



Spring 2015

Nebraska State Accountability (NeSA)

Reading, Mathematics, and Science

Alternate Assessment

Technical Report

September 2015

Prepared by Data Recognition Corporation





2015 NEBRASKA STATE ACCOUNTABILITY (NeSA) ALTERNATE ASSESSMENT 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 Alternate Assessments	
• Purpose of the NeSA	
• Phase-In Schedule for NeSA Alternate Assessment	
• Advisory Committees	
1.3. Administration.....	2

2. ITEM AND TEST DEVELOPMENT

2.1. Content Standards.....	3
2.2. Test Blueprints.....	4
2.3. Multiple-Choice Items.....	4
2.4. Item Development and Review	4
• Item Writer Training	
• Item Writing	
• Item Review	
• Editorial Review of Items	
• Universally Designed Assessments	
• Depth of Knowledge	
2.5. Item Banking	12
2.6. The Operational Form Construction Process	12
• Review of the Items and Test Forms	
2.7. Reading Assessment.....	14
• Test Design	
• Equating Design	
2.8. Mathematics Assessment.....	15
• Test Design	
• Equating Design	
2.9 Science Assessment.....	16
• Test Design	
• Equating Design	

3. STUDENT DEMOGRAPHICS AND ACCOMMODATIONS	17
4. CLASSICAL ITEM STATISTICS	
4.1. Item Difficulty	25
4.2. Item-Total Correlation	26
4.3. Percent Selecting Each Response Option	27
4.4. Point-Biserial Correlations of Response Options	28
4.5. Percent of Students Omitting an Item	28
5. RASCH ITEM CALIBRATION	
5.1. Description of the Rasch Model	29
5.2. Checking Rasch Assumptions	29
• Unidimensionality	
• Local Independence	
• Item Fit	
5.3. Rasch Item Statistics	38
6. EQUATING AND SCALING	
6.1. Equating	39
6.2. Scaling	41
7. FIELD TEST ITEM DATA SUMMARY	
7.1. Classical Item Statistics	46
8. RELIABILITY	
8.1. Coefficient Alpha	49
8.2. Standard Error of Measurement	50
8.3. Conditional Standard Error of Measurement (CSEM)	51
8.4. Decision Consistency and Accuracy	52
9. VALIDITY	
9.1. Evidence Based on Test Content	55
9.2. Evidence Based on Internal Structure	55
• Item-Test Correlations	
• Item Response Theory Dimensionality	
• Strand Correlations	
9.3. Evidence Related to the Use of the Rasch Model	61

10. REFERENCES	62
11. APPENDICES	
A. NeSA-AAR Test Blueprint.....	65
B. NeSA-AAM Test Blueprint.....	79
C. NeSA-AAS Test Blueprint	104
D. Confidentiality Agreement	123
E. Fairness in Testing Manual.....	124
F. Reading Key Verification and Foil Analysis.....	140
G. Mathematics Key Verification and Foil Analysis	149
H. Science Key Verification and Foil Analysis.....	159
I. Overview of Rasch Measurement.....	163
J. Reading, Mathematics, and Science Operational Form Calibration Summaries.....	167
K. Reading Item Bank Difficulties	174
L. Mathematics Item Bank Difficulties.....	181
M. Science Item Bank Difficulties.....	188
N. Reading Pre- and Post-Equating Summary	191
O. Mathematics Pre- and Post-Equating Summary	196
P. Science Pre- and Post-Equating Summary	202
Q. Reading Raw-to-Scale Conversion Tables and Distributions of Ability.....	205
R. Mathematics Raw-to-Scale Conversion Tables and Distributions of Ability	212
S. Science Raw-to-Scale Conversion Tables and Distributions of Ability.....	219
T. Reading, Mathematics, and Science Demographic Summary Sheets	222
U. Reading, Mathematics, and Science Strand Reliability and SEM.....	239



1. BACKGROUND

1.1 PURPOSE AND ORGANIZATION OF THIS REPORT

This report documents the technical aspects of the 2015 Nebraska Alternate Assessment Reading (NeSA-AAR), Mathematics (NeSA-AAM), and Nebraska Science (NeSA-AAS) operational tests, along with the NeSA-AAR, NeSA-AAM and NeSA-AAS embedded field tests, covering details of item and test development process, administration procedures, and psychometric methods and summaries.

1.2 BACKGROUND OF THE NEBRASKA STATE ACCOUNTABILITY (NE-SA)

Previous Nebraska Alternate Assessments: Prior to 2009, Alternate Assessments were not required. Districts had the ability to locally administer Alternate Assessments to students of their districts.

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-AAR for alternate reading assessments, NeSA-AAM for alternate mathematics, NeSA-AAS for alternate science. The alternate assessments in reading and mathematics were administered in grades 3-8 and 11; science was administered in grades 5, 8, and 11.

The NeSA-Alternate Assessment (NeSA-Alt) consists entirely of multiple choice items and are administered in a paper pencil format. 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 Alternate Assessment: The NDE prescribed the regular and the Alternate assessments starting in the 2009-2010 school year to be phased in as shown in Table 1.1. 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 1.1: NeSA Regular and Alternate Assessment 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

Alternate Assessment, and provided technical advice, guidance, and research to help the NDE make informed decisions regarding standards, assessment, and accountability.

1.3 ADMINISTRATION

The NeSA-Alt assessments are administered to students individually. The test administrator reads a prepared script for each item. As part of the assessment, the administrator may read the items multiple times and each student responds in their primary mode of communication. Test administrators record each response on the answer sheet. Students are able to utilize a full range of allowable accommodations that are detailed in documentation from the Nebraska Department of Education. If it becomes clear that a student is unable to respond to questions, the test administrator is required to record this on the answer sheet. Students who were administered the test but unable to respond count as participants but receive a zero score.

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, and science. 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-AAR. 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-AAR 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-AAM 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-AAM at grade 11. The assessable standards for each grade level are distributed among the four reporting categories: Number Sense Concepts, Geometric/Measurement Concepts, Algebraic Concepts, and Data Analysis/Probability Concepts.

The science component of the NeSA-AAS is composed of items that address indicators in grade-band strands 3–5, 6–8, and 9–12. The NeSA-AAS 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.

The NeSA-Alt are based on the same set of content standards that were extended by a team of special education specialists. The extended indicators detail underlying skills that students need

to master prior to attaining mastery of the full standard. The NeSA-Alt are aligned to the extended indicators.

2.2 TEST BLUEPRINTS (TABLE OF SPECIFICATIONS)

The test blueprints, or Table of Specifications (TOS), for each assessment include lists of all the standards, organized by reporting categories. The test blueprints also contain the Depth of Knowledge (DOK) level ranges assigned to each standard and the range of test items to be part of the assessment by extended indicator. The NeSA-AAR test blueprint (Appendix A) was originally developed and approved in fall 2009. The NeSA-AAM test blueprint (Appendix B) was originally developed and approved in fall 2010. The NeSA-AAS test blueprint (Appendix C) was originally developed and approved in fall 2011.

As part of the maturation of the NeSA-Alt program, NDE undertook to clarify the TOS in fall 2013 based on a careful examination of the overall pool of items within the NeSA-Alt item bank and the characteristics of the previous successful operational administrations. As a result, clarifications were made to all three TOS to better reflect the historical content of the NeSA-Alt program, and the clarified TOS were posted to NDE's website in advance of the 2013-2014 school year. It is important to point out that the clarifications made to the TOS bring the NeSA-Alt TOS into alignment with the actual historical NeSA-Alt test blueprints but did not change the breadth or depth of the content assessed within the actual NeSA-Alt program.

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 three 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 ITEM DEVELOPMENT AND REVIEW

The most significant considerations in the item and test development process are: aligning the items to the grade level extended 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 extended 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.

Nebraska State Accountability Alternate Assessment 2015 Technical Report

- 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.
- 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-Alt, 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 extended indicator each item was measuring.
- **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 four DOK categories:

recall, skill/concept, strategic thinking, and extended thinking (Webb, 2002). The NeSA-Alt items are classified based on DOK stages, subsets of the four categories. The stages include: responding, reproducing, recalling and basic reasoning.

- **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 and special education specialists 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.

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 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 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 alternate 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 for the alternate test. 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. 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 cognitive 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-Alt assessments include items written at levels 1 and 2. Levels 3 and 4 items are not included due to the test being comprised of only MC items and the cognitive level of students taking the alternate assessments. In addition, the NeSA-Alt items are classified based on DOK stages—subsets of the four DOK levels. The stages include responding, reproducing, recalling at DOK 1, and basic reasoning at DOK 2.

Reading Level 1-Stage 1: Responding to Discourse Materials

Level 1-Stage 1 requires students to display the ability to respond to or indicate, or acknowledge text or discourse related features. Some examples that represent, but do not constitute all of, Level 1-Stage 1 performance are:

- Student demonstrates the ability to attend to pictures/symbols/objects pertinent to a story
- Students display attention to people, surroundings, or materials.
- Student attends while teacher reads.

Reading Level 1-Stage 2: Reproduce Discourse Related Materials

Level 1-Stage 2 requires students to display the ability to copy, replicate, repeat, re-enact, mirror, or match text or discourse related features. Some examples that represent, but do not constitute all of, Level 1-Stage 2 performance are:

- Students match pictures and/or words that depict emotions such happy, sad, or angry.
- Students match printed words to objects.

Reading Level 1-Stage 3: Recalls Information about Discourse Related Materials

Level 1-Stage 3 requires the ability to recite or recall facts or information. Involves the ability to distinguish between text-based or discourse features. Some examples that represent, but do not constitute all of, Level 1-Stage3 performance are:

- Students demonstrate understanding or new words or passages by making connections with personal experience via speech, writing, signs, or assistive device.
- Students retell information taken from printed materials.
- Students answer who, what and where questions about a story.

Reading Level 2-Stage 4: Basic Reasoning

Level 2-Stage 4 requires processing beyond recall and observation. This requires both comprehension and subsequent processing of text. It also involves ordering, classifying text as well as identifying patterns, relationships, and main points. Some examples that represent, but do not constitute all of, Level 2-Stage 4 performance are:

- Students correct grammar mistakes in a reading selection.
- Students summarize the main idea of paragraph.
- Students identify the author's purpose for writing a brief passage.

Mathematics Level 1-Stage 1: Responding to Mathematical Materials

Level 1-Stage 1 requires the ability to respond to, indicate, or acknowledge mathematical features. Some examples that represent, but do not constitute all of, Level1-Stage 1 performance are:

- Students are able to recognize that there is a difference in patterns.
- Students respond to math ideas using appropriate vocabulary.

Mathematics Level 1-Stage 2: Reproduce Mathematical Features

Level 1-Stage 2 requires the ability to copy, replicate, repeat, re-enact, mirror, or match mathematical features. Some examples that represent, but do not constitute all of, Level 1-Stage 2 performance are:

- Students will write numbers accurately in a variety of contexts.
- Student accurately sort basic shapes into groups
- Student is able to accurately identify location terms when prompted (i.e., next to, between, over, under).

Mathematics Level 1-Stage 3: Recalls Information about Mathematical Features

Level 1-Stage 3 requires students to recall or observe facts, definitions, terms. It also involves simple one-step procedures. The stage also includes computing simple algorithms (e.g., sum, quotient). Some examples that represent, but do not constitute all of, Level 1-Stage3 performance are:

- Students locate a pattern in order to solve a problem
- Students measures using feet and yards.
- Students use a calculator or concrete objects to add and subtract.

Mathematics Level 2-Stage 4: Basic Reasoning

Level 2-Stage 4 requires students to make decisions of how to approach a problem. This may require students to compare, classify, organize, estimate or order data. This also typically involves two-step procedures. Some examples that represent, but do not constitute all of, Level 2-Stage 4 performance are:

- Student reads problem and determines operation to solve the problem.
- Student selects geometric figure from group of figures based on the definition of the geometric figure.
- Student determines how to solve for unknown value in equation or inequality and then selects solution.

Science Level 1-Stage 1: Responding to Scientific Features

Level 1-Stage 1 requires the ability to respond to or indicate or acknowledge scientific features. Some examples that represent, but do not constitute all of, Level1-Stage 1 performance are:

- Students attend to a teacher conducting scientific inquiry.
- Students respond to science ideas using appropriate vocabulary.

Science Level 1-Stage 2: Reproduce Scientific Features

Level 1-Stage 2 requires the ability to copy, replicate, repeat, re-enact, mirror, or match scientific ideas. Some examples that represent, but do not constitute all of, Level 1-Stage 2 performance are:

- Students copy figure of animal with distinguishing features.
- Student matches numbers on measuring devices.
- Student is able to accurately match descriptions of living and nonliving objects to visual representations.

Science Level 1-Stage 3: Recalls Information about Scientific Features

Level 1-Stage 3 requires students to recall or observe facts, definitions, terms. It also involves simple one-step procedures. The stage also requires a demonstration of a rote response, use of a well-known formula, or follow a set procedure (like a recipe), or preform a clearly defined series of steps. Some examples that represent, but do not constitute all of, Level 1-Stage3 performance are:

- Students recall or recognize a fact, term, or property.
- Students identify the correct measuring device to perform a task.
- Students perform a routine safety procedure.

Science Level 2-Stage 4: Basic Reasoning

Level 2-Stage 4 requires students to make decisions of how to approach a question or problem. This may require students to classify, organize, estimate, make observations or collect and order data. This also typically involves two-step procedures. Some examples that represent, but do not constitute all of, Level 2-Stage 4 performance are:

- Students make observations and collect data.
- Students organize and display data in tables, graphs, and charts.
- Students describe and explain examples and non-examples of science concepts.

2.5 ITEM BANKING

Prior to 2013, NDE exclusively maintained an item bank that provided a repository of item image, history, statistics, and usage. The item bank included a record of all newly created items together with item data from each item field test. It also included all data from the operational administration of the items. Within the item bank, NDE:

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

In 2014 NDE transitioned the item bank to DRC. DRC now maintains the alternate item bank in their system known as IDEAS, and it now functions as a repository of item image, history, statistics, and usage for the NeSA-Alt. 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.

2.6 THE OPERATIONAL FORM CONSTRUCTION PROCESS

The Spring 2015 operational forms were constructed in Lincoln, Nebraska in early September of 2014. The forms were constructed by a team of specialists representing special education, the Nebraska

Department of Education, and DRC testing experts. Training was provided collaboratively by NDE and DRC for the forms construction process.

Prior to arrival in Lincoln, DRC Test Development content specialists reviewed the test blueprints and the item pool to ensure that there was alignment between the items and the indicators, including the number of items per standard for each content-area test.

The specialists were provided 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 0.35 (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. The overall summary of the actual approved p -value and biserial of the forms is provided in the summary table later in this document. Below is the psychometric guidelines followed for item selection.

Psychometric Guidelines for Item Selection for a New Assessment

The main headings are more or less in order of precedence. This effectively means that content and reliability (*Ila and Iib*) define the pool of eligible items, from which items are selected based in p -value to match a target. *Guideline* is used here in the sense of *guiding principle*, not in the sense of *strict rule*. It is often, perhaps typically, necessary to deviate from these principles for a few items. There is no guideline for what a *few items* means.

- I. ***Item content: match the blue print.***
- II. Item-Total Correlation: (for MC items, point-biserial correlation)
 - a. Absolutely no correlations less than zero. This is a requirement, not a guideline.
 - b. Ideally, for MC items, point-biserial correlation should be greater than 0.35.**
 - i. A low correlation indicates there is a *smart* way to get the item wrong or *not-smart* way to get it right.
 - ii. The lower the value, the less discriminating the item.
- III. p -Value for correct response on MC
 - a. Target **mean percent correct about 65%** plus or minus a couple percent.
 - b. Ideally, all items greater than 40% and less than 85%**
 - c. For an existing assessment, the target mean percent correct should approximate past forms.

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. The item cards were compiled in binders and sorted 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, 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 using software into a form building template (Perform) which contained the item pool with statistics and item characteristics. The template automatically calculated the *p*-value, biserial, number of items per indicator and standard, number of items per DOK level, 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.

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. NDE 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.

2.7 READING ASSESSMENT

Test Design: The NeSA-AAR operational test includes operational items and field test items. The form pools contained 25 operational items and 16 field test items.

Table 2.7.1 Reading 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	25	8	33	2	25	16
4	25	8	33	2	25	16
5	25	8	33	2	25	16
6	25	8	33	2	25	16
7	25	8	33	2	25	16
8	25	8	33	2	25	16
11	25	8	33	2	25	16

Equating Design: Spring 2015 was the sixth operational administration of the NeSA-AAR. Approximately 20–40% of the assessment was constructed from items field tested from Spring 2009–2014. The approximate remaining 60–80% of the assessment was constructed from an overlap of items from the 2014 operational (core) item positions from the Spring 2014 operational forms.

In addition to the operational items, each student received 8 selected field test items. Equating was accomplished by anchoring on the operational items and calibrating the field test items concurrently.

2.8 MATHEMATICS ASSESSMENT

Test Design: The NeSA-AAM operational test includes operational items and field test items. The form pools contained 25 or 30 operational items (depending on the grade) with 16 field test items.

Table 2.8.1 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	25	8	33	2	25	16
4	30	8	38	2	30	16
5	30	8	38	2	30	16
6	30	8	38	2	30	16
7	30	8	38	2	30	16
8	30	8	38	2	30	16
11	30	8	38	2	30	16

Equating Design: Spring 2015 was the fifth operational administration of the NeSA-AAM. Approximately 20–40% of the assessment was constructed from items field tested from Spring 2010–

2014. The approximate remaining 60–80% 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 8 selected field test items. Equating was accomplished by anchoring on the operational items and calibrating the field test items concurrently.

2.9 SCIENCE ASSESSMENT

Test Design: The NeSA-AAS operational test includes operational and field test items. Depending on grade, the form pools contained 25 or 30 operational items (depending on the grade) with 16 field test items.

Table 2.9.1 Science 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
5	25	8	33	2	25	16
8	25	8	33	2	25	16
11	30	8	38	2	30	16

Equating Design: Spring 2015 was the fourth operational administration of the NeSA-AAS.

Approximately 20–40% of the assessment was constructed from items field tested in Spring 2011–2014. The approximate remaining 60–80% 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 8 field test items. Equating was accomplished by anchoring on the operational items and calibrating the field test items concurrently.

3. STUDENT DEMOGRAPHICS AND ACCOMMODATIONS

Gender, ethnicity, food program status (FRL), Limited English Proficiency/English Language Learners (LEP/ELL) status, and accommodation status data was collected for all students who participated and attempted the 2015 NeSA-Alt. This summary of student demographics by grade and content area is provided in Tables 3.1.1– 3.1.7. These tables show that for each grade, around 300 students took the assessment. Of those students across grades, approximately two-thirds are males, over half are white, and less than one fifth are Hispanic. Among the students across grades, over half are eligible for FRL, and almost all are non-LEP/ELL. In terms of the test accommodations, there are over half of the students across grade and content area that report at least one type of accommodation (see row ‘Total’ for ‘Accommodation’ in the table). Across all grades, the ‘Timing/Schedule/Setting’ is the most utilized accommodation, followed by the ‘Response’ and ‘Content Presentation’.

Table 3.1.1 Grade 3 NeSA-Alt Summary Data: Demographics and Accommodations

Grade 3		Reading		Mathematics	
		Count	%	Count	%
All Students		265	100.0	256	100.0
Gender	Female	88	33.2	87	34.0
	Male	177	66.8	169	66.0
Race/Ethnicity	American Indian/Alaska Native	8	3.0	8	3.1
	Asian	7	2.6	6	2.3
	Black	26	9.8	26	10.2
	Hispanic	44	16.6	43	16.8
	Native Hawaiian or other Pacific Islander	0	0.0	0	0.0
	White	170	64.2	163	63.7
	Two or More Races	10	3.8	10	3.9
Food Program	Yes	154	58.1	148	57.8
	No	111	41.9	108	42.2
LEP/ELL	Yes	6	2.3	5	2.0
	No	259	97.7	251	98.0
Accommodations	Content Presentation	141	53.2	134	52.3
	Response	160	60.4	152	59.4
	Timing/Schedule/Setting	184	69.4	175	68.4
	Direct Linguistic Support with Test Directions	3	1.1	2	0.8
	Direct Linguistic Support with Content and Test items	2	0.8	0	0.0
	Indirect Linguistic Support	2	0.8	2	0.8
	Total	186	70.2	177	69.1

Table 3.1.2 Grade 4 NeSA-Alt Summary Data: Demographics and Accommodations

Grade 4		Reading		Mathematics	
		Count	%	Count	%
All Students		291	100.0	289	100.0
Gender	Female	106	36.4	106	36.7
	Male	185	63.6	183	63.3
Race/Ethnicity	American Indian/Alaska Native	10	3.4	10	3.5
	Asian	6	2.1	6	2.1
	Black	22	7.6	22	7.6
	Hispanic	62	21.3	60	20.8
	Native Hawaiian or other Pacific Islander	0	0.0	0	0.0
	White	176	60.5	176	60.9
	Two or More Races	15	5.2	15	5.2
Food Program	Yes	186	63.9	186	64.4
	No	105	36.1	103	35.6
LEP/ELL	Yes	3	1.0	2	0.7
	No	288	99.0	287	99.3
Accommodations	Content Presentation	156	53.6	154	53.3
	Response	168	57.7	162	56.1
	Timing/Schedule/Setting	206	70.8	201	69.6
	Direct Linguistic Support with Test Directions	3	1.0	3	1.0
	Direct Linguistic Support with Content and Test items	3	1.0	1	0.3
	Indirect Linguistic Support	2	0.7	2	0.7
	Total	211	72.5	207	71.6

Table 3.1.3 Grade 5 NeSA-Alt Summary Data: Demographics and Accommodations

Grade 5		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		332	100.0	334	100.0	325	100.0
Gender	Female	109	32.8	115	34.4	109	33.5
	Male	223	67.2	219	65.6	216	66.5
Race/Ethnicity	American Indian/Alaska Native	9	2.7	9	2.7	8	2.5
	Asian	3	0.9	3	0.9	3	0.9
	Black	31	9.3	32	9.6	30	9.2
	Hispanic	69	20.8	69	20.7	71	21.8
	Native Hawaiian or other Pacific Islander	0	0.0	0	0.0	0	0.0
	White	211	63.6	211	63.2	204	62.8
	Two or More Races	9	2.7	10	3.0	9	2.8
Food Program	Yes	195	58.7	193	57.8	189	58.2
	No	137	41.3	141	42.2	136	41.8
LEP/ELL	Yes	3	0.9	3	0.9	3	0.9
	No	329	99.1	331	99.1	322	99.1
Accommodations	Content Presentation	195	58.7	191	57.2	184	56.6
	Response	191	57.5	192	57.5	184	56.6
	Timing/Schedule/Setting	230	69.3	225	67.4	218	67.1
	Direct Linguistic Support with Test Directions	2	0.6	1	0.3	3	0.9
	Direct Linguistic Support with Content and Test items	3	0.9	3	0.9	2	0.6
	Indirect Linguistic Support	2	0.6	1	0.3	2	0.6
	Total	237	71.4	233	69.8	225	69.2

Table 3.1.4 Grade 6 NeSA-Alt Summary Data: Demographics and Accommodations

Grade 6		Reading		Mathematics	
		Count	%	Count	%
All Students		331	100.0	339	100.0
Gender	Female	112	33.8	120	35.4
	Male	219	66.2	219	64.6
Race/Ethnicity	American Indian/Alaska Native	4	1.2	4	1.2
	Asian	4	1.2	4	1.2
	Black	40	12.1	39	11.5
	Hispanic	56	16.9	57	16.8
	Native Hawaiian or other Pacific Islander	0	0.0	0	0.0
	White	213	64.4	222	65.5
	Two or More Races	14	4.2	13	3.8
Food Program	Yes	192	58.0	193	56.9
	No	139	42.0	146	43.1
LEP/ELL	Yes	3	0.9	3	0.9
	No	328	99.1	336	99.1
Accommodations	Content Presentation	179	54.1	189	55.8
	Response	178	53.8	188	55.5
	Timing/Schedule/Setting	223	67.4	231	68.1
	Direct Linguistic Support with Test Directions	0	0.0	0	0.0
	Direct Linguistic Support with Content and Test items	0	0.0	0	0.0
	Indirect Linguistic Support	1	0.3	1	0.3
	Total	227	68.6	239	70.5

Table 3.1.5 Grade 7 NeSA-Alt Summary Data: Demographics and Accommodations

Grade 7		Reading		Mathematics	
		Count	%	Count	%
All Students		327	100.0	329	100.0
Gender	Female	124	37.9	129	39.2
	Male	203	62.1	200	60.8
Race/Ethnicity	American Indian/Alaska Native	5	1.5	5	1.5
	Asian	4	1.2	4	1.2
	Black	37	11.3	37	11.2
	Hispanic	61	18.7	58	17.6
	Native Hawaiian or other Pacific Islander	0	0.0	0	0.0
	White	204	62.4	209	63.5
	Two or More Races	16	4.9	16	4.9
Food Program	Yes	181	55.4	175	53.2
	No	146	44.6	154	46.8
LEP/ELL	Yes	6	1.8	6	1.8
	No	321	98.2	323	98.2
Accommodations	Content Presentation	179	54.7	177	53.8
	Response	184	56.3	185	56.2
	Timing/Schedule/Setting	228	69.7	225	68.4
	Direct Linguistic Support with Test Directions	0	0.0	0	0.0
	Direct Linguistic Support with Content and Test items	1	0.3	1	0.3
	Indirect Linguistic Support	0	0.0	0	0.0
	Total	228	69.7	226	68.7

Table 3.1.6 Grade 8 NeSA-Alt Summary Data: Demographics and Accommodations

Grade 8		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		334	100.0	338	100.0	327	100.0
Gender	Female	126	37.7	127	37.6	125	38.2
	Male	208	62.3	211	62.4	202	61.8
Race/Ethnicity	American Indian/Alaska Native	7	2.1	7	2.1	7	2.1
	Asian	8	2.4	8	2.4	8	2.4
	Black	41	12.3	41	12.1	42	12.8
	Hispanic	52	15.6	54	16.0	51	15.6
	Native Hawaiian or other Pacific Islander	0	0.0	0	0.0	0	0.0
	White	212	63.5	213	63.0	206	63.0
	Two or More Races	14	4.2	15	4.4	13	4.0
Food Program	Yes	180	53.9	186	55.0	178	54.4
	No	154	46.1	152	45.0	149	45.6
LEP/ELL	Yes	3	0.9	3	0.9	2	0.6
	No	331	99.1	335	99.1	325	99.4
Accommodations	Content Presentation	161	48.2	168	49.7	158	48.3
	Response	167	50.0	172	50.9	160	48.9
	Timing/Schedule/Setting	199	59.6	201	59.5	191	58.4
	Direct Linguistic Support with Test Directions	0	0.0	0	0.0	0	0.0
	Direct Linguistic Support with Content and Test items	0	0.0	0	0.0	0	0.0
	Indirect Linguistic Support	0	0.0	0	0.0	0	0.0
	Total	202	60.5	207	61.2	194	59.3

Table 3.1.7 Grade 11 NeSA-Alt Summary Data: Demographics and Accommodations

Grade 11		Reading		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		309	100.0	318	100.0	307	100.0
Gender	Female	102	33.0	104	32.7	102	33.2
	Male	207	67.0	214	67.3	205	66.8
Race/Ethnicity	American Indian/Alaska Native	8	2.6	8	2.5	7	2.3
	Asian	9	2.9	9	2.8	9	2.9
	Black	32	10.4	33	10.4	32	10.4
	Hispanic	42	13.6	42	13.2	42	13.7
	Native Hawaiian or other Pacific Islander	0	0.0	0	0.0	0	0.0
	White	210	68.0	215	67.6	208	67.8
	Two or More Races	8	2.6	11	3.5	9	2.9
Food Program	Yes	174	56.3	177	55.7	173	56.4
	No	135	43.7	141	44.3	134	43.6
LEP/ELL	Yes	0	0.0	0	0.0	0	0.0
	No	309	100.0	318	100.0	307	100.0
Accommodations	Content Presentation	132	42.7	141	44.3	134	43.6
	Response	128	41.4	138	43.4	127	41.4
	Timing/Schedule/Setting	169	54.7	180	56.6	172	56.0
	Direct Linguistic Support with Test Directions	2	0.6	1	0.3	3	1.0
	Direct Linguistic Support with Content and Test items	1	0.3	0	0.0	1	0.3
	Indirect Linguistic Support	2	0.6	0	0.0	2	0.7
	Total	172	55.7	187	58.8	173	56.4

4. CLASSICAL ITEM STATISTICS

This chapter provides an overview of the most familiar item-level statistics obtained from classical (traditional) item analysis: item difficulty, item discrimination, distractor distribution, and omits or blanks. The following results pertain only to operational NeSA-Alt 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-Alt, a test development goal is to include a wide range of item difficulties. Typically, test developers target *p*-values in the range of 0.40 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-AAR 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	1	0	6	8	8	2	0	0.661	25
4	0	0	0	0	3	6	8	7	1	0	0.643	25
5	0	0	0	0	0	5	14	5	1	0	0.643	25
6	0	0	0	1	1	8	8	6	1	0	0.631	25
7	0	0	0	0	1	5	9	8	2	0	0.672	25
8	0	0	0	0	0	6	10	8	1	0	0.674	25
11	0	0	0	1	2	2	11	8	1	0	0.645	25

Table 4.1.2 Summary of Proportion Correct for NeSA-AAM 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	1	2	3	13	2	4	0	0.657	25
4	0	0	0	1	3	7	9	9	1	0	0.643	30
5	0	0	0	1	3	8	9	9	0	0	0.633	30
6	0	0	0	0	5	8	9	7	1	0	0.620	30
7	0	0	0	0	1	6	9	11	3	0	0.678	30
8	0	0	0	1	5	4	7	10	3	0	0.640	30
11	0	0	0	1	7	5	4	9	4	0	0.630	30

Table 4.1.3 Summary of Proportion Correct for NeSA-AAS 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	0	3	6	9	5	2	0	0.635	25
8	0	0	0	0	1	7	11	2	4	0	0.663	25
11	0	0	0	1	3	7	9	9	1	0	0.643	30

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-Alt tests, Pearson 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-Alt overall) would be more likely to answer any given NeSA-Alt item correctly, while students with low ability (i.e., those who perform poorly on the NeSA-Alt 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-Alt 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-AAR

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
3	0	0	0	3	1	10	11	25
4	0	0	0	0	5	7	13	25
5	0	0	0	3	6	12	4	25
6	0	0	0	1	1	11	12	25
7	0	1	1	2	6	12	3	25
8	0	0	0	3	8	6	8	25
11	0	0	1	3	7	11	3	25

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

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
3	0	0	0	2	2	6	15	25
4	0	0	1	1	4	8	16	30
5	0	0	0	2	8	10	10	30
6	0	0	0	1	3	17	9	30
7	0	0	3	1	5	10	11	30
8	0	0	0	5	9	12	4	30
11	0	1	2	3	7	9	8	30

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

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
5	0	0	1	2	6	8	8	25
8	0	0	0	2	9	9	5	25
11	0	0	1	1	4	8	16	30

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-AAR operational items are presented in Appendix F. Detailed results of the item analyses for the NeSA-AAM operational items are presented in Appendix G. Detailed results of the item analyses for the NeSA-AAS 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-Alt 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-Alt items in recent history. Rasch models have several advantages over true-score 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 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-AAR, NeSA-AAM, and NeSA-AAS assessment.

5.1 DESCRIPTION OF THE RASCH MODEL

The Rasch dichotomous model was used to calibrate the NeSA-Alt items. All NeSA-Alt assessment contains 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 provides 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-Alt, 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 up to five components with eigenvalues greater than one. As can be seen in Table 5.2.1, the primary dimension for NeSA-AAR explained about 25 percent to 30 percent of the total variance across Grades 3–8 and 11. The eigenvalues of the second dimension ranged from 1.8 to 2.3. This indicates that the second dimension accounted for only 1.8 to 2.3 units out of about 37 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-Alt assessment.

Table 5.2.1 NeSA-AAR Results from PCA

Grade	Contrast	Eigenvalue	Explained Variance
3	measures	11.6	31.7%
	1	2.2	8.8%
	2	1.6	6.2%
	3	1.4	5.7%
	4	1.3	5.4%
	5	1.3	5.4%
4*	measures	12.8	33.9%
	1	2.1	8.3%
	2	1.8	7.1%
	3	1.5	5.9%
	4		
	5		
5	measures	8.5	25.3%
	1	1.9	7.7%
	2	1.6	6.4%
	3	1.4	5.8%
	4	1.3	5.3%
	5	1.3	5.3%
6*	measures	11.2	30.8%
	1	2.0	8.0%
	2	1.5	6.2%
	3		
	4		
	5		
7	measures	9.5	27.5%
	1	2.3	9.0%
	2	1.6	6.5%
	3	1.4	5.7%
	4	1.3	5.3%
	5	1.3	5.2%
8	measures	8.9	26.2%
	1	1.9	7.6%
	2	1.7	6.8%
	3	1.5	6.0%
	4	1.4	5.6%
	5	1.3	5.3%
11	measures	10.1	28.8%
	1	1.8	7.1%
	2	1.6	6.2%
	3	1.4	5.7%
	4	1.3	5.3%
	5	1.2	5.0%

Table 5.2.2 NeSA-AAM Results from PCA

Grade	Contrast	Eigenvalue	Explained Variance
3*	measures	12.2	32.8%
	1	1.9	7.5%
	2	1.8	7.4%
	3	1.6	6.2%
	4	1.4	5.7%
	5		
4*	measures	13.4	30.9%
	1	2.5	8.2%
	2	1.9	6.3%
	3		
	4		
	5		
5	measures	11.8	28.2%
	1	2.4	8.1%
	2	2.0	6.7%
	3	1.7	5.5%
	4	1.5	5.1%
	5	1.4	4.6%
6	measures	11.8	28.3%
	1	2.4	8.1%
	2	1.9	6.2%
	3	1.7	5.6%
	4	1.6	5.2%
	5	1.3	4.5%
7	measures	12.6	29.6%
	1	2.7	8.9%
	2	2.0	6.5%
	3	1.4	4.8%
	4	1.4	4.8%
	5	1.3	4.4%
8	measures	11.4	27.6%
	1	2.2	7.4%
	2	2.0	6.6%
	3	1.7	5.6%
	4	1.4	4.6%
	5	1.3	4.4%
11	measures	12.8	30.0%
	1	2.4	8.1%
	2	2.0	6.5%
	3	1.6	5.4%
	4	1.5	4.9%
	5	1.4	4.6%

Table 5.2.3 NeSA-AAS Results from PCA

Grade	Contrast	Eigenvalue	Explained Variance
5	measures	9.3	27.2%
	1	2.8	11.2%
	2	1.6	6.3%
	3	1.5	5.9%
	4	1.4	5.5%
	5	1.3	5.1%
8	measures	9.4	27.3%
	1	2.2	8.9%
	2	1.6	6.4%
	3	1.5	6.2%
	4	1.4	5.5%
	5	1.3	5.0%
11	measures	13.4	30.8%
	1	1.8	6.1%
	2	1.7	5.8%
	3	1.6	5.5%
	4	1.6	5.2%
	5	1.4	4.6%

*Only contrasts with eigenvalues greater than one were extracted.

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-Alt 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 -0.04 for the NeSA-Alt tests (since most of the NeSA-Alt tests had more than 25 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 -0.04 . The vast majority of the correlations were very

small, suggesting local item independence generally holds for the NeSA-Alt reading, mathematics, and science assessments.

Table 5.2.4 Summary of Item Residual Correlations for NeSA-AAR

Statistics	3	4	5	6	7	8	11
<i>N</i>	300	300	300	300	300	300	300
Mean	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04	-0.04
<i>SD</i>	0.07	0.08	0.06	0.07	0.07	0.07	0.06
Minimum	-0.23	-0.22	-0.22	-0.22	-0.21	-0.24	-0.21
P10	-0.14	-0.13	-0.12	-0.12	-0.14	-0.12	-0.12
P25	-0.09	-0.09	-0.08	-0.09	-0.10	-0.09	-0.08
P50	-0.04	-0.04	-0.04	-0.04	-0.04	-0.05	-0.04
P75	0.01	0.01	0.00	0.01	0.01	0.00	0.00
P90	0.05	0.07	0.04	0.05	0.06	0.05	0.04
Maximum	0.16	0.19	0.20	0.23	0.18	0.22	0.15
>0.20	0	0	0	1	0	1	0

Table 5.2.5 Summary of Item Residual Correlations for NeSA-AAM

	Mathematics						
Statistics	3	4	5	6	7	8	11
<i>N</i>	300	435	435	435	435	435	435
Mean	-0.04	-0.03	-0.03	-0.03	-0.03	-0.03	-0.03
<i>SD</i>	0.08	0.09	0.08	0.08	0.08	0.07	0.08
Minimum	-0.23	-0.26	-0.23	-0.20	-0.24	-0.24	-0.27
P10	-0.13	-0.14	-0.13	-0.13	-0.13	-0.12	-0.13
P25	-0.09	-0.09	-0.09	-0.09	-0.09	-0.08	-0.08
P50	-0.04	-0.03	-0.04	-0.04	-0.04	-0.04	-0.03
P75	0.01	0.03	0.02	0.03	0.02	0.02	0.01
P90	0.06	0.08	0.08	0.07	0.09	0.06	0.07
Maximum	0.21	0.27	0.22	0.23	0.25	0.21	0.21
>0.20	2	4	2	2	4	1	2

Table 5.2.6 Summary of Item Residual Correlations for NeSA-AAS

Statistics	Science		
	5	8	11
<i>N</i>	300	300	435
Mean	-0.04	-0.04	-0.03
<i>SD</i>	0.10	0.08	0.07
Minimum	-0.32	-0.23	-0.18
P10	-0.15	-0.13	-0.12
P25	-0.10	-0.09	-0.08
P50	-0.04	-0.04	-0.03
P75	0.02	0.01	0.01
P90	0.09	0.06	0.05
Maximum	0.32	0.19	0.21
>0.20	5	0	1

Item Fit: WINSTEPS provides two item fit statistics (infit and outfit) 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 a different variance or as a standardized statistic (Zstd with mean = 0 and variance = 1).

MnSq values are more difficult to interpret due to an asymmetrical distribution, 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-AA, the sample sizes can be considered small. 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. Deviation in excess of the expected value can be interpreted as noise or lack of fit between the responses and the model. Values lower than the expected value can be interpreted as item redundancy or overfitting items (too predictable and/or too much redundancy), and values greater than the expected value indicate underfitting items (too unpredictable and/or too much noise). Rules of thumb regarding “practically significant” MnSq values vary. 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. Deviation in excess of the expected value can be interpreted as noise or lack of fit between the items and the model. Values lower than the expected value can be interpreted as item redundancy or overfitting items (too predictable and/or too much redundancy), and values greater than the expected value indicate underfitting items (too unpredictable and/or too much noise). Rules of thumb regarding “practically significant” Zstd values vary. 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 -2 to $+2$ are given practical importance.

Table 5.2.7 lists the summary statistics of infit and outfit mean square statistics for the NeSA-Alt 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-Alt reading, mathematics, and science tests, including the mean, *SD*, and minimum and maximum values. The number of items within the range of $[-2, +2]$ is also reported in Table 5.2.8. As can be seen, the mean values for both fit statistics were close to 0.00 for all tests. Most of the items had infit values falling in the range of $[-2, +2]$. 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. Overall, these results indicate that the NeSA-Alt item data fits Rasch model well.

Table 5.2.7 Summary of Infit and Outfit Mean Square Statistics for 2015 NeSA-Alt Tests

		Infit Mean Square					Outfit Mean Square				
		Mean	SD	MIN	MAX	[0.7, 1.3]	Mean	SD	MIN	MAX	[0.7, 1.3]
Reading	3	1.00	0.18	0.75	1.39	22/25	0.97	0.32	0.47	1.72	16/25
	4	1.00	0.14	0.74	1.32	24/25	0.99	0.35	0.49	2.21	18/25
	5	1.00	0.13	0.81	1.29	25/25	1.01	0.31	0.53	1.95	19/25
	6	1.00	0.13	0.84	1.33	24/25	1.05	0.50	0.63	3.24	21/25
	7	1.00	0.17	0.79	1.47	24/25	1.02	0.33	0.56	1.89	19/25
	8	1.00	0.15	0.77	1.31	24/25	0.99	0.29	0.54	1.55	15/25
	11	1.00	0.14	0.76	1.33	24/25	0.98	0.27	0.57	1.51	18/25
Mathematics	3	0.99	0.16	0.70	1.38	22/25	0.94	0.43	0.29	2.41	14/25
	4	0.99	0.17	0.80	1.51	28/30	0.93	0.32	0.50	1.98	20/30
	5	1.00	0.11	0.81	1.21	30/30	0.96	0.24	0.53	1.40	22/30
	6	1.00	0.13	0.79	1.33	29/30	0.97	0.24	0.55	1.49	22/30
	7	0.99	0.22	0.71	1.53	27/30	0.94	0.36	0.41	1.76	18/30
	8	1.00	0.12	0.78	1.19	30/30	0.96	0.21	0.60	1.26	26/30
	11	0.99	0.15	0.80	1.41	28/30	0.95	0.29	0.49	1.71	20/30
Science	5	1.00	0.20	0.67	1.56	21/25	0.97	0.32	0.48	1.72	16/25
	8	0.99	0.13	0.73	1.19	25/25	0.96	0.26	0.45	1.38	19/25
	11	1.00	0.16	0.80	1.45	29/30	0.95	0.30	0.37	1.60	22/30

Table 5.2.8 Summary of Infit and Outfit Z STD Statistics for 2015 NeSA-Alt Tests

		Infit Z STD					Outfit Z STD				
		Mean	SD	MIN	MAX	[-2.0, 2.0]	Mean	SD	MIN	MAX	[-2.0, 2.0]
Reading	3	-0.05	2.18	-3.39	4.68	15/25	-0.26	1.69	-2.67	3.56	18/25
	4	0.03	1.79	-3.38	4.26	21/25	-0.16	1.63	-3.28	3.04	19/25
	5	-0.09	2.10	-3.09	4.78	16/25	-0.05	2.05	-3.03	3.85	16/25
	6	-0.08	1.72	-2.53	4.47	18/25	-0.05	1.75	-2.66	4.67	21/25
	7	-0.10	2.58	-2.95	6.80	12/25	-0.07	2.18	-2.83	6.36	16/25
	8	-0.01	2.19	-3.73	4.84	15/25	-0.06	2.20	-3.13	4.74	15/25
	11	0.01	1.95	-2.66	4.80	15/25	-0.16	1.70	-2.32	3.42	19/25
Mathematics	3	0.12	1.73	-2.27	4.90	21/25	-0.23	1.81	-2.31	5.22	18/25
	4	0.07	2.13	-2.07	6.13	23/30	-0.23	1.66	-2.15	4.97	27/30
	5	0.12	1.65	-2.70	3.54	24/30	-0.10	1.67	-2.98	2.64	21/30
	6	-0.02	1.96	-3.66	4.87	20/30	-0.14	1.81	-3.71	4.88	24/30
	7	-0.08	3.04	-4.26	7.08	13/30	-0.16	2.43	-3.45	5.99	18/30
	8	0.07	1.88	-3.55	3.31	21/30	-0.13	1.62	-3.22	2.63	21/30
	11	0.10	2.26	-3.88	6.32	22/30	-0.11	1.94	-2.69	4.71	20/30
Science	5	-0.02	2.93	-4.59	8.32	16/25	-0.15	2.44	-3.96	5.96	17/25
	8	0.12	1.67	-2.37	3.01	15/25	-0.06	1.59	-2.49	2.58	17/25
	11	0.03	2.05	-2.34	5.92	20/30	-0.20	1.75	-2.34	4.10	21/30

5.3 RASCH ITEM STATISTICS

Item calibration was implemented via WINSTEPS 3.90.0 program (Linacre, 2015). 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. Logits have an interval scale, meaning that two items with logits of 0.0 and +1.0 (respectively) are the same distance apart (in difficulty) as two items with logits 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-Alt items covered a relatively wide range of difficulties. The range describes the spread of the items. Some tests are narrower than others. 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-AAR, NeSA-AAM, and NeSA-AAS

	Grade	N	Mean	SD	Min	Max	Range
Reading	3	25	0.01	0.72	-1.43	1.85	3.29
	4	25	0.01	0.71	-1.38	1.37	2.75
	5	25	0.13	0.55	-1.11	0.92	2.03
	6	25	0.34	0.63	-0.82	1.75	2.56
	7	25	0.06	0.65	-1.35	1.40	2.76
	8	25	0.59	0.56	-1.43	1.45	2.88
	11	25	0.14	0.73	-1.18	1.76	2.94
Mathematics	3	25	-0.25	0.92	-2.21	1.24	3.45
	4	30	0.00	0.77	-1.48	1.46	2.95
	5	30	-0.11	0.75	-1.38	1.60	2.98
	6	30	0.11	0.66	-1.22	1.19	2.41
	7	30	-0.18	0.77	-1.53	1.13	2.66
	8	30	-0.14	0.79	-1.43	1.39	2.82
	11	30	-0.32	0.95	-1.88	1.15	3.03
Science	5	25	-1.15	0.79	-3.21	-0.11	3.10
	8	25	-1.24	0.72	-2.82	0.03	2.85
	11	30	-1.18	0.73	-3.20	-0.04	3.15

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.

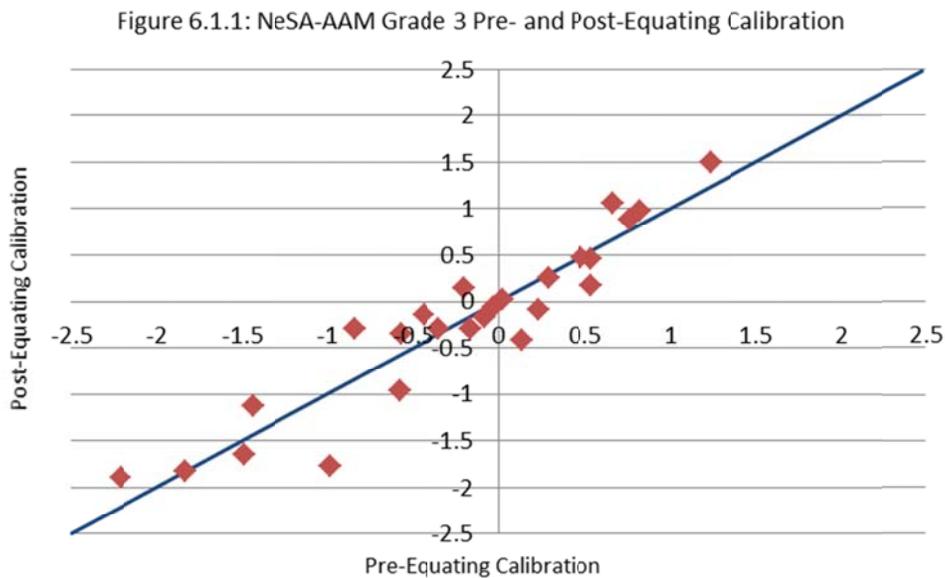
6.1 EQUATING

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-AAM Pre- and Post-Equating Comparison

	Grade						
	3	4	5	6	7	8	11
Correlation	.94*	.92*	.94*	.94*	.94*	.92*	.94*
SD pre	.87	.76	.72	.68	.74	.78	.95
SD post	.92	.76	.74	.69	.72	.77	.95
SD Ratio	1.05	1.00	1.03	1.02	.97	.99	1.00

*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 calculated 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 score scale. 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-AAR for 2015, Appendix R for NeSA-AAM for 2015 and Appendix S for NeSA-AAS 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-Alt 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 State Board of Education per the standard setting.

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-AAR Standard Setting Technical Report,

http://www.education.ne.gov/Assessment/pdfs/2010_NeSA_AAR_Standard_Setting_Report.pdf

2011 NeSA-AAR and NeSA-AAM Standard Setting Technical Report,

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

and *2012 NeSA-AAS Standard Setting Technical Report,*

<http://www.education.ne.gov/Assessment/pdfs/NeSA-AAS%20Standard%20Setting%20Results.pdf>

Given the scale score and the logit standards defining the performance level, it is sufficient to define the final scale score metric. To ensure proper rounding on all future forms, the calculations used 84.501 and 134.501 as the scale score performance standards. The arithmetic was done using logits rounded to four decimals and the final constants for the slope and intercept of the transformation were rounded to five. Scale scores are rounded to whole numbers.

The transformation to scale scores is:

$$SS = a + b * \text{logit},$$

where:

$$b = \frac{134.501 - 84.501}{x_E - x_M},$$

and where x_E is the logit for *Exceeds Standards* and x_M is the logit for *Meets Standards*.

Therefore:

$$a = 84.501 - bx_M \text{ or,}$$

$$a = 134.501 - bx_E .$$

Calculations of the slopes and intercepts for all grades of the NeSA-AAR scale score conversion are given in Table 6.2.1, for NeSA-AAM 6.2.2, and for NeSA-AAS 6.2.3. The raw-to-scale conversions are provided in Appendices Q, R, and S.

Table 6.2.1 NeSA-AAR 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	0.2501	1.8426	1 to 84	85 to 134	135 to 200	31.39720	76.64840
4	0.2536	1.8106	1 to 84	85 to 134	135 to 200	32.11300	76.35520
5	0.2612	1.5392	1 to 84	85 to 134	135 to 200	39.12360	74.28010
6	0.4202	2.0909	1 to 84	85 to 134	135 to 200	29.92760	71.92460
7	0.4169	1.7456	1 to 84	85 to 134	135 to 200	37.63080	68.81120
8	0.6792	2.3138	1 to 84	85 to 134	135 to 200	30.58680	63.72690
11	0.2362	1.8139	1 to 84	85 to 134	135 to 200	31.69170	77.01370

Table 6.2.2 NeSA-AAM 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 <i>b</i>	Intercept <i>a</i>
3	-0.0819	1.6006	1 to 84	85 to 134	135 to 200	29.71770	86.93460
4	0.4250	1.7728	1 to 84	85 to 134	135 to 200	37.09750	68.73270
5	-0.0108	1.3462	1 to 84	85 to 134	135 to 200	36.84600	84.89680
6	0.2970	2.0591	1 to 84	85 to 134	135 to 200	28.37520	76.07320
7	0.2953	1.7471	1 to 84	85 to 134	135 to 200	34.44000	74.33050
8	0.4528	1.7661	1 to 84	85 to 134	135 to 200	38.07200	67.26220
11	0.2976	1.2809	1 to 84	85 to 134	135 to 200	50.84920	69.36900

Table 6.2.3 NeSA-AAS 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 <i>b</i>	Intercept <i>a</i>
5	-1.0631	0.3571	1 to 84	85 to 134	135 to 200	35.20631	121.93783
8	-0.7286	0.5524	1 to 84	85 to 134	135 to 200	39.03201	112.94872
11	-0.8043	0.6780	1 to 84	85 to 134	135 to 200	33.73136	111.64013

Complete frequency distributions of the state scale scores for the NeSA-AAR, NeSA-AAM, and NeSA-AAS are provided in Appendices Q, R, and S as part of the raw-to-scale-score conversion tables. In addition, descriptive statistics of the state raw scores, scale scores, and performance levels are computed for subgroups based on gender, ethnicity, special education status, limited English proficiency status, and food program eligibility status in Appendix T. 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 2015 NeSA-AAR State Scale Score Summary, All Students

Grade	Count	Scale Score		Quartile		
		Mean	S.D.	First	Second	Third
3	265	110.0	53.9	74	117	146
4	291	107.9	56.7	68	104	147
5	332	113.3	55.5	77	119	161
6	331	108.0	54.6	69	112	146
7	327	111.9	54.8	75	110	152
8	334	113.8	48.3	84	119	146
11	309	109.6	48.8	78	108	139

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

Grade	Count	Scale Score		Quartile		
		Mean	S.D.	First	Second	Third
3	256	110.1	52.1	77	120	147
4	289	105.7	59.9	66	104	146
5	334	111.2	54.3	75	116	147
6	339	101.1	48.8	67	101	137
7	329	108.0	52.4	63	107	151
8	338	96.2	51.7	61	98	131
11	318	97.9	59.7	53	101	149

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

Grade	Count	Scale Score		Quartile		
		Mean	S.D.	First	Second	Third
5	325	110.9	53.4	73	118	145
8	327	105.7	54.4	76	106	135
11	307	112.9	57.2	74	111	152

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-AAR field test item in Appendix F, for NeSA-AAM Appendix G and NeSA-AAS 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-AAR 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	0	1	0	8	3	3	1	0	0.603	16
4	0	0	1	1	3	3	3	3	2	0	0.586	16
5	0	1	1	1	3	5	2	2	1	0	0.536	16
6	0	0	1	1	4	4	4	1	1	0	0.552	16
7	0	0	1	3	1	6	1	4	0	0	0.549	16
8	0	0	2	4	3	2	4	1	0	0	0.485	16
11	0	0	2	1	5	3	4	1	0	0	0.501	16

Nebraska State Accountability Alternate Assessment 2015 Technical Report

	Item Point-biserial Correlation							
Grade	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	Total
3	0	0	0	4	5	2	5	16
4	0	0	0	3	4	3	6	16
5	2	0	2	2	4	2	4	16
6	0	2	0	3	3	5	3	16
7	2	0	5	1	4	2	2	16
8	0	0	4	5	3	2	2	16
11	0	3	3	4	1	3	2	16

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

	Item Proportion Correct											
Grade	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	≤0.7	≤0.8	≤0.9	>0.9	Mean	Total
3	0	0	1	2	1	5	5	0	2	0	0.575	16
4	0	0	0	1	3	4	4	3	1	0	0.599	16
5	0	0	1	2	4	2	3	3	1	0	0.557	16
6	0	0	0	2	3	6	3	2	0	0	0.540	16
7	0	0	1	2	5	6	1	1	0	0	0.500	16
8	0	1	0	2	3	4	3	3	0	0	0.537	16
11	0	0	0	6	6	2	2	0	0	0	0.451	16

	Item Point-biserial Correlation							
Grade	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	Total
3	0	0	2	0	5	1	8	16
4	0	0	1	4	3	2	6	16
5	0	0	3	2	3	3	5	16
6	0	0	1	2	3	9	1	16
7	0	1	2	3	6	4	0	16
8	0	1	1	3	4	4	3	16
11	1	2	3	2	5	3	0	16

Table 7.1.3 Summary of Statistics for NeSA-AAS 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	0	2	1	2	2	1	8	0	0	0.598	16
8	0	0	2	0	1	5	3	4	1	0	0.604	16
11	0	0	0	1	3	4	4	3	1	0	0.599	16

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
5	0	0	2	2	3	3	6	16
8	1	1	3	1	3	3	4	16
11	0	0	1	4	3	2	6	16

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:

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

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

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

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-Alt 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-Alt. 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 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	25	0.93	1.9
	4	25	0.93	1.9
	5	25	0.91	2.0
	6	25	0.94	1.9
	7	25	0.90	2.0
	8	25	0.91	2.0
	11	25	0.91	2.0
Mathematics	3	25	0.94	1.8
	4	30	0.94	2.1
	5	30	0.93	2.1
	6	30	0.94	2.1
	7	30	0.92	2.1
	8	30	0.92	2.2
	11	30	0.91	2.2
Science	5	25	0.92	2.0
	8	25	0.91	2.0
	11	30	0.94	2.0

Appendix U 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. Reliability estimates for subgroups based on gender, ethnicity, special education status, limited English proficiency status, and food program eligibility status are not computed for the NeSA-Alt tests due to the small sample size of some subgroups.

8.2 STANDARD ERROR OF MEASUREMENT

The *SEM* in the true score model 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-Alt 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 true score 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 *CSEM* s across the score scale seemed reasonable for most NeSA-Alt 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-Alt Tests

		Min	Max	CSEM	CSEM
	Grade	CSEM	CSEM	B/M	M/E
Reading	3	13	58	13	18
	4	14	59	14	18
	5	16	72	16	20
	6	12	55	13	17
	7	16	69	16	21
	8	13	56	13	17
	11	13	58	13	18
Mathematics	3	13	55	13	17
	4	14	68	15	19
	5	14	68	14	18
	6	11	52	11	16
	7	13	63	14	19
	8	15	70	15	21
	11	20	94	21	25
Science	5	15	65	15	18
	8	17	72	17	22
	11	13	62	13	18

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-Alt 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-Alt 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 & Hanson, 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-AAR 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	0.91	0.92	0.88	0.94	0.91	0.92	0.88
	4	0.94	0.92	0.92	0.89	0.94	0.92	0.92	0.90
	5	0.93	0.92	0.91	0.89	0.93	0.91	0.91	0.89
	6	0.94	0.93	0.92	0.90	0.94	0.93	0.92	0.91
	7	0.93	0.92	0.90	0.89	0.93	0.92	0.90	0.89
	8	0.94	0.90	0.91	0.87	0.93	0.90	0.91	0.87
	11	0.93	0.91	0.90	0.87	0.93	0.90	0.90	0.87

Table 8.4.2 NeSA-AAM 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	0.95	0.89	0.93	0.84	0.95	0.89	0.93	0.86
	4	0.93	0.92	0.91	0.89	0.93	0.92	0.91	0.89
	5	0.93	0.90	0.91	0.87	0.93	0.90	0.91	0.87
	6	0.94	0.93	0.92	0.90	0.94	0.93	0.92	0.90
	7	0.93	0.89	0.91	0.84	0.93	0.90	0.91	0.86
	8	0.91	0.91	0.88	0.88	0.91	0.91	0.88	0.88
	11	0.91	0.90	0.88	0.86	0.91	0.90	0.88	0.86

Table 8.4.3 NeSA-AAS 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	0.94	0.91	0.91	0.87	0.94	0.91	0.92	0.87
	8	0.92	0.90	0.89	0.85	0.92	0.89	0.89	0.86
	11	0.94	0.93	0.92	0.90	0.94	0.92	0.92	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-Alt 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-Alt for grades 3 to 8 and 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 alternate 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.

9.2 EVIDENCE BASED ON INTERNAL STRUCTURE

As described in the *Standards for Educational and Psychological Testing* (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-Alt reading, mathematics, and science tests were essentially unidimensional,

providing evidence supporting interpretations based on the total scores for the respective NeSA-Alt 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-AAR tests have two strands (denoted by R.1 and R.2), the NeSA-AAM tests have four strands (denoted by M.1, M.2, M.3, and M.4), and the NeSA-AAS have four strands (denoted by S.1, S.2, S.3, and S.4) for each grade and content area.

For each grade, Pearson 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-Alt 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	0.77	0.80	—			
M.2	0.79	0.85	0.83	—		
M.3	0.66	0.75	0.71	0.72	—	
M.4	0.67	0.76	0.75	0.71	0.69	—

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	0.87	—				
M.1	0.83	0.86	—			
M.2	0.85	0.88	0.85	—		
M.3	0.79	0.80	0.81	0.79	—	
M.4	0.72	0.75	0.70	0.75	0.71	—

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	0.84	—								
M.1	0.81	0.83	—							
M.2	0.80	0.77	0.85	—						
M.3	0.67	0.69	0.72	0.70	—					
M.4	0.67	0.66	0.67	0.70	0.66	—				
S.1	0.70	0.73	0.72	0.74	0.62	0.65	—			
S.2	0.81	0.80	0.82	0.77	0.65	0.61	0.70	—		
S.3	0.74	0.75	0.79	0.79	0.68	0.70	0.72	0.75	—	
S.4	0.77	0.76	0.79	0.78	0.65	0.67	0.71	0.77	0.76	—

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	0.87	—				
M.1	0.79	0.85	—			
M.2	0.76	0.82	0.83	—		
M.3	0.76	0.77	0.79	0.79	—	
M.4	0.77	0.76	0.78	0.75	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	0.83	—				
M.1	0.72	0.75	—			
M.2	0.73	0.72	0.71	—		
M.3	0.74	0.79	0.76	0.72	—	
M.4	0.76	0.76	0.77	0.74	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	0.83	—								
M.1	0.71	0.75	—							
M.2	0.74	0.77	0.73	—						
M.3	0.72	0.75	0.76	0.74	—					
M.4	0.69	0.73	0.70	0.71	0.68	—				
S.1	0.68	0.69	0.64	0.68	0.65	0.60	—			
S.2	0.75	0.78	0.69	0.68	0.73	0.66	0.68	—		
S.3	0.80	0.80	0.70	0.76	0.71	0.74	0.69	0.77	—	
S.4	0.74	0.73	0.66	0.71	0.67	0.66	0.68	0.69	0.73	—

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	0.81	—								
M.1	0.58	0.66	—							
M.2	0.70	0.78	0.64	—						
M.3	0.74	0.80	0.64	0.77	—					
M.4	0.63	0.67	0.53	0.65	0.64	—				
S.1	0.74	0.80	0.58	0.77	0.79	0.64	—			
S.2	0.75	0.80	0.58	0.72	0.75	0.64	0.80	—		
S.3	0.72	0.80	0.58	0.73	0.77	0.64	0.77	0.82	—	
S.4	0.75	0.81	0.60	0.76	0.79	0.71	0.78	0.81	0.81	—

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-Alt 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 close to 1.00, while across-content-area strand correlations generally ranged from 0.83 to 1.00. 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	0.92	—				
M.1	0.94	0.92	—			
M.2	0.95	0.96	0.98	—		
M.3	0.89	0.95	0.93	0.94	—	
M.4	0.94	1.00	1.00	0.96	1.00	—

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	1.00	—				
M.1	0.98	0.98	—			
M.2	1.00	1.00	0.98	—		
M.3	1.00	1.00	1.00	1.00	—	
M.4	0.93	0.94	0.89	0.96	1.00	—

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	1.00	—								
M.1	0.99	0.96	—							
M.2	0.98	0.91	1.00	—						
M.3	1.00	1.00	1.00	1.00	—					
M.4	0.91	0.85	0.88	0.93	1.00	—				
S.1	0.94	0.94	0.93	0.98	0.91	0.94	—			
S.2	1.00	0.97	1.00	0.96	0.96	0.83	0.95	—		
S.3	1.00	0.96	1.00	1.00	1.00	1.00	1.00	1.00	—	
S.4	0.98	0.92	0.97	0.96	1.00	0.91	0.97	0.98	1.00	—

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	1.00	—				
M.1	0.96	0.99	—			
M.2	0.92	0.96	1.00	—		
M.3	0.93	0.90	0.97	0.98	—	
M.4	0.97	0.93	1.00	0.96	0.92	—

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	1.00	—				
M.1	0.97	0.97	—			
M.2	0.93	0.89	0.95	—		
M.3	0.97	1.00	1.00	0.94	—	
M.4	0.98	0.94	1.00	0.94	0.92	—

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	1.00	—								
M.1	0.92	0.94	—							
M.2	0.93	0.93	0.96	—						
M.3	0.95	0.96	1.00	0.98	—					
M.4	0.93	0.96	0.99	0.97	0.98	—				
S.1	0.97	0.96	0.96	0.98	0.91	0.95	—			
S.2	1.00	1.00	0.98	0.93	0.97	0.97	1.00	—		
S.3	1.00	0.97	0.91	0.96	1.00	1.00	1.00	1.00	—	
S.4	0.97	0.93	0.90	0.94	0.97	0.95	1.00	0.98	0.96	—

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	1.00	—								
M.1	0.96	1.00	—							
M.2	0.88	0.92	1.00	—						
M.3	0.96	0.98	1.00	0.96	—					
M.4	0.97	0.97	1.00	0.96	0.98	—				
S.1	1.00	1.00	1.00	1.00	1.00	1.00	—			
S.2	0.95	0.94	0.92	0.86	0.93	0.95	1.00	—		
S.3	0.90	0.94	0.91	0.86	0.99	0.93	1.00	0.98	—	
S.4	0.96	0.97	0.97	0.92	1.00	1.00	1.00	0.99	0.98	—

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-Alt, 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-Alt data, indicating the appropriateness of using the Rasch models to analyze the NeSA-Alt data.

In addition, the Rasch model was also used to link different operational NeSA-Alt 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.