDE representatives and DRC content specialists. Training was provided by DRC for the forms construction process.

Prior to the construction of the operational forms, DRC Test Development content specialists reviewed the test blueprints to ensure that there was alignment between the items and the indicators, including the number of items per standard for each content-area test.

DRC psychometricians provided Test Development specialists with an overview of the psychometric guidelines and targets for operational forms construction. The foremost guideline was for item content to match the test blueprint (Table of Specifications) for the given content. The point-biserial correlation guideline was to be greater than 0.3 (with a requirement for no point-biserial correlation less than zero). In addition, the average target p -value for each test was to be about 0.65. A Differential Item Functioning (DIF) code of C was to be avoided (unless no other items were available to fulfill a blueprint requirement). The overall summary of the actual approved p -value and biserial of the forms is provided in the summary table later in this document.

DRC Test Development specialists printed a copy of each item card, with accompanying item characteristics, image, and psychometric data. Test Development specialists verified the accuracy of each item card, making sure that the item image has its correct item characteristics. Test Development specialists carefully reviewed each item card's psychometric data to ensure it is complete and reasonable. For Reading, the item cards (items and passages) were compiled in binders and sorted by p -values from highest to lowest by passage with associated items. For Mathematics and science, the item cards were compiled in binders and sorted by p -values from highest to lowest by standard and indicator.

The NDE and DRC also checked to see that each item met technical quality for well-crafted items, including:

- only one correct answer,
- wording that is clear and concise,
- grammatical correctness,
- appropriate item complexity and cognitive demand,
 - appropriate range of difficulty,
 - appropriate depth-of-knowledge alignment,
- aligned with principles of Universal Design, and
- free of any content that might be offensive, inappropriate, or biased (content bias).

NDE representatives and DRC Test Development specialists made initial grade-level selections of the items (passages and items for Reading), known as the “pull list,” to be included on the 2015 operational forms. The goal was for the first pull of the items to meet the Table of Specification (TOS) guidelines and psychometric guidelines specific to each content area. As items were selected, the unique item codes were entered into a form building template (Perform) which contained the item pool with statistics and item characteristics. Perform automatically calculated the p -value, biserial, number of items per indicator and standard, number of items per DOK level (1, 2, or 3), and distribution of answer key as items were selected for each grade. As items were selected, the item characteristics (key, DOK, and alignment to indicator) were verified.

Differential Item Functioning in Operational Form Construction: DIF is present when the likelihood of success on an item is influenced by group membership. A pattern of such results may suggest the presence of, but does not prove, *item bias*. Actual item bias may present negative group stereotypes, may use language that is more familiar to one subpopulation than to another, or may present information in a format that disadvantages certain learning styles. While the source of item bias is often clear to trained judges, many instances of DIF may have no identifiable cause (resulting in false positives). As such, DIF is not used as a substitute for rigorous, hands-on reviews by content and bias specialists. Instead, DIF helps to organize the review of the instances in which bias is suggested. No items are automatically rejected simply because a statistical method flagged them or automatically accepted because they were not flagged.

During the operational form-pull process, the DIF code for every item proposed for use in the operational (core) is examined. To the greatest extent possible, the blueprint is met through the use of items with statistical DIF codes of A. Although DIF codes of B and C are not desirable and are deliberately avoided, the combination of the required blueprint and the depth of the available operational-ready item pool occasionally require that items with B and C DIF are considered for operational use. In addition, for passage-based tests like reading (in which each item available in the item pool is linked to a set of passage-based items), the ability to use a minimum number of items associated with a passage may require the use of an item with a B or C DIF code. In any case, prior to allowing exceptions of this nature, every attempt is made to re-craft the core to avoid the use of the item with B or C DIF. Before allowing any exception to be made, the item in question is examined to determine whether the suggested bias is identifiable. If the suggested bias is determined to be valid, the item is not used.

Review of the Items and Test Forms: At every stage of the test development process the match of the item to the content standard was reviewed and verified, since establishing content validity is one of the most important aspects in the legal defensibility of a test. As a result, it is essential that an item selected for a form link directly to the content curriculum standard and performance standard to which it is measuring. Test Development specialists verified all items against their classification codes and item maps, both to evaluate the correctness of the classification and to ensure that the given task measures what it purports to measure.

2.8 READING ASSESSMENT

Test Design: The NeSA-R operational test includes operational passages with associated items and one field test passage with associated items. This test was administered online via the test engine developed and managed by DRC (INSIGHT Online Learning System). One form of the test was also published in a printed test booklet for students needing accommodation provided by paper/pencil test. Depending on grade, the forms contained 45 to 50 operational items.

Table 2.8.1 Reading 2015 Operational Test

Grade	Total No. of MC Core Items	No. of Embedded FT Items per Form (1 passage)	Total Items per Form	Total No. of Equivalent FT Forms	Total Core Points	Total No. of MC Items Added to the Bank
3	45	10	55	5	45	50
4	45	10	55	5	45	50
5	48	10	58	5	48	50
6	48	10	58	5	48	50
7	48	10	58	5	48	50
8	50	10	60	5	50	50
11	50	10	60	5	50	50

Psychometric Targets: The goal for the operational forms was to meet a mean p -value of approximately 0.65 with values restricted to the range of 0.30 to 0.90 and point-biserial correlations greater than 0.25, based on previous field test results. However, these targets are secondary to constructing the best test possible. Some compromises were allowed when necessary to best meet the objective of the assessment, to conform to the test specifications, and to operate within the limitations of the item bank.

Equating Design: Spring 2015 was the sixth operational administration of NeSA-R. Approximately 70% of the assessment was constructed from passages and related items field tested from Spring 2009–2014. The approximate remaining 30% of the assessment was constructed from an overlap of items and passages from the 2014 operational (core) item positions from the Spring 2014 operational forms.

In addition to the operational passage sets, each student received one randomly selected field test passage with 10 associated field test items. The passages and items taken by each student were administered in two testing sessions each intended to be administered in a single class period. The operational passages were administered to the student in a random order, but the field test passage was maintained in a fixed position. Items within a passage were administered in a fixed order for the passage. Equating was accomplished by anchoring on the operational passage items and calibrating the field test items concurrently.

2.9 MATHEMATICS ASSESSMENT

Test Design: The NeSA-M operational test includes operational and field test items. This test was administered online via the test engine developed and managed by DRC (INSIGHT Online Learning System). One form of the test was also published in a printed test booklet for students needing

accommodation provided by paper/pencil test. Depending on grade, the forms contained 50 to 60 operational items.

Table 2.9.1 Mathematics 2015 Operational Test

Grade	Total No. of MC Core Items	No. of Embedded FT Items per Form	Total Items per Form	Total No. of Equivalent FT Forms	Total Core Points	Total No. of MC Items Added to the Bank
3	50	10	60	5	50	50
4	55	10	65	5	55	50
5	55	10	65	5	55	50
6	58	10	68	5	58	50
7	58	10	68	5	58	50
8	60	10	70	5	60	50
11	60	10	70	5	60	50

Psychometric Targets: The goal for the operational forms was to meet a mean *p*-value of approximately 0.65 with values restricted to the range of 0.3 to 0.9 and point-biserial correlations greater than 0.25, based on previous field test results. However, these targets are secondary to constructing the best test possible. Some compromises were allowed when necessary to best meet the objective of the assessment, to conform to the test specifications, and to operate within the limitations of the item bank.

Equating Design: Spring 2015 was the fifth operational administration of NeSA-M. Approximately 70% of the assessment was constructed from items field tested from Spring 2010–2014. The approximate remaining 30% of the assessment was constructed from an overlap of items from the 2014 operational (core) item positions from the 2014 operational forms.

In addition to the operational items, each student received 10 randomly selected field test items. The items taken by each student were administered in two testing sessions each intended to be administered in a single class period. The operational items were administered to the student in a random order, but the field test items were maintained in fixed positions. Equating was accomplished by anchoring on the operational items and calibrating the field test items concurrently.

2.10 SCIENCE ASSESSMENT

Test Design: The NeSA-S operational test includes operational and field test items. This test was administered online via the test engine developed and managed by DRC (INSIGHT Online Learning System). One form of the test was also published in a printed test booklet for students needing

accommodation provided by paper/pencil test. Depending on grade, the forms contained 50 or 60 operational items.

Table 2.10.1 Science 2015 Operational Test

Grade	No. Operational Items	No. of Embedded FT Items per Form	Total Items	Total No. of FT Forms	Total No. of Items Field Tested
5	50	10	60	5	50
8	60	10	70	5	50
11	60	10	70	5	50

Psychometric Targets: The goal for the operational forms was to meet a mean p -value of approximately 0.65 with values restricted to the range of 0.3 to 0.9 and point-biserial correlations greater than 0.25, based on previous field test results. However, these targets are secondary to constructing the best test possible. Some compromises were allowed when necessary to best meet the objective of the assessment, to conform to the test specifications, and to operate within the limitations of the item bank.

Equating Design: Spring 2015 was the fourth operational administration of NeSA-S. Approximately 70% of the assessment was constructed from items field tested in Spring 2011–2014. The approximate remaining 30% of the assessment was constructed from an overlap of items from the 2014 operational (core) item positions from the 2014 operational forms.

In addition to the operational items, each student received 10 randomly selected field test items. The items taken by each student were administered in two testing sessions each intended to be administered in a single class period. The operational items were administered to the student in a random order, but the field test items were maintained in fixed positions. Equating was accomplished by anchoring on the operational items and calibrating the field test items concurrently.

3. STUDENT DEMOGRAPHICS

Three areas of student demographics are discussed below, summary demographics and accommodations, summary information on the number of students tested with breakdowns by mode, and summary information on testing times.

3.1 DEMOGRAPHICS AND ACCOMMODATIONS

Gender, ethnicity, food program status (FRL), Limited English Proficiency/English Language Learners (LEP/ELL) status, Special Education status (SPED), and accommodation status data was collected for all students who participated and attempted the 2015 NeSA assessments. This summary of student demographics by grade and content area is provided in Tables 3.1.1– 3.1.7. These tables show around 22,000 students took the assessment in each grade. Of those students across grades, half are males, half are females, over two thirds white, and less than one fifth are Hispanic. Among the students across grades, about 37% to 47% are eligible for FRL, 2% to 9% are LEP/ELL, and 11% to 16% belong to at least one SPED category. For all three of these programs/categories, the participation rate is slightly lower for upper grade students. In terms of the test accommodations, there are about 6% to 16% of the students across grade and content area that report at least one type of accommodation (see row ‘Total’ for ‘Accommodation’ in the table). Similar to the rate for FRL, LEP/ELL, and SPED across grades, the rate for accommodation is lower for high school students (Grade 11). Across all grades, the ‘Timing/Schedule/Setting’ is the most utilized accommodation (about 6-10% for Grade 3-8, and 4% for Grade 11), followed by the ‘Content Presentation’ (about 6-9% for Grade 3-7, and 2-5% for Grade 8 and 11).

Table 3.1.1 Grade 3 NeSA Summary Data: Demographics and Accommodations

Grade 3		Reading		Mathematics	
		Count	%	Count	%
All Students		23013	100.00	23130	100.00
Gender	Female	11258	48.92	11314	48.91
	Male	11755	51.08	11816	51.09
Race/Ethnicity	American Indian/Alaska Native	324	1.41	323	1.40
	Asian	567	2.46	594	2.57
	Black	1619	7.04	1625	7.03
	Hispanic	4237	18.41	4300	18.59
	Native Hawaiian or other Pacific Islander	34	0.15	34	0.15
	White	15384	66.85	15407	66.61
	Two or More Races	848	3.68	847	3.66

Grade 3		Reading		Mathematics	
		Count	%	Count	%
Food Program	Yes	11087	48.18	11175	48.31
	No	11926	51.82	11955	51.69
LEP/ELL	Yes	2199	9.56	2310	9.99
	No	20814	90.44	20820	90.01
Special Education	Yes	3688	16.03	3697	15.98
	No	19325	83.97	19433	84.02
Accommodations	Content Presentation	1780	7.73	1829	7.91
	Response	901	3.92	1109	4.79
	Timing/Schedule/Setting	2054	8.93	2066	8.93
	Direct Linguistic Support with Test Directions	1469	6.38	1569	6.78
	Direct Linguistic Support with Content and Test items	1834	7.97	1901	8.22
	Indirect Linguistic Support	1602	6.96	1698	7.34
	Spanish	24	0.10	67	0.29
	Braille*	1	0.00	1	0.00
	Large Print*	10	0.04	11	0.05
	Audio	208	0.90	226	0.98
	Total	4025	17.49	4168	18.02

*Count represents the number of booklets ordered. This is not tracked.

Table 3.1.2 Grade 4 NeSA Summary Data: Demographics and Accommodations

Grade 4		Reading		Mathematics	
		Count	%	Count	%
All Students		22590	100.00	22685	100.00
Gender	Female	10992	48.66	11033	48.64
	Male	11598	51.34	11652	51.36
Race/Ethnicity	American Indian/Alaska Native	268	1.19	268	1.18
	Asian	551	2.44	570	2.51

Nebraska State Accountability 2015 Technical Report

Grade 4		Reading		Mathematics	
		Count	%	Count	%
	Black	1567	6.94	1579	6.96
	Hispanic	4050	17.93	4105	18.10
	Native Hawaiian or other Pacific Islander	29	0.13	29	0.13
	White	15284	67.66	15296	67.43
	Two or More Races	841	3.72	838	3.69
Food Program	Yes	10524	46.59	10600	46.73
	No	12066	53.41	12085	53.27
LEP/ELL	Yes	1619	7.17	1726	7.61
	No	20971	92.83	20959	92.39
Special Education	Yes	3680	16.29	3681	16.23
	No	18910	83.71	19004	83.77
Accommodations	Content Presentation	2040	9.03	2046	9.02
	Response	1034	4.58	1339	5.90
	Timing/Schedule/Setting	2318	10.26	2322	10.24
	Direct Linguistic Support with Test Directions	1257	5.56	1344	5.92
	Direct Linguistic Support with Content and Test items	1432	6.34	1523	6.71
	Indirect Linguistic Support	1298	5.75	1368	6.03
	Spanish	29	0.13	67	0.30
	Braille*	1	0.00	0	0.00
	Large Print*	12	0.05	13	0.06
	Audio	217	0.96	231	1.02
	Total	3986	17.64	4087	18.02

*Count represents the number of booklets ordered. This is not tracked.

Table 3.1.3 Grade 5 NeSA Summary Data: Demographics and Accommodations

Grade 5		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		22878	100.00	22946	100.00	22949	100.00
Gender	Female	11207	48.99	11231	48.95	11238	48.97
	Male	11671	51.01	11715	51.05	11711	51.03
Race/Ethnicity	American Indian/Alaska Native	318	1.39	318	1.39	318	1.39
	Asian	550	2.40	566	2.47	567	2.47
	Black	1476	6.45	1479	6.45	1476	6.43
	Hispanic	4061	17.75	4108	17.90	4110	17.91
	Native Hawaiian or other Pacific Islander	29	0.13	29	0.13	30	0.13
	White	15648	68.40	15650	68.20	15655	68.22
	Two or More Races	796	3.48	796	3.47	793	3.46
Food Program	Yes	10577	46.23	10630	46.33	10627	46.31
	No	12301	53.77	12316	53.67	12322	53.69
LEP/ELL	Yes	1133	4.95	1210	5.27	1211	5.28
	No	21745	95.05	21736	94.73	21738	94.72
Special Education	Yes	3672	16.05	3667	15.98	3672	16.00
	No	19206	83.95	19279	84.02	19277	84.00
Accommodations	Content Presentation	2163	9.45	2209	9.63	2144	9.34
	Response	1071	4.68	1429	6.23	1058	4.61
	Timing/Schedule/Setting	2441	10.67	2433	10.60	2358	10.27
	Direct Linguistic Support with Test Directions	810	3.54	811	3.53	812	3.54
	Direct Linguistic Support with Content and Test items	989	4.32	1022	4.45	1049	4.57
	Indirect Linguistic Support	765	3.34	866	3.77	791	3.45
	Spanish	30	0.13	70	0.31	69	0.30
	Braille*	2	0.01	2	0.01	2	0.01

Nebraska State Accountability 2015 Technical Report

Grade 5		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
	Large Print*	10	0.04	9	0.04	8	0.03
	Audio	218	0.95	223	0.97	241	1.05
	Total	3649	15.95	3682	16.05	3628	15.81

*Count represents the number of booklets ordered. This is not tracked.

Table 3.1.4 Grade 6 NeSA Summary Data: Demographics and Accommodations

Grade 6		Reading		Mathematics	
		Count	%	Count	%
All Students		22377	100.00	22445	100.00
Gender	Female	10859	48.53	10883	48.49
	Male	11518	51.47	11562	51.51
Race/Ethnicity	American Indian/Alaska Native	294	1.31	294	1.31
	Asian	536	2.40	550	2.45
	Black	1466	6.55	1471	6.55
	Hispanic	3966	17.72	4009	17.86
	Native Hawaiian or other Pacific Islander	22	0.10	22	0.10
	White	15299	68.37	15304	68.18
	Two or More Races	794	3.55	795	3.54
Food Program	Yes	10047	44.90	10109	45.04
	No	12330	55.10	12336	54.96
LEP/ELL	Yes	694	3.10	774	3.45
	No	21683	96.90	21671	96.55
Special Education	Yes	3402	15.20	3393	15.12
	No	18975	84.80	19052	84.88
Accommodations	Content Presentation	1994	8.91	1942	8.65
	Response	889	3.97	1479	6.59
	Timing/Schedule/Setting	2172	9.71	2131	9.49

Grade 6		Reading		Mathematics	
		Count	%	Count	%
	Direct Linguistic Support with Test Directions	529	2.36	524	2.33
	Direct Linguistic Support with Content and Test items	580	2.59	626	2.79
	Indirect Linguistic Support	518	2.31	504	2.25
	Spanish	32	0.14	69	0.31
	Braille*	4	0.02	3	0.01
	Large Print*	8	0.04	8	0.04
	Audio	308	1.38	315	1.40
	Total	3058	13.67	3125	13.92

*Count represents the number of booklets ordered. This is not tracked.

Table 3.1.5 Grade 7 NeSA Summary Data: Demographics and Accommodations

Grade 7		Reading		Mathematics	
		Count	%	Count	%
	All Students	22212	100.00	22285	100.00
Gender	Female	10923	49.18	10949	49.13
	Male	11289	50.82	11336	50.87
Race/Ethnicity	American Indian/Alaska Native	294	1.32	291	1.31
	Asian	516	2.32	534	2.40
	Black	1369	6.16	1369	6.14
	Hispanic	3927	17.68	3989	17.90
	Native Hawaiian or other Pacific Islander	26	0.12	26	0.12
	White	15343	69.08	15339	68.83
	Two or More Races	737	3.32	737	3.31
Food Program	Yes	9982	44.94	10049	45.09
	No	12230	55.06	12236	54.91
LEP/ELL	Yes	609	2.74	696	3.12
	No	21603	97.26	21589	96.88

Nebraska State Accountability 2015 Technical Report

Grade 7		Reading		Mathematics	
		Count	%	Count	%
Special Education	Yes	3137	14.12	3135	14.07
	No	19075	85.88	19150	85.93
Accommodations	Content Presentation	1618	7.28	1644	7.38
	Response	739	3.33	1381	6.20
	Timing/Schedule/Setting	1728	7.78	1714	7.69
	Direct Linguistic Support with Test Directions	363	1.63	423	1.90
	Direct Linguistic Support with Content and Test items	342	1.54	430	1.93
	Indirect Linguistic Support	310	1.40	341	1.53
	Spanish	31	0.14	82	0.37
	Braille*	3	0.01	3	0.01
	Large Print*	11	0.05	12	0.05
	Audio	478	2.15	515	2.31
	Total	2516	11.33	2682	12.04

*Count represents the number of booklets ordered. This is not tracked.

Table 3.1.6 Grade 8 NeSA Summary Data: Demographics and Accommodations

Grade 8		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		21904	100.00	21985	100.00	21993	100.00
Gender	Female	10743	49.05	10763	48.96	10762	48.93
	Male	11161	50.95	11222	51.04	11231	51.07
Race/Ethnicity	American Indian/Alaska Native	343	1.57	342	1.56	342	1.56
	Asian	484	2.21	504	2.29	504	2.29
	Black	1453	6.63	1458	6.63	1456	6.62
	Hispanic	3717	16.97	3775	17.17	3774	17.16
	Native Hawaiian or other Pacific Islander	24	0.11	24	0.11	24	0.11
	White	15216	69.47	15216	69.21	15225	69.23

Nebraska State Accountability 2015 Technical Report

Grade 8		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
	Two or More Races	667	3.05	666	3.03	668	3.04
Food Program	Yes	9477	43.27	9532	43.36	9542	43.39
	No	12427	56.73	12453	56.64	12451	56.61
LEP/ELL	Yes	485	2.21	586	2.67	588	2.67
	No	21419	97.79	21399	97.33	21405	97.33
Special Education	Yes	2935	13.40	2925	13.30	2933	13.34
	No	18969	86.60	19060	86.70	19060	86.66
Accommodations	Content Presentation	1422	6.49	1417	6.45	1460	6.64
	Response	568	2.59	1195	5.44	657	2.99
	Timing/Schedule/Setting	1562	7.13	1555	7.07	1493	6.79
	Direct Linguistic Support with Test Directions	277	1.26	344	1.56	325	1.48
	Direct Linguistic Support with Content and Test items	272	1.24	377	1.71	349	1.59
	Indirect Linguistic Support	235	1.07	283	1.29	276	1.25
	Spanish	56	0.26	96	0.44	98	0.45
	Braille*	1	0.00	1	0.00	1	0.00
	Large Print*	14	0.06	13	0.06	13	0.06
	Audio	469	2.14	470	2.14	485	2.21
	Total	2285	10.43	2470	11.23	2354	10.70

*Count represents the number of booklets ordered. This is not tracked.

Table 3.1.7 Grade 11 NeSA Summary Data: Demographics and Accommodations

Grade 11		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		21125	100.00	21107	100.00	21102	100.00
Gender	Female	10312	48.81	10308	48.84	10299	48.81
	Male	10813	51.19	10799	51.16	10803	51.19
Race/Ethnicity	American Indian/Alaska Native	243	1.15	239	1.13	241	1.14

Nebraska State Accountability 2015 Technical Report

Grade 11		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
	Asian	466	2.21	466	2.21	464	2.20
	Black	1289	6.10	1292	6.12	1293	6.13
	Hispanic	3300	15.62	3303	15.65	3301	15.64
	Native Hawaiian or other Pacific Islander	31	0.15	31	0.15	31	0.15
	White	15213	72.01	15196	72.00	15190	71.98
	Two or More Races	583	2.76	580	2.75	582	2.76
Food Program	Yes	7993	37.84	7997	37.89	7988	37.85
	No	13132	62.16	13110	62.11	13114	62.15
LEP/ELL	Yes	374	1.77	388	1.84	389	1.84
	No	20751	98.23	20719	98.16	20713	98.16
Special Education	Yes	2369	11.21	2353	11.15	2364	11.20
	No	18756	88.79	18754	88.85	18738	88.80
Accommodations	Content Presentation	688	3.26	647	3.07	683	3.24
	Response	345	1.63	819	3.88	473	2.24
	Timing/Schedule/Setting	984	4.66	990	4.69	993	4.71
	Direct Linguistic Support with Test Directions	130	0.62	131	0.62	137	0.65
	Direct Linguistic Support with Content and Test items	119	0.56	144	0.68	152	0.72
	Indirect Linguistic Support	136	0.64	138	0.65	143	0.68
	Spanish	82	0.39	88	0.42	95	0.45
	Braille*	1	0.00	1	0.00	1	0.00
	Large Print*	7	0.03	7	0.03	7	0.03
	Audio	258	1.22	260	1.23	273	1.29
	Total	1315	6.22	1492	7.07	1389	6.58

*Count represents the number of booklets ordered. This is not tracked.

3.2 STUDENTS TESTED AND MODE SUMMARY DATA

As noted in Chapters One and Two, the 2015 NeSA assessments were administered online to the extent practical. One form of the test was also published in a printed test booklet for students needing accommodation of a paper/pencil test. Tables 3.2.1 – 3.2.3 report the number of students in each test mode. For NeSA-R, between 2% and 8% of students took the assessment in the paper-based version with the lower percentages occurring in middle and high schools.

Table 3.2.1 NeSA-R Number of Students Tested

Grade	Total	Online	Paper	Percent Paper
3	23013	21330	1683	7
4	22590	20921	1669	7
5	22878	21371	1507	7
6	22377	21135	1242	6
7	22212	21265	947	4
8	21904	21076	828	4
11	21125	20648	477	2

For NeSA-M, between 2% and 8% of students took the assessment in the paper-based version.

Table 3.2.2 NeSA-M Number of Students Tested

Grade	Total	Online	Paper	Percent Paper
3	23130	21373	1757	8
4	22685	20951	1734	8
5	22946	21370	1576	7
6	22445	21146	1299	6
7	22285	21332	953	4
8	21985	21096	889	4
11	21107	20612	495	2

For NeSA-S, between 2% and 7% of students took the assessment in the paper version.

Table 3.2.3 NeSA-S Number of Students Tested

Grade	Total	Online	Paper	Percent Paper
5	22949	21446	1503	7
8	21993	21155	838	4
11	21102	20618	484	2

The number of students, across content area and grade level, who took the 2015 NeSA tests online instead of paper pencil is similar to that of the 2015 NeSA tests.

3.3 TESTING TIME

Online testing time for the 2015 NeSA assessments was examined for each grade and content area. The data in Tables 3.3.1, 3.3.2, and 3.3.3 were compiled based on students who had a *single login, a single logout, and responded to all the items*. Similar to 2014, students from upper grade levels, on average, spent slightly less time for all content areas, and there was a slightly bigger difference in the time spent in sessions 1 and 2, indicating a tendency toward less time in the second session. As compared to 2014, the average 2015 online testing time slightly increased, especially among the lower grade students. This is probably because there are more first-time online test takers among them. The rest of the distribution of times was comparable. The outliers on the other end, greater than 90 minutes, are also interesting because this data does not include students who *paused out*, had the test ended due to inactivity, or were reactivated. It appears that they were actively involved with the test for the full time between the login and logout, but it raises the question of how fully engaged those students may have been for that amount of time.

Table 3.3.1 Duration of Reading Online Testing Sessions

Grade	3		4		5		6		7		8		11	
	1	2	1	2	1	2	1	2	1	2	1	2	1	2
<5	0	1	0	5	2	0	2	1	7	18	2	18	88	149
5-10	9	15	8	56	1	11	6	16	14	70	42	50	184	335
10-15	108	157	43	422	60	125	36	107	69	459	404	257	775	1286
15-20	508	613	286	1211	367	643	199	467	368	1753	1511	976	2549	3280
20-25	1154	1379	839	2155	1033	1543	556	1407	1210	2956	2874	2143	4126	4336
25-30	1946	2044	1621	2778	1852	2387	1266	2311	1948	3610	3784	3027	4229	3810
30-35	2439	2426	2351	2806	2374	2751	1840	2730	2748	3343	3636	3303	3189	2831
35-40	2531	2479	2560	2496	2599	2729	2269	2849	2857	2587	2817	3086	2122	1702
40-45	2283	2265	2450	2096	2300	2301	2372	2398	2736	1982	1935	2344	1201	1072
45-50	1967	1856	2147	1648	2164	1867	2365	1944	2293	1333	1374	1744	753	623
50-55	1733	1666	1741	1189	1725	1520	1977	1575	1857	897	826	1245	472	416
55-60	1395	1322	1481	969	1407	1256	1638	1149	1369	563	561	802	323	217
60-65	1206	1135	1148	718	1104	923	1325	890	961	417	372	596	201	144
65-70	872	880	904	543	880	721	1017	640	730	286	251	406	129	117
70-75	661	637	684	383	726	541	805	511	499	249	153	245	75	74
75-80	509	489	506	278	580	436	618	376	352	134	115	194	42	55
80-85	416	359	415	193	412	322	477	272	272	104	76	136	34	38
85-90	258	258	335	170	314	205	400	260	219	99	70	89	28	30
>90	1236	1279	1343	698	1424	1047	1907	1180	702	355	241	372	103	95
Total	21231	21260	20862	20814	21324	21328	21075	21083	21211	21215	21044	21033	20623	2061
Mean	88.5	66.3	74.5	62.6	58.8	67.0	76.5	52.6	58.7	51.2	48.6	51.2	47.4	40.5

Table 3.3.2 Duration of Mathematics Online Testing Sessions

Grade	3		4		5		6		7		8		11	
	1	2	1	2	1	2	1	2	1	2	1	2	1	2
<5	3	1	1	0	1	2	3	4	11	17	13	29	114	188
5-10	10	29	6	13	4	15	5	10	19	27	25	42	182	368
10-15	222	407	121	187	51	175	17	74	72	195	78	173	432	851
15-20	1059	1516	843	945	340	1023	137	385	381	1170	400	880	1315	2241
20-25	2350	2635	1984	2074	1157	2299	574	1341	1203	2754	1286	2345	2827	3898
25-30	2828	3069	2915	2768	1962	3000	1378	2518	2212	3774	2560	3452	3946	4310
30-35	2817	2787	2904	2840	2522	2870	2291	3036	3033	3586	3178	3472	3768	3393
35-40	2477	2333	2586	2534	2480	2410	2644	2957	3035	2780	3232	2990	2884	2123
40-45	1940	1861	2141	2129	2294	2026	2546	2359	2678	2134	2709	2199	1892	1185
45-50	1659	1454	1608	1591	2070	1687	2317	1940	2219	1576	2153	1662	1165	752
50-55	1250	1104	1270	1246	1673	1250	1889	1478	1812	1071	1514	1129	707	475
55-60	911	881	971	959	1421	1020	1490	1162	1312	679	1077	776	455	260
60-65	800	710	776	808	1034	783	1190	878	945	445	802	532	300	187
65-70	588	497	602	581	886	627	916	603	669	303	571	352	190	85
70-75	457	416	455	465	694	463	691	460	462	167	396	254	118	70
75-80	403	319	349	336	542	332	512	369	329	158	273	164	76	52
80-85	287	270	287	291	406	272	445	251	212	90	188	137	53	35
85-90	208	182	171	220	338	227	326	239	180	72	131	93	41	24
>90	1034	846	909	915	1427	840	1722	1032	509	301	470	394	123	72
Total	21303	21317	20899	20902	21302	21321	21093	21096	21293	21299	21056	21075	20588	20569
Mean	76.7	62.6	68.3	56.1	72.2	53.4	63.6	57.7	58.4	40.8	53.7	44.2	45.5	38.0

Table 3.3.3 Duration of Science Online Testing Sessions

Grade	5		8		11	
	1	2	1	2	1	2
<5	3	4	25	30	119	229
5-10	202	396	248	638	1568	2655
10-15	2186	3302	3213	5424	7766	8599
15-20	4308	4865	6125	6468	6199	5305
20-25	4156	3957	4677	3694	2659	2005
25-30	2969	2666	2641	1927	1065	826
30-35	2186	1929	1574	1120	565	383
35-40	1515	1260	933	715	249	192
40-45	1095	887	525	361	152	115
45-50	840	616	331	221	82	91
50-55	524	424	220	155	47	64
55-60	377	321	160	97	35	49
60-65	266	209	115	93	26	15
65-70	222	132	79	46	16	19
70-75	128	110	53	17	8	10
75-80	89	82	47	24	8	8
80-85	83	59	18	15	7	6
85-90	56	44	29	16	3	6
>90	213	142	118	67	21	26
Total	21418	21405	21131	21128	20595	20603
Mean	35.3	38.3	32.0	27.1	25.7	22.5

4. CLASSICAL ITEM STATISTICS

This chapter provides an overview of the most familiar item-level statistics obtained from classical (true-score model) item analysis: item difficulty, item discrimination, distractor distribution, and omits or blanks. The following results pertain only to operational NeSA items (i.e., those items that contributed to a student’s total test score). Rasch item statistics are discussed in Chapter Five, and test-level statistics are found in Chapter Six. The statistics provide information about the quality of the items based on student responses in an operational setting. The following sections provide descriptions of the item summary statistics found in Appendices F, G, and H.

4.1 ITEM DIFFICULTY

Item difficulty (p -value) is the proportion of examinees in the sample who answered the item correctly. For example, if an item has a p -value of 0.89, it means 89 percent of the students answered the item correctly. Relatively lower values correspond to more difficult items and those that have relatively higher values correspond to easier items. Items that are either very hard or very easy provide little information about student differences in achievement. On a standards-referenced test like the NeSA, a test development goal is to include a wide range of item difficulties. Typically, test developers target p -values in the range of 0.30 to 0.90. Mathematically, information is maximized and standard errors minimized when the p -value equals 0.50. Experience suggests that multiple choice items are effective when the student is more likely to succeed than fail and it is important to include a range of difficulties matching the distribution of student abilities (Wright & Stone, 1979). Occasionally, items that fall outside the desired range can be justified for inclusion when the educational importance of the item content or the desire to measure students with very high or low achievement override the statistical considerations. Summary p -value information across all grades for each content area is shown in Tables 4.1.1 – 4.1.3. In general, most of the items fall into the p -value range of 0.4 to 0.9, which is appropriate for a criterion-referenced assessment.

Table 4.1.1 Summary of Proportion Correct for NeSA-R Operational Items

Grade	Item Proportion Correct										Mean	Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	≤0.7	≤0.8	≤0.9	>0.9		
3	0	0	0	0	2	7	17	14	5	0	0.679	45
4	0	0	0	0	4	7	12	17	4	1	0.680	45
5	0	0	0	0	4	5	11	19	6	3	0.706	48
6	0	0	0	0	3	6	14	13	11	1	0.703	48
7	0	0	0	1	2	6	15	15	9	0	0.693	48
8	0	0	0	0	5	7	10	21	6	1	0.688	50
11	0	0	0	0	3	9	14	10	12	2	0.694	50

Table 4.1.2 Summary of Proportion Correct for NeSA-M Operational Items

Grade	Item Proportion Correct										Mean	Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	≤0.7	≤0.8	≤0.9	>0.9		
3	0	0	0	0	4	11	12	15	7	1	0.676	50
4	0	0	0	0	3	13	14	16	9	0	0.677	55
5	0	0	0	0	3	10	14	18	9	1	0.690	55
6	0	0	0	0	4	11	19	12	12	0	0.673	58
7	0	0	0	1	3	11	19	20	4	0	0.667	58
8	0	0	0	0	7	9	19	20	5	0	0.661	60
11	0	0	0	3	4	11	22	16	4	0	0.643	60

Table 4.1.3 Summary of Proportion Correct for NeSA-S Operational Items

Grade	Item Proportion Correct										Mean	Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	≤0.7	≤0.8	≤0.9	>0.9		
5	0	0	0	1	4	9	9	19	8	0	0.679	50
8	0	0	0	0	6	14	15	17	8	0	0.661	60
11	0	0	0	1	5	12	17	16	9	0	0.666	60

4.2 ITEM-TOTAL CORRELATION

Item-total correlation describes the relationship between performance on the specific item and performance on the entire form. For the NeSA tests, Pearson’s product-moment correlation coefficient between item scores and test scores is used to indicate this relationship. For MC items, the statistic is typically referred to as point-biserial correlation. This index indicates an item’s ability to differentiate between high and low achievers (i.e., item discrimination power). It is expected that students with high ability (i.e., those who perform well on the NeSA overall) would be more likely to answer any given NeSA item correctly, while students with low ability (i.e., those who perform poorly on the NeSA overall) would be more likely to answer the same item incorrectly. However, an interaction can exist between item discrimination and item difficulty. Items answered correctly (or incorrectly) by a large proportion of examinees (i.e., the items have extreme *p*-values) can have reduced power to discriminate and thus can have lower correlations.

The correlation coefficient can range from -1.0 to +1.0. If the aforementioned expectation is met (high-scoring students tend to get the item right while low-scoring students do not), the correlation between the item score and the total test score will be both positive and noticeably large in its magnitude (i.e., well above zero), meaning the item is a good discriminator between high- and low-ability students. Items with negative correlations are flagged and referred to Test Development as possible mis-keys. Mis-keyed items will be corrected and rescored prior to computing the final item statistics. Negative

correlations can also indicate problems with the item content, structure, or students’ opportunity to learn. Items with point-biserial values of less than 0.2 are flagged and referred to content specialists for review before being considered for use on future forms. As seen below in Tables 4.2.1 – 4.2.3, no items in the 2015 NeSA tests have negative point-biserial correlations and most are above 0.30, indicating good item discrimination.

Table 4.2.1 Summary of Point-biserial Correlations for NeSA-R

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
3	0	0	3	16	23	3	0	45
4	0	0	6	26	13	0	0	45
5	0	0	6	19	23	0	0	48
6	0	0	5	21	20	2	0	48
7	0	0	1	19	24	4	0	48
8	0	0	8	22	18	2	0	50
11	0	0	2	26	17	5	0	50

Table 4.2.2 Summary of Point-biserial Correlations for NeSA-M

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
3	0	1	1	23	23	2	0	50
4	0	0	3	22	24	6	0	55
5	0	0	5	18	27	5	0	55
6	0	0	4	18	26	10	0	58
7	0	0	3	16	24	15	0	58
8	0	0	2	16	33	9	0	60
11	0	0	3	16	26	14	1	60

Table 4.2.3 Summary of Point-biserial Correlations for NeSA-S

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
5	0	2	11	22	15	0	0	50
8	0	1	6	29	24	0	0	60
11	0	0	7	28	23	2	0	60

4.3 PERCENT SELECTING EACH RESPONSE OPTION

This index indicates the effectiveness of each distractor. In general, one expects the correct response to be the most attractive, although this need not hold for unusually challenging items. This statistic for the correct response option is identical to the p -value when considering MC items with a single correct response. Please see the detailed summary statistics for each grade and content area in Appendices F, G, and H.

4.4 POINT-BISERIAL CORRELATIONS OF RESPONSE OPTIONS

This index describes the relationship between selecting a response option for a specific item and performance on the entire test. The correlation between an incorrect answer and total test performance should be negative. The desired pattern is strong positive values for the correct option and strong negative values for the incorrect options. Any other pattern indicates a problem with the item or with the key. These patterns would imply a high ability way to answer incorrectly or a low ability way to answer correctly. Examples of these situations could be an item with an ambiguous or misleading distractor that was attractive to high-performing examinees or an item that depended on experience outside of instruction that was unrelated to ability. This statistic for the correct option is identical to the item-total correlation for MC items. Please see the detailed summary statistics for each grade and content area in Appendices F, G, and H.

4.5 PERCENT OF STUDENTS OMITTING AN ITEM

This statistic is useful for identifying problems with testing time and test layout. If the omit percentage is large for a single item, it could indicate a problem with the layout or content of an item. For example, students tend to skip items with wordy stems or that otherwise appear difficult or time consuming. While there is no hard and fast rule for what *large* means, and it varies with groups and ages of students, five percent omits is often used as a preliminary screening value.

Detailed results of the item analyses for the NeSA-R operational items are presented in Appendix F. Detailed results of the item analyses for the NeSA-M operational items are presented in Appendix G. Detailed results of the item analyses for the NeSA-S operational items are presented in Appendix H. Based on these analyses, items were selected for review if the p -value was less than 0.25 and the item-total correlation was less than 0.2. Items were identified as probable mis-keys if the p -value for the correct response was less than one of the incorrect responses and the item-total correlation was negative.

5. RASCH ITEM CALIBRATION

The psychometric model used for the NeSA is based on the work of Georg Rasch (1960). Rasch models have had a long-standing presence in applied testing programs and have been the methodology used to calibrate NeSA items in recent history. Rasch models have several advantages over true-score test theory, so it has become the standard procedure for analyzing item response data in large-scale assessments. However, Rasch models have a number of strong requirements related to dimensionality, local independence, and model-data fit. Resulting inferences derived from any application of Rasch models rests strongly on the degree to which the underlying requirements are met.

Generally, item calibration is the process of estimating a difficulty-parameter estimate to each item on an assessment so that all items are placed onto a common scale. This chapter briefly introduces the Rasch model, reports the results from evaluations of the adequacy of the Rasch requirements, and summarizes Rasch item statistics for the 2015 NeSA Reading, Mathematics, and Science assessments.

5.1 DESCRIPTION OF THE RASCH MODEL

The Rasch dichotomous model was used to calibrate the NeSA items. All NeSA assessments contain only MC items. According to the Rasch model, the probability of answering an item correctly is based on the difference between the ability of the student and the difficulty of the item. The Rasch model places both student ability and item difficulty (estimated in terms of log-odds, or logits) on the same continuum. When the model assumptions are met, the Rasch model provides estimates of a person's ability that are independent of the items employed in the assessment and conversely, estimates item difficulty independently of the sample of examinees (Rasch, 1960; Wright & Panchapakesan, 1969). (As noted in Chapter Four, interpretation of item p -values confounds item difficulty and student ability.) Appendix I contains a more detailed overview of Rasch measurement.

5.2 CHECKING RASCH ASSUMPTIONS

Since the Rasch model was the basis of all calibration, scoring, and scaling analyses associated with the NeSA, the validity of the inferences from these results depends on the degree to which the assumptions of the model were met and how well the model fits the test data. Therefore, it is important to check these assumptions. This section evaluates the dimensionality of the data, local item independence, and item fit. It should be noted that only operational items were analyzed since they are the basis of student scores.

Unidimensionality: Rasch models assume that one dominant dimension determines the difference among students' performances. Principal components analysis (PCA) can be used to assess the unidimensionality assumption. The purpose of the analysis is to verify whether any other dominant component(s) exist among the items. If any other dimensions are found, the unidimensionality assumption would be violated.

Tables 5.2.1, 5.2.2, and 5.2.3 present the PCA results for the reading, mathematics, and science assessments, respectively. The results include the eigenvalues and the percentage of variance explained for the first five components. As can be seen in Table 5.2.1, the primary dimension for NeSA-R explained about 22 percent to 26 percent of the total variance across Grades 3–8 and 11. The eigenvalues of the second dimension ranged from 1.4 to 1.7. This indicates that the second dimension accounted for only 1.4 to 1.7 units out of 66 - 84 units of total variance. Similar patterns are observed for the Mathematics and the Science test. Overall, the PCA suggests that there is one clearly dominant dimension for each NeSA assessment.

Table 5.2.1 Results from PCA – Reading

Grade	Component	Eigenvalue	Explained Variance
3	measures	14.1	23.9%
	1	1.6	3.6%
	2	1.3	2.9%
	3	1.3	2.9%
	4	1.2	2.8%
	5	1.2	2.6%
4	measures	12.9	22.2%
	1	1.4	3.2%
	2	1.3	2.9%
	3	1.3	2.8%
	4	1.2	2.7%
	5	1.2	2.7%
5	measures	15.2	24.0%
	1	1.8	3.7%
	2	1.3	2.8%
	3	1.3	2.7%
	4	1.3	2.6%
	5	1.2	2.5%
6*	measures	14.4	23.1%
	1	1.5	3.2%
	2	1.3	2.8%
	3	1.3	2.7%
	4	1.2	2.5%
	5		
7*	measures	16.4	25.4%
	1	1.5	3.1%
	2	1.4	2.9%
	3	1.3	2.6%
	4	1.2	2.4%
	5		
8	measures	14.6	22.7%
	1	1.7	3.3%
	2	1.4	2.8%
	3	1.3	2.6%
	4	1.3	2.5%
	5	1.3	2.5%

Grade	Component	Eigenvalue	Explained Variance
11	measures	17.0	25.3%
	1	1.8	3.5%
	2	1.4	2.8%
	3	1.2	2.5%
	4	1.2	2.4%
	5	1.2	2.3%

*Only four components with eigenvalues greater than one were extracted.

Table 5.2.2 Results from PCA – Mathematics

Grade	Component	Eigenvalue	Explained Variance
3*	measures	17.0	25.4%
	1	1.6	3.2%
	2	1.4	2.8%
	3		
	4		
	5		
4	measures	18.7	25.4%
	1	1.7	3.1%
	2	1.6	2.9%
	3	1.5	2.7%
	4	1.5	2.7%
	5	1.4	2.5%
5*	measures	17.3	24.0%
	1	1.7	3.1%
	2	1.5	2.8%
	3	1.3	2.4%
	4		
	5		
6	measures	19.5	25.2%
	1	1.9	3.2%
	2	1.6	2.7%
	3	1.4	2.5%
	4	1.3	2.2%
	5	1.2	2.1%

Grade	Component	Eigenvalue	Explained Variance
7	measures	20.5	26.1%
	1	1.9	3.3%
	2	1.6	2.7%
	3	1.4	2.4%
	4	1.3	2.3%
	5	1.3	2.2%
8	measures	20.8	25.8%
	1	1.7	2.8%
	2	1.6	2.7%
	3	1.5	2.6%
	4	1.4	2.4%
	5	1.4	2.3%
11	measures	21.2	26.1%
	1	2.0	3.3%
	2	1.7	2.9%
	3	1.6	2.6%
	4	1.4	2.3%
	5	1.2	2.1%

Table 5.2.3 Results from PCA – Science

Grade	Component	Eigenvalue	Explained Variance
5*	measures	14.0	21.9%
	1	1.7	3.3%
	2	1.5	2.9%
	3	1.3	2.5%
	4		
	5		
8*	measures	17.4	22.5%
	1	1.6	2.6%
	2	1.4	2.4%
	3	1.3	2.2%
	4	1.3	2.1%
	5		

Grade	Component	Eigenvalue	Explained Variance
11	measures	17.4	22.5%
	1	1.7	2.9%
	2	1.5	2.4%
	3	1.4	2.3%
	4	1.3	2.2%
	5	1.2	2.1%

Local Independence: Local independence (LI) is a fundamental assumption of IRT. No relationship should exist between examinees’ responses to different items after accounting for the abilities measured by a test. Many indicators of LI are framed by the form of local independence proposed by McDonald (1979) that the conditional covariances of all pairs of item responses, conditioned on the abilities, are required to be equal to zero.

Residual item correlations provided in WINSTEPS for each item pair were used to assess local dependence among the NeSA items. Three types of residual correlations are available in WINSTEPS: raw, standardized, and logit. It should be noted that the raw score residual correlation essentially corresponds to Yen’s $Q3$ index, a popular LI statistic. The expected value for the $Q3$ statistic is approximately $-1/(k-1)$ when no local dependence exists, where k is test length (Yen, 1993). Thus, the expected $Q3$ values should be approximately -0.02 for the NeSA tests (since most of the NeSA tests had more than 50 core items). Index values that are greater than 0.20 indicate a degree of local dependence that probably should be examined by test developers (Chen & Thissen, 1997).

Since the three residual correlations are very similar, the default “standardized residual correlation” in WINSTEPS was used for these analyses. Tables 5.2.4 – 5.2.6 show the summary statistics—mean, SD, minimum, maximum, and several percentiles (P10, P25, P50, P75, P90)—for all the residual correlations for each test. The total number of item pairs (N) and the number of pairs with the residual correlations greater than 0.20 are also reported in this table. The mean residual correlations were slightly negative and the values were close to -0.02 . The vast majority of the correlations were very small; suggesting local item independence generally holds for the NeSA reading, mathematics, and science assessments.

Table 5.2.4 Summary of Item Residual Correlations for NeSA-R

Statistics	3	4	5	6	7	8	11
<i>N</i>	990	990	1128	1128	1128	1225	1225
Mean	-0.02	-0.02	-0.02	-0.02	-0.02	-0.02	-0.02
<i>SD</i>	0.02	0.02	0.02	0.02	0.02	0.02	0.02
Minimum	-0.07	-0.08	-0.08	-0.08	-0.08	-0.09	-0.09
P10	-0.05	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
P25	-0.04	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03
P50	-0.02	-0.02	-0.02	-0.02	-0.02	-0.02	-0.02
P75	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01	-0.01
P90	0.00	0.00	0.01	0.00	0.00	0.00	0.01
Maximum	0.12	0.15	0.19	0.13	0.08	0.25	0.09
>0.20	0	0	0	0	0	1	0

Table 5.2.5 Summary of Item Residual Correlations for NeSA-M

	Mathematics						
Statistics	3	4	5	6	7	8	11
<i>N</i>	1225	1485	1485	1653	1653	1770	1770
Mean	-0.02	-0.02	-0.02	-0.02	-0.02	-0.02	-0.02
<i>SD</i>	0.03	0.03	0.02	0.02	0.03	0.03	0.03
Minimum	-0.07	-0.08	-0.08	-0.11	-0.10	-0.08	-0.10
P10	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
P25	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03
P50	-0.02	-0.02	-0.02	-0.02	-0.02	-0.02	-0.02
P75	-0.01	-0.01	-0.01	0.00	-0.01	-0.01	0.00
P90	0.00	0.01	0.01	0.01	0.01	0.01	0.01
Maximum	0.39	0.45	0.23	0.15	0.31	0.36	0.37
>0.20	2	3	1	0	2	4	4

Table 5.2.6 Summary of Item Residual Correlations for NeSA-S

Statistics	Science		
	5	8	11
<i>N</i>	1225	1770	1770
Mean	-0.02	-0.02	-0.02
<i>SD</i>	0.02	0.02	0.02
Minimum	-0.08	-0.06	-0.10
P10	-0.04	-0.04	-0.04
P25	-0.03	-0.03	-0.03
P50	-0.02	-0.02	-0.02
P75	-0.01	-0.01	-0.01
P90	0.00	0.00	0.00
Maximum	0.45	0.27	0.26
>0.20	1	2	2

Item Fit: WINSTEPS provides two item fit statistics (infit/weighted and outfit/unweighted) for evaluating the degree to which the Rasch model predicts the observed item responses. Each fit statistic can be expressed as a mean square (MnSq) statistic with each statistic having an expected value of 1 and a different variance for each mean square or as a standardized statistic (Zstd with an expected mean = 0 and expected variance = 1).

MnSq values are more difficult to interpret due to an asymmetrical distribution and unique variance, while Zstd values are more oriented toward standardized statistical significance. Though both are informative, the Zstd values are less likely to be sensitive to the large sample sizes and have better distributional properties (Smith, Schumacker, & Bush, 1998). In the case of the NeSA-RMS, the sample sizes can be considered large ($n > 5,000$). The outfit statistic tends to be affected more by unexpected responses far from the person, item, or rating scale category measure (i.e., it is more sensitive to outlying, off-target, and low information responses that are very informative with regard to fit). The infit statistic tends to be affected more by unexpected responses close to the person, item, or rating scale category measure (i.e., with more information, but contributing little to the understanding of fit).

The expected MnSq value is 1.0 and can range from 0 to positive infinity. Values greater than 1.0 can be interpreted as indicating the presence of noise or lack of fit between the responses and the measurement model. Values less than 1.0 can be interpreted as item consistency or overfitting (too predictable and/or too much redundancy). Rules of thumb regarding “practically significant” MnSq values vary from author to author. More conservative users might prefer items with MnSq values that range from 0.8 to 1.2. Others believe reasonable test results can be achieved with values from 0.5 to 1.5. In the results below, values outside of 0.7 to 1.3 are given practical importance.

The expected Zstd value is 0.0 with an expected *SD* of 1.0 and can effectively range from -9.99 to +9.99 in WINSTEPS. Fit values greater than 0.0 can be interpreted as indicating the presence of noise or lack of fit between the items and the model (underfitting). Fit values lower than 0.0 can be interpreted as item redundancy or overfitting items (too predictable and/or too much redundancy). Rules of thumb regarding “practically significant” Zstd values vary from author to author. More conservative users might prefer items with Zstd values that range from -2 to +2. Others believe reasonable test results can be achieved with values from -3 to +3. In the results below, values outside of -3 to +3 are given practical importance.

Table 5.2.7 lists the summary statistics of infit and outfit mean square statistics for the NeSA reading, mathematics, and science tests, including the mean, *SD*, and minimum and maximum values. The number of items within the range of [0.7, 1.3] is also reported in Table 5.2.7. As can be seen, the mean values for both fit statistics were close to 1.00 for all tests. Most of the items had infit values falling in the range of [0.7, 1.3]. Though more outfit values fell outside this range than infit values, it is not surprising given that the infit statistic mutes the effects of anomalous response by extreme students.

Table 5.2.8 lists the summary statistics of infit and outfit Zstd statistics for the NeSA reading, mathematics, and science tests, including the mean, *SD*, and minimum and maximum values. The number of items within the range of [-3, +3] is also reported in Table 5.2.8. As can be seen, the mean values for both fit statistics were variable, ranging from -1.22 to 0.65. The fact that 16 of the 17 infit means were negative and 15 of the 17 outfit means were negative suggests that on average the data overfit the Rasch model, i.e. the data were a bit more consistent than expected by the probabilistic model.

Table 5.2.7 Summary of Infit and Outfit Mean Square Statistics for 2014 NeSA Tests

		Infit Mean Square					Outfit Mean Square				
		Mean	<i>SD</i>	MIN	MAX	[0.7, 1.3]	Mean	<i>SD</i>	MIN	MAX	[0.7, 1.3]
Reading	3	1.00	0.08	0.87	1.19	45/45	0.98	0.14	0.76	1.30	44/45
	4	1.00	0.07	0.87	1.15	45/45	0.99	0.13	0.66	1.28	44/45
	5	1.00	0.09	0.85	1.15	48/48	0.96	0.19	0.53	1.28	45/48
	6	1.00	0.09	0.82	1.22	48/48	0.97	0.17	0.60	1.38	40/48
	7	0.99	0.09	0.84	1.28	48/48	0.97	0.19	0.54	1.42	40/48
	8	1.00	0.08	0.83	1.20	50/50	0.98	0.15	0.62	1.29	49/50
	11	0.99	0.09	0.80	1.16	50/50	0.97	0.20	0.51	1.38	41/50
Mathematics	3	1.00	0.07	0.87	1.13	50/50	1.00	0.15	0.70	1.53	48/50
	4	1.00	0.08	0.84	1.18	55/55	0.99	0.16	0.68	1.36	52/55
	5	1.00	0.08	0.86	1.19	55/55	1.00	0.16	0.72	1.36	53/55
	6	1.00	0.10	0.86	1.30	57/58	0.99	0.17	0.66	1.48	56/58
	7	1.00	0.10	0.82	1.26	58/58	0.98	0.19	0.61	1.37	52/58
	8	0.99	0.10	0.82	1.28	60/60	0.99	0.17	0.65	1.43	54/60
	11	1.00	0.11	0.82	1.27	60/60	0.99	0.19	0.62	1.42	53/60

		Infit Mean Square					Outfit Mean Square				
		Mean	SD	MIN	MAX	[0.7, 1.3]	Mean	SD	MIN	MAX	[0.7, 1.3]
Science	5	1.00	0.08	0.88	1.19	50/50	0.98	0.15	0.65	1.29	49/50
	8	1.00	0.08	0.89	1.20	60/60	0.99	0.14	0.73	1.46	59/60
	11	1.00	0.08	0.86	1.15	60/60	0.98	0.14	0.71	1.27	60/60

Table 5.2.8 Summary of Infit and Outfit Z STD Statistics for 2014 NeSA Tests

		Infit Z STD					Outfit Z STD				
		Mean	SD	MIN	MAX	[-3.0, 3.0]	Mean	SD	MIN	MAX	[-3.0, 3.0]
Reading	3	-0.78	7.29	-9.99	9.99	8/45	-1.08	7.49	-9.99	9.99	9/45
	4	-0.59	6.35	-9.99	9.99	14/45	-0.50	7.01	-9.99	9.99	11/45
	5	-0.45	7.92	-9.99	9.99	6/48	-1.22	8.29	-9.99	9.99	8/48
	6	-0.31	7.24	-9.99	9.99	6/48	-1.04	7.24	-9.99	9.99	8/48
	7	-0.35	7.41	-9.99	9.99	13/48	-0.39	7.25	-9.99	9.99	11/48
	8	-0.75	7.14	-9.99	9.99	13/50	-0.83	7.49	-9.99	9.99	15/50
	11	0.16	7.29	-9.99	9.99	16/50	0.65	7.42	-9.99	9.99	12/50
Mathematics	3	-0.29	7.41	-9.99	9.99	11/50	0.25	7.15	-9.99	9.99	12/50
	4	-0.23	7.53	-9.99	9.99	9/55	-0.08	7.40	-9.99	9.99	12/55
	5	-0.89	7.34	-9.99	9.99	14/55	-0.52	7.38	-9.99	9.99	12/55
	6	-1.14	7.90	-9.99	9.99	11/58	-1.09	7.82	-9.99	9.99	10/58
	7	-0.52	8.30	-9.99	9.99	8/58	-0.59	7.83	-9.99	9.99	7/58
	8	-0.71	7.93	-9.99	9.99	9/60	-0.26	7.82	-9.99	9.99	6/60
	11	-0.54	7.91	-9.99	9.99	9/60	-0.67	7.42	-9.99	9.99	15/60
Science	5	-0.21	7.76	-9.99	9.99	11/50	-0.31	7.87	-9.99	9.99	11/50
	8	-0.92	7.53	-9.99	9.99	15/60	-0.96	7.66	-9.99	9.99	11/60
	11	-0.17	7.71	-9.99	9.99	7/60	-0.57	7.70	-9.99	9.99	6/60

5.3 RASCH ITEM STATISTICS

WINSTEPS 3.90.0 program (Linacre, 2015) was used for item calibration. The characteristics of calibration samples are reported in Chapter Three. These samples only include the students who attempted the tests. All omits (no response) and multiple responses (more than one response selected) were scored as incorrect answers (coded as 0s) for calibration.

As noted earlier, the Rasch model expresses item difficulty (and student ability) in units referred to as *logits* rather than on the proportion-correct metric. Large negative logits represent easier items while large positive logits represent more difficult items. The logit metrics is an interval scale,

meaning that two items with logit difficulties of 0.0 and +1.0 have the same difference in difficulty as two items with logit difficulties of +3.0 and +4.0.

Appendices J, K, L, and M report the Rasch calibration summaries and logit difficulties for all the operational items. Table 5.3.1 summarizes the Rasch logit difficulties of the operational items on each test. The minimum and maximum values and standard deviations suggest that the NeSA items covered a relatively wide range of difficulties. It is important to note that the logit difficulty values presented have not been linked to a common scale of measurement. Therefore, the relative magnitude of the statistics across subject areas and grades cannot be compared. The item pool was then updated with the item statistics.

Table 5.3.1 Summary of Rasch Item Difficulties for NeSA-R, NeSA-M, and NeSA-S

	Grade	N	Mean	SD	Min	Max	Range
Reading	3	45	-0.36	0.59	-1.72	0.70	2.42
	4	45	-0.50	0.69	-2.22	1.02	3.25
	5	48	-0.45	0.72	-1.87	1.14	3.01
	6	48	-0.59	0.73	-2.30	0.77	3.07
	7	48	-0.48	0.71	-1.98	1.03	3.01
	8	50	-0.62	0.65	-1.43	1.05	2.49
	11	50	-0.89	0.79	-2.60	0.34	2.94
Mathematics	3	50	-0.67	0.76	-2.34	0.74	3.09
	4	55	-0.65	0.76	-2.11	0.73	2.83
	5	55	-0.70	0.71	-2.36	0.95	3.31
	6	58	-0.63	0.66	-2.27	0.76	3.03
	7	58	-0.57	0.64	-1.96	1.07	3.03
	8	60	-0.67	0.69	-1.48	0.65	2.12
	11	60	-0.63	0.64	-2.09	1.17	3.26
Science	5	50	-0.78	0.73	-2.15	0.61	2.75
	8	60	-0.70	0.70	-2.14	0.53	2.67
	11	60	-0.70	0.65	-2.09	0.74	2.83

6. EQUATING AND SCALING

As discussed earlier in Chapter 2, the 2015 test forms were constructed with items that were either field tested, or used operationally on a previously administered NeSA test. NeSA assessments are constructed each year allowing each NeSA assessment to be different from the previous year's assessment. To ensure that all forms for a given grade and content area provide comparable scores, and to ensure the passing standards across different administrations are equivalent, the new operational items need to be placed on the bank scale via equating to bring the 2015 NeSA raw-score-to-Rasch-ability scale to the previous operational scale. When the new 2015 NeSA tests are placed on the bank's scale, the resulting scale scores for the new test form will be the same as the scale scores of the previous operational form such that students performing at the same level of (underlying) achievement should receive the same score (i.e., scale score). The resulting scale scores will be used for score reporting and performance level classification. Once operational items are equated, field test items are then placed on the bank scale and are then ready for future operational use.

This chapter begins with a summary of the entire NeSA equating procedures. This is followed by a scaling analysis that transforms raw scores to scale scores that represent the same skill level on every test form. Some summary results of the state scale score performance are also provided.

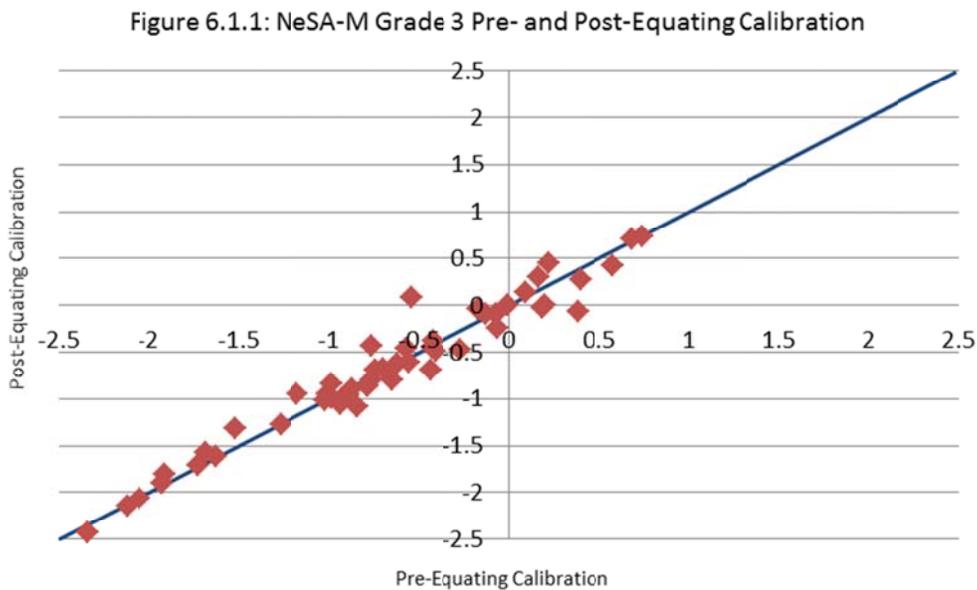
6.1 EQUATING

The equating design employed for NeSA is often referred to as a common-item non-equivalent groups (CINEG) design, which uses a set of anchor items that appear on two forms to adjust for differences in test difficulty across years. As discussed earlier, the 2015 NeSA test forms were constructed with items from previous administrations. The items were previously either field-test or operational items. If the item difficulty estimated from the previous administrations are within estimation error for the current administration, the entire set of the 2015 NeSA operational items can serve as the linking set. This means that the raw to scale score conversion tables can be established prior to the operational administration. This is often referred to as the pre-equating process because it is conducted before the operational test is administered. The most appealing feature of the pre-equating process, when applicable, is its ability to facilitate immediate score reporting for tests which have tight reporting windows.

However, it may not be appropriate to assume that the operational items will maintain their relative difficulty across administrations. The same item can perform differently across administrations due to changes in the item's position or changes in the students' experiences. Once the 2015 operational test data was available, DRC Psychometric Services staff, together with NDE, evaluated the item difficulty equivalence using a post-equating check procedure (Robust Z) to identify items that show significant difficulty changes from the bank values. If no unstable items are identified, the 2015 equating process would result in the pre-equating solution. On the other hand, if an item or items are found to be outside the normal estimation error, a post-equated solution would be used. The sub-set of 2015 operational

items, with those identified items excluded, was used as the set to estimate the link constant to map the 2015 test to the bank scale. This equating process is known as the post-equating because the equating occurs after the administration of the operation test and the raw-to-scale-score conversion is generated based on the operational test data.

As part of the post-equating check procedures, DRC Psychometric Services staff evaluated the item difficulty equivalence by comparing the old banked item calibration (called pre-calibration) with a new unanchored calibration of the 2015 data (called post-calibration). The evaluations were conducted for each grade and content area, using both visual graphing and statistical methods. The post-calibrated item difficulties (logits) were plotted against the pre-calibration for each grade and content area (see Appendices N – P). Ideally, these scatter plots should have a strong linear trend, closely clustered about the unit slope line. Items straying from the trend line did not perform in the same way in both administrations. The figure below illustrates an example of pre- and post-calibration plots for the 2015 NeSA-M test (Grade 3). Graphically, there is one apparent outlier item significantly above the unit slope line. It is located at the center of the scale, above -0.5 on the x -axis. This item has a Robust Z value of -5.782 , which is above the critical value of ± 1.645 . This item is harder for the 2015 population. All the other items fall, more or less, on the unit slope line, indicating consistent performance (within estimation error) in both years.



DRC Psychometric Services employed the Robust Z statistic (Huynh, 2000; Huynh & Rawls, 2009) for the post-equating check. This method focuses on the correlations between the pre- and post-calibrated item difficulties, and the ratio of standard deviations (SD) between the two calibrations. The correlation between the two estimates of item difficulty should be 0.95 or higher and the ratio of

standard deviations between the two sets of estimates of the item difficulty should range between 0.90 and 1.10 (Huynh & Meyer, 2010). To detect inconsistent item difficulty estimates, a critical value for the Robust Z statistic of ± 1.645 was used. The outlier identified in Figure 6.1.1 was detected using the Robust Z statistic.

Table 6.1.1 contains these statistics of correlation and SD ratio for the 2015 NeSA-M test. The item difficulty correlation for Grade 5 is the only statistic that falls below the criteria defined above. Appendices N – P contain the same statistics for each grade and content combination.

Table 6.1.1 NeSA-M Pre- and Post-Equating Comparison

	Grade						
	3	4	5	6	7	8	11
Correlation	0.98	0.97	0.93*	0.95	0.97	0.96	0.97
SD pre	0.76	0.76	0.67	0.66	0.61	0.69	0.64
SD post	0.75	0.74	0.67	0.68	0.61	0.64	0.64
SD Ratio	0.98	0.98	1.00	1.03	1.01	0.92	0.99

*The Grade 5 correlation was the only value that didn't meet the Robust Z criteria

Across all three content areas, the test forms with values below the ideal ranges of Robust Z correlation, or SD ratio values were further evaluated by the NDE in determining whether to include items that exceeded the Robust Z critical value of ± 1.645 in the linking set used for the post-equating. Items that exceeded the Robust Z critical value were then deleted, one item at a time, until both the item difficulty correlation and the SD ratio fell within the prescribed limits.

To summarize the 2015 NeSA test equating solutions, NDE decided to adopt a post-equating results for NeSA-M Grade 5 and all NeSA-R grades. For these tests, test equating was adjusted by excluding the items exceeding the critical value until the Robust Z criteria were met. A new raw-to-scale-score conversion table was created for these tests. For the other grades and content areas, NDE decided to use a pre-equating solution, keep the whole set of operational items in the linking set and then apply to the existing raw-to-scale-score conversion table.

6.2 SCALING

The purpose of a scaling analysis is to create a scale score for test reporting. The basic score on any test is the raw score, which is the number of items answered correctly or the total score points earned. However, the raw score alone does not present a wide-ranging picture of test performance because it is not on an equal-interval scale and can be interpreted only in terms of a particular set of items. Since a given raw score may not represent the same skill level on every test form, scale scores were assigned

to each raw score point to adjust for slight shifts in item difficulties and permit valid comparison across all test administrations within a particular content area.

Defining the scale score metric is an important, albeit arbitrary, step. Mathematically, scale scores are a linear transformation of the logit scores and thus do not alter the relationships or the displays. Scale scores are the numbers that will be reported to describe the performance of the students, schools, and systems. They will define the ranges of the performance levels, appear on individual student reports and school accountability analyses, and be dissected in newspaper accounts.

Appendix Q contains the detailed raw-score-to-scale-score conversion tables that were used to assign scale scores to students based on the total number correct scores from the NeSA-R for 2015, Appendix R for NeSA-M for 2015 and Appendix S for NeSA-S 2015. Because the relationship between raw and scale scores depends on the difficulties of the specific items on the form, these tables will change for every operational form.

There are two primary considerations when establishing the metric:

- Multiply the logit by a value large enough to make decimal points unnecessary for student scores, and
- Shift the scale enough to avoid negative values for low scale scores.

The scale chosen, for all grades and content areas of the NeSA assessment, range from 0 to 200. The value of 0 is reserved for students who were not tested or were otherwise invalidated. Thus, any student who attempted the test will receive a scale score equal to 1 even if the student gave no correct responses. No student tested will receive a scale score higher than 200 or lower than 1, even if this requires constraining the scale score calculation. It is possible that a future form will be easy enough that the upper limit of 200 is not invoked even for a perfect paper or could be difficult enough that the lower limit is not invoked.

As part of its deliberations concerning defining the performance levels, the State Board of Education specified that the *Meets the Standards* performance level have a scale score of 85 and that the *Exceeds the Standards* level have a scale score of 135. The logit standards defining the performance levels were adopted by the SBE per the standard setting and standard validation completed in 2010 for NeSA-R, in 2011 for NeSA-M, and in 2012 for NeSA-S.

Complete documentation of all standard setting events are presented in separate documents and are placed on the Nebraska State Department of Education website labeled:

2010 NeSA-Reading Standard Setting Technical Report,

http://www.education.ne.gov/Assessment/pdfs/2010_NeSA_Reading_Standard_Setting_Tech_%20Report.pdf ,

2011 NeSA-Mathematics Standard Setting Technical Report,

http://www.education.ne.gov/Assessment/pdfs/2011_NeSA_Math_Standard_Setting_Tech_Report.pdf

Table 6.2.2 NeSA-M Conversion of Logits to Scale Scores

Grade	Logit Cut Points		Scale Score Ranges by Performance Level			Conversion	
	B/M	M/E	Below	Meets	Exceeds	Slope b	Intercept a
3	-0.6	1.1000	1 to 84	85-134	135 to 200	29.41176	102.15706
4	-0.6	1.2000	1 to 84	85-134	135 to 200	27.77778	101.17667
5	-0.57	1.1597	1 to 84	85-134	135 to 200	28.90675	100.98685
6	-0.47	1.1816	1 to 84	85-134	135 to 200	30.27367	98.73862
7	-0.45	1.2500	1 to 84	85-134	135 to 200	29.41176	97.74529
8	-0.4	1.3000	1 to 84	85-134	135 to 200	29.41176	96.2747
11	-0.29	1.1000	1 to 84	85-134	135 to 200	35.97122	94.94165

Table 6.2.3 NeSA-S Conversion of Logits to Scale Scores

Grade	Logit Cut Points		Scale Score Ranges by Performance Level			Conversion	
	B/M	M/E	Below	Meets	Exceeds	Slope b	Intercept a
5	-0.4971	1.0580	1 to 84	85-134	135 to 200	32.15095	100.49331
8	-0.4543	1.0378	1 to 84	85-134	135 to 200	33.50958	99.73252
11	-0.5407	1.3130	1 to 84	85-134	135 to 200	26.97256	99.09502

Complete frequency distributions of the state scale scores for the NeSA-R, NeSA-M, and NeSA-S are provided in Appendices Q, R, and S as part of the raw-to-scale-score conversion tables. A simple summary of the reading, mathematics, and science distributions can be found in Tables 6.2.4, 6.2.5, and 6.2.6.

Table 6.2.4 2014 NeSA-R State Scale Score Summary, All Students

Grade	Count	Scale Score		Quartile		
		Mean	S.D.	First	Second	Third
3	23013	118.6	35.6	93	117	141
4	22590	121.0	39.2	94	121	146
5	22878	129.0	43.3	98	132	161
6	22377	121.7	41.6	93	120	152
7	22212	128.0	44.3	97	129	161
8	21904	117.4	39.5	89	118	146
11	21125	109.7	45.8	80	113	139

Table 6.2.5 2015 NeSA-M State Scale Score Summary, All Students

Grade	Count	Scale Score		Quartile		
		Mean	S.D.	First	Second	Third
3	23130	113.1	36.0	88	111	137
4	22685	112.0	33.9	87	111	133
5	22946	113.4	36.2	87	109	136
6	22445	110.0	37.3	82	109	135
7	22285	110.0	38.8	81	106	135
8	21985	105.0	38.9	77	102	128
11	21107	102.5	46.0	67	99	135

Table 6.2.6 2015 NeSA-S State Scale Score Summary, All Students

Grade	Count	Scale Score		Quartile		
		Mean	S.D.	First	Second	Third
5	22949	107.1	33.3	84	106	129
8	21993	106.1	36.3	79	105	131
11	21102	104.7	29.3	84	105	124

7. FIELD TEST ITEM DATA SUMMARY

As noted in Chapter Two, in addition to the operational items, field test items were embedded in all content areas and grade level assessments in order to expand the item pool for future form development. Field test items are items being administered for the first time to gather statistical information. These items do not count toward an individual student’s score. All field tested items were analyzed statistically following classical item analysis methods including proportion correct, point-biserial correlation, and DIF.

7.1 CLASSICAL ITEM STATISTICS

Indices known as classical item statistics included the item p -value and the point-biserial correlations for MC items. For MC items, the p -value reflects the proportion of students who answered the item correctly. In general, more capable students are expected to respond correctly to easy items and less capable students are expected to respond incorrectly to difficult items. The primary way of detecting such conditions is through the point-biserial correlation coefficient for dichotomous (MC) items. The point-biserial correlation will be positive if the total test mean score is higher for the students who respond correctly to MC items and negative when the reverse is true.

The traditional statistics are computed for each NeSA-R field test item in Appendix F, for NeSA-M in Appendix G and for NeSA-S in Appendix H. Tables 7.1.1, 7.1.2, and 7.1.3 provide summaries of the distributions of item proportion correct and point-biserial correlations. For future form construction, items with negative point-biserial correlations are never considered for operational use. Items with correlations less than 0.2 or proportion correct less than 0.3 or greater 0.9 are avoided when possible.

Table 7.1.1 Summary of Statistics for NeSA-R 2015 Field Test Items

Grade	Item Proportion Correct										Mean	Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	≤0.7	≤0.8	≤0.9	>0.9		
3	0	0	5	0	12	12	10	5	6	0	0.573	50
4	0	0	1	6	7	9	12	11	4	0	0.601	50
5	0	0	2	1	6	7	11	10	10	3	0.661	50
6	0	1	0	3	5	8	12	15	5	1	0.648	50
7	0	0	0	6	7	9	9	11	8	0	0.619	50
8	0	1	1	4	6	4	8	18	5	3	0.644	50
11	0	0	1	1	6	9	19	9	4	1	0.637	50

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
3	1	4	7	25	11	2	0	50
4	1	4	8	24	12	1	0	50
5	0	5	7	20	17	1	0	50
6	1	4	7	29	8	1	0	50
7	0	5	14	13	15	3	0	50
8	4	1	13	12	18	2	0	50
11	4	2	9	13	20	2	0	50

Table 7.1.2 Summary of Statistics for NeSA-M 2015 Field Test Items

Grade	Item Proportion Correct										Mean	Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	≤0.7	≤0.8	≤0.9	>0.9		
3	0	0	2	2	3	7	9	10	13	4	0.694	50
4	1	0	0	0	2	2	11	15	10	9	0.745	50
5	0	0	1	0	4	10	10	11	12	2	0.687	50
6	0	0	0	2	4	6	13	10	10	5	0.699	50
7	0	0	2	5	8	9	7	15	4	0	0.595	50
8	1	0	3	2	9	14	11	6	4	0	0.574	50
11	0	1	3	3	11	12	14	5	1	0	0.544	50

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
3	1	2	7	19	16	5	0	50
4	1	2	9	21	15	2	0	50
5	0	1	9	12	19	9	0	50
6	0	2	4	18	15	11	0	50
7	2	1	3	13	20	11	0	50
8	2	1	5	18	19	5	0	50
11	0	2	6	13	20	8	1	50

Table 7.1.3 Summary of Statistics for NeSA-S 2015 Field Test Items

Grade	Item Proportion Correct										Mean	Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	≤0.7	≤0.8	≤0.9	>0.9		
5	0	1	1	0	3	6	3	14	14	8	0.737	50
8	0	2	2	3	5	5	8	12	9	4	0.651	50
11	0	3	5	9	6	2	12	9	4	0	0.532	50

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
5	1	5	12	21	9	2	0	50
8	7	2	11	17	12	1	0	50
11	4	10	14	17	4	1	0	50

7.2 DIFFERENTIAL ITEM FUNCTIONING

DIF occurs when examinees with the same ability level but different group memberships do not have the same probability of answering an item correctly. This pattern of results may suggest the presence of *item bias*. Items exhibiting DIF were referred to content specialists to determine possible bias. No statistical procedure should be used as a substitute for rigorous, hands-on reviews by content and bias specialists. The statistical results can help organize the review so the effort is concentrated on the most problematic cases. Further, no items should be automatically rejected simply because a statistical method flagged them or accepted because they were not flagged.

For MC items, the Mantel-Haenszel procedure (Mantel & Haenszel, 1959) for detecting DIF is a commonly used technique in educational testing. The procedure as implemented by DRC contrasts a focal group with a reference group. While it makes no practical difference in the analysis which group is defined as the focal group, the group most apt to be disadvantaged by a biased measurement is typically defined as the focal group. In these analyses, the focal group was female for gender-based DIF and minority for ethnicity-based DIF; reference groups were male and white, respectively.

To assist the review committees in interpreting the analyses, the items are assigned a severity code based on the magnitude of the MH statistic. Items classified as A+ or A- have little or no statistical indication of DIF. Items classified as B+ or B- have some indication of DIF but may be judged to be acceptable for future use. Items classified as C+ or C- have strong evidence of DIF and should be reviewed and possibly rejected from the eligible item pool. The plus sign indicates that the item favors the focal group and a minus sign indicates that the item favors the reference group. Tables 7.2.1 – 7.2.3 show summaries of the DIF statistics. The first column defines the focal group. Appendices T, U, and V provide more summary information on DIF analysis.

Table 7.2.1 Summary of DIF by Code for NeSA-R 2015 Field Test

Grade 3	A+	A-	B+	B-	C+	C-	FT Items
Female	27	23	0	0	0	0	50
Black	10	36	0	4	0	0	50
Hispanic	20	25	0	4	0	1	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	6	3	0	1	0	0	50
2 or more Races	4	5	0	1	0	0	50

Grade 4	A+	A-	B+	B-	C+	C-	FT Items
Female	21	28	0	1	0	0	50
Black	12	31	0	4	0	3	50
Hispanic	11	37	0	2	0	0	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	6	4	0	0	0	0	50

Grade 5	A+	A-	B+	B-	C+	C-	FT Items
Female	23	27	0	0	0	0	50
Black	5	34	0	7	0	4	50
Hispanic	13	33	0	3	0	1	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	4	6	0	0	0	0	50

Nebraska State Accountability 2015 Technical Report

Grade 6	A+	A-	B+	B-	C+	C-	FT Items
Female	31	18	0	1	0	0	50
Black	13	29	0	7	0	1	50
Hispanic	17	30	0	2	0	1	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	1	8	0	1	0	0	50

Grade 7	A+	A-	B+	B-	C+	C-	FT Items
Female	26	21	2	1	0	0	50
Black	8	34	1	7	0	0	50
Hispanic	15	31	0	3	0	1	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	0	0	0	0	0	0	50

Grade 8	A+	A-	B+	B-	C+	C-	FT Items
Female	32	17	1	0	0	0	50
Black	8	30	0	8	0	4	50
Hispanic	10	34	0	4	0	2	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	0	0	0	0	0	0	50

Grade 11	A+	A-	B+	B-	C+	C-	FT Items
Female	34	10	2	3	0	1	50
Black	12	32	0	4	0	2	50
Hispanic	20	21	0	8	0	1	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	0	0	0	0	0	0	50

Table 7.2.2 Summary of DIF by Code for NeSA-M 2015 Field Test

Grade 3	A+	A-	B+	B-	C+	C-	FT Items
Female	23	23	1	3	0	0	50
Black	8	29	0	10	0	3	50
Hispanic	14	26	1	9	0	0	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	2	6	1	1	0	0	50
2 or more Races	5	4	0	0	0	0	50

Grade 4	A+	A-	B+	B-	C+	C-	FT Items
Female	22	27	0	0	1	0	50
Black	6	30	0	12	0	2	50
Hispanic	18	30	0	1	0	1	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	1	2	0	0	0	0	50
2 or more Races	5	5	0	0	0	0	50

Nebraska State Accountability 2015 Technical Report

Grade 5	A+	A-	B+	B-	C+	C-	FT Items
Female	23	24	0	3	0	0	50
Black	10	25	2	9	0	4	50
Hispanic	20	24	0	3	0	3	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	1	0	1	0	0	50
2 or more Races	2	7	0	1	0	0	50

Grade 6	A+	A-	B+	B-	C+	C-	FT Items
Female	29	19	0	2	0	0	50
Black	12	23	0	11	0	4	50
Hispanic	15	28	0	7	0	0	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	3	3	0	0	0	0	50

Grade 7	A+	A-	B+	B-	C+	C-	FT Items
Female	20	28	1	0	0	1	50
Black	8	33	0	8	0	1	50
Hispanic	8	40	0	2	0	0	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	4	1	0	0	0	0	50

Grade 8	A+	A-	B+	B-	C+	C-	FT Items
Female	25	25	0	0	0	0	50
Black	12	30	0	5	0	2	50
Hispanic	13	35	0	2	0	0	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	1	0	0	0	0	0	50

Grade 11	A+	A-	B+	B-	C+	C-	FT Items
Female	21	28	1	0	0	0	50
Black	11	34	0	4	0	0	50
Hispanic	9	40	0	1	0	0	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	0	0	0	0	0	0	50

Table 7.2.3 Summary of DIF by Code for NeSA-S 2015 Field Test

Grade 5	A+	A-	B+	B-	C+	C-	FT Items
Female	15	32	2	1	0	0	50
Black	10	24	0	13	0	3	50
Hispanic	12	34	0	3	0	1	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	2	0	0	50
2 or more Races	3	4	0	1	0	0	50

Nebraska State Accountability 2015 Technical Report

Grade 8	A+	A-	B+	B-	C+	C-	FT Items
Female	17	27	3	3	0	0	50
Black	11	29	0	8	0	2	50
Hispanic	9	38	0	2	0	1	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	0	0	0	0	0	0	50

Grade 11	A+	A-	B+	B-	C+	C-	FT Items
Female	23	24	0	3	0	0	50
Black	14	27	1	7	0	1	50
Hispanic	17	28	0	5	0	0	50
American Indian/Alaskan Native	0	0	0	0	0	0	50
Asian	0	0	0	0	0	0	50
2 or more Races	0	0	0	0	0	0	50

8. RELIABILITY

This chapter addresses the reliability of NeSA-Alt test scores. According to Mehrens and Lehmann (1975) reliability is defined as:

... the degree of consistency between two measures of the same thing. (p. 88).

8.1 COEFFICIENT ALPHA

The ability to measure consistently is a necessary prerequisite for making appropriate interpretations (i.e., showing evidence of valid use of results). Conceptually, reliability can be referred to as the consistency of the results between two measures of the same thing. This consistency can be seen in the degree of agreement between two measures on two occasions. Operationally, such comparisons are the essence of the mathematically defined reliability indices.

All measures consist of an accurate, or true, component and an inaccurate, or error, component. Errors occur as a natural part of the measurement process and can never be eliminated entirely. For example, uncontrollable factors such as differences in the physical environment and changes in examinee disposition may increase error and decrease reliability. This is the fundamental premise of traditional reliability analysis and measurement theory. Stated explicitly, this relationship can be seen as the following:

$$\text{Observed Score} = \text{True Score} + \text{Error} \quad (8.1)$$

To facilitate a mathematical definition of reliability, these components can be rearranged to form the following ratio:

$$\text{Reliability} = \frac{\text{TrueScoreVariance}}{\text{ObservedScoreVariance}} = \frac{\text{TrueScoreVariance}}{\text{TrueScoreVariance} + \text{ErrorScoreVariance}} \quad (8.2)$$

When there is no error, the reliability is true score variance divided by true score variance, which equals 1. However, as more error influences the measure, the error component in the denominator of the ratio increases. As a result, the reliability decreases.

The reliability index used for the 2015 administration of the NeSA was the Coefficient Alpha α (Cronbach, 1951). Acceptable α values generally range in the mid to high 0.80s to low 0.90s. The total test Coefficient Alpha reliabilities of the whole population are presented in Table 8.1.1 for each grade and content area of the NeSA. The table contains test length in total number of items (L), test reliabilities, and traditional standard errors of measurement (SEM). As can be seen in the table, all reading, mathematics, and science forms for grades 3-11 have Coefficient Alphas in the high 0.80s or low 0.90s. Overall, these α values provide evidence of good reliability.

Table 8.1.1 Reliabilities and Standard Errors of Measurement

	Grade	<i>L</i>	Reliability	<i>SEM</i>
Reading	3	45	0.91	2.8
	4	45	0.89	2.8
	5	48	0.90	2.8
	6	48	0.90	2.8
	7	48	0.91	2.8
	8	50	0.90	2.9
	11	50	0.91	2.9
Mathematics	3	50	0.91	2.9
	4	55	0.92	3.0
	5	55	0.92	3.0
	6	58	0.93	3.1
	7	58	0.94	3.1
	8	60	0.94	3.2
	11	60	0.94	3.2
Science	5	50	0.89	3.0
	8	60	0.92	3.3
	11	60	0.92	3.3

Reliability estimates for subgroups based on gender, ethnicity, special education status, limited English proficiency status, and food program eligibility status are also computed and reported in Appendix W. Results show fairly high reliability indices for all subpopulations in the high 0.80s to low 0.90s across grades and content areas, which indicates that the NeSA is not only reliable for the population as a whole, but it is also reliable for subpopulations of interest under NCLB. Appendix X present α for the content strands. Given that α is a function of test length, the smaller item counts for the content standards result in lower values of α which is to be expected. Overall, these two sets of values provide evidence of good reliability.

8.2 STANDARD ERROR OF MEASUREMENT

The traditional *SEM* uses the information from the test along with an estimate of reliability to make statements about the degree to which error influences individual scores. The *SEM* is based on the premise that underlying traits, such as academic achievement, cannot be measured exactly without a perfectly precise measuring instrument. The standard error expresses unreliability in terms of the raw-score metric. The *SEM* formula is provided below:

$$SEM = SD\sqrt{1 - reliability}. \tag{8.3}$$

This formula indicates that the value of the *SEM* depends on both the reliability coefficient and the standard deviation of test scores. If the reliability were equal to 0.00 (the lowest possible value), the

SEM would be equal to the standard deviation of the test scores. If test reliability were equal to 1.00 (the highest possible value), the *SEM* would be 0.0. In other words, a perfectly reliable test has no measurement error (Harvill, 1991). *SEMs* were calculated for each NeSA grade and content area using raw scores and displayed in Table 8.1.1.

8.3 CONDITIONAL STANDARD ERROR OF MEASUREMENT (CSEM)

The preceding discussion reviews the traditional approach to judging a test's consistency. This approach is useful for making overall comparisons between alternate forms. However, it is not very useful for judging the precision with which a specific student's score is known. The Rasch measurement models provide "conditional standard errors" that pertain to each unique ability estimate. Therefore, the *CSEM* may be especially useful in characterizing measurement precision in the neighborhood of a score level used for decision-making—such as cut scores for identifying students who meet a performance standard.

The complete set of conditional standard errors for every obtainable score can be found in Appendices Q, R and S as part of the raw-to-scale-score conversions for each grade and content area. Values were derived using the calibration data file described in Chapter Six and are on the scaled score metric. The magnitudes of *CSEMs* across the score scale seemed reasonable for most NeSA tests that the values are lower in the middle of the score range and increase at both extremes (i.e., at smaller and larger scale scores). This is because ability estimates from scores near the center of the test scoring range are known much more precisely than abilities associated with extremely high or extremely low scores. Table 8.3.1 reports the minimum *CSEM* of the scale score associated with the zero total test score (Min *CSEM*), the maximum *CSEM* of the scale score associated with the perfect total test score (Max *CSEM*), *CSEM* at the cuts of Below and Meets performance levels (*CSEM* B/M), and *CSEM* at the cuts of Meets and Exceeds performance levels (*CSEM* M/E) for each grade and content area. *CSEM* values at the cut score were generally associated with smaller *CSEM* values, indicating that more precise measurement occurs at these cuts.

Table 8.3.1 CSEM of the Scale Scores for 2015 NeSA Tests

		Min	Max	CSEM	CSEM
	Grade	CSEM	CSEM	B/M	M/E
Reading	3	9	52	9	11
	4	11	67	11	14
	5	12	72	12	14
	6	12	69	12	14
	7	12	71	12	14
	8	11	68	11	14
	11	12	73	12	15
Mathematics	3	9	54	9	12
	4	8	51	8	11
	5	8	53	8	11
	6	8	55	8	11
	7	8	54	8	11
	8	8	54	8	11
	11	10	66	10	13
Science	5	10	59	10	14
	8	9	61	9	12
	11	7	49	7	11

8.4 DECISION CONSISTENCY AND ACCURACY

When criterion-referenced tests are used to place the examinees into two or more performance classifications, it is useful to have some indication of how accurate or consistent such classifications are. Decision consistency refers to the degree to which the achievement level for each student can be replicated upon retesting using an equivalent form (Huynh, 1976). Decision accuracy describes the extent to which achievement-level classification decisions based on the administered test form would agree with the decisions that would be made on the basis of a perfectly reliable test. In a standards-based testing program there should be great interest in knowing how consistently and accurately students are classified into performance categories.

Since it is not feasible to repeat NeSA testing in order to estimate the proportion of students who would be reclassified in the same achievement levels, a statistical model needs to be imposed on the data to project the consistency or accuracy of classifications solely using data from the available administration (Hambleton & Novick, 1973). Although a number of procedures are available, two well-known methods were developed by Hanson and Brennan (1990) and Livingston and Lewis (1995) utilizing specific true score models. These approaches are fairly complex, and the cited sources contain

details regarding the statistical models used to calculate decision consistency from the single NeSA administration.

Several factors might affect decision consistency. One important factor is the reliability of the scores. All other things being equal, more reliable test scores tend to result in more similar reclassifications. Another factor is the location of the cutscore in the score distribution. More consistent classifications are observed when the cutscores are located away from the mass of the score distribution. The number of performance levels is also a consideration. Consistency indices for four performance levels should be lower than those based on three categories because classification using four levels would allow more opportunity to change achievement levels. Finally, some research has found that results from the Hanson and Brennan (1990) method on a dichotomized version of a complex assessment yield similar results to the Livingston and Lewis method (1995) and the method by Stearns and Smith (2007).

The results for the overall consistency across all three achievement levels are presented in Tables 8.4.1 – 8.4.3. The tabled values, derived using the program *BB-Class* (Brennan, 2004), show that consistency values across the two methods are generally very similar. Across all content areas, the overall decision consistency ranged from the mid 0.80s to the low 0.90s while the decision accuracy ranged from the high 0.80s to the mid 0.90s. If a parallel test were administered, at least 85% or more of students would be classified in the same way. Dichotomous decisions using the Meets cuts (Below/Meets) generally have the highest consistency values and exceeded 0.90 in all cases. The pattern of decision accuracy across different cuts is similar to that of decision consistency.

Table 8.4.1 NeSA-R Decision Consistency Results

Content Area	Grade	Livingston & Lewis				Hanson & Brennan			
		Decision Accuracy		Decision Consistency		Decision Accuracy		Decision Consistency	
		Meets	Exceeds	Meets	Exceeds	Meets	Exceeds	Meets	Exceeds
Reading	3	0.94	0.92	0.91	0.88	0.94	0.92	0.92	0.89
	4	0.93	0.90	0.90	0.86	0.93	0.90	0.90	0.86
	5	0.94	0.91	0.92	0.87	0.94	0.91	0.92	0.87
	6	0.94	0.90	0.91	0.86	0.94	0.90	0.91	0.86
	7	0.94	0.91	0.92	0.88	0.94	0.91	0.92	0.88
	8	0.94	0.90	0.91	0.86	0.94	0.90	0.91	0.86
	11	0.93	0.91	0.90	0.87	0.93	0.91	0.90	0.87

Table 8.4.2 NeSA-M Decision Consistency Results

Content Area	Grade	Livingston & Lewis				Hanson & Brennan			
		Decision Accuracy		Decision Consistency		Decision Accuracy		Decision Consistency	
		Meets	Exceeds	Meets	Exceeds	Meets	Exceeds	Meets	Exceeds
Math	3	0.93	0.92	0.93	0.92	0.93	0.92	0.91	0.89
	4	0.94	0.93	0.94	0.93	0.94	0.93	0.91	0.90
	5	0.93	0.93	0.93	0.93	0.93	0.93	0.91	0.90
	6	0.94	0.93	0.94	0.93	0.94	0.93	0.91	0.90
	7	0.94	0.93	0.94	0.93	0.94	0.93	0.92	0.91
	8	0.94	0.94	0.94	0.94	0.93	0.94	0.91	0.92
	11	0.94	0.94	0.94	0.94	0.94	0.94	0.91	0.92

Table 8.4.3 NeSA-S Decision Consistency Results

Content Area	Grade	Livingston & Lewis				Hanson & Brennan			
		Decision Accuracy		Decision Consistency		Decision Accuracy		Decision Consistency	
		Meets	Exceeds	Meets	Exceeds	Meets	Exceeds	Meets	Exceeds
Science	5	0.92	0.92	0.88	0.89	0.92	0.92	0.88	0.89
	8	0.93	0.93	0.90	0.90	0.93	0.93	0.90	0.90
	11	0.93	0.93	0.90	0.90	0.93	0.93	0.91	0.90

9. VALIDITY

As defined in the *Standards for Educational and Psychological Testing* (AERA, APA, & NCME, 2014), “Validity refers to the degree to which evidence and theory support the interpretation of test scores for proposed uses of tests” (p. 11). The validity process involves the collection of a variety of evidence to support the proposed test score interpretations and uses. This entire technical report describes the technical aspects of the NeSA tests in support of their score interpretations and uses. Each of the previous chapters contributes important evidence components that pertain to score validation: test development, test scoring, item analysis, Rasch calibration, scaling, and reliability. This chapter summarizes and synthesizes the evidence based on the framework presented in *The Standards*.

9.1 EVIDENCE BASED ON TEST CONTENT

Content validity addresses whether the test adequately samples the relevant material it purports to cover. The NeSA for grades 3 through 11 is a criterion-referenced assessment. The criteria referenced are the Nebraska reading and mathematics content standards. Each assessment was based on and was directly aligned to the Nebraska statewide content standards to ensure good content validity.

For criterion-referenced, standards-based assessment, the strong content validity evidence is derived directly from the test construction process and the item scaling. The item development and test construction process, described above, ensures that every item aligns directly to one of the content standards. This alignment is foremost in the minds of the item writers and editors. As a routine part of item selection prior to an item appearing on a test form, the review committees check the alignment of the items with the standards and make any adjustments necessary. The result is consensus among the content specialists and teachers that the assessment does in fact assess what was intended.

The empirical item scaling, which indicates where each item falls on the logit ability-difficulty continuum, should be consistent with what theory suggests about the items. Items that require more knowledge, more advanced skills, and more complex behaviors should be empirically more difficult than those requiring less. Evidence of this agreement is contained in the item summary tables in Appendices K, L, and M, as well as the success of the Bookmark and Contrasting Groups standard setting processes (in the separate *2010 NeSA-R Standard Setting Technical Report*, *2011 NeSA-M Standard Setting Technical Report* and *2012 NeSA-S Standard Setting Technical Report*). Panelists participating in the Bookmark process work from an item booklet in which items are ordered by their empirical difficulties. Discussions about placement of the bookmarks almost invariably focus on the knowledge, skills, and behaviors required of each item, and, overall, panelists were comfortable with the item ordering and spacing. Contrasting Groups participants, using their knowledge and experience with their students, placed their students in a corresponding Performance Level.

9.2 EVIDENCE BASED ON INTERNAL STRUCTURE

As described in the *Standards* (2014), internal-structure evidence refers to the degree to which the relationships between test items and test components conform to the construct on which the proposed test interpretations are based.

Item-Test Correlations: Item-test correlations are reviewed in Chapter Four. All values are positive and of acceptable magnitude.

Item Response Theory Dimensionality: Results from principle components analyses are presented in Chapter Five. The NeSA reading, mathematics, and science tests were essentially unidimensional, providing evidence supporting interpretations based on the total scores for the respective NeSA tests.

Strand Correlations: Correlations and disattenuated correlations between strand scores within each content area are presented below. This data can also provide information on score dimensionality that is part of internal-structure evidence. As noted in Chapter Two and also in Table 9.2.1, the NeSA-R tests have two strands (denoted by R.1 and R.2), the NeSA-M tests have four strands (denoted by M.1, M.2, M.3, and M.4), and the NeSA-S have four strands (denoted by S.1, S.2, S.3, and S.4) for each grade and content area.

For each grade, Pearson’s correlation coefficients between these strands are reported in Tables 9.2.2.a through 9.2.2.g. The intercorrelations between the strands within the content areas are positive and generally range from moderate to high in value.

Table 9.2.1 NeSA Content Strands

Content	Code	Strand
Reading	R.1	Vocabulary
	R.2	Comprehension
Mathematics	M.1	Number Sense
	M.2	Geometric/Measurement
	M.3	Algebraic
	M.4	Data Analysis/Probability
Science	S.1	Inquiry, the Nature of Science, and Technology
	S.2	Physical Science
	S.3	Life Science
	S.4	Earth and Space Science

Table 9.2.2.a Correlations between Reading and Mathematics Strands for Grade 3

Grade 3	R.1	R.2	M.1	M.2	M.3	M.4
R.1	—					
R.2	0.79	—				
M.1	0.63	0.66	—			
M.2	0.61	0.63	0.70	—		
M.3	0.56	0.60	0.71	0.62	—	
M.4	0.59	0.63	0.69	0.60	0.62	—

Table 9.2.2.b Correlations between Reading and Mathematics Strands for Grade 4

Grade 4	R.1	R.2	M.1	M.2	M.3	M.4
R.1	—					
R.2	0.73	—				
M.1	0.60	0.67	—			
M.2	0.60	0.66	0.77	—		
M.3	0.53	0.58	0.72	0.67	—	
M.4	0.55	0.63	0.66	0.63	0.57	—

Table 9.2.2.c Correlations between Reading, Mathematics, and Science Strands for Grade 5

Grade 5	R.1	R.2	M.1	M.2	M.3	M.4	S.1	S.2	S.3	S.4
R.1	—									
R.2	0.76	—								
M.1	0.64	0.71	—							
M.2	0.57	0.61	0.72	—						
M.3	0.55	0.60	0.72	0.59	—					
M.4	0.61	0.67	0.72	0.64	0.64	—				
S.1	0.62	0.70	0.63	0.55	0.56	0.64	—			
S.2	0.60	0.64	0.60	0.56	0.52	0.59	0.63	—		
S.3	0.60	0.66	0.57	0.53	0.50	0.58	0.62	0.65	—	
S.4	0.57	0.62	0.57	0.53	0.50	0.56	0.60	0.64	0.65	—

Table 9.2.2.d Correlations between Reading and Mathematics Strands for Grade 6

Grade 6	R.1	R.2	M.1	M.2	M.3	M.4
R.1	—					
R.2	0.75	—				
M.1	0.59	0.67	—			
M.2	0.55	0.62	0.75	—		
M.3	0.60	0.68	0.77	0.70	—	
M.4	0.57	0.66	0.74	0.67	0.71	—

Table 9.2.2.e Correlations between Reading and Mathematics Strands for Grade 7

Grade 7	R.1	R.2	M.1	M.2	M.3	M.4
R.1	—					
R.2	0.79	—				
M.1	0.63	0.69	—			
M.2	0.56	0.63	0.75	—		
M.3	0.62	0.70	0.80	0.73	—	
M.4	0.56	0.63	0.71	0.65	0.71	—

Table 9.2.2.f Correlations between Reading, Mathematics, and Science Strands for Grade 8

Grade 8	R.1	R.2	M.1	M.2	M.3	M.4	S.1	S.2	S.3	S.4
R.1	—									
R.2	0.75	—								
M.1	0.61	0.67	—							
M.2	0.58	0.62	0.74	—						
M.3	0.61	0.68	0.81	0.72	—					
M.4	0.62	0.69	0.76	0.72	0.76	—				
S.1	0.64	0.69	0.63	0.60	0.63	0.64	—			
S.2	0.63	0.66	0.62	0.60	0.62	0.64	0.68	—		
S.3	0.68	0.72	0.63	0.61	0.62	0.64	0.70	0.71	—	
S.4	0.64	0.69	0.62	0.61	0.62	0.63	0.66	0.69	0.74	—

Table 9.2.2.g Correlations between Reading, Mathematics, and Science Strands for Grade 11

Grade 11	R.1	R.2	M.1	M.2	M.3	M.4	S.1	S.2	S.3	S.4
R.1	—									
R.2	0.77	—								
M.1	0.52	0.59	—							
M.2	0.62	0.68	0.69	—						
M.3	0.63	0.71	0.72	0.82	—					
M.4	0.56	0.62	0.60	0.69	0.71	—				
S.1	0.65	0.72	0.56	0.67	0.67	0.59	—			
S.2	0.67	0.71	0.57	0.68	0.68	0.60	0.70	—		
S.3	0.70	0.73	0.56	0.66	0.68	0.59	0.71	0.74	—	
S.4	0.65	0.66	0.51	0.63	0.61	0.55	0.66	0.72	0.71	—

The correlations in Tables 9.2.2.a through 9.2.2.g are based on the observed strand scores. These observed-score correlations are weakened by existing measurement error contained within each strand. As a result, disattenuating the observed correlations can provide an estimate of the relationships between strands if there is no measurement error. The disattenuated correlation coefficients can be computed from the observed correlations (reported in Tables 9.2.2.a – 9.2.2.g) and the reliabilities for

each strand (Spearman, 1904, 1910). Disattenuated correlations very near 1.00 might suggest that the same or very similar constructs are being measured. Values somewhat less than 1.00 might suggest that different strands are measuring slightly different aspects of the same construct. Values markedly less than 1.00 might suggest the strands reflect different constructs.

Tables 9.2.3.a through 9.2.3.g show the corresponding disattenuated correlations for the 2015 NeSA tests for each grade. Given that none of these strands has perfect reliabilities (see Chapter Eight), the disattenuated strand correlations are higher than their observed score counterparts. Some within-content-area correlations are very high (e.g., above 0.95), suggesting that the within-content-area strands might be measuring essentially the same construct. This, in turn, suggests that some strand scores might not provide unique information about the strengths or weaknesses of students.

On a fairly consistent basis, the correlations between the strands within each content area were higher than the correlations between strands across different content areas. In general, within-content-area strand correlations were mostly greater than 0.90, while across-content-area strand correlations generally ranged from 0.75 to 0.92. Such a pattern is expected since the two content area tests were designed to measure different constructs.

Table 9.2.3.a Disattenuated Strand Correlations for Reading and Mathematics: Grade 3

Grade 3	R.1	R.2	M.1	M.2	M.3	M.4
R.1	—					
R.2	0.97	—				
M.1	0.81	0.77	—			
M.2	0.82	0.78	0.89	—		
M.3	0.80	0.78	0.95	0.89	—	
M.4	0.86	0.84	0.94	0.87	0.95	—

Table 9.2.3.b Disattenuated Strand Correlations for Reading and Mathematics: Grade 4

Grade 4	R.1	R.2	M.1	M.2	M.3	M.4
R.1	—					
R.2	0.96	—				
M.1	0.80	0.79	—			
M.2	0.83	0.80	0.95	—		
M.3	0.79	0.77	0.96	0.92	—	
M.4	0.88	0.90	0.94	0.93	0.93	—

Table 9.2.3.c Disattenuated Strand Correlations for Reading, Mathematics and Science: Grade 5

Grade 5	R.1	R.2	M.1	M.2	M.3	M.4	S.1	S.2	S.3	S.4
R.1	—									
R.2	0.98	—								
M.1	0.83	0.81	—							
M.2	0.83	0.79	0.93	—						
M.3	0.84	0.82	0.99	0.92	—					
M.4	0.87	0.85	0.92	0.92	0.96	—				
S.1	0.93	0.91	0.84	0.82	0.88	0.93	—			
S.2	0.87	0.84	0.78	0.82	0.81	0.84	0.94	—		
S.3	0.85	0.84	0.73	0.76	0.76	0.81	0.91	0.94	—	
S.4	0.87	0.83	0.77	0.80	0.80	0.84	0.93	0.97	0.97	—

Table 9.2.3.d Disattenuated Strand Correlations for Reading and Mathematics: Grade 6

Grade 6	R.1	R.2	M.1	M.2	M.3	M.4
R.1	—					
R.2	0.98	—				
M.1	0.79	0.78	—			
M.2	0.79	0.77	0.95	—		
M.3	0.83	0.82	0.96	0.93	—	
M.4	0.82	0.82	0.94	0.92	0.95	—

Table 9.2.3.e Disattenuated Strand Correlations for Reading and Mathematics: Grade 7

Grade 7	R.1	R.2	M.1	M.2	M.3	M.4
R.1	—					
R.2	0.97	—				
M.1	0.79	0.79	—			
M.2	0.76	0.78	0.93	—		
M.3	0.79	0.81	0.94	0.92	—	
M.4	0.79	0.82	0.94	0.93	0.95	—

Table 9.2.3.f Disattenuated Strand Correlations for Reading, Mathematics and Science: Grade 8

Grade 8	R.1	R.2	M.1	M.2	M.3	M.4	S.1	S.2	S.3	S.4
R.1	—									
R.2	0.96	—								
M.1	0.81	0.79	—							
M.2	0.80	0.77	0.95	—						
M.3	0.81	0.80	0.98	0.93	—					
M.4	0.84	0.82	0.94	0.94	0.94	—				
S.1	0.91	0.87	0.82	0.82	0.83	0.85	—			
S.2	0.88	0.82	0.80	0.81	0.79	0.83	0.94	—		
S.3	0.92	0.87	0.79	0.81	0.78	0.82	0.94	0.94	—	
S.4	0.89	0.86	0.81	0.83	0.80	0.83	0.92	0.94	0.98	—

Table 9.2.3.g Disattenuated Strand Correlations for Reading, Mathematics and Science: Grade 11

Grade 11	R.1	R.2	M.1	M.2	M.3	M.4	S.1	S.2	S.3	S.4
R.1	—									
R.2	0.96	—								
M.1	0.79	0.80	—							
M.2	0.79	0.79	0.97	—						
M.3	0.79	0.80	0.98	0.95	—					
M.4	0.82	0.82	0.95	0.94	0.93	—				
S.1	0.92	0.91	0.86	0.87	0.86	0.88	—			
S.2	0.91	0.85	0.84	0.84	0.83	0.85	0.95	—		
S.3	0.94	0.89	0.82	0.82	0.83	0.84	0.97	0.97	—	
S.4	0.90	0.82	0.77	0.81	0.77	0.81	0.93	0.96	0.96	—

9.3 EVIDENCE RELATED TO THE USE OF THE RASCH MODEL

Since the Rasch model is the basis of all calibration, scaling, and linking analyses associated with the NeSA, the validity of the inferences from these results depends on the degree to which the assumptions of the model are met as well as the fit between the model and test data. As discussed at length in Chapter Five, the underlying assumptions of Rasch models were essentially met for all the NeSA data, indicating the appropriateness of using the Rasch models to analyze the NeSA data.

In addition, the Rasch model was also used to link different operational NeSA tests across years. The accuracy of the linking also affects the accuracy of student scores and the validity of score uses. DRC Psychometric Services staff conducted verifications to check the accuracy of the procedures, including item calibration, conversions from the raw score to the Rasch ability estimate, and conversions from the Rasch ability estimates to the scale scores.

10. REFERENCES

- 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.
- Andersen, E. B. (1977). Sufficient statistics and latent trait models. *Psychometrika*, *42*, 69-81.
- Andrich, D. (1988). *Rasch models for measurement*. Newberry Park, CA: Sage.
- Brennan, R. L. (2004). BB-Class (Version 1.0). [Computer software] Iowa City, IA: University of Iowa, Center for Advanced Studies in Measurement & Assessment. Retrieved from www.education.uiowa.edu/casma.
- Chen, W., & Thissen, D. (1997). Local dependence indexes for item pairs using item response theory. *Journal of Educational and Behavioral Statistics*, *22*, 265-289.
- Cizek, G. J., & Bunch, M. B. (2007). *Standard setting: A guide to establishing and evaluating performance standards on tests*. Thousand Oaks, CA: Sage.
- Cronbach, L. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, *16*, 297-334.
- Fischer, G., & Molenaar, I. (1995). *Rasch models : Foundations, recent developments, and applications*. New York, NY: Springer.
- Hambleton, R. K., & Novick, M. R. (1973). Toward an integration of theory and method for criterion-referenced tests. *Journal of Educational Measurement*, *10*, 159-170.
- Hanson, B. A., & Brennan, R. L. (1990). An investigation of classification consistency indexes estimated under alternative strong true score theory models. *Journal of Educational Measurement*, *27*, 345-359.
- Harvill, L. M. (1991). Standard error of measurement. *Educational Measurement: Issues and Practices*, *10*(2), 33-41.
- Huynh, H. (1976). On the reliability of decisions in domain-referenced testing. *Journal of Educational Measurement*, *13*, 253-264.
- Huynh, H. (2000). Guidelines for Rasch linking for PACT. Memorandum to Paul Sandifer on June 18, 2000. Columbia, SC: Available from Author.
- Huynh, H., & Rawls, A. (2009). A comparison between Robust z and 0.3-logit difference procedures in assessing stability of linking items for the Rasch model. In E. V. Smith, Jr., & G. E. Stone (Eds.) *Applications of Rasch measurement in criterion-referenced testing*. (pp. 429-442). Maple Grove, MN: JAM Press.

- Huynh, H., & Meyer, P. (2010). Use of Robust z in detecting unstable items in item response theory models. *Practical Research, Assessment, and Evaluation*, 15, (2). Retrieved from ????
- Lewis, D. M., Mitzel, H. C., & Green, D. R. (1996, June). Standard setting: A bookmark approach. In D. R. Green (Chair), *IRT-Based standard-setting procedures utilizing behavioral anchoring*. Symposium conducted at the Council of Chief State School Officers National Conference on Large-Scale Assessment, Phoenix, AZ.
- Linacre, J. M. (2009). *A user's guide to WINSTEPS MINISTEP Rasch-model computer programs*. Chicago, IL: Winsteps.com
- Linacre, J. M. (2015). Winsteps® Rasch measurement computer program (V3.90). Beaverton, OR: Winsteps.com.
- Livingston, S., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32, 179-197.
- Ludlow, L. H. (1986). Graphical analysis of item response theory residuals. *Applied Psychological Measurement*, 10, 217-229.
- Mantel, N., & Haenszel, W. (1959). Statistical aspects of the analysis of data from retrospective studies of disease. *Journal of the National Cancer Institute*, 22, 719-748.
- McDonald, R. P. (1979). The structural analysis of multivariate data: A sketch of general theory. *Multivariate Behavioral Research*, 14, 21-38.
- Mead, R. J. (1976). *Assessing the fit of data to the Rasch model through the analysis of residuals*. Unpublished doctoral dissertation. Chicago, IL: University of Chicago.
- Mead, R. J. (2008). *A Rasch primer: The measurement theory of Georg Rasch*. (Psychometrics Services Research Memorandum 2008–001). Maple Grove, MN: Data Recognition Corporation.
- Mehrens, W. A., & Lehmann, I. J. (1975) *Standardized tests in education* (2nd ed.). New York, NY: Holt, Rinehart, and Winston.
- Mogilner, A. (1992). *Children's writer's world book*. Cincinnati, OH: Writer's Digest Books.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Copenhagen, Denmark: Danish Institute for Educational Research.
- Rasch, G. (1977). On specific objectivity: An attempt at formalizing the request for generality and validity of scientific statements. *Danish Yearbook of Philosophy*, 14, 58-94.
- Spearman C. (1904). The proof and measurement of association between two things. *American Journal of Psychology*, 15, 72-101.

- Spearman C. (1910). Correlation calculated from faulty data. *British Journal of Psychology*, 3, 271-295.
- Smith, E. V., Jr., & Smith, R. M. (Eds.). (2004). *Introduction to Rasch measurement*. Maple Grove, MN: JAM Press.
- Smith, R. M. (2000). Fit analysis in latent trait measurement models. *Journal of Applied Measurement*, 1, 199-218.
- Smith, R. M., Schumacker, R. E., & Bush, M. J. (1998). Using item mean squares to evaluate fit to the Rasch model. *Journal of Outcome Measurement*, 2, 66-78.
- Stearns, M., & Smith R. M. (2008). Estimation of classification consistency indices for complex assessments: Model based approaches. *Journal of Applied Measurement*, 9, 305-315.
- Taylor, S. E., Frackenpohl, H., White, C. E., Nieroroda, B. W., Browning, C. L., & Brisner, E. P. (1989). *EDL Core Vocabularies in Reading, Mathematics, Science, and Social Studies*. Orlando, FL: Steck-Vaughn Company.
- Thompson, S., Johnston, C. J., & Thurlow, M. L. (2002). *Universal design applied to large scale assessments*. National Center on Educational Outcomes Synthesis Report 44. Minneapolis, MN: University of Minnesota.
- Webb, N. L. (2002). *Alignment study in language arts, mathematics, science, and social studies of state standards and assessment for four states*. Washington, DC: Council of Chief State School Officers.
- Wright, B. D. (1968). Sample-free test calibration and person measurement. In *Proceedings of the 1967 Invitational Conference on Testing Problem* (pp. 85-101). Princeton, NJ: Educational Testing Service.
- Wright, B. D., & Mok, M. M. C. (2004). An overview of the family of Rasch measurement models. In E. V. Smith, Jr., & R. M. Smith (Eds.) *Introduction to Rasch measurement* (pp. 25-47). Maple Grove, MN: JAM Press.
- Wright, B. D., & Panchapakesan, N. (1969). A procedure of sample-free item analysis. *Educational and Psychological Measurement*, 29, 23-48.
- Wright, B. D., & Stone, M. H. (1979). *Best test design*. Chicago, IL: MESA Press.
- Yen, W. M. (1993). Scaling performance assessments: Strategies for managing local item dependence. *Journal of Educational Measurement*, 30, 187-213.

