



Spring 2023

Nebraska Student-Centered Assessment System (NSCAS)

Alternate Assessment

ELA • Mathematics • Science

TECHNICAL REPORT

December 2023

Prepared by Data Recognition Corporation



2023 Nebraska Student-Centered Assessment System (NSCAS) Alternate Assessment

Technical Report

Table of Contents

Chapter 1: Background	1
1.1. Purpose and organization of this report.....	1
1.2. Background of the Nebraska Student-Centered Assessment System Alternate Tests (NSCAS)	1
1.2.1. Purpose of the NSCAS Alternate:	1
1.2.2. History of Nebraska Alternate Assessments:	1
1.2.3. Phase-In Schedule for NSCAS Alternate Assessment:	1
1.2.4. Advisory Committees:	2
1.2.5. College and Career Ready Standards for Mathematics:	2
1.2.6. New College and Career Ready Standards for Science:	2
1.2.7. College and Career Ready Standards for English Language Arts:	3
1.3. Administration	3
Chapter 2: Item and Test Development.....	4
2.1. Content Standards.....	4
2.2. Test Blueprints (Table of Specifications).....	5
2.3. Multiple-Choice Items.....	6
2.4. Item Development and Review	6
2.4.1. Item Writer Training:	6
2.4.2. Item Writing:	7
2.4.3. Item Review:	8
2.4.4. Editorial Review of Items:	8
2.4.5. Universally Designed Assessments:	8
2.4.6. Depth of Knowledge (DOK):	11
2.5. Item Banking.....	13
2.6. The Operational Form Construction Process	13
2.6.1. Review of the Items and Test Forms:	15
2.7. English Language Arts Assessment	15
2.7.1. Test Design:	15
2.7.2. Equating Design:	16

2.8. Mathematics Assessment	17
2.8.1. Test Design:	17
2.8.2. Equating Design:	17
2.9. Science Assessment	18
2.9.1. Test Design:	18
2.9.2. Equating Design:	18
Chapter 3: STUDENT DEMOGRAPHICS AND ACCOMMODATIONS	19
Chapter 4: CLASSICAL ITEM STATISTICS	27
4.1. ITEM DIFFICULTY	27
4.2. ITEM-TOTAL CORRELATION	28
4.3. PERCENT SELECTING EACH RESPONSE OPTION	30
4.4. POINT-BISERIAL CORRELATIONS OF RESPONSE OPTIONS	30
4.5. PERCENT OF STUDENTS OMITTING AN ITEM	30
Chapter 5: RASCH ITEM CALIBRATION	32
5.1. DESCRIPTION OF THE RASCH MODEL	32
5.2. CHECKING RASCH ASSUMPTIONS	32
5.2.1. Unidimensionality:	32
5.2.2. Local Independence:	36
5.2.3. Item Fit:	38
5.3. RASCH ITEM STATISTICS	40
Chapter 6: EQUATING AND SCALING	42
6.1. Equating	42
6.2. SCALING	44
Chapter 7: FIELD TEST ITEM DATA SUMMARY	50
7.1. CLASSICAL ITEM STATISTICS	50
Chapter 8: RELIABILITY	52
8.1. COEFFICIENT ALPHA	52
8.2. STANDARD ERROR OF MEASUREMENT	53
8.3. CONDITIONAL STANDARD ERROR OF MEASUREMENT (CSEM)	54
8.4. DECISION CONSISTENCY AND ACCURACY	55
Chapter 9: VALIDITY	58
9.1. EVIDENCE BASED ON TEST CONTENT	58
9.2. EVIDENCE BASED ON INTERNAL STRUCTURE	59
9.2.1. Item-Test Correlations:	59
9.2.2. Item Response Theory Dimensionality:	59

9.2.3. Strand Correlations:	59
9.3. EVIDENCE RELATED TO THE USE OF THE RASCH MODEL.....	65
Chapter 10: References.....	66

Appendices

A. NSCAS-AAELA Test Blueprint.....	71
B. NSCAS-AAM Test Blueprint.....	93
C. NSCAS-AAS Test Blueprint	130
D. Confidentiality and Security Agreements	162
E. Fairness in Testing Manual.....	165
F. ELA Key Verification and Foil Analysis.....	183
G. Mathematics Key Verification and Foil Analysis.....	191
H. Science Key Verification and Foil Analysis	199
I. Overview of Rasch Measurement.....	203
J. ELA Item Bank Difficulties.....	207
K. Mathematics Item Bank Difficulties.....	215
L. Science Item Bank Difficulties	223
M. ELA Raw-to-Scale Conversion Tables and Distributions of Ability.....	227
N. Mathematics Raw-to-Scale Conversion Tables and Distributions of Ability.....	235
O. Science Raw-to-Scale Conversion Tables and Distributions of Ability.....	243
P. ELA, Mathematics, and Science Demographic Summary Sheets	247
Q. ELA, Mathematics, and Science Strand Reliability and SEM.....	257

Chapter 1: Background

1.1. Purpose and organization of this report

This report documents the technical aspects of the 2023 Nebraska Student-Centered Assessment System Alternate tests of English Language Arts (NSCAS-AAELA), Mathematics (NSCAS-AAM), and Science (NSCAS-AAS)—including operational tests and embedded field tests—covering details of item and test development, administration procedures, and psychometric methods and summaries.

1.2. Background of the Nebraska Student-Centered Assessment System Alternate Tests (NSCAS)

1.2.1. Purpose of the NSCAS Alternate:

The NSCAS Alternate is for students with the most significant cognitive disabilities who are assessed against Nebraska’s College and Career Ready Extended Indicators.

1.2.2. History of 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.

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 Nebraska academic content standards for reading, mathematics, science, and writing in Nebraska’s K-12 public schools. The assessment system was named NeSA (Nebraska State Accountability). The NeSA-Alternate Assessment (NeSA-AA) consisted of multiple-choice items administered in a paper-pencil format. Initially delivered directly by the Nebraska Department of Education (NDE), NDE contracted with Data Recognition Corporation (DRC) starting in 2011-2012 to support administration and reporting of its statewide alternate assessments.

Nebraska assessment designs were revised following the adoption of new Nebraska College and Career Ready Standards. In 2017-2018, Nebraska replaced Nebraska State Accountability (NeSA) tests with the Nebraska Student-Centered Assessment System (NSCAS). 2017-2018 also marked the first online administration of alternate assessments via DRC INSIGHT.

1.2.3. Phase-In Schedule for NSCAS Alternate Assessment:

The NDE prescribed alternate assessments starting in the 2009-2010 school year, phased in as shown in Table 1.2.1.

Table 1.2.1: 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
Mathematics aligned to CCR Standards	2017	2018	3 through 8 plus high school
Science aligned to CCR Standards*	2021	2022	5, 8, and 11
ELA aligned to new CCR Standards	2023	2023	3 through 8 plus high school

*Note that the transition of science to the CCR standards was delayed by one year due to the interruptions to testing in spring 2020 caused by the Covid 19 Pandemic.

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

1.2.5. College and Career Ready Standards for Mathematics:

The existing College and Career Ready Standards for Mathematics were adopted by the State Board of Education in September of 2015. Districts were required to adopt these standards within one year.

Student scores for the 2022-2023 NSCAS-AAM were calculated using only operational items aligned to 2015 College and Career Ready Standards for Mathematics. 2017-2018 was the first operational assessment of items aligned to the revised standards. This report includes technical information about embedded field test items, which were also aligned to the 2015 standards.

1.2.6. New College and Career Ready Standards for Science:

New College and Career Ready Standards for Science were adopted by the State Board of Education in September of 2017. Typically, districts are required to adopt new standards within one year. However, operational delays and delays due to the Covid 19 Pandemic resulted in a timeline shift. The first operational alternate science assessment aligned to the new standards took place in spring 2022.

Student scores for the 2022-2023 NSCAS-AAS were calculated using only operational items aligned to the 2017 College and Career Ready Standards for Science. This report includes technical information about embedded field test items, which were also aligned to the new standards.

1.2.7. College and Career Ready Standards for English Language Arts:

New College and Career Ready Standards for English Language Arts were adopted by the State Board of Education in September 2021. Districts were required to adopt these standards within one year. Spring 2023 was an operational field test aligned to the new 2021 College and Career Ready Standards for English Language Arts. Student scores for the 2022-2023 NSCAS-AAELA were calculated after completing the standard setting process in the summer of 2023. This report includes technical information about embedded field test items, which were also aligned to the new 2021 standards.

1.3. Administration

NSCAS Alternate tests 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. Students capable of responding to items via a touch-enabled device or computer may respond directly to items online using DRC INSIGHT. For other students, test administrators record each response and transcribe them into DRC INSIGHT for scoring. 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 in the DRC INSIGHT Portal—an online data collection system. Students who were administered the test but unable to respond count as participants but receive a zero score.

Chapter 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 seven 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.
- College and Career Ready Standards for English Language Arts were adopted by the State Board of Education in September of 2014. Spring 2016 was the final administration of the NeSA-AAR and spring 2017 marked the first administration of the NeSA-English Language Arts (ELA) Alternate Assessment. Adjustments were made to the College and Career Ready Standards for English Language Arts, which were adopted by the State Board of Education in September 2021. The 2022-2023 NSCAS-AAELA was the first assessment aligned to the 2021 standards.
- College and Career Ready Standards for Mathematics were adopted by the State Board of Education in September of 2015. Spring 2017 was the final administration of the NeSA-AAM and spring 2018 marked the first operational administration of Mathematics alternate assessments aligned to the College and Career Ready Standards.
- New College and Career Ready Standards for Science were adopted by the State Board of Education in September of 2017. Spring 2019 was the final administration of the science alternate assessment aligned to the previous standards. The original plan for testing in Spring 2020 was to have a standalone field test with items aligned to the new College and Career Ready Science standards. Due to the cancellation of state assessments in Spring 2020 because of the Covid 19 Pandemic, the Spring 2021 assessment was made up of 100% field test items. Spring 2022 was the first operational administration of the NSCAS-AAS, which was aligned to the new standards.
- Spring 2018 marked the transition from the Nebraska State Accountability (NeSA) program to the Nebraska Student-Centered Assessment System (NSCAS).

The Nebraska Language Arts College and Career Ready Standards are the foundation for NSCAS-AAELA. 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 NSCAS-AAELA tests at grade 11. The ELA standards for each grade are represented in items that are distributed between three reporting categories: Vocabulary, Comprehension, and Writing.

- The Vocabulary standards include vocabulary acquisition and use as well as vocabulary context and connotation.
- The Comprehension standards include both reading prose and poetry and reading informational text and include central ideas and details as well as knowledge and ideas (supporting evidence).
- The Writing standards include production and modes of writing.

The mathematics component of the NSCAS-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 NSCAS-AAM at grade 11. The assessable standards for each grade level are distributed among the four reporting categories: Number, Algebra, Geometry, and Data.

The science component of the NSCAS-AAS is composed of items that address the College-and Career Ready Science Standards extended indicators and access points in grade levels 5, 8 and 11. The NSCAS-AAS assesses the standards for each grade level: grade 5, grade 8, and grade 11. The assessable standards for grade 5, 8 and 11 are distributed among three reporting categories: Physical Science, Life Science, and Earth and Space Science.

NSCAS Alternate tests 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. NSCAS Alternate tests 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 NSCAS-AAELA test blueprint (Appendix A) was developed and approved in the fall of 2022. The NSCAS-AAM test blueprint (Appendix B) was originally developed and approved in fall 2016. The NSCAS-AAS test blueprint was developed and approved in fall 2021 to reflect the new science standards.

Since all three content areas are part of the maturation of the NSCAS Alternate program, NDE revised the TOS appropriately based on careful examination of the overall pool of items within the alternate assessment item bank and the characteristics of the previous successful operational administrations. As

a result, all three TOS reflect the current reporting categories and information as expected of the College and Career Ready Standards for English language arts, mathematics and science.

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 up to 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.

With the three-dimensional expectations of the new College and Career Ready science standards, clusters have been developed for the alternate science assessment that are using the same stimulus text to support student responses and understanding of the extended indicators and access points. The NSCAS-AAS still uses the multiple-choice format and the items within the clusters do stand alone or are not dependent on the student's response from one item to the next, but the items are connected by a common topic or text.

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.
- 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.

2.4.1. 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 NSCAS Alternate, a cadre of item writers was trained in the following areas:

- 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.

2.4.2. 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 NSCAS Alternate 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.

2.4.3. Item Review:

Throughout the item development process, and since the items are developed by Nebraska educators, content experts and special education specialists have reviewed the items using the following guidelines. 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.

2.4.4. Editorial Review of Items:

After items were written and reviewed, 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 NDE to review and revise the items prior to field testing to ensure highest level of quality possible.

2.4.5. 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). 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. 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, 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. 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 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

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 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 NDE's editing process:

- Reduce excessive length.
- Use clear and common words.
- Avoid ambiguous or multiple meaning words.
- Avoid irregularly spelled words.
- Avoid proper names.
- Avoid inconsistent naming and graphic conventions.
- Avoid unclear signals about how to direct attention.
- Avoid "mark all" questions.
- Maximize 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 NDE to develop a style guide that includes dimensions of style that are consistent with universal design.

2.4.6. Depth of Knowledge (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 NSCAS Alternate 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 NSCAS Alternate items are classified based on DOK stages—subsets of the four DOK levels. The stages include reproducing, recalling at DOK 1, and basic reasoning at DOK 2.

ELA 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:

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

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

Level 1-Stage 3 requires the ability to recite or recall facts or information. It also 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:

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

ELA 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:

- Student corrects grammar mistakes in a reading selection.
- Student summarizes the main idea of paragraph.
- Student identifies the author’s purpose for writing a brief passage.

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:

- Student writes numbers accurately in a variety of contexts.
- Student accurately sort basic shapes into groups
- Student accurately identifies location terms when prompted (e.g., 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:

- Student locates a pattern in order to solve a problem
- Student measures using feet and yards.
- Student uses 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 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:

- Student copies figure of animal with distinguishing features.
- Student matches numbers on measuring devices.
- Student accurately matches 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 perform a clearly defined series of steps. Some examples that represent, but do not constitute all of, Level 1-Stage3 performance are:

- Student recalls or recognizes a fact, term, or property.
- Student identifies the correct measuring device to perform a task.

- Student performs 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:

- Student makes observations and collect data.
- Student organizes and displays data in tables, graphs, and charts.
- Student describes and explains 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 NSCAS Alternate. 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 2023 operational forms were constructed by DRC and NDE in Lincoln, Nebraska in September of 2022. 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 form 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 are 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, item 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.

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 2023 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) were verified.

2.6.1. 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. English Language Arts Assessment

2.7.1. Test Design:

With the adoption of revised ELA standards in 2021, the spring 2023 NSCAS-AAELA (English Language Arts) test was the first operational test and included field test items. The operational test items were used for reporting purposes after the standard setting process in the summer of 2023. The form pools contained 28 operational test items and 16 field test items aligned to the new ELA standards, as shown in Table 2.7.1. Statistics for field-test items can be found in Chapter 7 of this Technical Report.

Table 2.7.1 NSCAS-AAELA 2023 Operational Field Test

Grade	Total No. of 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 Items Added to the Bank
3	28	8	36	2	28	16
4	28	8	36	2	28	16
5	28	8	36	2	28	16
6	28	8	36	2	28	16
7	28	8	36	2	28	16
8	28	8	36	2	28	16
11	28	8	36	2	28	16

Table 2.7.2 NSCAS-AAELA 2023 Points per Content Domain

Grade	Number of Points	Reading Vocabulary Points	Reading Comprehension Points		Writing Points
3	28	6-8	13-17		6-7
			RP 7-9	RI 6-8	
4	28	6-8	13-17		6-7
			RP 7-9	RI 6-8	
5	28	6-8	12-16		6-7
			RP 6-8	RI 6-8	
6	28	6-8	12-16		6-7
			RP 6-8	RI 6-8	
7	28	6-7	13-17		6-7
			RP 6-8	RI 7-9	
8	28	6-7	13-17		6-7
			RP 6-8	RI 7-9	
11	28	6-7	13-17		6-7
			RP 6-8	RI 7-9	

2.7.2. Equating Design:

Refer to Chapter 6 for information about the equating design.

2.8. Mathematics Assessment

2.8.1. Test Design:

The NSCAS-AAM test included operational items and field test items. The form pools contained 25 or 30 operational items (depending on the grade) with 16 field-test items, as shown in Table 2.8.1.

Table 2.8.1 NSCAS-AAM 2023 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

Table 2.8.2 NSCAS-AAM 2023 Points per Content Domain

Grade	Number of Items	Number Items	Algebra Items	Geometry Items	Data Items
3	25	8-11	3-6	6-9	2-4
4	30	8-12	3-6	5-8	2-4
5	30	6-12	3-6	5-8	2-4
6	30	8-11	6-8	6-8	4-6
7	30	4-8	8-12	4-6	2-4
8	30	4-8	8-10	5-8	2-4
11	30	4-6	8-10	6-10	3-6

2.8.2. Equating Design:

Refer to chapter 6 for information about the equating design.

2.9. Science Assessment

2.9.1. Test Design:

The 2023 NSCAS-AAS was the second operational assessment aligned to the 2017 standards. The test included both operational and embedded field test items. Depending on grade, the forms consisted of the same number of items as on the 2022 operational test design, as shown in Table 2.9.1.

2.9.1 NSCAS-AAS 2023 Operational Test Table

Grade	Total No. of MC Core Items	No. of 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	28	8	36	2	28	16

Table 2.9.2 NSCAS-AAS 2023 Points per Content Domain

Grade	Number of Points	Physical Science Points	Life Science Points	Earth and Space Sciences Points
5	25	6-10	6-10	8-12
8	25	8-12	8-12	6-10
HS	28	8-14	10-14	8-14

Table 2.9.2 NSCAS-AAS 2023 DOK Percentages per Stage

Grade	Number of Points	DOK Level 1,2	DOK Level 1,3	DOK Level 2,4
5	25	5%-15%	40%-60%	30%-50%
8	25	5%-15%	40%-60%	30%-50%
HS	28	5%-15%	40%-60%	25%-45%

2.9.2. Equating Design:

Refer to chapter 6 for information about the equating design.

Chapter 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 2023 NSCAS-AA. This summary of student demographics by grade and content area is provided in Tables 3.1 through 3.8. These tables show that for each grade, approximately 235 to 260 students took the assessment in grades 3 through 8, and approximately 225 with High School. Of those students across grades, approximately two-thirds were males, over half were white, and approximately one fifth were Hispanic. Among the students across grades, over half were eligible for the food program, and almost all were non-LEP/ELL. The proportion of students taking the paper and pencil accommodation was approximately 20% in grades 3 through 8 and slightly over 10% with High School.

Table 3.1 Number of Alternate Tests Administered

Grade	ELA	Mathematics	Science
3	238	235	—
4	258	255	—
5	253	251	251
6	234	235	—
7	262	260	—
8	239	235	237
HS	229	227	226

Table 3.2 Grade 3 NSCAS-AA Summary Data: Demographics and Accommodations

Grade 3		ELA		Mathematics	
		Count	%	Count	%
All Students		238	100.00	235	100.00
Gender	Female	96	40.34	94	40.00
	Male	142	59.66	141	60.00
Race/Ethnicity	American Indian/Alaska Native	4	1.68	4	1.70
	Asian	9	3.78	9	3.83
	Black	26	10.92	26	11.06
	Hispanic	50	21.01	50	21.28
	Native Hawaiian or Other Pacific Islander	1	0.42	1	0.43
	White	134	56.30	131	55.74
	Two or More Races	14	5.88	14	5.96
Food Program	Yes	143	60.08	142	60.43
	No	95	39.92	93	39.57
LEP/ELL	Yes	6	2.52	6	2.55
	No	232	97.48	229	97.45
Accommo- dations	Paper and Pencil	57	23.95	54	22.98
	Computation Supports	23	9.66	63	26.81
	Assistive Technology	23	9.66	23	9.79
	Specialized Presentation	27	11.34	26	11.06

Table 3.3 Grade 4 NSCAS-AA Summary Data: Demographics and Accommodations

Grade 4		ELA		Mathematics	
		Count	%	Count	%
All Students		258	100.00	255	100.00
Gender	Female	94	36.43	93	36.47
	Male	164	63.57	162	63.53
Race/Ethnicity	American Indian/Alaska Native	4	1.55	3	1.18
	Asian	7	2.71	7	2.75
	Black	19	7.36	18	7.06
	Hispanic	54	20.93	54	21.18
	Native Hawaiian or Other Pacific Islander	1	0.39	1	0.39
	White	162	62.79	161	63.14
	Two or More Races	11	4.26	11	4.31
Food Program	Yes	150	58.14	148	58.04
	No	108	41.86	107	41.96
LEP/ELL	Yes	3	1.16	3	1.18
	No	255	98.84	252	98.82
Accommo- dations	Paper and Pencil	61	23.64	62	24.31
	Computation Supports	23	8.91	86	33.73
	Assistive Technology	26	10.08	26	10.20
	Specialized Presentation	23	8.91	23	9.02

Table 3.4 Grade 5 NSCAS-AA Summary Data: Demographics and Accommodations

Grade 5		ELA		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		253	100.00	251	100.00	251	100.00
Gender	Female	85	33.60	85	33.86	85	33.86
	Male	168	66.40	166	66.14	166	66.14
Race/Ethnicity	American Indian/Alaska Native	4	1.58	4	1.59	4	1.59
	Asian	5	1.98	5	1.99	5	1.99
	Black	23	9.09	23	9.16	23	9.16
	Hispanic	62	24.51	61	24.30	61	24.30
	Native Hawaiian or Other Pacific Islander	0	0.00	0	0.00	0	0.00
	White	148	58.50	147	58.57	147	58.57
	Two or More Races	11	4.35	11	4.38	11	4.38
Food Program	Yes	144	56.92	142	56.57	142	56.57
	No	109	43.08	109	43.43	109	43.43
LEP/ELL	Yes	6	2.37	6	2.39	6	2.39
	No	247	97.63	245	97.61	245	97.61
Accommodations	Paper and Pencil	52	20.55	51	20.32	49	19.52
	Computation Supports	17	6.72	77	30.68	30	11.95
	Assistive Technology	23	9.09	22	8.76	20	7.97
	Specialized Presentation	32	12.65	30	11.95	31	12.35

Table 3.5 Grade 6 NSCAS-AA Summary Data: Demographics and Accommodations

Grade 6		ELA		Mathematics	
		Count	%	Count	%
All Students		234	100.00	235	100.00
Gender	Female	82	35.04	83	35.32
	Male	152	64.96	152	64.68
Race/Ethnicity	American Indian/Alaska Native	2	0.85	2	0.85
	Asian	10	4.27	10	4.26
	Black	21	8.97	21	8.94
	Hispanic	38	16.24	38	16.17
	Native Hawaiian or Other Pacific Islander	2	0.85	2	0.85
	White	147	62.82	148	62.98
	Two or More Races	14	5.98	14	5.96
Food Program	Yes	121	51.71	122	51.91
	No	113	48.29	113	48.09
LEP/ELL	Yes	4	1.71	4	1.70
	No	230	98.29	231	98.30
Accommo- dations	Paper and Pencil	39	16.67	40	17.02
	Computation Supports	25	10.68	90	38.30
	Assistive Technology	24	10.26	25	10.64
	Specialized Presentation	35	14.96	36	15.32

Table 3.6 Grade 7 NSCAS-AA Summary Data: Demographics and Accommodations

Grade 7		ELA		Mathematics	
		Count	%	Count	%
All Students		262	100.00	260	100.00
Gender	Female	94	35.88	93	35.77
	Male	168	64.12	167	64.23
Race/Ethnicity	American Indian/Alaska Native	7	2.67	7	2.69
	Asian	7	2.67	7	2.69
	Black	20	7.63	20	7.69
	Hispanic	54	20.61	53	20.38
	Native Hawaiian or Other Pacific Islander	1	0.38	1	0.38
	White	154	58.78	153	58.85
	Two or More Races	19	7.25	19	7.31
Food Program	Yes	151	57.63	150	57.69
	No	111	42.37	110	42.31
LEP/ELL	Yes	4	1.53	4	1.54
	No	258	98.47	256	98.46
Accommo- dations	Paper and Pencil	43	16.41	43	16.54
	Computation Supports	30	11.45	122	46.92
	Assistive Technology	39	14.89	40	15.38
	Specialized Presentation	36	13.74	36	13.85

Table 3.7 Grade 8 NSCAS-AA Summary Data: Demographics and Accommodations

Grade 8		ELA		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		239	100.00	235	100.00	237	100.00
Gender	Female	74	30.96	73	31.06	74	31.22
	Male	165	69.04	162	68.94	163	68.78
Race/Ethnicity	American Indian/Alaska Native	2	0.84	2	0.85	2	0.84
	Asian	7	2.93	7	2.98	7	2.95
	Black	20	8.37	20	8.51	20	8.44
	Hispanic	55	23.01	55	23.40	55	23.21
	Native Hawaiian or Other Pacific Islander	0	0.00	0	0.00	0	0.00
	White	137	57.32	133	56.60	135	56.96
	Two or More Races	18	7.53	18	7.66	18	7.59
Food Program	Yes	147	61.51	145	61.70	147	62.03
	No	92	38.49	90	38.30	90	37.97
LEP/ELL	Yes	1	0.42	1	0.43	1	0.42
	No	238	99.58	234	99.57	236	99.58
Accommo- dations	Paper and Pencil	43	17.99	41	17.45	41	17.30
	Computation Supports	15	6.28	87	37.02	32	13.50
	Assistive Technology	32	13.39	32	13.62	31	13.08
	Specialized Presentation	34	14.23	33	14.04	34	14.35

Table 3.8 High School NSCAS-AA Summary Data: Demographics and Accommodations

HS		ELA		Mathematics		Science	
		Count	%	Count	%	Count	%
All Students		229	100.00	227	100.00	226	100.00
Gender	Female	87	37.99	86	37.89	86	38.05
	Male	142	62.01	141	62.11	140	61.95
Race/Ethnicity	American Indian/Alaska Native	5	2.18	5	2.20	5	2.21
	Asian	5	2.18	5	2.20	5	2.21
	Black	23	10.04	23	10.13	23	10.18
	Hispanic	50	21.83	50	22.03	50	22.12
	Native Hawaiian or Other Pacific Islander	2	0.87	2	0.88	2	0.88
	White	136	59.39	134	59.03	133	58.85
	Two or More Races	8	3.49	8	3.52	8	3.54
Food Program	Yes	126	55.02	125	55.07	125	55.31
	No	103	44.98	102	44.93	101	44.69
LEP/ELL	Yes	2	0.87	2	0.88	2	0.88
	No	227	99.13	225	99.12	224	99.12
Accommo- dations	Paper and Pencil	33	14.41	42	18.50	31	13.72
	Computation Supports	18	7.86	62	27.31	24	10.62
	Assistive Technology	9	3.93	7	3.08	9	3.98
	Specialized Presentation	10	4.37	9	3.96	9	3.98

Chapter 4: CLASSICAL ITEM STATISTICS

This chapter provides an overview of the most familiar item-level statistics obtained from classical item analysis: item difficulty, item discrimination, distractor distribution, and omits or blanks. The following results pertain only to operational NSCAS-AA 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 operational 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 0.79, it means 79 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 NSCAS-AA, 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 0.80. 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 through 4.1.3. In general, most of the items fall into the p -value range of 0.4 to 0.8, which is appropriate for a criterion-referenced assessment. In reading the following tables, the heading ≤ 0.1 describes items between 0.0 and 0.1, and the heading ≤ 0.2 describes items between 0.1 and 0.2, etc.

Table 4.1.1 Summary of Proportion Correct for NSCAS-AAELA Operational 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	1	9	7	8	2	0	0	0.550	28
4	0	0	0	4	7	7	6	4	0	0	0.559	28
5	0	0	0	2	4	13	4	5	0	0	0.570	28
6	0	0	0	1	11	9	6	1	0	0	0.531	28
7	0	0	0	5	5	7	8	3	0	0	0.546	28
8	0	0	0	3	9	9	6	1	0	0	0.531	28
HS	0	0	0	1	4	6	11	5	1	0	0.615	28

Table 4.1.2 Summary of Proportion Correct for NSCAS-AAM Operational 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	0	2	4	14	4	1	0	0	0.535	25
4	0	0	0	3	9	10	5	3	0	0	0.536	30
5	0	0	2	2	6	10	7	3	0	0	0.547	30
6	0	0	1	6	6	10	5	2	0	0	0.507	30
7	0	0	0	5	7	11	5	2	0	0	0.528	30
8	0	0	0	8	5	10	4	1	2	0	0.515	30
HS	0	0	0	3	2	11	9	4	1	0	0.583	30

Table 4.1.3 Summary of Proportion Correct for NSCAS-AAS Operational 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
5	0	0	1	5	3	5	6	4	1	0	0.560	25
8	0	0	0	3	9	3	5	5	0	0	0.557	25
HS	0	0	0	3	4	7	7	6	1	0	0.599	28

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 NSCAS-AA tests, point-biserial 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 NSCAS-AA overall) would be more likely to answer any given NSCAS-AA item correctly, while students with low ability (i.e., those who perform poorly on the NSCAS-AA 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 are 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.

No items in the 2023 NSCAS-AA tests have negative point-biserial correlations and most are above 0.30 as seen below in Tables 4.2.1 through 4.2.3, indicating good item discrimination.

Table 4.2.1 Summary of Point-biserial Correlations for NSCAS-AAELA

	Item Point-biserial Correlation							
Grade	≤ 0.1	≤ 0.2	≤ 0.3	≤ 0.4	≤ 0.5	≤ 0.6	> 0.6	Total
3	1	0	6	11	9	1	0	28
4	0	2	4	5	11	6	0	28
5	0	0	6	7	12	3	0	28
6	1	0	4	9	9	4	1	28
7	0	1	6	4	13	4	0	28
8	0	1	2	5	12	8	0	28
HS	0	1	3	6	10	8	0	28

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

	Item Point-biserial Correlation							
Grade	≤ 0.1	≤ 0.2	≤ 0.3	≤ 0.4	≤ 0.5	≤ 0.6	> 0.6	Total
3	0	4	4	5	9	3	0	25
4	0	2	6	11	10	1	0	30
5	0	2	9	5	9	5	0	30
6	0	5	7	9	7	2	0	30
7	0	1	5	11	10	3	0	30
8	0	2	5	14	8	1	0	30
HS	0	2	6	8	7	6	1	30

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

Grade	Item Point-biserial Correlation							Total
	≤0.1	≤0.2	≤0.3	≤0.4	≤0.5	≤0.6	>0.6	
5	0	2	3	8	6	6	0	25
8	0	1	3	5	13	3	0	25
HS	0	1	7	9	5	5	1	28

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 NSCAS-AAELA, NSCAS-AAM, NSCAS-AAS operational items are presented in Appendix F, G, and H. Based on these analyses, items were selected for review if the *p*-value was less than 0.25 and the point-biserial 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.

Chapter 5: RASCH ITEM CALIBRATION

The psychometric model used for the NSCAS-AA 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 NSCAS-AA 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 2023 NSCAS-AAELA, NSCAS-AAM, and NSCAS-AAS assessment.

5.1. DESCRIPTION OF THE RASCH MODEL

The Rasch dichotomous model was used to calibrate the NSCAS-AA items. All NSCAS-AA 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 NSCAS-AA, 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.

5.2.1. Unidimensionality:

Rasch models assume that one dominant dimension determines the difference among students' performances. Principal components analysis of residuals (PCAR) 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 PCAR results for the ELA, 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 NSCAS-AAELA explained about 27 to 32 percent of the total variance across Grades 3–8 and High School. The eigenvalues of the second dimension ranged from 2.1 to 3.0. This indicates that the second dimension accounted for only 2.1 to 2.6 units out of 25 to 30 units of total variance (the total variance is the same as the number of items on a form). Similar patterns are observed for the mathematics, and science assessments. Overall, the PCAR suggests that there is one clearly dominant dimension for each NSCAS-AA assessment.

Table 5.2.1 NSCAS-AAELA Results from PCAR

Grade	Contrast	Eigenvalue	Explained Variance
3	Measures	10.6	27.4%
	1	2.6	9.4%
	2	1.7	6.1%
	3	1.7	5.9%
	4	1.6	5.7%
	5	1.5	5.3%
4	Measures	12.1	30.1%
	1	2.4	8.7%
	2	1.8	6.5%
	3	1.6	5.6%
	4	1.5	5.4%
	5	1.4	5.1%
5	measures	11.2	28.7%
	1	2.5	8.8%
	2	1.7	6.1%
	3	1.7	5.9%
	4	1.6	5.5%
	5	1.4	5.1%
6	Measures	10.2	26.7%
	1	2.2	7.8%
	2	1.9	6.9%
	3	1.8	6.6%
	4	1.5	5.2%
	5	1.4	5.2%
7	Measures	10.9	28.0%
	1	2.6	9.3%
	2	1.9	6.7%
	3	1.7	6.1%
	4	1.4	5.1%
	5	1.4	4.8%
8	Measures	12.1	30.1%
	1	2.3	8.2%
	2	1.9	6.8%
	3	1.8	6.3%
	4	1.4	5.1%
	5	1.4	5.1%
HS	Measures	12.9	31.6%
	1	2.6	9.4%
	2	2.1	7.4%
	3	1.7	6.1%
	4	1.5	5.5%
	5	1.5	5.4%

Table 5.2.2 NSCAS-AAM Results from PCAR

Grade	Contrast	Eigenvalue	Explained Variance
3	Measures	8.2	24.8%
	1	3.0	11.9%
	2	1.9	7.6%
	3	1.5	5.9%
	4	1.4	5.5%
	5	1.3	5.0%
4*	Measures	10.2	25.3%
	1	3.2	10.5%
	2	2.0	6.7%
	3	1.5	5.1%
	4	--	--
	5	--	--
5	Measures	11.3	27.4%
	1	4.2	14.0%
	2	2.3	7.6%
	3	1.6	5.3%
	4	1.4	4.6%
	5	1.4	4.5%
6	Measures	9.4	23.9%
	1	3.9	12.9%
	2	2.3	7.8%
	3	1.6	5.2%
	4	1.5	5.1%
	5	1.4	4.8%
7	Measures	10.6	26.1%
	1	3.0	9.9%
	2	2.1	7.2%
	3	1.7	5.7%
	4	1.5	4.9%
	5	1.4	4.5%
8	Measures	11.7	28.0%
	1	3.2	10.8%
	2	2.0	6.8%
	3	1.7	5.7%
	4	1.7	5.6%
	5	1.4	4.6%
HS	Measures	11.1	27.1%
	1	3.2	10.5%
	2	1.9	6.5%
	3	1.7	5.7%
	4	1.5	5.0%
	5	1.4	4.7%

*Only eigenvalues greater than 1 reported

Table 5.2.3 NSCAS-AAS Results from PCAR

Grade	Contrast	Eigenvalue	Explained Variance
5	Measures	11.7	31.9%
	1	2.5	10.0%
	2	1.9	7.6%
	3	1.7	6.7%
	4	1.4	5.8%
	5	1.4	5.5%
8	Measures	12.0	32.4%
	1	2.4	9.5%
	2	1.8	7.1%
	3	1.5	6.0%
	4	1.5	6.0%
	5	1.4	5.6%
HS	Measures	12.1	30.2%
	1	3.2	11.4%
	2	2.2	7.7%
	3	1.8	6.5%
	4	1.5	5.3%
	5	1.4	5.1%

5.2.2. Local Independence:

Local independence (LI) is a fundamental assumption of Rasch models. 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 NSCAS-AA 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 NSCAS-AA tests (since all of the NSCAS-AA 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, 5.2.5, and 5.2.6 show the summary statistics—median, interquartile range (IQR), 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 median

residual correlations were slightly negative, and the values were around -0.04 . Most of the correlations were very small, suggesting local item independence generally holds for the NSCAS-AA ELA, mathematics, and science assessments.

Table 5.2.4 Summary of Item Residual Correlations for NSCAS-AAELA

Statistics	3	4	5	6	7	8	HS
N	378	378	378	378	378	378	378
Median	-0.04	-0.04	-0.03	-0.04	-0.04	-0.04	-0.04
IQR	0.12	0.11	0.11	0.12	0.12	0.12	0.13
Minimum	-0.24	-0.28	-0.27	-0.23	-0.28	-0.25	-0.27
P10	-0.14	-0.14	-0.13	-0.15	-0.14	-0.14	-0.16
P25	-0.10	-0.09	-0.10	-0.10	-0.10	-0.10	-0.10
P50	-0.04	-0.04	-0.03	-0.04	-0.04	-0.04	-0.04
P75	0.02	0.02	0.02	0.02	0.02	0.02	0.03
P90	0.09	0.08	0.08	0.08	0.09	0.07	0.09
Maximum	0.26	0.19	0.24	0.27	0.21	0.32	0.25
>0.20	3	0	2	2	1	2	4

Table 5.2.5 Summary of Item Residual Correlations for NSCAS-AAM

Statistics	3	4	5	6	7	8	HS
N	300	435	435	435	435	435	435
Median	-0.05	-0.04	-0.05	-0.04	-0.05	-0.04	-0.04
IQR	0.15	0.13	0.20	0.16	0.15	0.16	0.13
Minimum	-0.30	-0.29	-0.31	-0.34	-0.24	-0.27	-0.31
P10	-0.17	-0.15	-0.20	-0.19	-0.15	-0.16	-0.16
P25	-0.12	-0.10	-0.13	-0.12	-0.11	-0.11	-0.10
P50	-0.05	-0.04	-0.05	-0.04	-0.05	-0.04	-0.04
P75	0.03	0.02	0.07	0.04	0.04	0.05	0.03
P90	0.11	0.11	0.15	0.14	0.10	0.11	0.10
Maximum	0.26	0.32	0.38	0.44	0.30	0.24	0.32
>0.20	3	8	19	18	7	9	4

Table 5.2.6 Summary of Item Residual Correlations for NSCAS-AAS

Statistics	5	8	HS
<i>N</i>	300	300	378
Median	-0.05	-0.04	-0.04
IQR	0.13	0.13	0.15
Minimum	-0.28	-0.24	-0.34
P10	-0.16	-0.15	-0.17
P25	-0.11	-0.11	-0.11
P50	-0.05	-0.04	-0.04
P75	0.02	0.02	0.03
P90	0.09	0.07	0.12
Maximum	0.27	0.24	0.30
>0.20	5	2	8

5.2.3. 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 NSCAS-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 NSCAS-AA ELA, 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 NSCAS-AA ELA 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 NSCAS-AA item data fits Rasch model well.

Table 5.2.7 Summary of Infit and Outfit Mean Square Statistics for 2023 NSCAS-AA Tests

		Infit Mean Square					Outfit Mean Square				
		Mean	SD	MIN	MAX	[0.7, 1.3]	Mean	SD	MIN	MAX	[0.7, 1.3]
ELA	3	0.98	0.09	0.84	1.24	28/28	0.98	0.15	0.70	1.46	27/28
	4	0.97	0.12	0.79	1.26	28/28	0.91	0.16	0.66	1.27	26/28
	5	0.98	0.09	0.85	1.15	28/28	0.95	0.12	0.73	1.15	28/28
	6	0.99	0.11	0.78	1.34	27/28	0.97	0.18	0.69	1.52	25/28
	7	0.99	0.12	0.80	1.21	28/28	0.96	0.18	0.69	1.45	25/28
	8	0.99	0.13	0.82	1.43	27/28	0.96	0.25	0.62	1.80	24/28
	HS	0.99	0.14	0.77	1.30	28/28	0.93	0.21	0.48	1.42	24/28
Mathematics	3	0.97	0.12	0.79	1.21	25/25	0.96	0.18	0.69	1.36	24/25
	4	0.99	0.09	0.80	1.21	30/30	0.97	0.13	0.72	1.36	29/30
	5	0.98	0.10	0.81	1.17	30/30	0.98	0.18	0.70	1.41	29/30
	6	0.99	0.11	0.80	1.20	30/30	0.96	0.14	0.74	1.26	30/30
	7	0.98	0.09	0.84	1.21	30/30	0.96	0.12	0.75	1.30	30/30
	8	0.98	0.08	0.82	1.18	30/30	0.96	0.14	0.61	1.30	29/30
	HS	0.96	0.08	0.79	1.14	30/30	0.95	0.12	0.69	1.16	29/30
Science	5	0.99	0.14	0.79	1.34	24/25	0.99	0.26	0.60	1.60	17/25
	8	0.99	0.12	0.77	1.32	24/25	0.98	0.19	0.61	1.32	21/25
	HS	0.99	0.14	0.67	1.23	27/28	0.94	0.22	0.53	1.30	22/28

Table 5.2.8 Summary of Infit and Outfit Z STD Statistics for 2023 NSCAS-AA Tests

		Infit Z STD					Outfit Z STD				
		Mean	SD	MIN	MAX	[-2.0, 2.0]	Mean	SD	MIN	MAX	[-2.0, 2.0]
ELA	3	-0.42	1.53	-3.31	2.57	22/28	-0.24	1.13	-2.00	2.42	27/28
	4	-0.50	2.06	-3.53	4.69	17/28	-0.56	1.09	-2.30	1.78	25/28
	5	-0.42	1.57	-3.00	2.38	20/28	-0.41	1.00	-2.11	1.26	26/28
	6	-0.27	1.93	-4.22	5.34	22/28	-0.27	1.43	-2.95	3.66	23/28
	7	-0.23	2.12	-3.88	3.98	17/28	-0.34	1.29	-2.44	2.74	24/28
	8	-0.24	1.92	-3.22	4.64	18/28	-0.30	1.38	-2.14	4.26	25/28
	HS	-0.19	1.87	-2.99	4.17	18/28	-0.31	1.25	-2.22	2.85	23/28
Mathematics	3	-0.57	2.34	-4.30	3.82	11/25	-0.47	1.80	-3.54	2.82	18/25
	4	-0.29	1.65	-4.11	3.91	23/30	-0.33	1.21	-3.12	2.82	27/30
	5	-0.35	1.75	-3.12	3.28	22/30	-0.15	1.17	-2.08	2.20	28/30
	6	-0.30	2.24	-4.84	4.21	17/30	-0.24	0.99	-2.08	1.84	28/30
	7	-0.51	1.73	-3.47	3.84	22/30	-0.31	1.01	-2.10	2.35	28/30
	8	-0.38	1.39	-3.21	3.69	25/30	-0.26	0.96	-2.09	2.24	28/30
	HS	-0.68	1.48	-3.46	2.50	23/30	-0.48	1.11	-2.85	1.52	27/30
Science	5	-0.10	2.16	-3.79	5.02	16/25	-0.17	1.57	-2.73	3.07	21/25
	8	-0.23	1.81	-4.19	5.04	19/25	-0.13	1.14	-2.20	1.90	23/25
	HS	-0.16	2.10	-5.31	3.51	16/28	-0.29	1.48	-3.45	2.23	23/28

5.3. RASCH ITEM STATISTICS

Item calibration was implemented via Winsteps 5.4.3 program (Linacre, 2023). 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 I, J, and K report the 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 NSCAS-AA 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.

Table 5.3.1 Summary of Rasch Item Difficulties for 2023 NSCAS-AA Tests

	Grade	N	Mean	SD	Min	Max	Range
ELA	3	28	0.00	0.60	-1.19	1.64	2.83
	4	28	0.00	0.61	-1.18	1.05	2.22
	5	28	0.00	0.62	-1.23	1.23	2.46
	6	28	0.00	0.44	-1.06	0.70	1.77
	7	28	0.00	0.63	-1.21	1.19	2.41
	8	28	0.00	0.52	-0.95	1.13	2.08
	HS	28	0.00	0.62	-1.33	1.39	2.73
Mathematics	3	25	0.04	0.50	-1.09	0.92	2.01
	4	30	0.06	0.54	-1.04	0.99	2.02
	5	30	-0.06	0.61	-1.13	1.20	2.33
	6	30	-0.05	0.54	-1.06	1.04	2.10
	7	30	-0.08	0.53	-1.08	0.80	1.88
	8	30	-0.04	0.71	-1.75	0.90	2.65
	HS	30	-0.02	0.59	-1.43	1.09	2.52
Science	5	25	-0.06	0.80	-1.60	1.41	3.00
	8	25	-0.06	0.74	-1.22	1.19	2.42
	HS	28	-0.25	0.77	-1.74	1.25	3.00

Chapter 6: EQUATING AND SCALING

As discussed earlier in Chapter 2, the 2023 test forms were constructed with items that were either field tested or used operationally on a previously administered NSCAS-AA test. NSCAS-AA assessments are constructed each year allowing each NSCAS-AA assessment to be different from the previous year's assessment.

Typically, 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 new year's NSCAS-AA raw-score-to-Rasch-ability scale to the previous operational scale. When the new NSCAS-AA 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 were used for score reporting and performance level classification. Once operational items were equated, field test items were then placed on the bank scale and are ready for future operational use.

This chapter begins with a summary of the 2023 NSCAS-AA 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. Summary results of the state scale score performance is also provided.

6.1. Equating

The equating design employed for NSCAS-AA 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 2023 NSCAS-AA test forms were constructed with items from previous administrations. The items were previously either field-test or operational items. Rescaling happened in the first operational administration with new College and Career Ready Standards in 2018 for mathematics, 2021 for science, and 2023 for ELA. All items from 2018, 2019, 2021, and 2022 in mathematics and items from 2021, and 2022 in science, were used as equating items in the post-equating solution or the items' bank parameters were used with the pre-equating solution. The tables below show the number of operational items in the 2023 forms for each grade and subject and how many were from previous administrations where the scales are the same. There was one item in grade 5 mathematics prior to 2018. These items were calibrated after the 2023 assessment to place them onto the new bank scale. Note that the 2023 ELA assessment was the first operational assessment aligned to new College and Career Ready ELA Standards and, therefore, it is not equated to prior years. The scaling constant obtained after the standard setting is shown in Table 6.2.1.

Table 6.1.1 Items and Previous Usage for 2023 NSCAS-AAM Tests

Grade	Total No. of Core Items	Equating Items								Non-Equating Items 2017 and Earlier
		2022		2021		2019		2018		
		OP	FT	OP	FT	OP	FT	OP	FT	
3	25	11	11	2	0	0	0	1	0	0
4	30	16	8	4	0	0	0	2	0	0
5	30	18	8	2	0	1	0	0	0	1
6	30	21	6	1	0	2	0	0	0	0
7	30	15	11	3	0	1	0	0	0	0
8	30	17	8	3	0	1	0	1	0	0
HS	30	12	11	1	0	5	0	1	0	0

Table 6.1.2 Items and Previous Usage for 2023 NSCAS-AAS Tests

Grade	Total No. of Core Items	Equating Items				Non- Equating Items 2020 and Earlier
		2022		2021		
		OP	FT	OP	FT	
5	25	13	6	0	6	0
8	25	15	6	0	4	0
HS	28	19	7	0	2	0

The equating steps for 2023 administration were as follows.

1. Equating items are items previously used that are on a common scale. For mathematics, the equating items were previously administered in 2022, 2021, 2019 and 2018. For Science, the equating items were in 2022, and 2021.
2. Calibrate the operational form without anchoring (i.e., free calibration using 2023 data).
3. Evaluate the across-year stability of the item difficulty parameter estimates.
 - a. Compute the correlation and ratio of standard deviations (SD) between the calibrated item parameters and banked item parameters. Examine whether the correlation is lower than 0.95 or the ratio of standard deviations is outside the range between 0.90 and 1.10. However, these are not used in the item removal process.
 - b. Robust-Z statistics are computed. Robust Z exceeding absolute value of 2.576 are considered extreme value, equivalent to 1% of normal distribution.
 - c. Visually inspect the item parameter differences with plots of banked item parameters and item parameters from Step 2. The information from the plots, as well as content expert's opinions, shall be considered before dropping any equating items.

4. Compute the mean shift equating adjustment using the remaining equating items and apply the adjustment to the item difficulty in Step 2.
5. Winsteps is run anchoring all of the operational item parameters and the Raw-to-Logit table is obtained from the Winsteps output.
6. Scaling transformation parameters are applied to the logit values to obtain the scale scores. The scaling transformation parameters are shown in Table 6.2.2. and Table 6.2.3.
 - a. The transformation to scale score formula is: $SS = a + b \times \text{Logit}$

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 provide the basis for school accountability analyses.

Appendix M, N, and O 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 2023 NSCAS-AA ELA, mathematics, and science. 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 scale score 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 of ELA, mathematics, and science range from 100 to 300. The value of 100 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 101 even if the student gave no correct responses. No student tested will receive a scale score higher than 300 or lower than 101, 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 300 is not invoked even for a perfect paper, or a future form could be difficult enough that the lower limit is not invoked.

As part of its deliberations concerning defining the performance levels, the Nebraska State Department of Education specified that the *On Track* performance level cut score have a scale score of 200. 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. Given the scale score and the logit standards defining the performance level, it is sufficient to define the final scale score metric.

The mathematics tests have a scale score of 200 for the *On Track* performance level cut score, while the *Advanced* performance level for mathematics cut scores varies per grade. 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 mathematics scale scores were initially calculated using the following formula:

$$SS = 200 + (\text{logit} - x_{L2}) * \frac{33.33}{\sigma}$$
, where x_{L2} is the logit for the *On Track* cut score for the given grade, and σ is the standard deviation of the students with that grade.

Calculations of the slopes and intercepts for all grades of the NSCAS-AAM scale score conversion are given in Table 6.2.2. The raw-to-scale conversions are provided in Appendices N.

The ELA and Science tests have a scale score of 200 for the *On Track* performance level cut score, and a scale score of 250 for the *Advanced* performance level cut scores.

The transformation to scale scores is:

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

where:

$$b = \frac{249.501 - 199.501}{x_{L3} - x_{L2}},$$

and where x_{L3} is the logit for *Advanced* cut score and x_{L2} is the logit for *On Track* cut scores.

Therefore:

$$a = 199.501 - bx_{L2} \text{ or,}$$

$$a = 249.501 - bx_{L3}.$$

Calculations of the slopes and intercepts for all grades of the NSCAS-AAELA scale conversion are given in Table 6.2.1 and NSCAS-AAS scale score conversion are given in Table 6.2.3. The raw-to-scale conversions are provided in Appendices M and O.

Table 6.2.1 NSCAS-AAELA Conversion of Logits to Scale Scores

Grade	Logit Cut Points		Scale Score Ranges by Performance Level			Conversion	
	Dev/OT	OT/Adv	Developing	On Track	Advanced	Slope b	Intercept a
3	-0.3170	1.6351	101 to 199	200 to 249	250 to 300	25.61344	207.62046
4	-0.3116	1.9223	101 to 199	200 to 249	250 to 300	22.38238	206.47535
5	0.0058	1.9211	101 to 199	200 to 249	250 to 300	26.10557	199.34959
6	-0.1460	1.8556	101 to 199	200 to 249	250 to 300	24.98002	203.14808
7	-0.0009	1.9336	101 to 199	200 to 249	250 to 300	25.84647	199.52426
8	-0.1537	1.8856	101 to 199	200 to 249	250 to 300	24.51822	203.26945
HS	0.3096	2.7319	101 to 199	200 to 249	250 to 300	20.64154	193.11038

Table 6.2.2 NSCAS-AAM Conversion of Logits to Scale Scores

Grade	Logit Cut Points		Scale Score Ranges by Performance Level			Conversion	
	Dev/OT	OT/Adv	Developing	On Track	Advanced	Slope b	Intercept a
3	0.0910	2.1216	101 to 199	200 to 251	252 to 300	25.42528	197.68630
4	0.0005	1.7252	101 to 199	200 to 249	250 to 300	28.76996	199.98562
5	0.0297	1.8621	101 to 199	200 to 251	252 to 300	28.11709	199.16492
6	0.0051	2.3544	101 to 199	200 to 266	267 to 300	28.40706	199.85512
7	0.0780	2.8239	101 to 199	200 to 282	283 to 300	30.19295	197.64495
8	0.1596	2.4971	101 to 199	200 to 269	270 to 300	30.01080	195.21028
HS	0.4859	2.4491	101 to 199	200 to 255	256 to 300	28.31535	186.24157

Table 6.2.3 NSCAS-AAS Conversion of Logits to Scale Scores

Grade	Logit Cut Points		Scale Score Ranges by Performance Level			Conversion	
	Dev/OT	OT/Adv	Developing	On Track	Advanced	Slope b	Intercept a
5	0.4624	2.1662	101 to 199	200 to 249	250 to 300	29.34617	185.9313
8	0.1030	2.6209	101 to 199	200 to 249	250 to 300	19.85782	197.4556
HS	-0.0795	1.8508	101 to 199	200 to 249	250 to 300	25.90271	201.5603

Complete frequency distributions of the state scale scores for the NSCAS-AAELA, NSCAS-AAM, and NSCAS-AAS are provided in Appendices M, N and O 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 P. Historical student performance summary statistics after the most recent rescaling of the assessment are shown in Tables 6.2.4, 6.2.5,

and 6.2.6 for the ELA, mathematics, and science. Summary statistics from prior to the standard setting are on a different scale and not comparable to the current scale, thus they are not presented in the tables.

Table 6.2.4 2023 NSCAS-AAELA State Scale Score Summary, All Students

Year	Grade	Count	Scale Score		Quartile		
			Mean	SD	First	Second	Third
2023	3	238	209.9	37.3	195	211	228
	4	258	209.6	37.8	192	210	229
	5	253	204.7	37.9	187	204	225
	6	234	204.5	35.9	188	203	227
	7	262	203.2	36.3	183	204	221
	8	239	204.1	39.5	184	200	227
	HS	229	204.9	35.0	183	203	227

Table 6.2.5 2023 NSCAS-AAM State Scale Score Summary, All Students

Year	Grade	Count	Scale Score		Quartile		
			Mean	SD	First	Second	Third
2023	3	235	202.7	35.4	186	199	221
	4	255	203.8	37.6	185	206	223
	5	251	199.6	37.7	181	201	223
	6	235	194.6	35.9	182	198	215
	7	260	197.7	38.4	178	195	223
	8	235	193.3	38.5	174	190	218
	HS	227	196.8	35.8	173	194	222
2022	3	235	201.7	36.5	188	201	224
	4	243	197.2	37.8	186	199	216
	5	235	197.0	36.7	182	199	216
	6	254	197.3	32.5	183	197	220
	7	239	196.5	42.8	178	196	218
	8	236	196.4	36.3	177	193	218
	HS	209	185.1	33.2	169	187	204
2021	3	205	200.7	40.0	190	203	221
	4	203	200.4	39.4	181	198	219
	5	223	199.2	36.8	180	201	222
	6	210	196.6	35.4	184	200	217
	7	222	196.6	38.0	177	196	219
	8	213	199.7	34.8	181	199	224
	HS	213	190.6	30.2	171	188	210
2019	3	244	202.6	44.5	183	206	230
	4	245	201.4	38.5	186	202	228
	5	236	202.0	40.8	180	197	222
	6	234	194.6	34.2	180	197	217
	7	232	202.7	37.4	182	205	230
	8	241	197.7	44.9	172	190	225
	HS	239	193.8	33.8	173	194	218

Table 6.2.6 2023 NSCAS-AAS State Scale Score Summary, All Students

Year	Grade	Count	Scale Score		Quartile		
			Mean	<i>SD</i>	First	Second	Third
2023	5	251	190.5	39.6	171	192	216
	8	237	199.6	32.6	187	198	217
	11	226	207.2	36.0	182	208	233
2022	5	235	191.1	40.2	167	194	217
	8	235	200.6	30.2	185	200	216
	HS	209	199.9	40.6	184	202	219

Chapter 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 levels 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 and point-biserial correlation.

7.1. CLASSICAL ITEM STATISTICS

Indices known as classical item statistics include 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 NSCAS-AAELA field test item in Appendix F, for NSCAS-AAM in Appendix G, and NSCAS-AAS in Appendix H. Tables 7.1.1 and 7.1.2 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 0.8 are avoided when possible. In reading the following tables, the heading ≤ 0.1 describes items between 0.0 and 0.1, and the heading ≤ 0.2 describes items between 0.1 and 0.2, etc.

Table 7.1.1 Summary of Proportion Correct for NSCAS-AA 2023 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		
ELA	3	0	0	4	4	5	2	1	0	0	0	0.411	16
	4	0	0	0	3	7	5	1	0	0	0	0.463	16
	5	0	0	2	3	6	2	2	1	0	0	0.463	16
	6	0	0	1	4	7	3	1	0	0	0	0.437	16
	7	0	0	1	5	6	2	2	0	0	0	0.454	16
	8	0	0	1	2	4	5	4	0	0	0	0.499	16
	HS	0	0	0	1	3	6	4	2	0	0	0.562	16
Mathematics	3	0	0	2	3	1	2	0	0	0	0	0.395	8
	4	0	0	2	0	2	3	1	0	0	0	0.478	8
	5	0	0	1	2	0	4	1	0	0	0	0.483	8
	6	0	0	3	2	3	0	0	0	0	0	0.355	8
	7	0	0	2	0	5	1	0	0	0	0	0.417	8
	8	0	0	1	4	2	1	0	0	0	0	0.402	8
	HS	0	0	1	2	1	3	1	0	0	0	0.461	8
Science	5	0	0	1	0	2	8	4	1	0	0	0.565	16
	8	0	0	0	5	3	3	4	1	0	0	0.509	16
	HS	0	0	0	2	3	5	4	2	0	0	0.553	16

Table 7.1.2 Summary of Point-biserial Correlations for NSCAS-AA 2023 Field Test Items

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

Chapter 8: RELIABILITY

This chapter addresses the reliability of NSCAS-AA 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 prerequisite for making appropriate interpretations (i.e., showing evidence of valid use of results). Conceptually, reliability can be referred to as the consistency of the results between two measures of the same thing. This consistency can be seen in the degree of agreement between two measures on two occasions. Operationally, such comparisons are the essence of the mathematically-defined reliability indices.

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

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

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

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

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 2023 administration of the NSCAS-AA 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 NSCAS-AA. The table contains test length in total number of items (L), test reliabilities, and traditional standard errors of measurement (SEM) in raw score points. As can be seen in the table, ELA, mathematics, and science forms have Coefficient Alphas over 0.84. Overall, these reliability values provide evidence of good reliability.

Table 8.1.1 Reliabilities and Standard Errors of Measurement

	Grade	<i>L</i>	Reliability	<i>SEM</i>
ELA	3	28	0.87	2.31
	4	28	0.89	2.26
	5	28	0.89	2.28
	6	28	0.89	2.31
	7	28	0.88	2.30
	8	28	0.90	2.26
	HS	28	0.90	2.20
Mathematics	3	25	0.86	2.22
	4	30	0.87	2.42
	5	30	0.88	2.39
	6	30	0.86	2.45
	7	30	0.88	2.41
	8	30	0.88	2.38
	HS	30	0.88	2.37
Science	5	25	0.87	2.12
	8	25	0.88	2.12
	HS	28	0.88	2.23

Appendix Q presents α for the content strands. Given that the reliability is a function of test length, the smaller item counts for the content standards result in lower reliability, 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 NSCAS-AA 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:

$$SEM = SD\sqrt{1 - \text{reliability}}. \quad (8.3)$$

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 NSCAS-AA 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 about 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 M, N, and O 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 NSCAS-AA tests: 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 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 test score that has the smallest *CSEM* (Min *CSEM*), the maximum *CSEM* of the scale score associated with a zero/perfect total test score (Max *CSEM*), *CSEM* at the cuts of *Developing* and *On Track* performance levels (*CSEM Dev/OT*), and *CSEM* at the cuts of *On Track* and *Advanced* performance levels (*CSEM OT/ADV*) 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 2023 NSCAS-AA Tests

	Grade	Min CSEM	Max CSEM	CSEM Dev/OT	CSEM OT/Adv
ELA	3	10	47	10	13
	4	9	41	9	12
	5	10	48	10	14
	6	10	46	10	14
	7	10	47	10	14
	8	10	45	10	13
	HS	8	38	8	15
Mathematics	3	10	47	10	16
	4	11	53	11	14
	5	11	52	11	15
	6	11	52	11	21
	7	11	55	11	31
	8	12	55	12	22
	HS	11	52	11	21
Science	5	13	54	13	22
	8	8	37	9	20
	HS	10	48	11	16

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 NSCAS-AA 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 NSCAS-AA 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) but considerably lower than the method developed by Stearns and Smith (2007).

The results for the overall consistency across all three achievement levels are presented in Table 8.4.1, Table 8.4.2, and Table 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 high 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. With the 2023 administration, dichotomous decisions using the Advanced cuts (On Track/Advanced) generally have the higher consistency values and exceeded 0.90 in most cases. The pattern of decision accuracy across different cuts is like that of decision consistency.

Table 8.4.1 NSCAS-AAELA Decision Consistency Results

Grade	Livingston & Lewis				Hanson & Brennan			
	Decision Accuracy		Decision Consistency		Decision Accuracy		Decision Consistency	
	OT	ADV	OT	ADV	OT	ADV	OT	ADV
3	0.909	0.916	0.870	0.889	0.908	0.916	0.874	0.890
4	0.913	0.927	0.878	0.901	0.915	0.928	0.881	0.903
5	0.909	0.898	0.871	0.879	0.909	0.898	0.875	0.879
6	0.901	0.937	0.860	0.917	0.900	0.936	0.862	0.918
7	0.900	0.898	0.860	0.884	0.901	0.898	0.864	0.885
8	0.908	0.944	0.871	0.924	0.907	0.944	0.872	0.925
HS	0.908	0.939	0.870	0.919	0.908	0.939	0.873	0.921

Table 8.4.2 NSCAS-AAM Decision Consistency Results

Grade	Livingston & Lewis				Hanson & Brennan			
	Decision Accuracy		Decision Consistency		Decision Accuracy		Decision Consistency	
	OT	ADV	OT	ADV	OT	ADV	OT	ADV
3	0.884	0.964	0.838	0.953	0.884	0.963	0.837	0.952
4	0.894	0.931	0.852	0.909	0.895	0.930	0.853	0.909
5	0.903	0.940	0.863	0.924	0.906	0.940	0.868	0.918
6	0.886	0.987	0.839	0.987	0.890	0.987	0.845	0.987
7	0.894	0.988	0.852	0.985	0.895	0.988	0.854	0.983
8	0.893	0.979	0.850	0.975	0.894	0.979	0.852	0.973
HS	0.895	0.965	0.853	0.953	0.894	0.966	0.854	0.951

Table 8.4.3 NSCAS-AAS Decision Consistency Results

Grade	Livingston & Lewis				Hanson & Brennan			
	Decision Accuracy		Decision Consistency		Decision Accuracy		Decision Consistency	
	OT	ADV	OT	ADV	OT	ADV	OT	ADV
5	0.890	0.956	0.845	0.947	0.894	0.956	0.850	0.942
8	0.898	0.962	0.857	0.957	0.899	0.962	0.859	0.954
HS	0.899	0.925	0.858	0.904	0.900	0.925	0.861	0.904

Chapter 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 NSCAS-AA 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 NSCAS-AA for grades 3 to 8 and High School is a criterion-referenced assessment. The criteria referenced are the Nebraska ELA, mathematics, and science content standards. Each assessment was based on and was directly aligned to the Nebraska statewide alternate content standards (i.e., extended indicators) 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 in chapter 2, 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.

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.

9.2.1. Item-Test Correlations:

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

9.2.2. Item Response Theory Dimensionality:

Results from principal components analyses are presented in Chapter Five. The NSCAS-AA ELA, mathematics, and science tests were essentially unidimensional, providing evidence supporting interpretations based on the total scores for the respective NSCAS-AA tests.

9.2.3. 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 NSCAS-AAELA tests have three strands (denoted by E.1, E.2, and E.3), NSCAS-AAM tests have four strands (denoted by M.1, M.2, M.3, and M.4), and the NSCAS-AAS tests have three strands (denoted by S.1, S.2, and S.3) for each grade.

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 NSCAS-AA Content Strands

Content	Code	Strand
ELA	E.1	Vocabulary
	E.2	Comprehension
	E.3	Writing
Mathematics	M.1	Number Sense
	M.2	Algebraic
	M.3	Geometric/Measurement
	M.4	Data Analysis/Probability
Science	S.1	Physical Science
	S.2	Life Science
	S.3	Earth and Space Sciences

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

Grade 3	E.1	E.2	E.3	M.1	M.2	M.3	M.4
E.1	—						
E.2	0.60	—					
E.3	0.52	0.69	—				
M.1	0.60	0.62	0.58	—			
M.2	0.31	0.55	0.58	0.57	—		
M.3	0.48	0.65	0.57	0.67	0.59	—	
M.4	0.51	0.42	0.40	0.53	0.37	0.41	—

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

Grade 4	E.1	E.2	E.3	M.1	M.2	M.3	M.4
E.1	—						
E.2	0.82	—					
E.3	0.67	0.67	—				
M.1	0.64	0.71	0.58	—			
M.2	0.58	0.62	0.53	0.65	—		
M.3	0.68	0.70	0.55	0.73	0.61	—	
M.4	0.56	0.64	0.49	0.58	0.52	0.58	—

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

Grade 5	E.1	E.2	E.3	M.1	M.2	M.3	M.4	S.1	S.2	S.3
E.1	—									
E.2	0.71	—								
E.3	0.56	0.68	—							
M.1	0.61	0.68	0.57	—						
M.2	0.54	0.62	0.51	0.68	—					
M.3	0.68	0.74	0.62	0.73	0.64	—				
M.4	0.59	0.69	0.52	0.60	0.64	0.72	—			
S.1	0.58	0.76	0.63	0.64	0.56	0.69	0.63	—		
S.2	0.56	0.72	0.60	0.64	0.58	0.68	0.63	0.64	—	
S.3	0.68	0.79	0.61	0.67	0.59	0.75	0.71	0.71	0.73	—

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

Grade 6	E.1	E.2	E.3	M.1	M.2	M.3	M.4
E.1	—						
E.2	0.76	—					
E.3	0.55	0.65	—				
M.1	0.67	0.70	0.58	—			
M.2	0.51	0.59	0.51	0.63	—		
M.3	0.42	0.49	0.50	0.52	0.62	—	
M.4	0.50	0.53	0.49	0.55	0.46	0.37	—

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

Grade 7	E.1	E.2	E.3	M.1	M.2	M.3	M.4
E.1	—						
E.2	0.72	—					
E.3	0.65	0.66	—				
M.1	0.61	0.63	0.57	—			
M.2	0.62	0.72	0.61	0.65	—		
M.3	0.60	0.61	0.53	0.63	0.61	—	
M.4	0.60	0.70	0.59	0.62	0.61	0.53	—

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

Grade 8	E.1	E.2	E.3	M.1	M.2	M.3	M.4	S.1	S.2	S.3
E.1	—									
E.2	0.77	—								
E.3	0.68	0.72	—							
M.1	0.61	0.57	0.64	—						
M.2	0.60	0.69	0.63	0.62	—					
M.3	0.68	0.73	0.65	0.62	0.67	—				
M.4	*	*	*	*	*	*	—			
S.1	0.71	0.77	0.67	0.60	0.64	0.69	*	—		
S.2	0.74	0.74	0.70	0.65	0.63	0.68	*	0.75	—	
S.3	0.65	0.66	0.62	0.64	0.61	0.66	*	0.66	0.71	—

*M.4 (Data Analysis/Probability) had no items assessed.

Table 9.2.2.g Correlations between ELA, Mathematics, and Science Strands for Grade HS

Grade HS	E.1	E.2	E.3	M.1	M.2	M.3	M.4	S.1	S.2	S.3
E.1	—									
E.2	0.71	—								
E.3	0.60	0.76	—							
M.1	0.55	0.69	0.59	—						
M.2	0.56	0.64	0.57	0.66	—					
M.3	0.61	0.66	0.63	0.63	0.70	—				
M.4	0.63	0.65	0.55	0.58	0.58	0.59	—			
S.1	0.67	0.77	0.70	0.69	0.68	0.74	0.64	—		
S.2	0.64	0.69	0.67	0.57	0.62	0.65	0.64	0.74	—	
S.3	0.60	0.70	0.64	0.54	0.52	0.61	0.55	0.70	0.63	—

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 through 9.2.2.g) and the reliabilities for each strand (Spearman, 1904, 1910). Disattenuated correlations very close to 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 2023 NSCAS-AA 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 higher than across-content-area strand correlations. 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 ELA and Mathematics: Grade 3

Grade 3	E.1	E.2	E.3	M.1	M.2	M.3	M.4
E.1	—						
E.2	0.84	—					
E.3	0.88	1.00	—				
M.1	0.95	0.86	0.97	—			
M.2	0.49	0.76	0.96	0.86	—		
M.3	0.73	0.85	0.90	1.00	0.86	—	
M.4	1.00	0.75	0.87	1.00	0.73	0.80	—

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

Grade 4	E.1	E.2	E.3	M.1	M.2	M.3	M.4
E.1	—						
E.2	1.00	—					
E.3	1.00	1.00	—				
M.1	0.90	0.92	0.95	—			
M.2	0.97	0.97	1.00	1.00	—		
M.3	0.97	0.92	0.91	1.00	1.00	—	
M.4	0.96	1.00	0.97	0.98	1.00	1.00	—

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

Grade 5	E.1	E.2	E.3	M.1	M.2	M.3	M.4	S.1	S.2	S.3
E.1	—									
E.2	1.00	—								
E.3	1.00	1.00	—							
M.1	0.93	0.89	0.98	—						
M.2	1.00	1.00	1.00	1.00	—					
M.3	1.00	0.98	1.00	1.00	1.00	—				
M.4	0.93	0.92	0.90	0.88	1.00	1.00	—			
S.1	0.96	1.00	1.00	0.99	1.00	1.00	1.00	—		
S.2	0.92	0.99	1.00	0.97	1.00	1.00	0.97	1.00	—	
S.3	0.99	0.98	0.99	0.90	0.99	1.00	0.99	1.00	1.00	—

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

Grade 6	E.1	E.2	E.3	M.1	M.2	M.3	M.4
E.1	—						
E.2	1.00	—					
E.3	0.90	0.98	—				
M.1	0.94	0.90	0.92	—			
M.2	0.82	0.87	0.93	0.97	—		
M.3	0.64	0.68	0.86	0.77	1.00	—	
M.4	1.00	1.00	1.00	1.00	1.00	0.80	—

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

Grade 7	E.1	E.2	E.3	M.1	M.2	M.3	M.4
E.1	—						
E.2	0.98	—					
E.3	1.00	1.00	—				
M.1	0.88	0.82	0.95	—			
M.2	0.90	0.95	1.00	0.90	—		
M.3	1.00	0.92	1.00	1.00	0.96	—	
M.4	1.00	1.00	1.00	0.98	0.98	0.97	—

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

Grade 8	E.1	E.2	E.3	M.1	M.2	M.3	M.4	S.1	S.2	S.3
E.1	—									
E.2	1.00	—								
E.3	0.96	0.95	—							
M.1	0.87	0.77	0.95	—						
M.2	0.83	0.89	0.89	0.89	—					
M.3	0.91	0.92	0.90	0.87	0.90	—				
M.4	*	*	*	*	*	*	—			
S.1	0.95	0.97	0.93	0.84	0.86	0.90	*	—		
S.2	1.00	0.96	1.00	0.94	0.87	0.91	*	1.00	—	
S.3	1.00	0.97	1.00	1.00	0.96	1.00	*	1.00	1.00	—

*M.4 (Data Analysis/Probability) had no item assessed.

Table 9.2.3.g Disattenuated Strand Correlations for ELA, Mathematics, and Science: Grade HS

Grade HS	E.1	E.2	E.3	M.1	M.2	M.3	M.4	S.1	S.2	S.3
E.1	—									
E.2	0.94	—								
E.3	0.88	1.00	—							
M.1	0.85	0.96	0.91	—						
M.2	0.85	0.87	0.86	1.00	—					
M.3	0.85	0.83	0.88	0.92	1.00	—				
M.4	0.98	0.90	0.85	0.94	0.94	0.87	—			
S.1	0.92	0.95	0.95	0.98	0.96	0.96	0.92	—		
S.2	0.95	0.91	1.00	0.88	0.96	0.91	0.99	1.00	—	
S.3	0.88	0.92	0.93	0.84	0.79	0.85	0.85	0.96	0.94	—

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 NSCAS-AA, 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 NSCAS-AA data, indicating the appropriateness of using the Rasch models to analyze the NSCAS-AA data.

In addition, the Rasch model was also used to link different operational NSCAS-AA 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.

Chapter 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).
- 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. (2023). Winsteps® Rasch measurement computer program (V5.4.3). 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.