



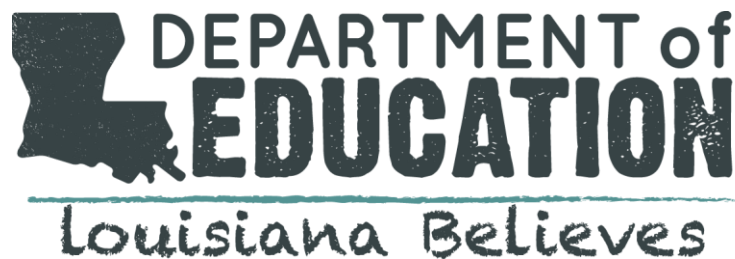
Pearson



# LEAP 2025 Social Studies Grades 3–8 Technical Report: 2020–2021

Prepared by DRC, Pearson, and WestEd

# LEAP 2025



## EXECUTIVE SUMMARY

---

The Louisiana Educational Assessment Program 2025 (LEAP 2025) is composed of tests that are carefully constructed to fairly assess the achievement of Louisiana students. This technical addendum provides information on the operational test administrations, scoring activities, analyses, and results of the spring 2021 administration of the LEAP 2025 Social Studies test, which used intact forms based on previously administered operational forms. For information on the development and construction processes for these forms, see the [2019 LEAP 2025 Social Studies Grades 3-8 Technical Report](#).

While this technical report and its associated materials have been produced in a way that can help educators understand the technical characteristics of the assessment used to measure student achievement, the information is primarily intended for use by those who evaluate tests, interpret scores, or use test results in making educational decisions. It is assumed that the reader has technical knowledge of test construction and measurement procedures, as stated in *Standards for Educational and Psychological Testing* (American Educational Research Association, American Psychological Association, and National Council on Measurement in Education, 2014).

The chapters of this technical report outline general information about the administration and scoring activities of the LEAP 2025 assessments, CTT (Classical Test Theory) and IRT (Item Response Theory) analysis results, and the interpretation of the scores on the tests. Additionally, because of conditions related to COVID-19, please use caution when making any inferences from the statistical results of the spring 2021 administration.

# Table of Contents

<b>EXECUTIVE SUMMARY .....</b>	<b>ii</b>
<b>1. Introduction .....</b>	<b>5</b>
<b>2. Test Administration.....</b>	<b>6</b>
Training of School Systems .....	6
Ancillary Materials.....	7
Return Material Forms and Guidelines.....	17
Security Checklists.....	17
Interpretive Guides .....	17
Time .....	17
Online Forms Administration, Grades 3–8 .....	17
Paper-Based Forms Administration, Grades 3 and 4 .....	18
Accessibility and Accommodations .....	18
Testing Windows .....	19
Test Security Procedures.....	19
Data Forensic Analyses.....	20
<b>3. Scoring Activities .....</b>	<b>22</b>
Constructed-Response and Extended-Response Scoring.....	24
<b>4. Data Analysis .....</b>	<b>33</b>
Classical Item Statistics.....	33
Differential Item Functioning.....	33

Pre-Equating for Intact Forms.....	38
Unidimensionality and Principal Component Analysis .....	39
Scaling .....	39
<b>5. Reliability and Validity.....</b>	<b>42</b>
Internal Consistency Reliability Estimation.....	42
Student Classification Accuracy and Consistency .....	43
Validity .....	45
<b>6. Statistical Summaries .....</b>	<b>47</b>
<b>References.....</b>	<b>54</b>
<b>Appendix A: Test Summary.....</b>	<b>57</b>
<b>Appendix B: Item Analysis Summary Report .....</b>	<b>66</b>
<b>Appendix C: Dimensionality.....</b>	<b>75</b>
<b>Appendix D: Scale Distribution and Statistical Report.....</b>	<b>80</b>
<b>Appendix E: Reliability and Classification Accuracy .....</b>	<b>93</b>

# 1. Introduction

The Louisiana Department of Education (LDOE) has a long and distinguished history in the development and administration of assessments that support its state accountability system and are aligned to its state content standards. Per state law, the LDOE is to administer statewide Social Studies assessments in grades 3–8 and in high school. Fulfilling the directive of the Louisiana State Board of Elementary and Secondary Education (BESE), the LDOE must deliver high-quality, Louisiana-specific standards-based assessments. Further, the LDOE and the BESE are committed to the development of rigorous assessments as one component of their comprehensive plan—Louisiana Believes—designed to ensure that every Louisiana student is on track to be successful in postsecondary education and the workforce.

The purpose of this technical addendum is to describe the processes for the spring 2021 administration of LEAP 2025 Social Studies. This report outlines the testing administrations, scoring activities, and psychometric analyses.

## 2. Test Administration

This chapter describes processes and activities implemented and information disseminated to help ensure standardized test administration procedures and, thus, uniform test administration conditions for students. According to the American Educational Research Association (AERA), American Psychological Association (APA), and National Council on Measurement in Education (NCME) (2014) *Standards for Educational and Psychological Testing* (hereafter the *Standards*), “The usefulness and interpretability of test scores require that a test be administered and scored according to the developer’s instructions” (111). This chapter examines how test administration procedures implemented for the Louisiana Education Assessment Program 2025 (LEAP 2025) strengthen and support the intended score interpretations and reduce construct-irrelevant variance that could threaten the validity of score interpretations.

### Training of School Systems

To ensure that the LEAP 2025 assessments are administered and scored in accordance with the department’s policies, the LDOE takes a primary role in communicating with and training school system personnel. The LDOE provides train-the-trainer opportunities for the district test coordinators, who in turn convey test administration training to schools within their school systems. The LDOE conducts quality-assurance visits during testing to ensure adherence to the standardized administration of the tests.

The district test coordinators are responsible for the schools within their systems. They disseminate information to each school, offer assistance with test administration, and serve as liaisons between the LDOE and their school systems. The LDOE also provides assistance with and interpretation of assessment data and test results.

## Ancillary Materials

Ancillary materials for LEAP 2025 test administration contribute to the body of evidence of the validity of score interpretation. This section examines how the test materials address the *Standards* related to test administration procedures.

For the spring test administration, Data Recognition Corporation (DRC) produces two administration manuals: the *LEAP 2025 Grades 3–4 Paper-Based Test Administration Manual (TAM)* and the *LEAP 2025 Grades 3–8 Computer-Based Test Administration Manual (TAM)*. The TAMs provide detailed instructions for administering the LEAP assessments. The manuals include information on test security, test administrator responsibilities, test preparation, administration of tests (computer-based or paper-based), and post-test procedures.

### Table of Contents for LEAP 2025 Paper-Based Test Administration Manual (TAM)

- Spring 2021 Notes and Reminders
- Test Administrator Pre-Administration Oath of Security and Confidentiality Statement
- Test Administrator Post-Administration Oath of Security and Confidentiality Statement
- Overview
- Test Security
  - Secure Test Materials
  - Testing Irregularities and Security Breaches
  - Testing Environment
  - Violations of Test Security
  - Answer Change Analysis
  - Voiding Student Tests
- Test Administrator Responsibilities
- Test Administration Checklists
  - Before Testing
  - During Testing
  - After Testing (Daily)

- After Testing (Last Day)
- Test Administrators' Frequently Asked Questions
- Test Materials
  - Receipt of Test Materials
- Testing Guidelines
  - Testing Eligibility
  - Test Schedule
  - Extended Time for Testing
- Testing Times for Grades 3 and 4
  - Makeup Testing
  - Testing Conditions
- Special Populations and Accommodations
  - IDEA Special Education Students
  - Students with One or More Disabilities According to Section 504
  - Gifted and Talented Special Education Students
  - Test Accommodations for Special Education and Section 504 Students
  - Special Considerations for Deaf and Hard-of-Hearing Students
  - English Learners (ELs)
- Hand-Coded Consumable Test Booklets
- Students Absent from Testing
- Consumable Test Booklet Coding
  - Coding the Demographic Section
- Sample Grade 3 English Language Arts Consumable Test Booklet
- General Instructions for LEAP 2025
  - Student Marking/Erasing on Consumable Test Booklet
  - Reading Directions to Students
  - Special Instructions
- Directions for Administering LEAP 2025 Tests
- Post-Test Procedures



- Test Administrator Oath of Security and Confidentiality Statement
- Used and Unused Consumable Test Booklets (Defined)
- Transferring Student Responses
- Returning Test Materials to the School Test Coordinator
- Index

## Table of Contents for LEAP 2025 Computer-Based Test Administration Manual (TAM)

- Spring 2021 Notes and Reminders
- Test Administrator Pre-Administration Oath of Security and Confidentiality Statement (CBT)
- Test Administrator Post-Administration Oath of Security and Confidentiality Statement (CBT)
- Overview
- Test Security
  - Secure Test Materials
  - Testing Irregularities and Security Breaches
  - Testing Environment
  - Violations of Test Security
  - Voiding Student Tests
- Test Administrator Responsibilities
  - Software Tools and Features for Test Administrators
- Test Administration Checklists
  - Before Testing
  - During Testing
  - After Testing (Daily)
  - After Testing (Last Day)
- Test Administrators' Frequently Asked Questions
- Testing Guidelines
  - Testing Eligibility

- Testing Schedule
  - Extended Time for Testing
- Testing Times for Grades 3 through 8
  - Makeup Testing
  - Testing Conditions
- Online Tools Training
- Student Tutorials
  - Student Tutorials
- Directions for Administering the Grades 3–8 LEAP 2025 Tests
- Special Populations and Accommodations
  - IDEA Special Education Students
  - Students with One or More Disabilities According to Section 504
  - Gifted and Talented Special Education Students
  - Test Accommodations for Special Education and Section 504 Students
  - Special Considerations for Deaf and Hard-of-Hearing Students
  - English Learners (ELs)
- Test Materials
  - Receipt of Test Materials
- General Instructions
  - Reading Directions to Students
- LEAP 2025: Grades 3–8 English Language Arts (All Sessions)
- LEAP 2025: Grades 3–8 Mathematics (All Sessions)
- LEAP 2025: Grades 3–8 Science (Sessions 1–3)
- LEAP 2025: Grades 3–8 Social Studies (Grades 3–4 Sessions 1–2, Grades 5–8 Sessions 1–3)
- Post-Test Procedures
  - Test Administrator Post-Administration Oath of Security and Confidentiality Statement
  - Returning Test Materials to the School Test Coordinator

- Index

DRC also produces Test Coordinators Manuals for paper-based test administration and for computer-based test administration. The TCMs provide detailed instructions for district and school test coordinators' responsibilities for distributing, collecting, and returning test materials to DRC for scoring.

#### Table of Contents for LEAP 2025 Paper-Based Testing Test Coordinators Manual (TCM)

- Key Dates
- Spring 2021 Alerts
- Pre-Administration Oath of Security and Confidentiality Statement
- Post-Administration Oath of Security and Confidentiality Statement
- General Information
- Test Security
  - Key Definitions
  - Violations of Test Security
  - Answer Change Analysis
  - Voiding Student Tests
- Testing Guidelines
  - Testing Eligibility
  - Testing Conditions
  - Test Schedule
  - Extended Time for Testing
  - Extended Breaks
  - Makeup Testing
  - Test Administration Resources
- Testing Times for Grades 3 and 4
- District Test Coordinator
  - Conduct Training Session
  - Receive Test Materials

- Large-Print and Braille Test Materials and Communication Assistance Scripts (CAS)
- Accommodated Materials
- Verify and Distribute Test Materials to School Test Coordinators
- Request Additional Test Materials and Bar-Code Labels
- Collect Materials from Schools After Testing
- Used and Unused Consumable Test Booklets (Defined)
- Unscorable Documents and Unscorable Document Labels
- Directions for Returning Test Materials to DRC in May
  - Pickup 1: ELA and Mathematics Scorable Test Materials
  - Pickup 2: Science and Social Studies Scorable Test Materials
  - Pickup 3: Nonscorable Test Materials
  - Final Checklist for Returning Test Materials to DRC
- School Test Coordinator
  - Receive and Verify Test Materials
  - Conduct Test Administration and Security Training Session
  - Supervise Application of Bar-Code Labels and Coding of Consumable Test Booklets
  - Soiled, Damaged, and Other Unscorable Consumable Test Booklets
  - Verify and Distribute Materials to Test Administrators
  - Supervise Test Administration
  - Collect Test Materials
  - Used and Unused Consumable Test Booklets (Defined)
  - Coding Responsibilities of Principals—Before Testing
  - Coding Responsibilities of Principals—Before or After Testing
  - Coding Responsibilities of Principals—After Testing
- Directions for Returning Test Materials to DTC
  - Pickup 1: ELA and Mathematics Scorable Test Materials
  - Pickup 2: Science and Social Studies Scorable Test Materials

- Pickup 3: Nonscorable Test Materials
- Final Checklist for Returning Test Materials to DTC
- Void Notification—Spring 2021
- Index

## Table of Contents for LEAP 2025 Computer-Based Testing Test Coordinators Manual (TCM)

- Key Dates Spring 2021
- Resources Available in DRC INSIGHT Portal (eDIRECT) Spring 2021
- Spring 2021 Alerts
- Pre-Administration Oath of Security and Confidentiality Statement
- Post-Administration Oath of Security and Confidentiality Statement
- General Information
  - DRC INSIGHT Portal (eDIRECT) and INSIGHT
- Test Security
  - Key Definitions
  - Violations of Test Security
- Testing Guidelines
  - Testing Eligibility
  - Testing Conditions
  - Testing Schedule
  - Extended Time for Testing
  - Extended Breaks
  - Accommodations
  - Makeup Testing
  - Test Administration Resources
- Testing Times for Grades 3 through 8
- Roles and Responsibilities
  - District Test Coordinator
  - School Test Coordinator

- Technology Coordinator
- Managing Test Tickets
  - Student Transfers
  - Locked Test Tickets
  - Technical Issues
  - Invalidating Test Tickets
- Resources for Online Testing
  - Test Administration Manuals
  - *DRC INSIGHT Portal (eDIRECT) User Guides*
  - *LEAP 2025 Accommodations and Accessibility Features User Guide*
  - *INSIGHT Technology User Guide*
  - Online Tools Training (OTT)
  - Student Tutorials
- Void Notification—Spring 2021

LDOE assessment staff review, provide feedback, and give final approval for these manuals. The manuals are inclusive of grades 3–8 English Language Arts (ELA), Mathematics, Social Studies, and Science.

The *Standards* contain multiple references relevant to test administration. Information in the LEAP 2025 test administration manuals addresses these in the following manner.

Directions for test administration found in the manual address Standard 4.15, which states:

The directions for test administration should be presented with sufficient clarity so that it is possible for others to replicate the administration conditions under which the data on reliability, validity, and (where appropriate) norms were obtained. Allowable variations in administration procedures should be clearly described. The process for reviewing requests for additional testing variations should also be documented (90).

The LEAP 2025 TAMs provide instructions for activities that happen before, during, and after testing with sufficient detail and clarity to support reliable test administrations by qualified test administrators. To ensure uniform administration conditions throughout the state, instructions in the TAMs describe the following: general rules of paper and online testing; assessment duration, timing, and sequencing information; and the materials required for testing.

Furthermore, the standardized procedures addressed in the TAMs need to be followed, as the *Standards* state in Standard 6.1: “Test administrators should follow carefully the standardized procedures for administration and scoring specified by the test developer and any instructions from the test user” (114). To ensure the usefulness and interpretability of test scores and to minimize sources of construct-irrelevant variance, it was essential that the LEAP 2025 was administered according to the prescribed test administration manual. It should be noted that adhering to the test schedule is also a critical component. The TCMs included instructions for scheduling the test within the state testing window. The TAMs and TCMs also contained the schedule for timing each test session.

**Standard 6.3.** Changes or disruptions to standardized test administration procedures or scoring should be documented and reported to the test user (115).

Department staff release annual test security reports that describe a wide range of improper activities that may occur during testing, including the following: copying and reviewing test questions with students; cueing students during testing, verbally or with written materials on the classroom walls; cueing students nonverbally, such as by tapping or nodding the head; allowing students to correct or complete answers after tests have been submitted; splitting sessions into two parts; ignoring the standardized directions in the online assessment; paraphrasing parts of the test to students; changing or completing (or allowing other school personnel to change or complete) student answers; allowing accommodations that are not written in the Individualized Education Program (IEP), Individual Accommodation Plan/504 Plan (IAP), or English Learner Plan (EL plan); allowing accommodations for students who do not have an IEP, IAP, EL plan; or defining terms on the test.

**Standard 6.4.** The testing environment should furnish reasonable comfort with minimal distractions to avoid construct-irrelevant variance (116).

The TAMs outline the steps that teachers should take to prepare the classroom testing environment for administering the LEAP 2025 test. These include the following:

- Determine the layout of the classroom environment.
- Plan seating arrangements. Allow enough space between students to prevent the sharing of answers.
- Eliminate distractions such as bells or telephones.
- Use a Do Not Disturb sign on the door of the testing room.
- Make sure classroom maps, charts, and any other materials that relate to the content and processes of the test are covered or removed or are out of the students' view.

**Standard 6.6.** Reasonable efforts should be made to ensure the integrity of test scores by eliminating opportunities for test takers to attain scores by fraudulent or deceptive means (116).

The TAMs present instructions for post-test activities to ensure that online tests are submitted and printed test materials are handled properly to maintain the integrity of student information and test scores. Detailed instructions guide test examiners in submitting all online test records. For students who were administered a large-print or braille version of the LEAP 2025 assessment, examiners are instructed to transcribe students' responses from the large-print or braille test book into the online testing system (INSIGHT) exactly as they responded in the large-print or braille test book.

**Standard 6.7.** Test users have the responsibility of protecting the security of test materials at all times (117).

Throughout the manuals, test coordinators and examiners are reminded of test security requirements and procedures to maintain test security. Specific actions that are direct violations of test security are so noted. Detailed information about test security procedures is presented under "Test Security" in the manuals.



## Return Material Forms and Guidelines

The paper-based TCM instructs test coordinators regarding procedures for organizing and packing materials and returning them to DRC for secure inventory purposes. LDOE assessment staff have opportunities to review, provide feedback, and give final approval of the guidelines. The purpose of the instructions is to ensure that secure test materials are properly accounted for and organized appropriately for return shipment.

## Security Checklists

As soon as printed test materials are received by a school system, the district test coordinator ensures that the first and last security barcodes on the tests match the packing list they received. The district test coordinator then packages the tests to be sent to schools. Upon returning test books to DRC, school and district test coordinators are required to complete and submit an accountability form that details the number of test books or printed test forms returned. This form also requires that systems/schools document nonstandard situations, including lost, damaged, destroyed, extra, or missing test books.

## Interpretive Guides

Essential to making valid interpretations of test scores is an understanding of what the test scores mean and how to interpret score reports. The Interpretive Guide is written for Louisiana teachers and administrators who receive the LEAP 2025 score reports.

<https://www.louisianabelieves.com/resources/library/assessment-guidance>

## Time

Each session of each content area test is timed to provide sufficient time for students to attempt all items. Only students with an extended time accommodation were permitted to exceed the established time limits of any given session. The manuals provide examiners with timing guidelines for the assessments.

## Online Forms Administration, Grades 3–8

The online forms are administered via DRC's INSIGHT online assessment system. School system and school personnel set up test sessions via DRC's INSIGHT portal (eDIRECT) and

print test tickets. Students enter their ticket information to access the test in INSIGHT. In addition, students have access to the Online Tools Training (OTT) before the testing window, which allows them to practice using tools and features within INSIGHT. Tutorials with online video clips that demonstrate features of the system are also available to students before testing.

## **Paper-Based Forms Administration, Grades 3 and 4**

Schools with testers at grades 3 and 4 had the option to participate in either paper-based or computer-based testing for the spring 2021 test. DRC prints and ships paper materials to the sites that choose paper-based testing. These materials are returned to DRC after testing for processing and scoring with the online tests.

## **Accessibility and Accommodations**

Accessibility features and accommodations include Access for All, Accessibility Features, and Accommodations.

- Access for All features are available to all students taking an assessment.
- Accessibility Features are available to students when deemed appropriate by a team of educators.
- Accommodations must appear in a student's IEP/IAP/EL plan.

Accommodations may be used with students who qualify under the Individuals with Disabilities Education Act (IDEA) and have an IEP or Section 504 of the Americans with Disabilities Act and have an IAP, or who are identified as English Learners (ELs) and have an EL plan.

Accommodations must be specified in the qualifying student's individual plan and must be consistent with accommodations used during daily classroom instruction and testing. The use of any accommodation must be indicated on the student information sheet at the time of test administration. AERA, APA, and NCME Standard 6.2 states:

When formal procedures have been established for requesting and receiving accommodations, test takers should be informed of these procedures in advance of testing (115).

In compliance with this standard, the TAM contains the list of Universal Tools, Designated Supports, and Accommodations permissible for the LEAP assessments. The following accommodations were provided by DRC for this administration:

- Braille
- Text-to-Speech
- Directions in Native Language

The following additional access and accommodation features were also available:

- Answers Recorded
- Extended Time
- Transferred Answers
- Individual/Small Group Administration
- Tests Read Aloud
- English/Native Language Word-to-Word Dictionary
- Directions Read Aloud/Clarified in Native Language
- Text-to-Speech for online testers
- Human Read Aloud
- Directions in Native Language

For more details about these accommodations, please refer to the [\*LEAP 2025 Accessibility and Accommodations Manual\*](#).

## **Testing Windows**

The computer-based testing window was available from April 26 through May 26, 2021. Paper-based testing occurred from April 28 through May 4, 2021.

## **Test Security Procedures**

Maintaining the security of all test materials is crucial to preventing the possibility of random or systematic errors, such as unauthorized exposure of test items that would affect the valid interpretation of test scores. Several test security measures are implemented for the LEAP 2025 assessments. Test security procedures are discussed throughout the TCMs and TAMs.

Test coordinators and administrators are instructed to keep all test materials in locked storage, except during actual test administration, and access to secure materials must be restricted to authorized individuals only (e.g., test administrators and the school test coordinator). During the testing sessions, test administrators are directly responsible for the security of the LEAP 2025 assessment and must account for all test materials and supervise the test administrators at all times.

## Data Forensic Analyses

Due to the importance of the LEAP 2025 assessment, it is prudent to ensure that the results from the assessments are based on effective instruction and true student achievement. To help ensure that scores are related to actual learning and that results are valid, data forensic analyses take place to assist in separating meaningful gains from spurious gains. It is important to note that although the results of the analyses may be used to identify potential problems within a school, the identification of a problem is not an accusation of misconduct.

Multiple methods are incorporated into the forensic analysis. The following methods are applied:

- Response Change Analysis
- Score Fluctuation Analysis
- Item Exposure Monitoring
- Web Monitoring
- Plagiarism Detection
- Alerts for Disturbing Content

**Response Change Analysis.** Students make changes to answer choices when taking the LEAP 2025 assessments, and this behavior is expected. Unfortunately, changes to student answers are sometimes influenced by school personnel who want to improve performance. Therefore, the response change analysis is conducted to identify school- and test administrator-level response change patterns that are statistically improbable when compared to the expected pattern at the state level.

**Score Fluctuation Analysis.** It is anticipated that performance on the LEAP 2025 assessments will improve over time for legitimate reasons such as changes in the

curriculum and improvement in instruction. However, large and unexpected score changes may be a sign of testing impropriety. The LDOE applied an approach where the state's level of change in performance from one year to the next is compared to schools' and test administrators' change in performance during the same time frame. Schools and test administrators are identified when the level of change is statistically unexpected.

**Item Exposure Monitoring.** Due to the re-use of the 2019 operational forms for the spring of 2021 administration, item performance was examined to ensure that item content had not been exposed. Frequently during the testing window, every item's moving  $p$ -value and point-biserial averages were produced. If an item's moving average  $p$ -value was larger than expected compared to the previous administrations, the item was flagged. Additionally, plots were produced for a visual inspection of the day-to-day patterns of item performance.

**Web Monitoring.** The content of the LEAP 2025 assessments should not appear outside the boundaries of the forms administered. To protect Louisiana test content, the internet is monitored for postings that contain, or appear to contain, potentially exposed and/or copied test content. When test content is verified, steps are taken to quickly remove the infringing content.

**Plagiarism Detection.** The LDOE monitors for two different plagiarism situations: copying from student to student and copying from an outside source, such as Wikipedia or another internet source. Instances of plagiarism are identified by human scorers and artificial intelligence. Alerts are set to identify responses that may indicate the possibility of teacher interference or plagiarism. Alerted responses are given additional review so that the appropriate response can be taken.

**Alerts for Disturbing Content.** Scorers for the LEAP 2025 assessments also have the ability to apply an alert flag to student responses that may indicate disturbing content (e.g., possible physical or emotional abuse, suicidal ideation, threats of harm to themselves or others). All alerted responses are automatically routed to the scoring director, who reviews and forwards appropriate responses to senior project staff for review. If it is concluded that a response warrants an alert, project management will contact the LDOE to take the necessary action. At no point during this process do scorers or staff have access to demographic information for any students participating in the assessment.

## 3. Scoring Activities

**Directory of Test Specification (DOTS) Process.** DRC creates a DOTS file, based on the approved test selection. The DOTS is a document containing information about each item on a test form, such as item identifier, item sequence, answer key, score points, subtype, session, alignment, and prior use of item. WestEd reviews and confirms the contents of the DOTS file as part of test review rounds. The DOTS file is then provided to the LDOE for review and final approval. Once approved, the information contained in the DOTS is used in scoring the test and in reporting.

**Selected-Response (SR) Item Keycheck.** SR items for Social Studies include multiple-choice (MC) and multiple-select (MS) questions. Pearson calculates MC and MS item statistics and flags items if item statistics fall outside expected ranges. For example, items are flagged if few students select the correct response ( $p$ -value less than 0.15), if the item does not discriminate well between students of lower and higher ability (point-biserial correlation less than 0.20), or if many students (more than 40%) select a certain incorrect response. Lists of flagged MC and MS items, with the reasons for flagging, are provided to LDOE and WestEd content staff for key verification. The staff reviews the list of flagged MC and MS items to confirm that the answer keys are accurate. Scoring of MC and MS items is also evaluated at data review.

**Scoring of Technology-Enhanced (TE) Items.** All TE items are processed through DRC's autoscoring engine and scored according to the assigned scoring rules established during content creation by WestEd in conjunction with the LDOE. DRC ensures that all rubrics and scoring rules are verified for accuracy before scoring any TE items. DRC has an established adjudication process for TE items to verify that correct answers are identified. DRC's TE scoring process includes the following procedures:

- A scoring rubric is created for each TE item. The rubric describes the one and only correct answer for dichotomously scored items (i.e., items scored as either right or wrong). If partial credit is possible, the rubric describes in detail the type of response that could receive credit for each score point.
- The information from each scoring rubric is entered into the scoring system within the item banking system so that the truth resides in one place along with the item image and other metadata. This scoring information designates specific

information that varies by item type. For example, for a drag-and-drop item, the information includes which objects are to be placed in each drop region to receive credit.

- The information is then verified by another autoscoring expert.
- After testing starts, reports are generated that show every response, how many students gave that response, and the score the scoring system provided for that response.
- The scoring is then checked against the scoring rubric using two levels of verification.
- If any discrepancies are found, the scoring information is modified and verified again. The scoring process is then rerun. This checking and modification process continues until no other issues are found.
- As a final check, a final report is generated that shows all student responses, their frequencies, and their received scores.

In the case of braille and accommodated print test forms, student responses to TE items are transcribed into the online system by a test administrator.

**Adjudication.** TE items and other eligible items identified in the test map are automatically scored as tests are processed. TE items are scored according to scoring rules in the DOTS, which includes scoring information for all item types.

The adjudication process focuses on detecting possible errors in scoring TE and MS items. DRC provides a report listing the frequency distributions of TE item responses and MS items. Members of the LDOE and WestEd content staff examine the TE and MS response distributions and the auto-frequency reports to evaluate whether the items are scored appropriately. In the event that scoring issues are identified, WestEd content staff and the LDOE recommend changes to the scoring algorithm. Any changes to the scoring algorithm are based on the LDOE's decisions. DRC, in turn, applies the approved scoring changes to any affected items.

## Constructed-Response and Extended-Response Scoring

Constructed-response items are scored by human raters trained by DRC. Extended-response items are scored by Project Essay Grade (PEG), an Artificial Intelligence (AI) scoring engine. Ten percent of the responses are scored twice to monitor and maintain inter-rater reliability. Scoring supervisors also conduct read-behinds and review all nonscores and alerts. Handscoring processing rules are detailed in the LEAP 2025 Spring 2021 Handscoring/AI Documentation document.

**Selection of Scoring Evaluators.** Standard 4.20 states the following:

The process for selecting, training, qualifying, and monitoring scorers should be specified by the test developer. The training materials, such as the scoring rubrics and examples of test takers' responses that illustrate the levels on the rubric score scale, and the procedures for training scorers should result in a degree of accuracy and agreement among scorers that allows the scores to be interpreted as originally intended by the test developer. Specifications should also describe processes for assessing scorer consistency and potential drift over time in raters' scoring (92).

The following sections explain how scorers are selected and trained for the LEAP 2025 handscoring process and describe how the scorers are monitored throughout the handscoring process.

**Recruitment and Interview Process.** DRC strives to develop a highly qualified, experienced core of evaluators to appropriately maintain the integrity of all projects. All readers hired by DRC to score 2020–2021 LEAP 2025 test responses had at least a four-year college degree.

DRC has a human resources director dedicated solely to recruiting and retaining the handscoring staff. Applications for reader positions are screened by the handscoring project manager, the human resources director, or recruiting staff to create a large pool of potential readers. In the screening process, preference is given to candidates with previous experience scoring large-scale assessments and with degrees emphasizing the appropriate content areas. At the personal interview, reader candidates are asked to demonstrate their proficiency in writing by responding to a DRC writing topic and their proficiency in mathematics by solving word problems with correct work shown. These



steps result in a highly qualified and diverse workforce. DRC personnel files for readers and team leaders include evaluations for each project completed. DRC uses these evaluations to place individuals on projects that best fit their professional backgrounds, their college degrees, and their performances on similar projects at DRC. Once placed, all readers go through rigorous training and qualifying procedures specific to the project on which they are placed. Any scorer who does not complete this training and does not demonstrate the ability to apply the scoring criteria by qualifying at the end of the process is not allowed to score live student responses.

**Security.** Whether training and scoring are conducted within a DRC facility or done remotely, security is essential to the handscoring process. When users log into DRC's secure, web-based scoring application, ScoreBoard, they are required to read and accept the security policy before they are allowed to access any project. For each project, scorers are also required to read and sign non-disclosure agreements, and during training emphasis is always given to what security means, the importance of maintaining security, and how this is accomplished.

Readers only have access to student responses they are qualified to score. Each scorer is assigned a unique username and password to access DRC's imaging system and must qualify before viewing any live student responses. DRC maintains full control of who may access the system and which item each scorer may score. No demographic data is available to scorers at any time.

Each DRC scoring center is a secure facility. Access to scoring centers is limited to badge-wearing staff and to visitors accompanied by authorized staff. All readers are made aware that no scoring materials may leave the scoring center. To prevent the unauthorized duplication of secure materials, cell phone/camera use within the scoring rooms is strictly forbidden. Readers only have access to student responses they are qualified to score.

In a remote environment, security reminders are given on a daily basis. Similar to the work that occurs within DRC scoring sites, in a remote environment, education about security expectations is the best way to maintain security of any project materials. DRC requires scorers working remotely to work in a private environment away from other people (including family members). Restrictions are in place that define the hours during the day scorers are able to log into the system. If any type of security breach were to

occur, immediate action would be taken to secure materials, and the employee would be terminated. DRC has the same policy within the scoring centers.

**Handscoring Training Process.** Standard 6.9 specifies:

Those responsible for test scoring should establish and document quality control processes and criteria. Adequate training should be provided. The quality of scoring should be monitored and documented. Any systematic source of scoring errors should be documented and corrected (118).

**Training Material Development.** DRC scoring supervisors train scorers using LDOE-approved training materials. These materials are developed by DRC and LDOE staff from a selection scored by Louisiana educators at rangefinding and include the following:

- Prompts and associated sources
- Rubrics
- Anchor sets
- Practice sets
- Qualifying sets

**Training and Qualifying Procedures.** Handscoring involves training and qualifying team leaders and evaluators, monitoring scoring accuracy and production, and ensuring security of both the test materials and the scoring facilities. The LDOE reviews training materials and oversees the training process.

**Qualifying Standards.** Scorers demonstrate their ability to apply the scoring criteria by qualifying (i.e., scoring with acceptable agreement with true scores on qualifying sets). After each qualifying set is scored, the DRC scoring director responsible for training leads the scorers in a discussion of the set.

Any scorer who does not qualify by the end of the qualifying process for an item is not allowed to score live student responses.

**Monitoring the Scoring Process.** Standard 6.8 states:

Those responsible for test scoring should establish scoring protocols. Test scoring that involves human judgment should include rubrics, procedures, and criteria for

scoring. When scoring of complex responses is done by computer, the accuracy of the algorithm and processes should be documented (118).

The following section explains the monitoring procedures that DRC uses to ensure that handscoring evaluators follow established scoring criteria while items are being scored. Detailed scoring rubrics, which specify the criteria for scoring, are available for all constructed- and extended-response items.

**Reader Monitoring Procedures.** Throughout the handscoring process, DRC project managers, scoring directors, and team leaders review the statistics that are generated daily. DRC uses one team leader for every 10 to 12 readers. If scoring concerns are apparent among individual scorers or if a scorer needs clarification on the scoring rules, team leaders address those issues on an individual basis. DRC supervisors typically monitor one out of five of the scorer's readings, making adjustments to that ratio as needed. If a supervisor disagrees with a reader's scores during monitoring, the supervisor provides retraining in the form of direct feedback to the reader, using rubric language and applicable training responses.

**Validity Sets and Inter-Rater Reliability.** In addition to the feedback that supervisors provide to readers during regular read-behinds and the continuous monitoring of inter-rater reliability and score point distributions, DRC also conducts validity scoring using LDOE-approved validity responses identified by DRC scoring supervisors during live scoring for newly operational items. Validity responses are inserted among the live student responses.

The validity responses are added to DRC's image handscoring system prior to the beginning of scoring. Validity reports compare readers' scores to predetermined scores and are used to help detect potential room drift as well as individual scorer drift. This data is used to make decisions regarding the retraining and/or release of scorers, as well as the rescoring of responses.

Approximately 10% of all live student responses are scored by a second reader to establish inter-rater reliability statistics for all constructed- and extended-response items. This procedure is called a "double-blind read" because the second reader does not know the first reader's score. DRC monitors inter-rater reliability based on the responses that are scored by two readers. If a scorer falls below the expected rate of agreement, the

team leader or scoring director retrain the scorer. If a scorer fails to improve after retraining and feedback, DRC removes the scorer from the project. In this situation, DRC removes all scores assigned by the scorer in question. The responses are then reassigned and rescored.

To monitor inter-rater reliability, DRC produces scoring summary reports daily. DRC's scoring summary reports display exact, adjacent, and nonadjacent agreement rates for each reader. These rates are calculated based on responses that are scored by two readers, and their definitions are included below.

- Percentage Exact (%EX)—total number of responses by reader where scores are the same, divided by the number of responses that were scored twice
- Percentage Adjacent (%AD)—total number of responses by reader where scores are one point apart, divided by the number of responses that were scored twice
- Percentage Nonadjacent (%NA)—total number of responses by reader where scores are more than one point apart, divided by the number of responses that were scored twice

Each reader is required to maintain a level of exact agreement on validity responses and on inter-rater reliability. Additionally, readers are required to maintain a low rate of nonadjacent agreement.

**Calibration Sets.** DRC pulls calibration responses for items. DRC uses these sets to perform calibration across the entire scorer population for an item if trends are detected (e.g., low agreement between certain score points if a certain type of response is missing from initial training). These calibrations are designed to help refocus scorers on how to properly use the scoring guidelines. They are selected to help illustrate particular points and familiarize scorers with the types of responses commonly seen during operational scoring. After readers score a calibration set, the scoring director reviews it from the front of the room, using rubric language and scoring concepts exemplified by the anchor responses to explain the reasoning behind each response's score.

**Reports and Reader Feedback.** Reader performance and intervention information are recorded in reader feedback logs. These logs track information about actions taken with individual readers to ensure scoring consistency in regard to reliability, score point distribution, and validity performance. In addition to the reader feedback logs, DRC

provides the LDOE with handscoring quality control reports for review throughout the scoring window.

**Inter-Rater Reliability.** A minimum of 10% of the responses for constructed-and extended-response items are scored independently by a second reader. This is the case regardless of whether the first reader is a human rater or AI. The statistics for inter-rater reliability are calculated for all items at all grades. To determine the reliability of scoring, the percentage of perfect agreement and adjacent agreement between the first and second scores is examined.

Tables 3.1–3.4 provide the inter-rater reliability and score point distributions by grade level for the constructed-response and extended-response items administered in the spring 2021 forms.

Table 3.1

*Inter-Rater Reliability for Operational Constructed-Response Items*

Grade	Item	Inter-Rater Reliability*			
		2x	Exact Agreement (%)	Adjacent Agreement (%)	Nonadjacent (%)
3	Item 1	≥17,320	93	7	0
	Item 2	≥14,300	94	5	0
4	Item 1	≥13,270	86	14	1
	Item 2	≥13,410	93	6	0
5	Item 1	≥36,450	97	2	0
	Item 2	≥21,190	92	8	0
6	Item 1	≥16,230	96	4	0
	Item 2	≥18,420	88	12	0
7	Item 1	≥15,680	88	12	0
	Item 2	≥16,400	90	10	0
8	Item 1	≥16,890	87	13	0
	Item 2	≥13,880	83	17	0

\* The percent may not add up to 100% due to rounding.

Table 3.2

*Score Point Distributions for Operational Constructed-Response Items*

Grade	Item	Score Point Distribution*					
		Total	"0" Rating (%)	"1" Rating (%)	"2" Rating (%)	Blank (%)	Nonscore Codes (%)**
3	Item 1	≥61,290	46	24	7	9	15
	Item 2	≥59,890	58	11	14	9	8
4	Item 1	≥58,760	37	30	19	7	6
	Item 2	≥58,980	71	12	2	7	7
5	Item 1	≥67,570	39	13	7	0	40
	Item 2	≥59,910	32	38	10	0	20
6	Item 1	≥59,050	51	36	4	0	9
	Item 2	≥60,040	34	44	9	0	13
7	Item 1	≥59,230	49	32	8	0	10
	Item 2	≥59,600	38	28	21	0	13
8	Item 1	≥59,450	35	41	12	0	12
	Item 2	≥57,950	35	46	12	0	7

\* The percent may not add up to 100% due to rounding.

\*\* Nonscore codes include Foreign language (F), Insufficient (I), Don't Understand (N), Refusal (R), Off Topic (T), and Unintelligible (U). Responses that cannot be assigned a score based on the rubric are assigned a nonscore code and count as zero points toward student scores.

Table 3.3

*Inter-Rater Reliability for Operational Extended-Response Items*

Grade	2x	Inter-Rater Reliability*			
		Dimension	Exact Agreement (%)	Adjacent Agreement (%)	Nonadjacent (%)
5	≥52,660	Content	95	5	0
		Claim	95	5	0
6	≥40,620	Content	94	6	0
		Claim	94	6	0
7	≥55,660	Content	96	4	0
		Claim	96	4	0
8	≥52,130	Content	90	10	0
		Claim	88	12	0

\* The percent may not add up to 100% due to rounding.

Table 3.4

*Score Point Distributions for Operational Extended-Response Items*

Grade	Score Point Distribution*								
	Total	Dimension	"0" (%)	"1" (%)	"2" (%)	"3" (%)	"4" (%)	Blank (%)	Nonscore Codes (%)**
5	≥75,700	Content	35	30	12	3	0	0	18
		Claim	44	24	10	2	0	0	18
6	≥71,180	Content	40	32	9	2	0	0	17
		Claim	51	24	6	1	0	0	17
7	≥79,130	Content	38	32	10	3	1	0	16
		Claim	45	25	9	3	1	0	16
8	≥77,030	Content	24	31	24	7	2	0	11
		Claim	23	31	25	8	3	0	11

\* The percent may not add up to 100% due to rounding.

\*\* Nonscore codes include Foreign language (F), Insufficient (I), Don't Understand (N), Refusal (R), Off Topic (T), and Unintelligible (U). Responses that cannot be assigned a score based on the rubric are assigned a nonscore code and count as zero points toward student scores.



## 4. Data Analysis

### Classical Item Statistics

This section describes the classical item analysis for data obtained from the operational LEAP 2025 Social Studies tests. The classical analysis includes statistical analysis based on the following types of items: multiple-choice/multiple-select items, rule-based machine-scored items such as technology-enhanced items, and handscored items such as constructed- and extended-response items. For each operational item, the statistical analysis produces item difficulty ( $p$ -value) and item discrimination (point-biserial).

Tables and figures that provide the information on classical item statistics for operational items for the spring 2021 test can be found in [Appendix B: Item Analysis Summary Report](#). Tables B.1.1–B.5.2 show summaries of classical item statistics. As a measure of item difficulty,  $p$  (or “the  $p$ -value”) indicates the average proportion of total points earned on an item. For example, if  $p = 0.50$  on an MC item, then half of the examinees earned a score of 1. If  $p = 0.50$  on a CR item, then examinees earned half of the possible points on average (e.g., 1 out of 2 possible points). A measure of point-biserial correlation indicates a measure of item discrimination. Items with higher item-total correlations provide better information about how well items discriminate between lower- and higher-performing students. Statistical analysis results for field-test (FT) items are stored in Pearson’s Assessment Banking and Building solutions for Interoperable assessment (ABBI) system. Placeholder (PH) items included on test forms are not part of any statistical analyses. Because the purpose of PH items is to maintain a consistent testing length and experience by occupying FT-item positions for administrations when no field testing takes place, these items do not require any statistical analysis.

### Differential Item Functioning

Differential item functioning (DIF) analyses are intended to statistically signal potential item bias. DIF is defined as a difference between similar ability groups’ (e.g., males or females that attain the same total test score) probability of getting an item correct. Because test scores can reflect many sources of variation, the test developers’ task is to create assessments that measure the intended knowledge and skills without introducing construct-irrelevant variance. When tests measure something other than what they are intended to measure, test scores may reflect those extraneous elements in addition to

what the test is purported to measure. If this occurs, these tests can be called biased (Angoff, 1993; Camilli & Shepard, 1994; Green, 1975; Zumbo, 1999). Different cultural and socioeconomic experiences are among some factors that can confound test scores intended to reflect the measured construct.

One DIF methodology applied to dichotomous items was the Mantel–Haenszel (*MH*) *DIF* statistic (Holland & Thayer, 1988; Mantel & Haenszel, 1959). The *MH* method is a frequently used method that offers efficient statistical power (Clauser & Mazor, 1998). The *MH* chi-square statistic is

$$MH_{\chi^2} = \frac{(\sum_k F_k - \sum_k E(F_k))^2}{\sum_k Var(F_k)},$$

where  $F_k$  is the sum of scores for the focal group at the  $k$ th level of the matching variable (Zwick, Donoghue, & Grima, 1993). Note that the *MH* statistic is sensitive to  $N$  such that larger sample sizes increase the value of chi-square.

In addition to the *MH* chi-square statistic, the *MH* delta statistic ( $\Delta MH$ ), first developed by the Educational Testing Service (ETS), is computed. To compute the  $\Delta MH$  *DIF*, the *MH* alpha (the odds ratio) is first calculated:

$$\alpha_{MH} = \frac{\sum_{k=1}^K N_{r1k} N_{f0k} / N_k}{\sum_{k=1}^K N_{f1k} N_{r0k} / N_k},$$

where  $N_{r1k}$  is the number of correct responses in the reference group at ability level  $k$ ,  $N_{f0k}$  is the number of incorrect responses in the focal group at ability level  $k$ ,  $N_k$  is the total number of responses,  $N_{f1k}$  is the number of correct responses in the focal group at ability level  $k$ , and  $N_{r0k}$  is the number of incorrect responses in the reference group at ability level  $k$ . The *MH* *DIF* statistic is based on a  $2 \times 2 \times M$  (2 groups  $\times$  2 item scores  $\times$   $M$  strata) frequency table, in which students in the reference (male or white) and focal (female or African American/Hispanic-Latino) groups are matched on their total raw scores.

The  $\Delta MH DIF$  is then computed as

$$\Delta MH DIF = -2.35 \ln(\alpha_{MH}).$$

Positive values of  $\Delta MH DIF$  indicate items that favor the focal group (i.e., positive DIF items are differentially easier for the focal group); negative values of  $\Delta MH DIF$  indicate items that favor the reference group (i.e., negative DIF items are differentially easier for the reference group). Ninety-five percent confidence intervals for  $\Delta MH DIF$  are used to conduct statistical tests.

The  $MH$  chi-square statistic and the  $\Delta MH DIF$  are used in combination to identify operational test items exhibiting strong, weak, or no DIF (Zieky, 1993). Table 4.1 defines the DIF categories for dichotomous items.

Table 4.1  
*DIF Categories for Dichotomous Items*

DIF Category	Criteria
A (negligible)	$\Delta MH DIF$   is not significantly different from 0.0 or is less than 1.0.
B (slight to moderate)	1.   $\Delta MH DIF$   is significantly different from 0.0 but not from 1.0, and is at least 1.0; OR 2.   $\Delta MH DIF$   is significantly different from 1.0, but is less than 1.5. Positive values are classified as "B+" and negative values as "B-."
C (moderate to large)	$\Delta MH DIF$   is significantly different from 1.0 and is at least 1.5. Positive values are classified as "C+" and negative values as "C-."

For polytomous items, the standardized mean difference ( $SMD$ ) (Dorans & Schmitt, 1991; Zwick, Thayer, & Mazzeo, 1997) and the Mantel  $\chi^2$  statistic (Mantel, 1963) are used to identify items with DIF.  $SMD$  estimates the average difference in performance between the reference group and the focal group while controlling for student ability. To calculate  $SMD$ , let  $M$  represent the matching variable (total test score). For all  $M = m$ , identify the students with raw score  $m$  and calculate the expected item score for the reference group ( $E_{rm}$ ) and the focal group ( $E_{fm}$ ).  $DIF$  is defined as  $D_m = E_{fm} - E_{rm}$ , and  $SMD$  is a weighted average of  $D_m$  using the weights  $w_m = N_{fm}$  (the number of students in the focal group with raw score  $m$ ), which gives the greatest weight at score levels most frequently attained by students in the focal group.

$$SMD = \frac{\sum_m w_m (E_{fm} - E_{rm})}{\sum_m w_m} = \frac{\sum_m w_m D_m}{\sum_m w_m}$$

*SMD* is converted to an effect-size metric by dividing it by the standard deviation of item scores for the total group. A negative *SMD* value indicates an item on which the focal group has a lower mean than the reference group, conditioned on the matching variable. On the other hand, a positive *SMD* value indicates an item on which the reference group has a lower mean than the focal group, conditioned on the matching variable.

The *MH DIF* statistic is based on a  $2 \times (T+1) \times M$  (2 groups  $\times$   $T+1$  item scores  $\times$   $M$  strata) frequency table, where students in the reference and focal groups are matched on their total raw scores ( $T$  = maximum score for the item). The Mantel  $\chi^2$  statistic is defined by the following equation:

$$\text{Mantel's } \chi^2 = \frac{(\sum_m \sum_t N_{rtm} Y_t - \sum_m \frac{N_{r+m}}{N_{++m}} \sum_t N_{+tm} Y_t)^2}{\sum_m \text{Var}(\sum_t N_{rtm} Y_t)}$$

The *p*-value associated with the Mantel  $\chi^2$  statistic and the *SMD* (on an effect-size metric) are used to determine DIF classifications. Table 4.2 defines the DIF categories for polytomous items.

Table 4.2

*DIF Categories for Polytomous Items*

DIF Category	Criteria
A (negligible)	Mantel $\chi^2$ <i>p</i> -value > 0.05 or $ SMD/SD  \leq 0.17$
B (slight to moderate)	Mantel $\chi^2$ <i>p</i> -value < 0.05 and $0.17 <  SMD/SD  < 0.25$
C (moderate to large)	Mantel $\chi^2$ <i>p</i> -value < 0.05 and $ SMD/SD  \geq 0.25$

Three DIF analyses are conducted for the operational test items only: female/male, African American/white, and Hispanic/white. That is, item score data are used to detect items on which female or male students performed unexpectedly well or unexpectedly poorly, given their performance on the full assessment. The same methods are used to detect items on which both African American/white and Hispanic-Latino/white students performed unexpectedly well or unexpectedly poorly, given their performance on the full assessment. The last two columns of Tables 4.3.1, 4.3.2 and 4.3.3 provide the number of

items flagged for DIF. Items flagged with A-DIF show negligible DIF, items flagged with B-DIF are said to exhibit slight to moderate DIF, and items with C-DIF are said to exhibit moderate to large DIF. Very few operational test items were flagged for C-DIF by either analysis.

Note that DIF flags for dichotomous items are based on the *MH* statistics while DIF flags for polytomous items are based on the combination of Mantel  $\chi^2$  *p*-value and *SMD* statistics. Tables 4.3.1, 4.3.2, and 4.3.3 summarize the operational-test DIF statistics for the operational items appearing on the spring 2021 test forms. Because the spring 2021 tests were administered during the COVID-19 pandemic, great caution should be applied when any statistical inference is drawn.

Table 4.3.1

*Summary of Female – Male DIF Flags for Operational Items: Spring 2021 Social Studies by Grade*

Grade	A	[B+],[B-]	[C+],[C-]
3	43	[0],[0]	[0],[0]
4	43	[0],[0]	[0],[0]
5	46	[0],[0]	[0],[0]
6	52	[0],[1]	[0],[0]
7	51	[0],[1]	[0],[0]
8	49	[2],[2]	[0],[0]

Table 4.3.2

*Summary of African American – White DIF Flags for Operational Items: Spring 2021 Social Studies by Grade*

Grade	A	[B+],[B-]	[C+],[C-]
3	43	[0],[0]	[0],[0]
4	41	[0],[2]	[0],[0]
5	46	[0],[0]	[0],[0]
6	52	[0],[1]	[0],[0]
7	50	[0],[1]	[0],[1]
8	52	[0],[1]	[0],[0]

Table 4.3.3

*Summary of Hispanic/Latino – White DIF Flags for Operational Items: Spring 2021 Social Studies by Grade*

Grade	A	[B+],[B-]	[C+],[C-]
3	42	[0],[1]	[0],[0]
4	43	[0],[0]	[0],[0]
5	46	[0],[0]	[0],[0]
6	53	[0],[0]	[0],[0]
7	50	[0],[2]	[0],[0]
8	52	[0],[1]	[0],[0]

## Pre-Equating for Intact Forms

Because the spring 2021 test administration used intact operational forms from spring 2019, the pre-equating method was applied. That is, the existing spring 2019 scoring tables were used for score report and performance classifications.

## Unidimensionality and Principal Component Analysis

[Appendix C: Dimensionality](#) provides information about principal component analysis of the LEAP 2025 Social Studies tests. Measurement implies order and magnitude along a single dimension (Andrich, 2004). Consequently, in the case of scholastic achievement, a one-dimensional scale is required to reflect this idea of measurement (Andrich, 1988, 1989). However, unidimensionality cannot be strictly met in a real testing situation because students' cognitive, personality, and test-taking factors usually have a unique influence on their test performance to some level (Andrich, 2004; Hambleton, Swaminathan, & Rogers, 1991). Consequently, what is required for unidimensionality to be met is an investigation of the presence of a dominant factor that influences test performance. This dominant factor is considered as the ability measured by the test (Andrich, 1988; Hambleton et al., 1991; Ryan, 1983).

To check the unidimensionality of the spring 2021 test, the relative sizes of the eigenvalues associated with a principal component analysis of the item set were examined using the Statistical Analysis System (SAS) program. The first and second principal component eigenvalues were compared *without rotation*. Tables C 1.1 and C 1.2 and Figures C 1.1 and C 1.2 summarize the results of the first and second principal component eigenvalues of the assessments. A general guideline in exploratory factor analysis suggests that a set of items may represent as many factors as there are eigenvalues greater than 1 because there is one unit of information per item and the eigenvalues sum to the total number of items. However, a set of items may have multiple eigenvalues greater than 1 and still be sufficiently unidimensional for analysis with IRT (Loehlin, 1987; Orlando, 2004). As seen from the tables and figures, the first component is substantially larger than the second eigenvalue for the spring 2021 test. Because the spring test was administered during the COVID-19 pandemic, great caution should be applied when any statistical inference is drawn.

### Scaling

Although the spring 2021 test used the preexisting scoring tables, general procedures for the scaling method are described here because scaling is directly associated with performance-level cuts. Based on the Standard Setting panelist recommendations and LDOE approval, the scale is set using two cut scores, Basic and Mastery, with fixed scale score points of 725 and 750, respectively. The scale scores for Approaching Basic and

Advanced are subsequently interpolated and vary by grades and subjects. The highest obtainable scale score (HOSS) and lowest obtainable scale score (LOSS) for the scale determined by the LDOE are 650 and 850.

IRT ability estimates ( $\theta$ s) are transformed to the reporting scale with a linear transformation equation of the form

$$SS = A\theta + B,$$

where  $SS$  is scale score,  $\theta$  is IRT ability,  $A$  is a slope coefficient, and  $B$  is an intercept. The slope can be calculated as

$$A = \frac{SS_{Mastery} - SS_{Basic}}{\theta_{Mastery} - \theta_{Basic}},$$

where  $\theta_{Mastery}$  is the Mastery cut score on the theta scale, and  $\theta_{Basic}$  is the Basic cut score on the theta scale.  $SS_{Mastery}$  and  $SS_{Basic}$  are the Mastery and Basic scale score cuts, respectively. With  $A$  calculated,  $B$  are derived from the equation

$$SS_{Mastery} = A\theta_{Mastery} + B,$$

which are rearranged as

$$B = SS_{Mastery} - A\theta_{Mastery} \text{ or } B = SS_{Mastery} - \frac{SS_{Mastery} - SS_{Basic}}{\theta_{Mastery} - \theta_{Basic}}\theta_{Mastery}.$$

Thus, the general equation for converting  $\theta$ s to scale scores is

$$SS = \left( \frac{SS_{Mastery} - SS_{Basic}}{\theta_{Mastery} - \theta_{Basic}} \right) \theta + \left( SS_{Mastery} - \frac{SS_{Mastery} - SS_{Basic}}{\theta_{Mastery} - \theta_{Basic}} \theta_{Mastery} \right).$$

The scaling constants  $A$  and  $B$  are calculated, and the Advanced cut score and the Approaching Basic cut score on the  $\theta$  scale are transformed to the reporting scale, rounded to the nearest integer. At this point, the score ranges associated with the five achievement levels are determined. The same scaling constants  $A$  and  $B$  are used to



convert student ability estimates to the reporting scale until new achievement-level standards are set. Descriptive statistics and frequency distribution of LEAP 2025 Social Studies scale scores can be found in [Appendix D: Scale Distribution and Statistical Report](#).

# 5. Reliability and Validity

## Internal Consistency Reliability Estimation

Internal consistency methods use data from a single administration to estimate test score reliability. For state assessments where student testing time is at a premium, internal consistency procedures have a practical advantage over reliability estimation procedures that require multiple test administrations. One of the most frequently used internal consistency reliability estimates is coefficient alpha (Cronbach, 1951). Coefficient alpha is based on the assumption that inter-item covariances constitute true-score variance and the fact that the average true score variance of items is greater than or equal to the average inter-item covariance. The formula for coefficient alpha is

$$\alpha = \left( \frac{N}{N-1} \right) \left( 1 - \frac{\sum_{i=1}^N s_{y_i}^2}{s_x^2} \right),$$

where  $N$  is the number of items on the test,  $s_{y_i}^2$  is the sample variance of the  $i$ th item (or component), and  $s_x^2$  is the observed score variance for the test. Coefficient alpha is appropriate for use when the items on the test are reasonably homogeneous. The homogeneity of LEAP 2025 Social Studies tests is evidenced through a dimensionality analysis. Dimensionality analyses results are discussed in “Chapter 4. Data Analysis.”

The reliability and classification accuracy reports in [Appendix E: Reliability and Classification Accuracy](#) provide Cronbach’s alpha for the total test. Cronbach’s alpha values for the spring 2021 tests were between 0.84 and 0.93. Because the spring test was administered during the COVID-19 pandemic, however, statistical inferences should be cautiously drawn from these results. Additional reliabilities were calculated on various demographic subgroups using the population of students (see [Appendix E: Reliability and Classification Accuracy](#)).

The subgroups are male/female, white/African American/Hispanic-Latino/Asian/American Indian or Alaska Native/Native Hawaiian or Other Pacific Islander/Two or More Races, Economically Disadvantaged, English Learners, Education Classification, and Section 504.

Cronbach's alpha estimates are computed for the entire test and each subscale by reporting category. Subscore reliability will generally be lower than total score reliability because reliability is influenced by the number of items as well as their covariation. In some cases, the number of items associated with a subscore is small (10 or fewer). Subscore results must be interpreted carefully when these measures reflect the limited number of items associated with the score.

## Student Classification Accuracy and Consistency

Students are classified into one of five performance levels based on their scale scores. It is important to know the reliability of student scores in any examination, but assessing the reliability of the classification decisions based on these scores is of even greater importance. Classification decision reliability is estimated by the probabilities of correct and consistent classification of students. Procedures were used from Livingston and Lewis (1995) and Lee, Hanson, and Brennan (2000) to derive accuracy and consistency classification measures.

**Accuracy of Classification.** According to Livingston and Lewis (1995, p. 180), the classification accuracy is "the extent to which the actual classifications of the test takers agree with those that would be made on the basis of their true scores, if their true scores could somehow be known." Accuracy estimates are calculated from cross-tabulations between "classifications based on an observable variable (scores on a test) and classifications based on an unobservable variable (the test takers' true scores)." True score is also referred to as a hypothetical mean of scores from all possible forms of the test if they could be somehow obtained (Young & Yoon, 1998).

**Consistency of Classification.** Classification consistency is "the agreement between classifications based on two non-overlapping, equally difficult forms of the test" (Livingston & Lewis, 1995, p. 180). Consistency is estimated using actual response data from a test and the test's reliability to statistically model two parallel forms of the test and compare the classifications on those alternate forms.

**Accuracy and Consistency Indices.** Three types of accuracy and consistency indices are generated: *overall*, *conditional-on-level*, and *cut point*, provided in [Appendix E: Reliability and Classification Accuracy](#). The *overall accuracy* of performance-level classifications is computed as a sum of the proportions on the diagonal of the joint distribution of true score and observed score levels. It is a proportion (or percentage) of correct classification across all the levels. While the overall accuracy index of the spring 2021 tests were between 0.625 and 0.716, the overall consistency values were between 0.530 and 0.618. Because the spring 2021 test was administered during the COVID-19 pandemic, however, great caution should be applied when any statistical inference is drawn.

Another way to express overall consistency is to use Cohen's Kappa ( $\kappa$ ) coefficient (Cohen, 1960). The overall coefficient Kappa when applying all cutoff scores together is

$$\kappa = \frac{P - P_c}{1 - P_c},$$

where  $P$  is the probability of consistent classification, and  $P_c$  is the probability of consistent classification by chance (Lee, Hanson, & Brennan, 2000).  $P$  is the sum of the diagonal elements, and  $P_c$  is the sum of the squared row totals. The PChance indices were between 0.218 and 0.246 for the spring 2021 tests.

Kappa is a measure of "how much agreement exists beyond chance alone" (Fleiss, 1973), which means that it provides the proportion of consistent classifications between two forms after removing the proportion of consistent classifications expected by chance alone. The Kappa indices were between 0.377 and 0.512 for the spring 2021 tests.

*Consistency conditional-on-level* is computed as the ratio between the proportion of correct classifications at the selected level (diagonal entry) and the proportion of all the students classified into that level (marginal entry).

*Accuracy conditional-on-level* is analogously computed. The only difference is that in the consistency table, both row and column marginal sums are the same, whereas in the accuracy table, the sum that is based on true status is used as a total for computing accuracy conditional on level.

Perhaps the most important indices for accountability systems are those for the accuracy and consistency of classification decisions made at specific cut points. To evaluate decisions at specific cut points, the joint distribution of all the performance levels is collapsed into a dichotomized distribution around that specific cut point.

## Validity

“Validity refers to the degree to which evidence and theory support the interpretations of test scores entailed by proposed users of tests. Validity is, therefore, the most fundamental consideration in developing and evaluating tests” (AERA/APA/NCME, 2014). The purpose of test score validation is not to validate the test itself but to validate interpretations of the test scores for particular purposes or uses. Test score validation is not a quantifiable property but an ongoing process, beginning at initial conceptualization and continuing throughout the entire assessment process.

The spring 2021 Social Studies tests were designed and developed to provide fair and accurate scores that support appropriate, meaningful, and useful educational decisions. As the technical addendum progresses, it details the procedures and processes applied to the LEAP 2025 Social Studies tests and their results. Validity evidence may be found in the following portions: Chapter 2 (Test Administration), Chapter 3 (Scoring Activities), Chapter 4 (Data Analysis), Chapter 5 (Reliability and Validity), and Chapter 6 (Statistical Summaries). For validity evidence related to the development and construction of the test form used in the spring 2021 administration, please refer to the [2018–2019 LEAP 2025 Social Studies Grades 3–8 Technical Report](#). Because the spring 2021 test was administered during the COVID-19 pandemic, any validity evidence associated with the spring test should be carefully interpreted.

The knowledge, expertise, and professional judgment offered by Louisiana educators ultimately ensure that the content for the LEAP 2025 Social Studies test is an adequate and representative sample of appropriate content, and that the content is a legitimate basis upon which to derive valid conclusions about student achievement. Participation by Louisiana educators throughout the process—from source selection, item development, and content and bias review to range-finding and standard setting—reinforces confidence in the content and design of the LEAP 2025 Social Studies test to derive valid inferences about Louisiana student performance.

Chapter 2 of the technical addendum describes the process, procedures, and policies that guide the administration of the LEAP 2025 assessments, including accommodations, test security, and detailed written procedures provided to test administrators and school personnel.

Chapter 3 describes scoring processes and activities for the LEAP 2025 Social Studies test. Although the spring 2021 test utilized a pre-equating method, Chapter 4 briefly describes classical data analysis, IRT, and scaling of the Social Studies tests, which derive scale scores from students' raw scores. In addition, Chapter 4 describes an analysis of DIF and includes gender and ethnicity DIF results. A summary of classical analysis and DIF results for the operational items is presented in [Appendix B: Item Analysis Summary Report](#).

Chapter 5 addresses Cronbach's alpha and marginal alpha as measures of internal consistency and also describes analysis procedures for classification consistency and classification accuracy.

Chapter 6 reports the statistical summaries of the spring 2021 Social Studies test.

## 6. Statistical Summaries

For the spring 2021 Social Studies test, the lowest obtainable scale score (LOSS) is 650 and the highest obtainable scale score (HOSS) is 850. Test results are provided in Tables 6.1.1 and 6.1.2. Scale score means and standard deviations as well as the percentages of students in each performance level are reported for the state and are disaggregated by demographic groups. In addition to the descriptive statistics presented in Tables 6.1.1–6.1.2, scale score frequency distributions are presented in [Appendix D: Scale Distribution and Statistical Report](#). Finally, because the spring 2021 test was administered during the COVID-19 pandemic, great caution should be applied when any statistical inference is drawn.

Table 6.1.1

## LEAP 2025 State Test Results: Spring 2021 Operational Social Studies Grade 3

	Scale Score			% at Performance Level**				
	N	Mean	SD	Unsatisfactory	Approaching Basic	Basic	Mastery	Advanced
TOTAL	≥49,530	714.18	40.40	30	30	20	15	5
Gender								
Female	≥24,250	714.78	39.47	29	31	21	15	5
Male	≥25,270	713.61	41.28	32	29	19	14	6
Ethnicity								
African American	≥20,920	700.49	36.16	42	33	16	8	2
American Indian or Alaska Native	≥310	715.28	39.49	29	31	19	16	5
Asian	≥810	740.19	40.24	12	21	25	27	16
Hispanic/Latino	≥4,830	708.90	38.83	35	30	19	12	3
Two or More Races	≥1,650	718.86	39.65	25	31	22	16	6
Native Hawaiian or Other Pacific Islander	≥40	723.61	29.15	9	41	32	18	NR
White	≥20,950	727.65	39.88	19	27	24	21	9
Economically Disadvantaged*								
No	≥12,610	737.45	39.31	13	22	25	26	13
Yes	≥36,710	706.24	37.61	36	32	18	11	3
English Learner								
No	≥47,060	715.31	40.45	29	30	20	15	6
Yes	≥2,470	692.75	32.81	50	33	13	4	1
Education Classification								
Regular	≥43,380	716.32	40.24	28	30	21	15	6
Special	≥6,140	699.09	38.28	45	31	13	8	3
Section 504 Status								
No	≥46,210	714.69	40.50	30	30	20	15	5
Yes	≥3,320	707.07	38.30	36	33	18	10	4

\* ≥210 students with no record of either No or Yes.

\*\* The percent may not add up to 100% due to rounding.



Table 6.1.2

LEAP 2025 State Test Results: Spring 2021 Operational Social Studies Grade 4

	Scale Score			% at Performance Level**				
	N	Mean	SD	Unsatisfactory	Approaching Basic	Basic	Mastery	Advanced
TOTAL	≥49,540	716.06	40.26	29	26	24	17	4
Gender								
Female	≥24,120	715.79	38.38	28	28	25	16	3
Male	≥25,410	716.33	41.97	30	25	23	17	5
Ethnicity								
African American	≥21,140	700.50	35.84	42	30	19	8	1
American Indian or Alaska Native	≥290	717.45	35.49	22	32	28	16	2
Asian	≥760	743.31	39.48	11	19	26	32	13
Hispanic/Latino	≥4,590	712.26	38.66	32	28	24	14	2
Two or More Races	≥1,630	722.41	38.90	23	25	28	21	4
Native Hawaiian or Other Pacific Islander	≥30	727.72	34.76	18	26	26	28	3
White	≥21,060	730.99	38.78	17	22	29	25	6
Economically Disadvantaged*								
No	≥12,990	740.01	37.69	11	18	30	31	10
Yes	≥36,300	707.60	37.62	35	29	22	12	2
English Learner								
No	≥47,410	717.10	40.29	28	26	25	17	4
Yes	≥2,130	693.08	31.93	50	32	15	3	NR
Education Classification								
Regular	≥43,430	718.89	39.83	26	26	26	18	4
Special	≥6,100	695.93	37.45	50	27	15	7	2
Section 504 Status								
No	≥45,100	716.91	40.48	29	26	24	17	4
Yes	≥4,430	707.52	36.91	35	31	22	11	2

\* ≥240 students with no record of either No or Yes.

\*\* The percent may not add up to 100% due to rounding.

Table 6.1.3

## LEAP 2025 State Test Results: Spring 2021 Operational Social Studies Grade 5

	Scale Score			% at Performance Level**				
	N	Mean	SD	Unsatisfactory	Approaching Basic	Basic	Mastery	Advanced
TOTAL	≥49,850	721.21	37.32	28	24	25	18	5
Gender								
Female	≥24,200	722.00	36.35	26	25	27	18	5
Male	≥25,650	720.46	38.20	29	23	24	19	5
Ethnicity								
African American	≥21,110	708.48	33.75	39	28	22	10	2
American Indian or Alaska Native	≥320	721.19	34.30	26	25	29	18	2
Asian	≥800	748.18	40.15	12	13	22	31	21
Hispanic/Latino	≥4,760	716.87	37.03	32	24	25	16	3
Two or More races	≥1,600	724.87	36.69	25	24	25	21	6
Native Hawaiian or Other Pacific Islander	≥20	731.36	44.54	25	18	29	14	14
White	≥21,210	733.53	36.09	17	20	29	26	8
Economically Disadvantaged*								
No	≥12,780	743.36	34.60	11	15	29	33	12
Yes	≥36,750	713.61	35.09	34	27	24	13	2
English Learner								
No	≥47,640	722.38	37.18	27	24	26	19	5
Yes	≥2,200	695.88	30.98	54	28	14	4	NR
Education Classification								
Regular	≥43,880	724.33	36.61	24	24	27	20	5
Special	≥5,970	698.26	34.40	53	25	14	6	1
Section 504 Status								
No	≥44,750	722.57	37.44	27	23	26	19	5
Yes	≥5,100	709.29	34.03	38	29	21	10	2

\* ≥310 students with no record of either No or Yes.

\*\* The percent may not add up to 100% due to rounding.

Table 6.1.4

## LEAP 2025 State Test Results: Spring 2021 Operational Social Studies Grade 6

	Scale Score			% at Performance Level**				
	N	Mean	SD	Unsatisfactory	Approaching Basic	Basic	Mastery	Advanced
TOTAL	≥51,530	718.11	38.38	30	23	25	14	7
Gender								
Female	≥25,330	719.25	37.10	28	25	26	15	7
Male	≥26,190	717.02	39.55	32	22	24	14	8
Ethnicity								
African American	≥22,290	703.50	34.29	44	27	20	7	2
American Indian or Alaska Native	≥310	721.32	34.88	23	27	28	16	6
Asian	≥760	750.03	37.92	9	14	24	27	26
Hispanic/Latino	≥4,600	714.15	38.83	34	23	23	13	6
Two or More Races	≥1,670	724.35	36.94	22	24	30	15	8
Native Hawaiian or Other Pacific Islander	≥40	726.13	37.20	27	19	25	19	10
White	≥21,830	732.21	36.39	17	20	30	21	12
Economically Disadvantaged*								
No	≥13,210	741.09	35.33	11	17	29	26	17
Yes	≥38,000	710.23	36.15	37	26	23	11	4
English Learner								
No	≥49,560	719.19	38.23	29	23	25	15	7
Yes	≥1,970	691.02	31.53	59	24	13	3	1
Education Classification								
Regular	≥45,680	721.75	37.54	26	24	27	16	8
Special	≥5,850	689.71	32.57	63	21	11	4	2
Section 504 Status								
No	≥46,060	719.84	38.42	28	23	26	15	8
Yes	≥5,460	703.56	34.78	45	27	19	8	3

\* ≥320 students with no record of either No or Yes.

\*\* The percent may not add up to 100% due to rounding.

Table 6.1.5

## LEAP 2025 State Test Results: Spring 2021 Operational Social Studies Grade 7

	Scale Score			% at Performance Level**				
	N	Mean	SD	Unsatisfactory	Approaching Basic	Basic	Mastery	Advanced
TOTAL	≥52,310	726.21	40.80	30	17	22	19	11
Gender								
Female	≥25,670	727.81	39.24	28	18	23	20	11
Male	≥26,630	724.66	42.20	33	16	21	19	12
Ethnicity								
African American	≥22,450	712.27	36.90	42	20	21	13	4
American Indian or Alaska Native	≥300	730.70	39.53	26	16	24	22	12
Asian	≥830	761.12	41.59	9	8	19	25	39
Hispanic/Latino	≥4,650	722.33	40.79	34	17	21	18	10
Two or More	≥1,590	730.70	39.40	26	16	24	21	13
Native Hawaiian or Other Pacific Islander	≥40	744.33	42.56	20	10	18	22	29
White	≥22,420	739.24	39.50	19	14	23	26	18
Economically Disadvantaged*								
No	≥14,090	749.56	38.02	12	11	22	30	25
Yes	≥37,870	717.64	38.35	37	19	22	15	7
English Learner								
No	≥50,390	727.32	40.65	29	17	22	20	12
Yes	≥1,920	696.97	33.32	60	19	14	6	1
Education Classification								
Regular	≥46,710	729.80	39.96	27	17	23	21	12
Special	≥5,600	696.20	34.96	63	17	11	6	3
Section 504 Status								
No	≥46,780	728.09	40.80	29	16	22	20	12
Yes	≥5,530	710.26	37.17	45	20	18	12	5

\* ≥340 students with no record of either No or Yes.

\*\* The percent may not add up to 100% due to rounding.

Table 6.1.6

## LEAP 2025 State Test Results: Spring 2021 Operational Social Studies Grade 8

	Scale Score			% at Performance Level**				
	N	Mean	SD	Unsatisfactory	Approaching Basic	Basic	Mastery	Advanced
TOTAL	≥51,830	730.88	39.70	24	18	23	25	10
Gender								
Female	≥25,530	733.78	37.76	20	19	25	26	10
Male	≥26,300	728.07	41.30	28	18	21	24	10
Ethnicity								
African American	≥22,160	716.71	36.60	34	23	23	17	4
American Indian or Alaska Native	≥310	732.81	38.88	22	14	29	27	8
Asian	≥800	764.35	39.76	8	7	15	34	36
Hispanic/Latino	≥4,210	724.72	40.81	30	17	22	23	8
Two or More	≥1,520	736.37	39.04	19	16	25	27	12
Native Hawaiian or Other Pacific Islander	≥40	741.93	33.20	11	25	18	30	16
White	≥22,760	744.21	37.06	13	14	24	34	15
Economically Disadvantaged*								
No	≥14,550	753.40	35.09	8	11	22	38	21
Yes	≥36,920	722.15	37.88	30	21	23	20	5
English Learner								
No	≥49,960	732.12	39.37	23	18	23	26	10
Yes	≥1,870	697.81	33.65	56	22	14	7	1
Education Classification								
Regular	≥46,640	734.69	38.40	20	18	24	27	11
Special	≥5,190	696.70	34.46	60	19	12	7	2
Section 504 Status								
No	≥46,520	732.67	39.57	22	18	23	26	10
Yes	≥5,300	715.24	37.30	37	23	21	15	4

\* ≥360 students with no record of either No or Yes.

\*\* The percent may not add up to 100% due to rounding.

# References

- AERA/APA/NCME. (2014). *Standards for educational and psychological testing*. Washington, DC: Author.
- Andrich, A. (1988). *Rasch models for measurement*. Newbury Park, CA: Sage Publications.
- Andrich, A. (1989). Distinctions between assumptions and requirements in measurement in the social sciences. In J. A. Keats, R. Taft, R. A. Heath, & H. H. Lovibond (Eds.), *Mathematical and theoretical systems*. North-Holland: Elsevier Science Publisher B.V.
- Andrich, A. (2004). *Modern measurement and analysis in social science*. Murdoch University, Perth, Western Australia.
- Angoff, W. (1993). Perspectives on differential item functioning methodology. In P. W. Holland & H. Warner (Eds.), *Differential item functioning* (pp. 3–24). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Camilli, G., & Shepard, L. A. (1994). *Methods for identifying biased test items*. Newbury Park, CA: Sage Publications.
- Clauser, B. E., & Mazor, K. M. (1998). Using statistical procedures to identify differentially functioning test items. *Educational Measurement: Issues and Practice*, 17, 31–44.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20(1), 37–47.
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16(3), 297–334.

- Dorans, N. J., & Schmitt, A. P. (1991). *Constructed response and differential item functioning: A pragmatic approach* (Research Report No. RR-91-47). Princeton, NJ: Educational Testing Service.
- Fleiss, J. L. (1973). *Statistical methods for rates and proportions*. New York: Wiley.
- Green, D. R. (1975, December). Procedures for assessing bias in achievement tests. Presented at the National Institute of Education Conference on Test Bias, Annapolis, MD.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Fundamentals of item response theory*. Newbury Park, CA: Sage Publications.
- Holland, P. W., & Thayer, D. T. (1988). Differential item performance and the Mantel-Haenszel procedure. In H. Wainer & H. I. Braun (Eds.), *Test validity* (pp. 129-145). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Lee, W., Hanson, B. A., & Brennan, R. L. (2000, October). Procedures for computing classification consistency and accuracy indices with multiple categories (ACT Research Report Series 2000-10). Iowa City: ACT, Inc.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32(2), 179-197.
- Loehlin, J. C. (1987). *Latent variable models*. Hillsdale, NJ: Lawrence Erlbaum Associates, Publishers.
- Mantel, N. (1963). Chi-square tests with one degree of freedom: Extensions of the Mantel-Haenszel procedure. *Journal of the American Statistical Association*, 58, 690-700.

- 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.
- Orlando, M. (2004, June). Critical issues to address when applying item response theory (IRT) models. Paper presented at the Drug Information Association, Bethesda, MD.
- Ryan, J. P. (1983). Introduction to latent trait analysis and item response theory. In W. E. Hathaway (Ed.), *Testing in the schools: New directions for testing and measurement* (p. 19). San Francisco: Jossey-Bass.
- Young, M. J., & Yoon, B. (1998, April). Estimating the consistency and accuracy of classifications in a standards-referenced assessment (CSE Technical Report 475). Center for the Study of Evaluation, National Center for Research on Evaluation, Standards, and Student Testing. Los Angeles: University of California, Los Angeles.
- Zieky, M. (1993). Practical questions in the use of DIF statistics in test development. In P. W. Holland & H. Wainer (Eds.), *Differential item functioning* (pp. 337–348). Hillsdale, NJ: Lawrence Erlbaum Associates.
- Zumbo, B. D. (1999). *A handbook on the theory and methods of differential item functioning (DIF): Logistic regression modeling as a unitary framework for binary and Likert-type (ordinal) item scores*. Ottawa: Directorate of Human Resources Research and Evaluation, Department of National Defense.
- Zwick, R., Donoghue, J. R., & Grima, A. (1993). Assessment of differential item functioning for performance tasks. *Journal of Educational Measurement*, 26, 44–66.
- Zwick, R., Thayer, D. T., & Mazzeo, J. (1997). Descriptive and inferential procedures for assessing differential item functioning in polytomous items. *Applied Measurement in Education*, 10(4), 321–344.



# Appendix A: Test Summary

## *Social Studies*

Contents
Table A.1 Item Type Summary: Spring 2021 Operational Social Studies
Table A.2 Raw Score Summary: Spring 2021 Operational Social Studies
Table A.3 Raw Score Summary by Reporting Category: Spring 2021 Operational Social Studies
Table A.4 Scale Score and Raw Score Summary: Spring 2021 Operational Social Studies

- Because the spring 2021 test was administered during the COVID-19 pandemic, great caution should be applied when any statistical inference is drawn.

Table A.1.1

*Item Type Summary: Spring 2021 Operational Social Studies*

Grade	MC	MS	TE	CR	ER*
3	38	3	0	2	0
4	37	4	0	2	0
5	38	2	3	2	1
6	42	4	4	2	1
7	43	2	4	2	1
8	38	8	4	2	1

\* Classical and IRT analyses are calculated and estimated separately for each dimension of the ER item, and the result summarizes both dimensions.

Table A.2.1

*Raw Score Summary: Spring 2021 Operational Social Studies*

Grade	N	Mean	SD	Min	Max	Mean_Pval	Mean_Pbis	Reliability*	SEM
3**	≥49,530	18.53	7.63	0	43	0.43	0.36	0.84	3.05
4**	≥49,540	20.09	8.07	1	45	0.47	0.38	0.86	3.02
5	≥49,850	20.50	9.61	0	55	0.40	0.39	0.88	3.33
6	≥51,530	28.45	11.27	0	63	0.48	0.40	0.91	3.38
7	≥52,310	26.83	11.90	1	65	0.47	0.41	0.91	3.57
8	≥51,830	32.32	13.17	1	65	0.53	0.44	0.93	3.48

\* Reliability is Cronbach's alpha.

\*\* The PBT (paper-based test) and CBT (computer-based test) are combined together for statistical analysis.

Table A.3.1

*Raw Score Summary by Reporting Category: Spring 2021 Operational Social Studies*

<b>Grade</b>	<b>Reporting Category</b>	<b>Mean</b>	<b>SD</b>	<b>Min</b>	<b>Max</b>	<b>Mean_Pval</b>	<b>Mean_Pbis</b>	<b>Reliability</b>	<b>SEM</b>
3	History	4.17	2.09	0	11	0.39	0.31	0.54	1.42
	Geography	4.17	2.17	0	11	0.41	0.35	0.58	1.41
	Civics	6.03	2.55	0	12	0.51	0.36	0.62	1.57
	Economics	4.16	2.54	0	11	0.40	0.41	0.53	1.74
4	History	5.22	2.60	0	12	0.45	0.39	0.49	1.86
	Geography	5.14	2.29	0	11	0.48	0.34	0.64	1.37
	Civics	5.24	2.32	0	11	0.49	0.35	0.60	1.47
	Economics	4.49	2.48	0	11	0.45	0.44	0.68	1.40
5	History	11.64	5.78	0	33	0.39	0.38	0.84	2.31
	Geography	2.97	1.66	0	7	0.44	0.36	0.45	1.23
	Civics	1.96	1.48	0	7	0.31	0.33	0.33	1.21
	Economics	3.93	2.16	0	9	0.50	0.48	0.53	1.48
6	History	13.06	5.04	0	28	0.48	0.38	0.83	2.08
	Geography	9.05	4.34	0	24	0.49	0.44	0.70	2.38
	Civics	2.99	1.64	0	7	0.46	0.44	0.37	1.30
	Economics	3.35	1.74	0	7	0.49	0.37	0.66	1.01
7	History	14.66	6.89	0	39	0.44	0.40	0.84	2.76
	Geography	2.50	1.53	0	6	0.46	0.43	0.72	0.81
	Civics	6.47	3.13	0	13	0.54	0.43	0.66	1.83
	Economics	3.21	1.76	0	7	0.46	0.42	0.53	1.21
8	History	19.32	8.64	0	42	0.51	0.46	0.89	2.87
	Geography	4.12	2.11	0	8	0.53	0.44	0.62	1.30
	Civics	4.68	2.00	0	8	0.60	0.40	0.57	1.31
	Economics	4.21	1.83	0	8	0.54	0.38	0.59	1.17

Table A.4.1.1

*Scale Score and Raw Score Summary: Spring 2021 Operational Social Studies Grade 3*

<b>Subgroup</b>	<b>N</b>	<b>Percent</b>	<b>Scale Score Mean</b>	<b>Scale Score SD</b>	<b>Raw Score Mean</b>	<b>Raw Score SD</b>
Total	≥49,530	100.00	714.18	40.40	18.53	7.63
Female	≥24,250	48.96	714.78	39.47	18.61	7.46
Male	≥25,270	51.02	713.61	41.28	18.46	7.79
African American	≥20,920	42.23	700.49	36.16	15.91	6.42
American Indian or Alaska Native	≥310	0.63	715.28	39.49	18.63	7.68
Asian	≥810	1.65	740.19	40.24	23.66	8.15
Hispanic/Latino	≥4,830	9.75	708.90	38.83	17.51	7.17
Two or More Races	≥1,650	3.34	718.86	39.65	19.38	7.57
Native Hawaiian or Other Pacific Islander	≥40	0.09	723.61	29.15	19.89	5.62
White	≥20,950	42.30	727.65	39.88	21.11	7.86
Economically Disadvantaged: No*	≥12,610	25.46	737.45	39.31	23.08	7.95
Economically Disadvantaged: Yes*	≥36,710	74.10	706.24	37.61	16.98	6.86
EL: No	≥47,060	95.00	715.31	40.45	18.74	7.67
EL: Yes	≥2,470	5.00	692.75	32.81	14.48	5.51
Regular Education	≥43,380	87.59	716.32	40.24	18.92	7.66
Special Education	≥6,140	12.41	699.09	38.28	15.79	6.85
Section 504 Status: No	≥46,210	93.29	714.69	40.50	18.63	7.66
Section 504 Status: Yes	≥3,320	6.71	707.07	38.30	17.14	7.04

\* Economic Status was not available for all students.

Table A.4.1.2

*Scale Score and Raw Score Summary: Spring 2021 Operational Social Studies Grade 4*

<b>Subgroup</b>	<b>N</b>	<b>Percent</b>	<b>Scale Score Mean</b>	<b>Scale Score SD</b>	<b>Raw Score Mean</b>	<b>Raw Score SD</b>
Total	≥49,540	100.00	716.06	40.26	20.09	8.07
Female	≥24,120	48.68	715.79	38.38	19.98	7.70
Male	≥25,410	51.30	716.33	41.97	20.19	8.41
African American	≥21,140	42.69	700.50	35.84	16.92	6.80
American Indian or Alaska Native	≥290	0.59	717.45	35.49	20.22	7.10
Asian	≥760	1.55	743.31	39.48	25.66	8.34
Hispanic/Latino	≥4,590	9.27	712.26	38.66	19.28	7.67
Two or More Races	≥1,630	3.30	722.41	38.90	21.34	7.96
Native Hawaiian or Other Pacific Islander	≥30	0.08	727.72	34.76	22.33	7.63
White	≥21,060	42.53	730.99	38.78	23.14	8.06
Economically Disadvantaged: No*	≥12,990	26.23	740.01	37.69	25.04	7.97
Economically Disadvantaged: Yes*	≥36,300	73.27	707.60	37.62	18.34	7.34
EL: No	≥47,410	95.69	717.10	40.29	20.30	8.10
EL: Yes	≥2,130	4.31	693.08	31.93	15.43	5.76
Regular Education	≥43,430	87.68	718.89	39.83	20.64	8.05
Special Education	≥6,100	12.32	695.93	37.45	16.15	7.09
Section 504 Status: No	≥45,100	91.04	716.91	40.48	20.27	8.13
Section 504 Status: Yes	≥4,430	8.96	707.52	36.91	18.28	7.22

\* Economic Status was not available for all students.

Table A.4.1.3

*Scale Score and Raw Score Summary: Spring 2021 Operational Social Studies Grade 5*

<b>Subgroup</b>	<b>N</b>	<b>Percent</b>	<b>Scale Score Mean</b>	<b>Scale Score SD</b>	<b>Raw Score Mean</b>	<b>Raw Score SD</b>
Total	≥49,850	100.00	721.21	37.32	20.50	9.61
Female	≥24,200	48.55	722.00	36.35	20.58	9.42
Male	≥25,650	51.45	720.46	38.20	20.42	9.79
African American	≥21,110	42.35	708.48	33.75	17.14	7.70
American Indian or Alaska Native	≥320	0.64	721.19	34.30	20.22	8.62
Asian	≥800	1.62	748.18	40.15	28.37	11.70
Hispanic/Latino	≥4,760	9.56	716.87	37.03	19.40	9.16
Two or More Races	≥1,600	3.22	724.87	36.69	21.38	9.70
Native Hawaiian or Other Pacific Islander	≥20	0.06	731.36	44.54	23.64	12.19
White	≥21,210	42.55	733.53	36.09	23.73	10.06
Economically Disadvantaged: No*	≥12,780	25.64	743.36	34.60	26.53	10.22
Economically Disadvantaged: Yes*	≥36,750	73.73	713.61	35.09	18.43	8.45
EL: No	≥47,640	95.57	722.38	37.18	20.78	9.65
EL: Yes	≥2,200	4.43	695.88	30.98	14.34	6.05
Regular Education	≥43,880	88.02	724.33	36.61	21.23	9.65
Special Education	≥5,970	11.98	698.26	34.40	15.14	7.41
Section 504 Status: No	≥44,750	89.76	722.57	37.44	20.86	9.72
Section 504 Status: Yes	≥5,100	10.24	709.29	34.03	17.33	7.92

\* Economic Status was not available for all students.

Table A.4.1.4

*Scale Score and Raw Score Summary: Spring 2021 Operational Social Studies Grade 6*

<b>Subgroup</b>	<b>N</b>	<b>Percent</b>	<b>Scale Score Mean</b>	<b>Scale Score SD</b>	<b>Raw Score Mean</b>	<b>Raw Score SD</b>
Total	≥51,530	100.00	718.11	38.38	28.45	11.27
Female	≥25,330	49.17	719.25	37.10	28.69	10.97
Male	≥26,190	50.83	717.02	39.55	28.21	11.54
African American	≥22,290	43.26	703.50	34.29	24.10	9.56
American Indian or Alaska Native	≥310	0.61	721.32	34.88	29.27	10.42
Asian	≥760	1.48	750.03	37.92	38.17	11.63
Hispanic/Latino	≥4,600	8.94	714.15	38.83	27.34	11.23
Two or More Races	≥1,670	3.25	724.35	36.94	30.21	11.00
Native Hawaiian or Other Pacific Islander	≥40	0.09	726.13	37.20	30.85	11.47
White	≥21,830	42.37	732.21	36.39	32.62	11.11
Economically Disadvantaged: No*	≥13,210	25.64	741.09	35.33	35.36	10.96
Economically Disadvantaged: Yes*	≥38,000	73.74	710.23	36.15	26.07	10.35
EL: No	≥49,560	96.17	719.19	38.23	28.75	11.26
EL: Yes	≥1,970	3.83	691.02	31.53	20.73	8.21
Regular Education	≥45,680	88.64	721.75	37.54	29.47	11.16
Special Education	≥5,850	11.36	689.71	32.57	20.46	8.64
Section 504 Status: No	≥46,060	89.39	719.84	38.42	28.96	11.33
Section 504 Status: Yes	≥5,460	10.61	703.56	34.78	24.14	9.74

\* Economic Status was not available for all students.

Table A.4.1.5

*Scale Score and Raw Score Summary: Spring 2021 Operational Social Studies Grade 7*

<b>Subgroup</b>	<b>N</b>	<b>Percent</b>	<b>Scale Score Mean</b>	<b>Scale Score SD</b>	<b>Raw Score Mean</b>	<b>Raw Score SD</b>
Total	≥52,310	100.00	726.21	40.80	26.83	11.90
Female	≥25,670	49.08	727.81	39.24	27.15	11.58
Male	≥26,630	50.92	724.66	42.20	26.53	12.19
African American	≥22,450	42.92	712.27	36.90	22.69	9.92
American Indian or Alaska Native	≥300	0.58	730.70	39.53	28.12	11.78
Asian	≥830	1.60	761.12	41.59	37.80	13.30
Hispanic/Latino	≥4,650	8.89	722.33	40.79	25.76	11.66
Two or More Races	≥1,590	3.05	730.70	39.40	28.03	11.80
Native Hawaiian or Other Pacific Islander	≥40	0.09	744.33	42.56	32.61	13.37
White	≥22,420	42.86	739.24	39.50	30.68	12.17
Economically Disadvantaged: No*	≥14,090	26.94	749.56	38.02	33.92	12.14
Economically Disadvantaged: Yes*	≥37,870	72.40	717.64	38.35	24.23	10.69
EL: No	≥50,390	96.33	727.32	40.65	27.14	11.92
EL: Yes	≥1,920	3.67	696.97	33.32	18.75	7.94
Regular Education	≥46,710	89.29	729.80	39.96	27.80	11.86
Special Education	≥5,600	10.71	696.20	34.96	18.75	8.70
Section 504 Status: No	≥46,780	89.43	728.09	40.80	27.38	11.99
Section 504 Status: Yes	≥5,530	10.57	710.26	37.17	22.19	10.00

\* Economic Status was not available for all students.



Table A.4.1.6

*Scale Score and Raw Score Summary: Spring 2021 Operational Social Studies Grade 8*

<b>Subgroup</b>	<b>N</b>	<b>Percent</b>	<b>Scale Score Mean</b>	<b>Scale Score SD</b>	<b>Raw Score Mean</b>	<b>Raw Score SD</b>
Total	≥51,830	100.00	730.88	39.70	32.32	13.17
Female	≥25,530	49.26	733.78	37.76	33.20	12.72
Male	≥26,300	50.74	728.07	41.30	31.48	13.55
African American	≥22,160	42.75	716.71	36.60	27.51	11.78
American Indian or Alaska Native	≥310	0.61	732.81	38.88	32.91	12.72
Asian	≥800	1.56	764.35	39.76	43.64	12.91
Hispanic/Latino	≥4,210	8.13	724.72	40.81	30.41	13.30
Two or More Races	≥1,520	2.95	736.37	39.04	34.17	13.02
Native Hawaiian or Other Pacific Islander	≥40	0.08	741.93	33.20	35.84	11.97
White	≥22,760	43.92	744.21	37.06	36.83	12.61
Economically Disadvantaged: No*	≥14,550	28.08	753.40	35.09	39.97	11.99
Economically Disadvantaged: Yes*	≥36,920	71.22	722.15	37.88	29.36	12.39
EL: No	≥49,960	96.38	732.12	39.37	32.73	13.11
EL: Yes	≥1,870	3.62	697.81	33.65	21.66	9.93
Regular Education	≥46,640	89.97	734.69	38.40	33.54	12.90
Special Education	≥5,190	10.03	696.70	34.46	21.41	10.26
Section 504 Status: No	≥46,520	89.76	732.67	39.57	32.93	13.17
Section 504 Status: Yes	≥5,300	10.24	715.24	37.30	27.02	11.98

\* Economic Status was not available for all students.

# Appendix B: Item Analysis Summary Report

## Summary Statistics Reports

### Social Studies

Contents
Table B.1.1 P-Value Summary by Item Type: Spring 2021 Operational Social Studies
Plot B.1.1 P-Value Summary by Item Type: Spring 2021 Operational Social Studies
Table B.2.1 Item-Total Correlation Summary by Item Type: Spring 2021 Operational Social Studies
Plot B.2.1 Item-Total Correlation Summary by Item Type: Spring 2021 Operational Social Studies
Table B.3.1 Corrected Point-Biserial Correlation Summary by Item Type: Spring 2021 Operational Social Studies
Plot B.3.1 Corrected Point-Biserial Correlation Summary by Item Type: Spring 2021 Operational Social Studies
Table B.4.1 Item-Total Correlation Summary by Reporting Category and Item Type: Spring 2021 Operational Social Studies
Table B.5.1 Statistically Flagged Items by Item Type: Spring 2021 Operational Social Studies

- Because the spring 2021 test was administered during the COVID-19 pandemic, great caution should be applied when any statistical inference is drawn.

Table B.1.1

*P-Value Summary by Grade: Spring 2021 Operational Social Studies*

<b>Grade</b>	<b>No. of Items</b>	<b>Minimum</b>	<b>25th Percentile</b>	<b>Median</b>	<b>75th Percentile</b>	<b>Maximum</b>
3	43	0.161	0.339	0.434	0.518	0.731
4	43	0.098	0.389	0.461	0.533	0.763
5	46	0.149	0.310	0.382	0.495	0.728
6	53	0.122	0.393	0.472	0.587	0.792
7	52	0.142	0.397	0.455	0.549	0.721
8	53	0.269	0.452	0.530	0.616	0.758

\* Classical and IRT analyses are calculated and estimated separately for each dimension of the ER item, and the result summarizes both dimensions.

Plot B.1.1

*P-Value Summary by Grade: Spring 2021 Operational Social Studies*

***Box and Whisker Plot***

**P-Value: Social Studies**

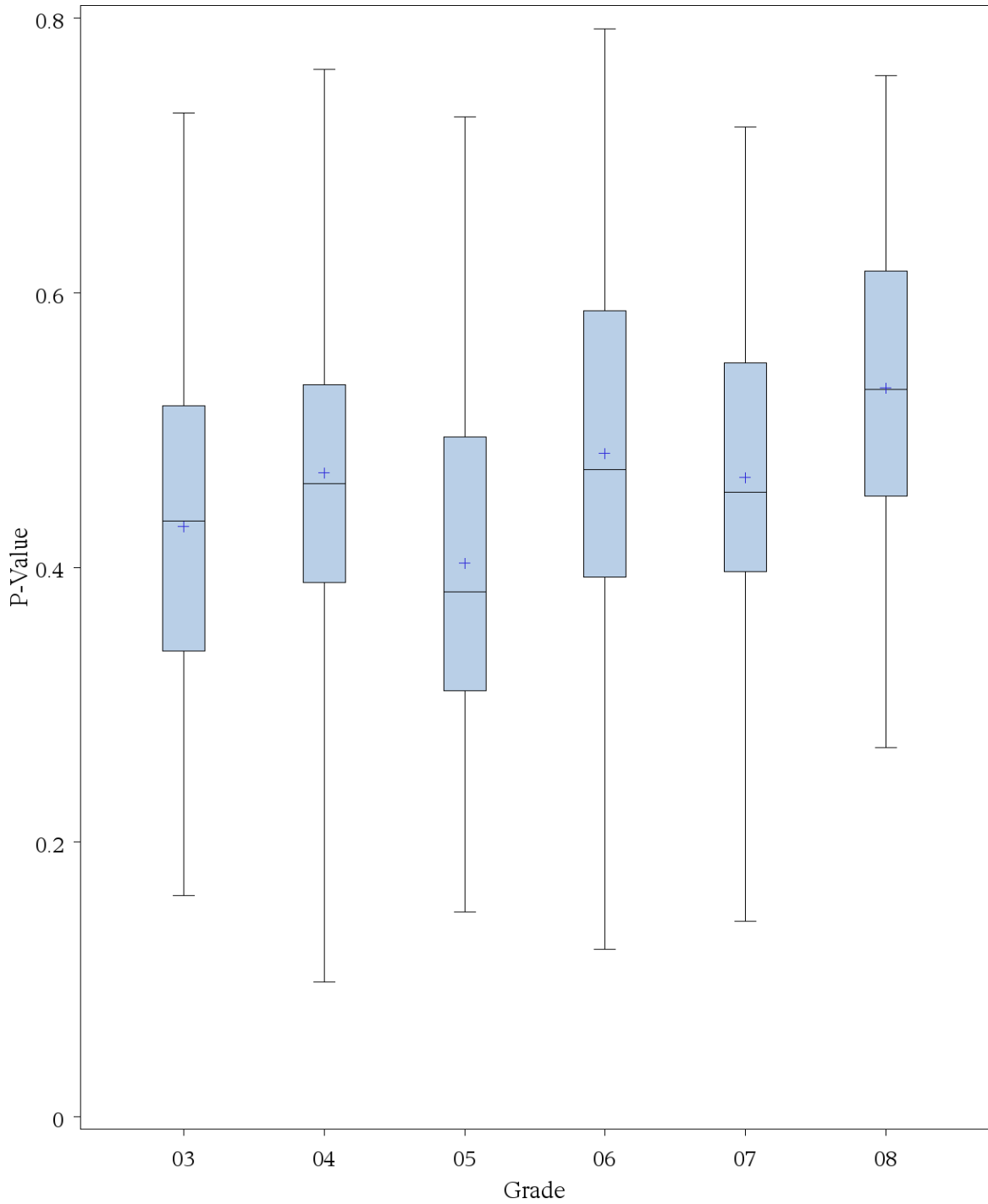


Table B.2.1

*Item-Total Correlation Summary by Grade: Spring 2021 Operational Social Studies*

<b>Grade</b>	<b>No. of Items</b>	<b>Minimum</b>	<b>25th Percentile</b>	<b>Median</b>	<b>75th Percentile</b>	<b>Maximum</b>
3	43	0.145	0.283	0.375	0.425	0.587
4	43	0.179	0.287	0.404	0.436	0.577
5	46	0.140	0.298	0.382	0.437	0.735
6	53	0.215	0.322	0.384	0.454	0.711
7	52	0.224	0.323	0.408	0.470	0.717
8	53	0.119	0.362	0.428	0.509	0.754

\* Classical and IRT analyses are calculated and estimated separately for each dimension of the ER item, and the result summarizes both dimensions.

Plot B.2.1

Item-Total Correlation Summary by Grade: Spring 2021 Operational Social Studies

***Box and Whisker Plot***

**Point-Biserial Correlation: Social Studies**

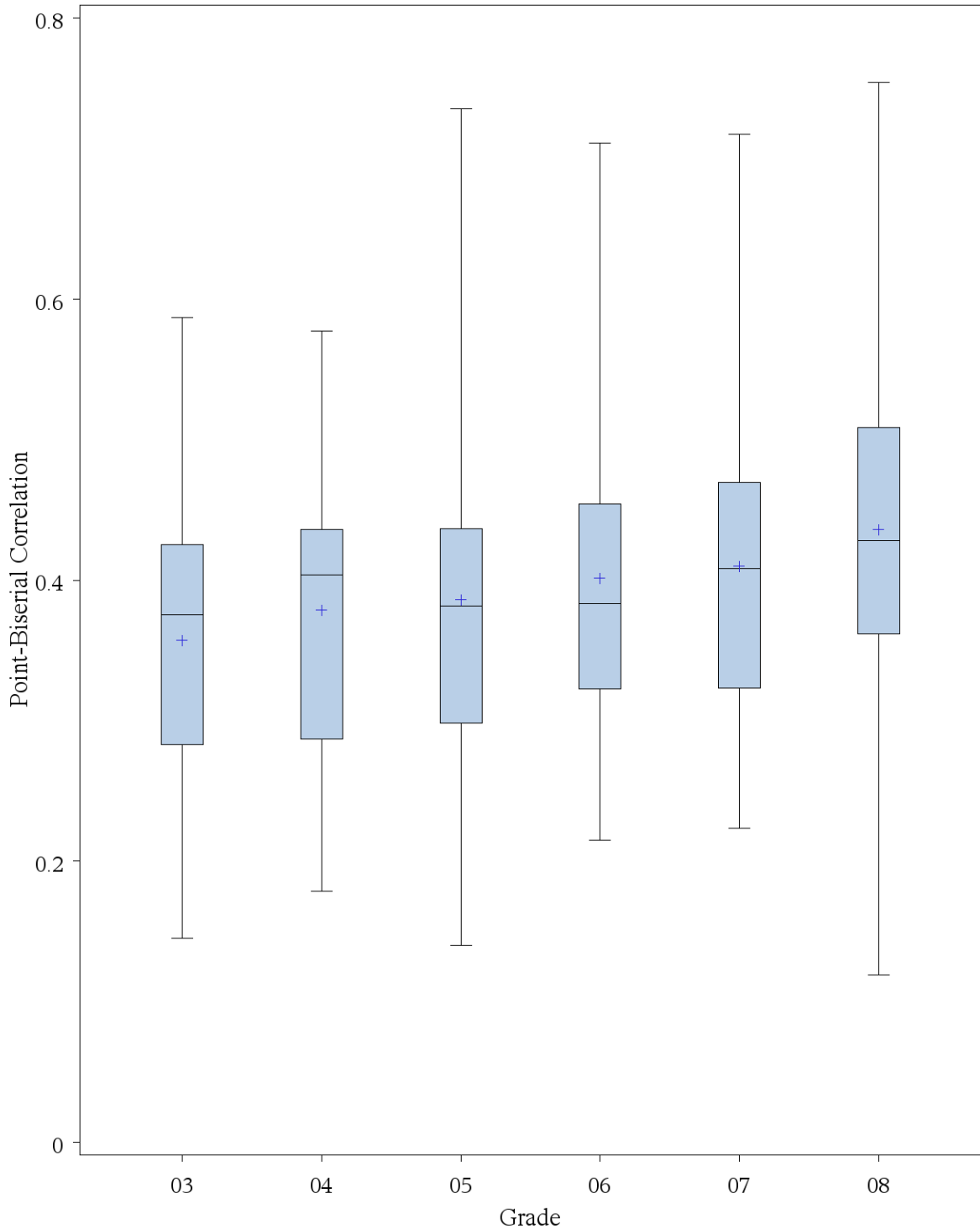


Table B.3.1

*Corrected Point-Biserial Correlation\* Summary by Grade: Spring 2021 Operational Social Studies*

<b>Grade</b>	<b>No. of Items</b>	<b>Minimum</b>	<b>25th Percentile</b>	<b>Median</b>	<b>75th Percentile</b>	<b>Maximum</b>
3	43	0.087	0.223	0.320	0.370	0.517
4	43	0.118	0.231	0.350	0.386	0.506
5	46	0.091	0.254	0.337	0.394	0.692
6	53	0.175	0.282	0.347	0.422	0.676
7	52	0.183	0.285	0.376	0.428	0.681
8	53	0.083	0.330	0.396	0.482	0.718

\* Corrected point-biserial correlation, which is slightly more robust than point-biserial correlation, calculates the relationship between the item score and the total test score after removing the item score from the total test score.

\*\* Classical and IRT analyses are calculated and estimated separately for each dimension of the ER item, and the result summarizes both dimensions.

Plot B.3.1

Corrected Point-Biserial Correlation by Grade: Spring 2021 Operational Social Studies

***Box and Whisker Plot***

**Corrected Point-Biserial Correlation: Social Studies**

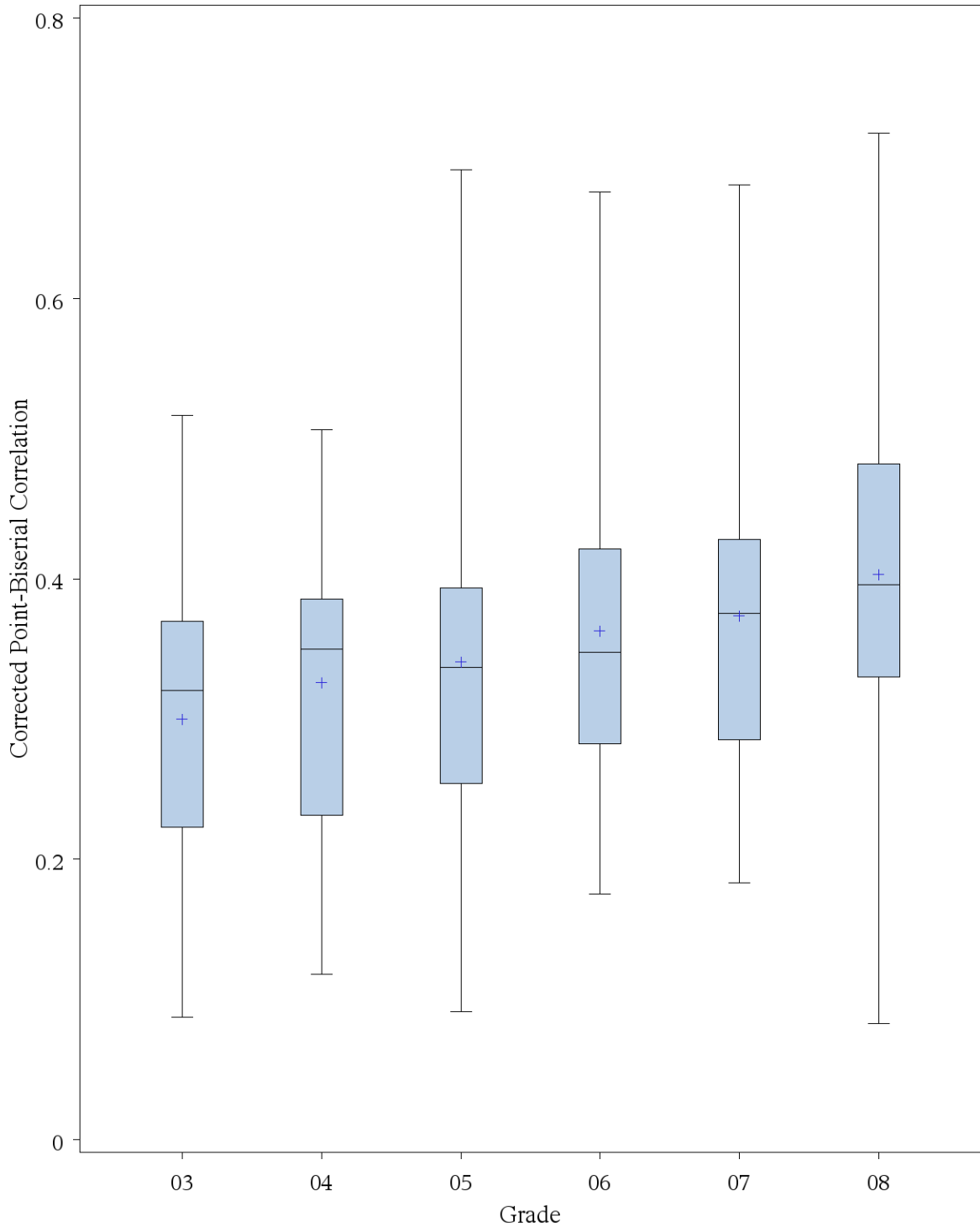




Table B.4.1

*Item-Total Correlation Summary by Reporting Category: Spring 2021 Operational Social Studies*

Grade	Reporting Category	No. of Items	Minimum	25th Percentile	Median	75th Percentile	Maximum
3	History	11	0.145	0.201	0.351	0.398	0.421
	Geography	10	0.164	0.283	0.348	0.415	0.510
	Civics	12	0.159	0.268	0.382	0.449	0.541
	Economics	10	0.285	0.364	0.388	0.425	0.587
4	History	11	0.225	0.302	0.398	0.468	0.577
	Geography	11	0.179	0.227	0.381	0.412	0.526
	Civics	11	0.220	0.265	0.326	0.428	0.539
	Economics	10	0.370	0.414	0.433	0.483	0.513
5	History	26	0.168	0.298	0.373	0.415	0.735
	Geography	7	0.231	0.279	0.382	0.431	0.445
	Civics	6	0.140	0.201	0.334	0.429	0.565
	Economics	7	0.352	0.414	0.498	0.532	0.568
6	History	25	0.215	0.321	0.373	0.437	0.550
	Geography	16	0.257	0.373	0.425	0.492	0.711
	Civics	5	0.316	0.322	0.454	0.521	0.608
	Economics	7	0.261	0.350	0.376	0.426	0.441
7	History	30	0.224	0.295	0.395	0.460	0.717
	Geography	5	0.239	0.363	0.450	0.462	0.628
	Civics	11	0.266	0.332	0.417	0.497	0.656
	Economics	6	0.285	0.323	0.407	0.504	0.577
8	History	31	0.301	0.363	0.447	0.525	0.754
	Geography	8	0.300	0.396	0.428	0.512	0.558
	Civics	7	0.119	0.324	0.425	0.509	0.598
	Economics	7	0.248	0.271	0.425	0.473	0.497

\* Classical and IRT analyses are calculated and estimated separately for each dimension of the ER item, and the result summarizes both dimensions.

Table B.5.1

*Statistically Flagged Items by Item Type: Spring 2021 Operational Social Studies*

Grade	Item Type	N OP Items	N Items Flagged for P-Value	N Items Flagged for Point-Biserial Correlation	N Items Flagged for DIF*	N Items Flagged for Omitting
3	CR	2	2	0	0	0
	MC	38	0	4	1	0
	MS	3	2	0	0	0
4	CR	2	1	0	0	0
	MC	37	0	2	1	0
	MS	4	1	0	1	0
5	CR	2	1	0	0	0
	ER**	1	1	0	0	0
	MC	38	3	2	0	0
	MS	2	0	0	0	0
	TE	3	1	0	0	0
6	CR	2	1	0	0	0
	ER**	1	1	0	0	0
	MC	42	0	0	2	0
	MS	4	1	0	0	0
	TE	4	0	0	0	0
7	CR	2	0	0	0	0
	ER**	1	1	0	0	0
	MC	43	0	0	0	0
	MS	2	0	0	2	0
	TE	4	1	0	3	0
8	CR	2	0	0	1	0
	ER**	1	0	0	1	0
	MC	38	0	0	4	0
	MS	8	0	1	0	0
	TE	4	0	0	0	0

\* The number of flagged DIF items includes both B and C DIF items.

\*\* Classical and IRT analyses are calculated and estimated separately for each dimension of the ER item, and the result summarizes both dimensions.

# Appendix C: Dimensionality

## *Dimensionality Reports*

### *Social Studies*

Contents
Table C.1.1 Reporting Category Intercorrelation Coefficients for Spring 2021 Operational Social Studies
Table C.2.1 First and Second Eigenvalues: Spring 2021 Operational Social Studies
Figure C.1.1 Principal Component Analysis Plot: Spring 2021 Operational Social Studies Grades 3 and 4
Figure C.1.2 Principal Component Analysis Plot: Spring 2021 Operational Social Studies Grades 5 through 8

- Because the spring 2021 test was administered during the COVID-19 pandemic, great caution should be applied when any statistical inference is drawn.

Table C.1.1

*Intercorrelation Coefficients among Reporting Categories: Spring 2021 Operational Social Studies Grade 3*

<b>Grade</b>	<b>Reporting Category</b>	<b>History</b>	<b>Geography</b>	<b>Civics</b>	<b>Economics</b>
3	History	1.00			
	Geography	0.52	1.00		
	Civics	0.52	0.56	1.00	
	Economics	0.55	0.57	0.58	1.00
4	History	1.00			
	Geography	0.59	1.00		
	Civics	0.60	0.53	1.00	
	Economics	0.64	0.58	0.59	1.00
5	History	1.00			
	Geography	0.62	1.00		
	Civics	0.57	0.44	1.00	
	Economics	0.73	0.57	0.50	1.00
6	History	1.00			
	Geography	0.77	1.00		
	Civics	0.65	0.66	1.00	
	Economics	0.62	0.62	0.51	1.00
7	History	1.00			
	Geography	0.67	1.00		
	Civics	0.76	0.61	1.00	
	Economics	0.68	0.53	0.62	1.00
8	History	1.00			
	Geography	0.75	1.00		
	Civics	0.71	0.60	1.00	
	Economics	0.68	0.58	0.55	1.00

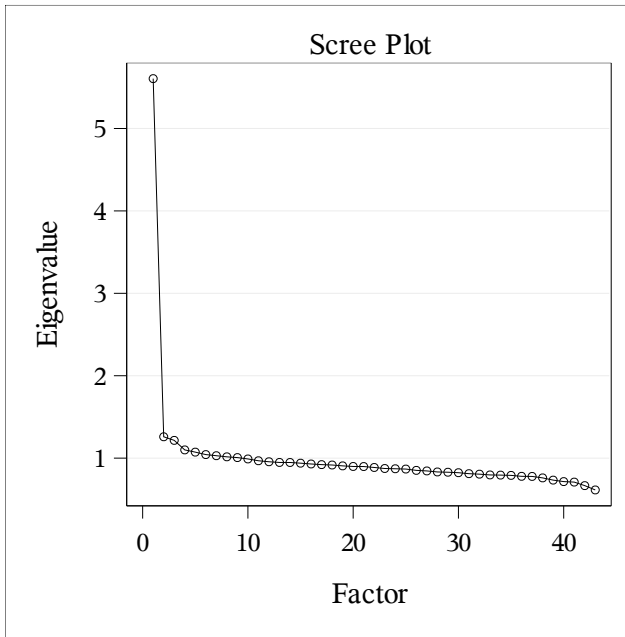
Table C.2.1

*First and Second Eigenvalues by Grade: Spring 2201 Operational Social Studies*

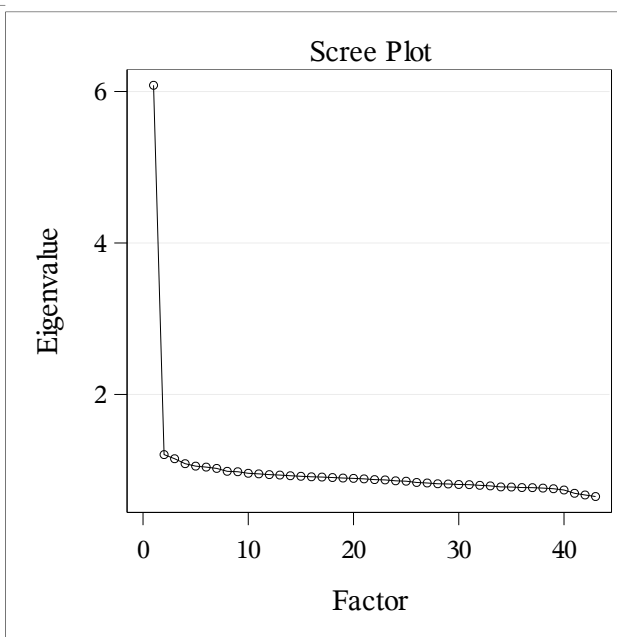
<b>Grade</b>	<b>Form</b>	<b>First Eigenvalue</b>	<b>Second Eigenvalue</b>	<b>Ratio</b>
3	Online	5.605	1.259	4.452
	Paper	6.082	1.206	5.043
4	Online	6.685	1.521	4.395
	Paper	6.689	1.475	4.535
5	Online	7.775	1.261	6.166
6	Online	9.588	1.264	7.585
7	Online	9.797	1.457	6.724
8	Online	11.384	1.320	8.624

Figure C.1.1

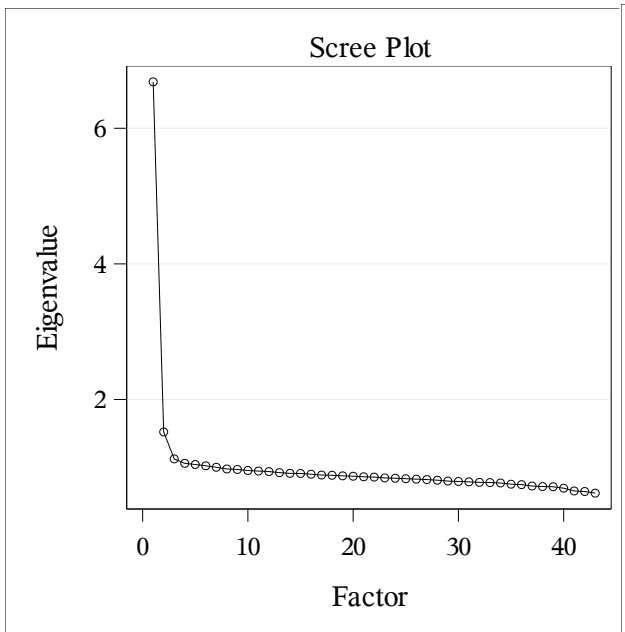
Principal Component Analysis Plot: Spring 2021 Operational Social Studies Grades 3 and 4



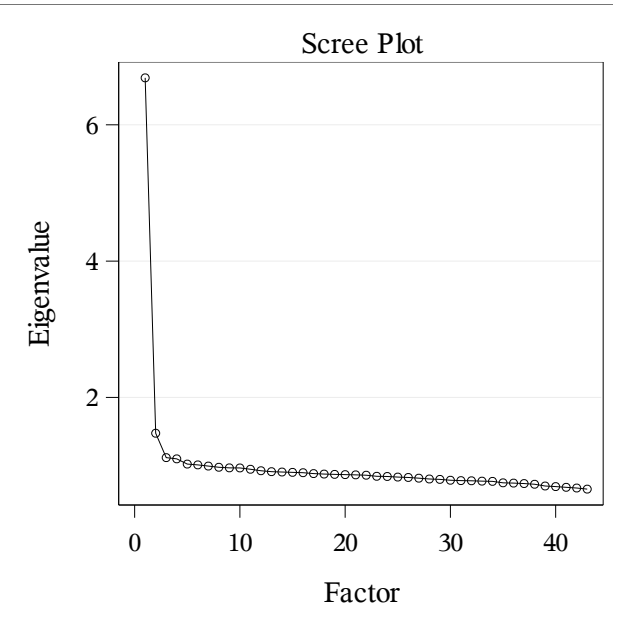
LEAP Social Studies Online: Grade 3



LEAP Social Studies Paper: Grade 3

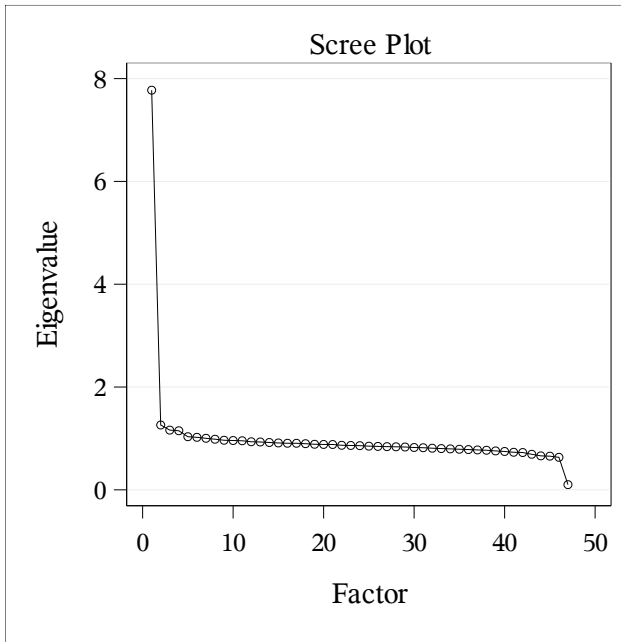


LEAP Social Studies Online: Grade 4

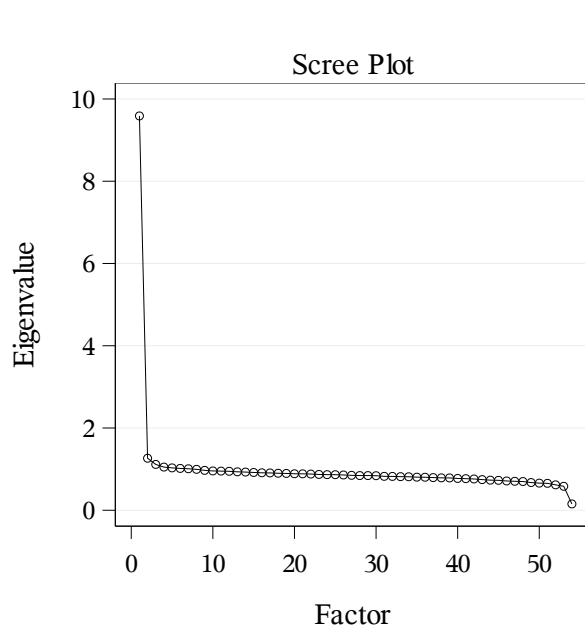


LEAP Social Studies Paper: Grade 4

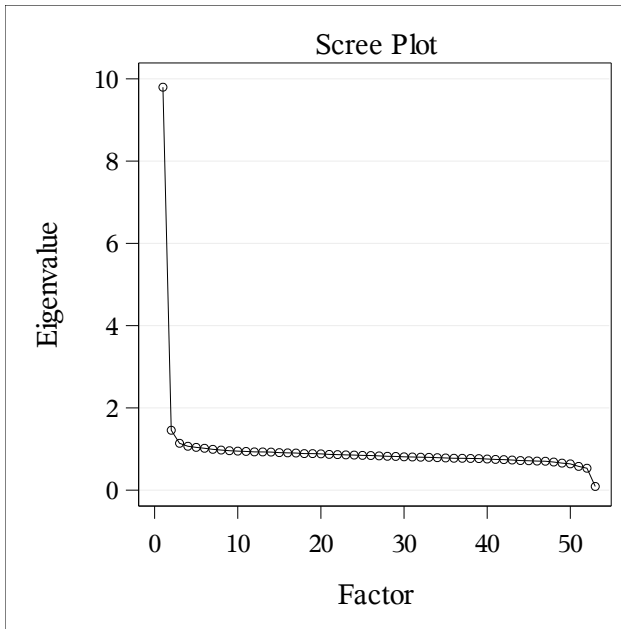
Figure C.1.2  
Principal Component Analysis Plot: Spring 2021 Operational Social Studies Grades 5–8



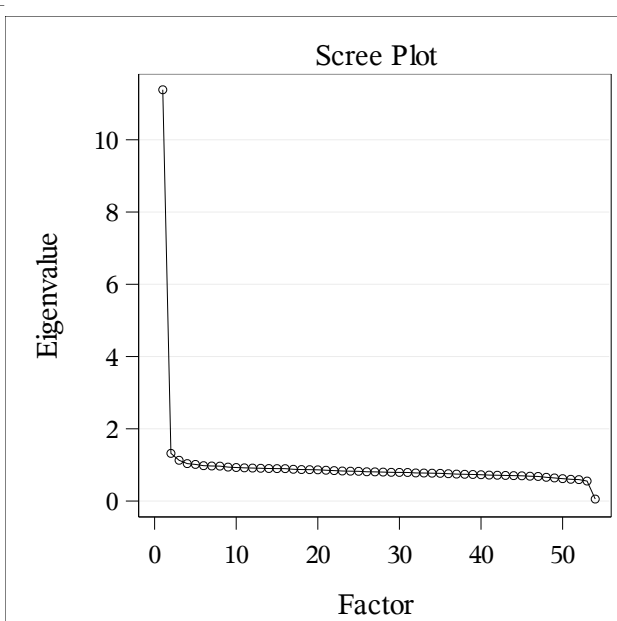
LEAP Social Studies Online: Grade 5



LEAP Social Studies Online: Grade 6



LEAP Social Studies Online: Grade 7



LEAP Social Studies Online: Grade 8

# Appendix D: Scale Distribution and Statistical Report

Contents
Table D.1.1 Scale Score Descriptive Statistics and Plots for Spring 2021 Social Studies
Table D.1.2 Frequency Distribution of Scale Scores for Spring 2021 Social Studies

- Because the spring 2021 test was administered during the COVID-19 pandemic, great caution should be applied when any statistical inference is drawn.



Table D.1.1 Scale Score Descriptive Statistics and Plots for Spring 2021 Social Studies Grade 3

DESCRIPTIVE STATISTICS - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 03

N	≥49530	Median	714.00
Mean	714.18	Variance	1632.46
Std deviation	40.40	Kurtosis	-0.5740
Skewness	0.1327	Std Error Mean	0.1815
Mode	650.00	Interquartile Range	56.00
Range	200.00		

Quantile	Estimate
100% Max	850
99%	810
95%	781
90%	768
75% Q3	744
50% Median	714
25% Q1	688
10%	653
5%	650
1%	650
0% Min	650

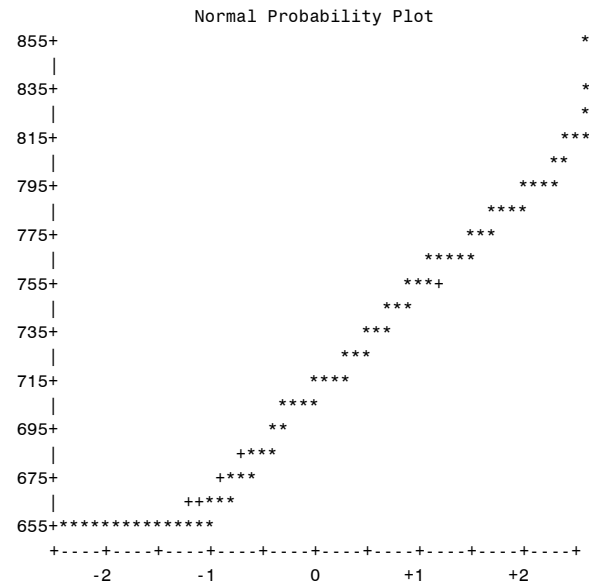
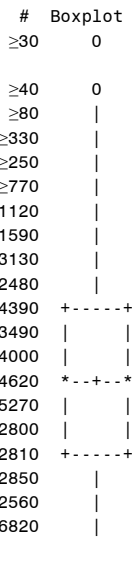
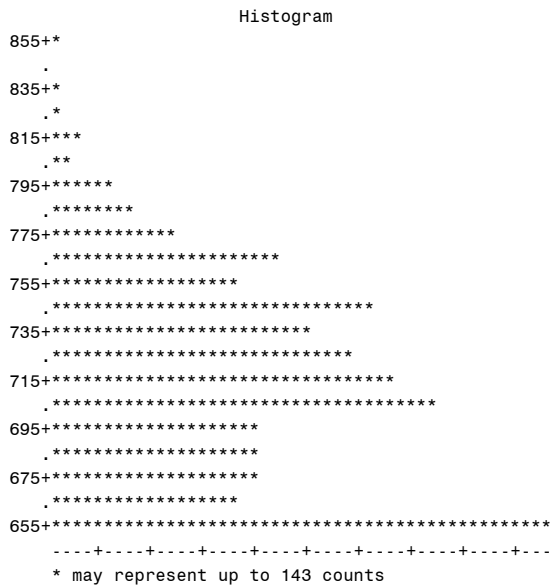


Table D.1.2 Frequency Distribution of Scale Scores for Spring 2021 Social Studies Grade 3

FREQUENCY DISTRIBUTION - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 03

SCALE_SCORE		Freq	Cum. Freq	Percent	Cum. Percent
650	*****	≥4600	≥4600	9.29	9.29
653	*****	≥2220	≥6820	4.50	13.78
668	*****	≥2560	≥9390	5.17	18.96
679	*****	≥2850	≥12240	5.77	24.72
688	*****	≥2810	≥15060	5.68	30.41
696	*****	≥2800	≥17860	5.66	36.07
703	*****	≥2720	≥20580	5.49	41.56
708	*****	≥2550	≥23130	5.15	46.71
714	*****	≥2350	≥25490	4.76	51.47
719	*****	≥2260	≥27760	4.57	56.04
724	*****	≥2060	≥29820	4.17	60.21
728	*****	≥1940	≥31770	3.92	64.13
732	*****	≥1770	≥33540	3.59	67.72
736	*****	≥1720	≥35270	3.47	71.20
740	*****	≥1610	≥36880	3.25	74.45
744	*****	≥1460	≥38340	2.95	77.40
748	*****	≥1320	≥39660	2.68	80.08
752	*****	≥1320	≥40980	2.66	82.74
756	*****	≥1160	≥42150	2.35	85.10
760	*****	≥1130	≥43290	2.30	87.39
764	*****	≥1030	≥44330	2.10	89.49
768	*****	≥950	≥45280	1.93	91.42
772	*****	≥870	≥46160	1.76	93.19
776	*****	≥720	≥46880	1.46	94.65
781	*****	≥600	≥47490	1.22	95.87
786	*****	≥520	≥48010	1.05	96.93
791	*****	≥420	≥48430	0.85	97.78
797	*****	≥350	≥48780	0.71	98.48
803	***	≥250	≥49030	0.51	98.99
810	***	≥190	≥49220	0.38	99.37
818	**	≥140	≥49360	0.28	99.66
827	*	≥80	≥49450	0.17	99.83
839	*	≥40	≥49500	0.09	99.93
850		≥30	≥49530	0.07	100.00

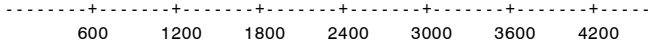


Table D.1.3 Scale Score Descriptive Statistics and Plots for Spring 2021 Social Studies Grade 4

DESCRIPTIVE STATISTICS - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 04

N	≥49540	Median	717.00
Mean	716.06	Variance	1620.96
Std deviation	40.26	Kurtosis	-0.5899
Skewness	0.0488	Std Error Mean	0.1809
Mode	650.00	Interquartile Range	57.00
Range	200.00		

Quantile	Estimate
100% Max	850
99%	808
95%	777
90%	768
75% Q3	746
50% Median	717
25% Q1	689
10%	650
5%	650
1%	650
0% Min	650

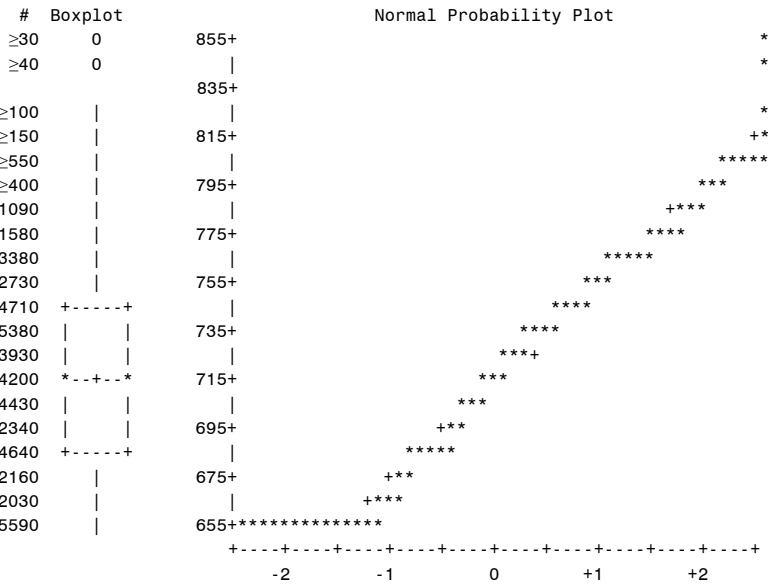
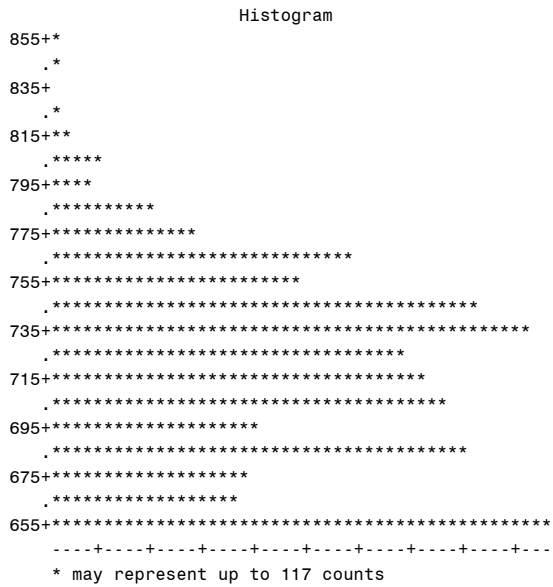


Table D.1.4 Frequency Distribution of Scale Scores for Spring 2021 Social Studies Grade 4

FREQUENCY DISTRIBUTION - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 04

SCALE_SCORE	Freq	Cum. Freq	Percent	Cum. Percent
650	≥5590	≥5590	11.29	11.29
662	≥2030	≥7620	4.10	15.39
673	≥2160	≥9780	4.36	19.76
681	≥2270	≥12060	4.59	24.34
689	≥2360	≥14420	4.78	29.12
696	≥2340	≥16760	4.72	33.85
702	≥2240	≥19010	4.53	38.38
707	≥2190	≥21200	4.43	42.80
712	≥2100	≥23310	4.25	47.05
717	≥2100	≥25410	4.24	51.29
722	≥2030	≥27450	4.11	55.41
726	≥1900	≥29350	3.84	59.24
730	≥1800	≥31150	3.64	62.88
734	≥1830	≥32990	3.71	66.59
738	≥1740	≥34730	3.51	70.10
742	≥1700	≥36430	3.43	73.53
746	≥1520	≥37950	3.07	76.60
749	≥1490	≥39440	3.02	79.62
753	≥1450	≥40890	2.93	82.55
757	≥1280	≥42170	2.59	85.13
761	≥1240	≥43420	2.51	87.65
764	≥1110	≥44530	2.24	89.89
768	≥1030	≥45560	2.08	91.97
773	≥890	≥46460	1.81	93.78
777	≥680	≥47150	1.39	95.17
782	≥570	≥47720	1.16	96.33
788	≥510	≥48240	1.04	97.37
794	≥400	≥48640	0.82	98.19
801	≥310	≥48960	0.63	98.82
808	≥230	≥49190	0.48	99.30
817	≥150	≥49350	0.32	99.62
828	≥100	≥49460	0.21	99.83
842	≥40	≥49500	0.10	99.93
850	≥30	≥49540	0.07	100.00

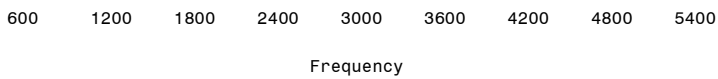
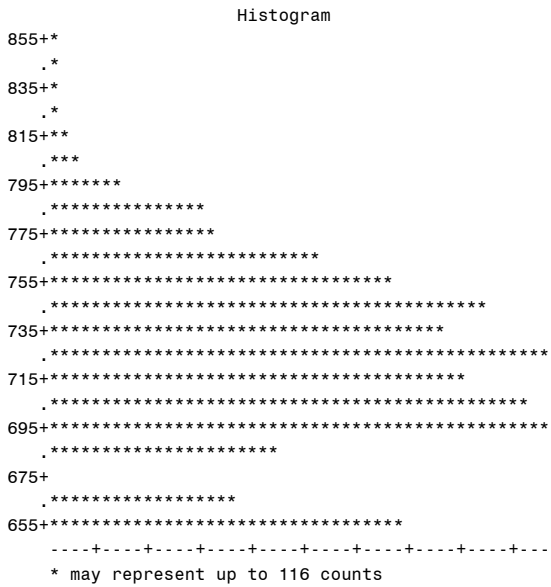


Table D.1.5 Scale Score Descriptive Statistics and Plots for Spring 2021 Social Studies Grade 5

DESCRIPTIVE STATISTICS - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 05

N	≥49850	Median	721.00
Mean	721.21	Variance	1393.00
Std deviation	37.32	Kurtosis	-0.4348
Skewness	-0.0574	Std Error Mean	0.1672
Mode	650.00	Interquartile Range	52.00
Range	200.00		

Quantile	Estimate
100% Max	850
99%	803
95%	780
90%	769
75% Q3	749
50% Median	721
25% Q1	697
10%	667
5%	650
1%	650
0% Min	650



#	Boxplot
<10	0
<10	0
≥20	0
≥70	
≥160	
≥290	
≥720	
≥1630	
≥1810	
≥2980	
≥3770	
≥4750	+---+
≥4340	
≥5560	*--*
≥4550	
≥5290	
≥5470	+---+
≥2480	
≥2020	
≥3860	

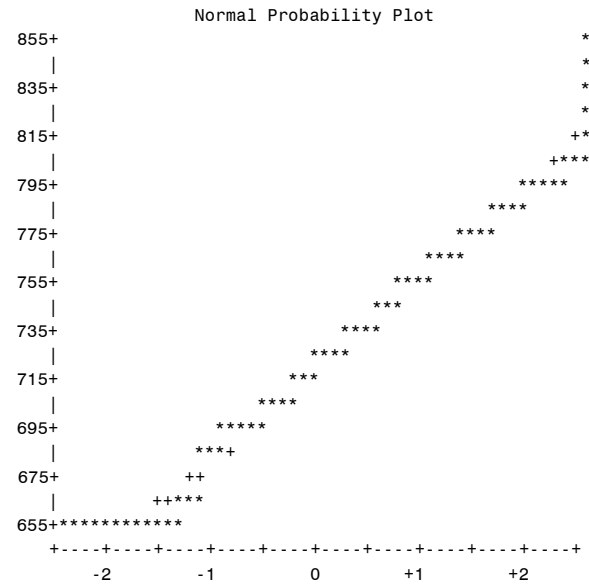


Table D.1.6 *Frequency Distribution of Scale Scores for Spring 2021 Social Studies Grade 5*

FREQUENCY DISTRIBUTION - SCALE SCORES  
Social Studies  
ALL STUDENTS  
GRADE 05

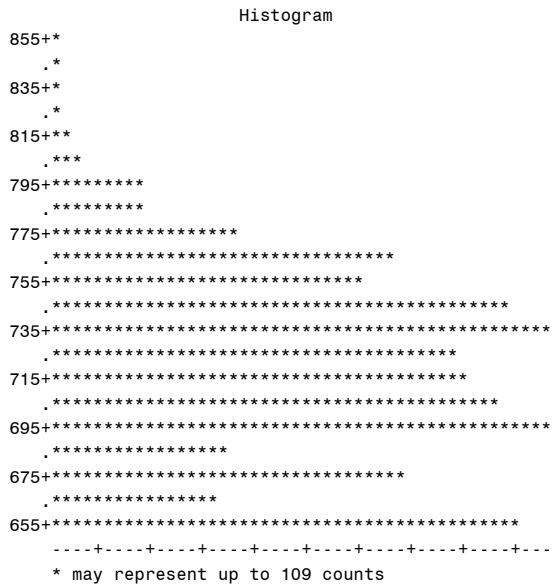
SCALE_SCORE		Freq	Cum. Freq	Percent	Cum. Percent
650	*****	≥3860	≥3860	7.75	7.75
667	*****	≥2020	≥5890	4.07	11.82
680	*****	≥2480	≥8370	4.98	16.80
690	*****	≥2710	≥11080	5.45	22.24
697	*****	≥2760	≥13840	5.54	27.78
703	*****	≥2750	≥16590	5.52	33.30
709	*****	≥2540	≥19140	5.11	38.41
713	*****	≥2350	≥21500	4.72	43.13
717	*****	≥2190	≥23690	4.41	47.54
721	*****	≥2030	≥25720	4.07	51.61
725	*****	≥1910	≥27640	3.84	55.45
728	*****	≥1610	≥29260	3.25	58.70
731	*****	≥1540	≥30800	3.09	61.79
735	*****	≥1420	≥32230	2.87	64.66
738	*****	≥1360	≥33600	2.75	67.40
740	*****	≥1320	≥34920	2.66	70.06
743	*****	≥1230	≥36160	2.48	72.54
746	*****	≥1110	≥37270	2.23	74.77
749	*****	≥1080	≥38360	2.17	76.95
751	*****	≥1010	≥39370	2.03	78.98
754	*****	≥970	≥40340	1.95	80.93
757	*****	≥940	≥41280	1.89	82.81
759	*****	≥850	≥42130	1.71	84.52
762	*****	≥840	≥42980	1.70	86.22
764	*****	≥750	≥43730	1.51	87.73
767	*****	≥750	≥44490	1.51	89.25
769	*****	≥620	≥45120	1.26	90.51
772	*****	≥640	≥45760	1.30	91.81
775	*****	≥590	≥46360	1.19	93.00
778	*****	≥570	≥46930	1.15	94.15
780	*****	≥450	≥47390	0.92	95.06
783	*****	≥440	≥47830	0.89	95.95
786	*****	≥400	≥48230	0.81	96.76
789	*****	≥330	≥48560	0.67	97.42
793	*****	≥290	≥48860	0.59	98.02
796	*****	≥230	≥49090	0.47	98.48
799	****	≥190	≥49290	0.39	98.88
803	****	≥170	≥49470	0.36	99.23
807	**	≥110	≥49580	0.24	99.47
811	**	≥90	≥49670	0.18	99.65
815	*	≥70	≥49750	0.15	99.80
820	*	≥40	≥49790	0.09	99.89
825	*	≥30	≥49820	0.06	99.95
831		≥10	≥49840	0.03	99.97
838		<10	≥49840	0.02	99.99
847		<10	≥49850	0.00	99.99
850		<10	≥49850	0.01	100.00

Table D.1.7 Scale Score Descriptive Statistics and Plots for Spring 2021 Social Studies Grade 6

DESCRIPTIVE STATISTICS - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 06

N	≥51530	Median	719.00
Mean	718.11	Variance	1473.21
Std deviation	38.38	Kurtosis	-0.6123
Skewness	0.0485	Std Error Mean	0.1691
Mode	650.00	Interquartile Range	56.00
Range	200.00		

Quantile	Estimate
100% Max	850
99%	803
95%	779
90%	769
75% Q3	746
50% Median	719
25% Q1	690
10%	666
5%	650
1%	650
0% Min	650



#	Boxplot
<10	0
<10	0
≥50	0
≥50	
≥180	
≥310	
≥920	
≥930	
≥1890	
≥3490	
≥3250	
≥4790	+----+
≥5200	
≥4200	
≥4270	*- -*
≥4670	
≥5180	+----+
≥1840	
≥3690	
≥1730	
≥4800	

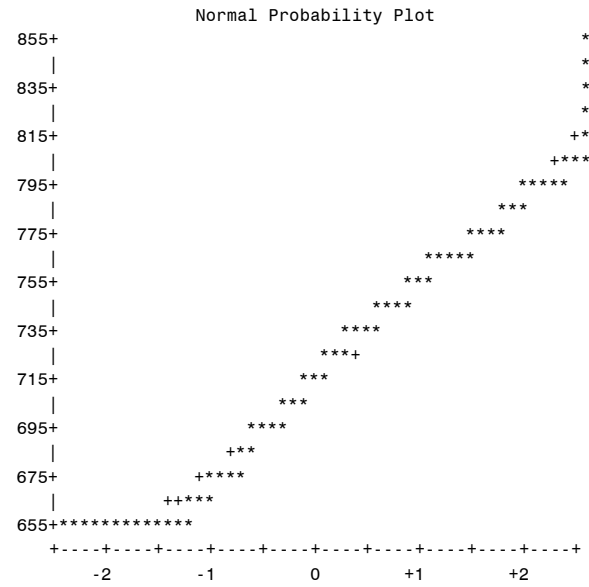


Table D.1.8 Frequency Distribution of Scale Scores for Spring 2021 Social Studies Grade 6

FREQUENCY DISTRIBUTION - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 06

SCALE_SCORE		Freq	Cum. Freq	Percent	Cum. Percent
650	*****	≥3330	≥3330	6.47	6.47
658	*****	≥1460	≥4800	2.85	9.32
666	*****	≥1730	≥6540	3.37	12.69
673	*****	≥1810	≥8350	3.53	16.22
679	*****	≥1870	≥10230	3.63	19.85
685	*****	≥1840	≥12070	3.57	23.42
690	*****	≥1720	≥13790	3.35	26.77
694	*****	≥1750	≥15540	3.40	30.17
698	*****	≥1700	≥17250	3.31	33.48
702	*****	≥1630	≥18890	3.18	36.66
706	*****	≥1580	≥20470	3.07	39.73
709	*****	≥1460	≥21930	2.84	42.56
713	*****	≥1470	≥23410	2.87	45.43
716	*****	≥1420	≥24840	2.77	48.20
719	*****	≥1370	≥26210	2.66	50.86
722	*****	≥1430	≥27640	2.77	53.64
725	*****	≥1380	≥29020	2.68	56.32
727	*****	≥1390	≥30410	2.71	59.03
730	*****	≥1340	≥31760	2.60	61.63
733	*****	≥1370	≥33130	2.67	64.30
735	*****	≥1280	≥34410	2.49	66.79
738	*****	≥1200	≥35620	2.34	69.13
741	*****	≥1230	≥36850	2.39	71.52
744	*****	≥1200	≥38060	2.33	73.85
746	*****	≥1200	≥39260	2.34	76.19
749	*****	≥1150	≥40410	2.24	78.43
752	*****	≥1120	≥41540	2.18	80.61
754	*****	≥1100	≥42640	2.14	82.75
757	*****	≥1020	≥43660	1.99	84.73
760	*****	≥950	≥44620	1.85	86.59
763	*****	≥910	≥45530	1.77	88.36
766	*****	≥810	≥46340	1.58	89.94
769	*****	≥810	≥47160	1.57	91.51
772	*****	≥700	≥47860	1.37	92.88
776	*****	≥630	≥48500	1.23	94.11
779	*****	≥550	≥49050	1.07	95.19
783	*****	≥480	≥49530	0.93	96.12
786	*****	≥440	≥49980	0.87	96.99
790	*****	≥390	≥50370	0.76	97.76
794	*****	≥310	≥50690	0.61	98.37
798	****	≥210	≥50910	0.42	98.79
803	****	≥170	≥51080	0.34	99.13
807	***	≥130	≥51220	0.27	99.40
812	**	≥100	≥51330	0.21	99.61
818	**	≥70	≥51410	0.15	99.76
823	*	≥50	≥51470	0.11	99.88
830	*	≥30	≥51500	0.06	99.94
837		≥20	≥51520	0.04	99.97
845		<10	≥51530	0.02	99.99
850		<10	≥51530	0.01	100.00

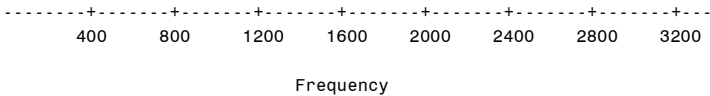


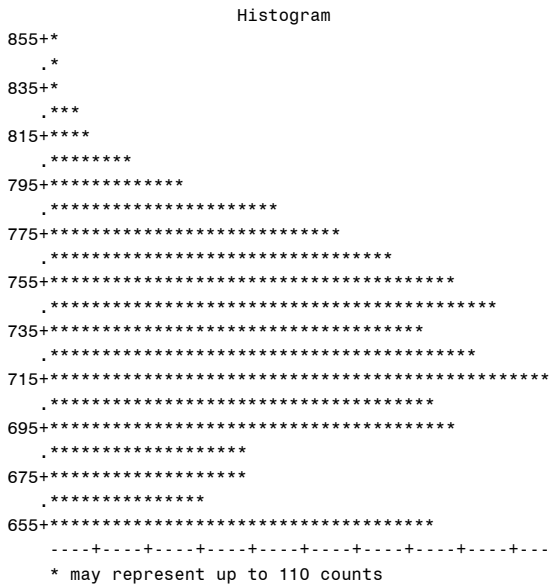


Table D.1.9 Scale Score Descriptive Statistics and Plots for Spring 2021 Social Studies Grade 7

DESCRIPTIVE STATISTICS - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 07

N	≥52310	Median	725.00
Mean	726.21	Variance	1664.81
Std deviation	40.80	Kurtosis	-0.5710
Skewness	0.0131	Std Error Mean	0.1784
Mode	650.00	Interquartile Range	58.00
Range	200.00		

Quantile	Estimate
100% Max	850
99%	816
95%	792
90%	781
75% Q3	755
50% Median	725
25% Q1	697
10%	667
5%	650
1%	650
0% Min	650



#	Boxplot
≥40	0
≥10	0
≥100	
≥230	
≥350	
≥850	
≥1330	
≥2330	
≥3030	
≥3540	
≥4210	+ - - - +
≥4690	
≥3920	
≥4490	* - - + *
≥5230	
≥3960	
≥4180	+ - - - +
≥2060	
≥2010	
≥1640	
≥4040	

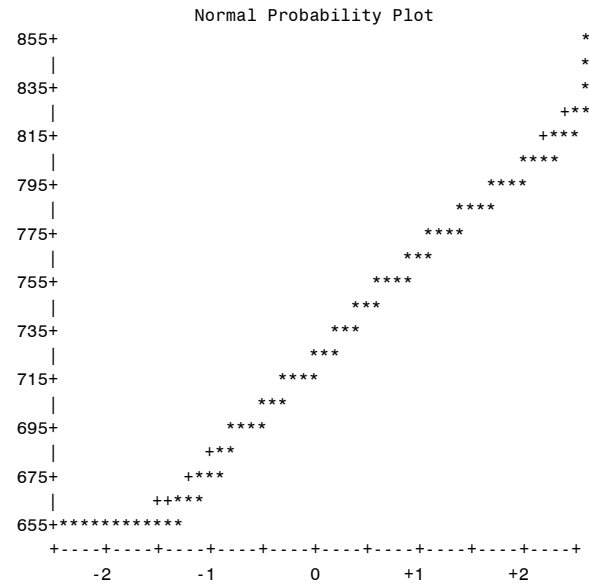


Table D.1.10 Frequency Distribution of Scale Scores for Spring 2021 Social Studies Grade 7

FREQUENCY DISTRIBUTION - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 07

SCALE_SCORE		Freq	Cum. Freq	Percent	Cum. Percent
650	*****	≥2660	≥2660	5.09	5.09
656	*****	≥1380	≥4040	2.64	7.73
667	*****	≥1640	≥5680	3.13	10.87
677	*****	≥2010	≥7690	3.85	14.71
684	*****	≥2060	≥9760	3.95	18.67
691	*****	≥2120	≥11880	4.05	22.72
697	*****	≥2060	≥13950	3.95	26.67
702	*****	≥2000	≥15950	3.83	30.50
706	*****	≥1960	≥17910	3.75	34.25
711	*****	≥1820	≥19740	3.49	37.73
715	*****	≥1760	≥21500	3.38	41.11
718	*****	≥1640	≥23140	3.14	44.25
722	*****	≥1580	≥24730	3.03	47.28
725	*****	≥1470	≥26200	2.82	50.10
728	*****	≥1430	≥27640	2.74	52.84
731	*****	≥1330	≥28970	2.54	55.38
734	*****	≥1270	≥30240	2.43	57.81
737	*****	≥1310	≥31560	2.52	60.33
740	*****	≥1250	≥32810	2.39	62.73
743	*****	≥1220	≥34030	2.33	65.06
745	*****	≥1130	≥35160	2.16	67.22
748	*****	≥1080	≥36250	2.08	69.30
750	*****	≥1070	≥37320	2.05	71.36
753	*****	≥1100	≥38420	2.10	73.46
755	*****	≥1050	≥39480	2.02	75.48
758	*****	≥980	≥40470	1.88	77.36
760	*****	≥940	≥41410	1.80	79.17
763	*****	≥910	≥42320	1.74	80.91
765	*****	≥870	≥43190	1.67	82.57
768	*****	≥810	≥44010	1.55	84.13
770	*****	≥830	≥44840	1.60	85.73
773	*****	≥800	≥45650	1.53	87.26
775	*****	≥700	≥46350	1.35	88.61
778	*****	≥680	≥47040	1.31	89.92
781	*****	≥670	≥47710	1.29	91.21
783	*****	≥610	≥48320	1.17	92.38
786	*****	≥520	≥48850	1.01	93.39
789	*****	≥510	≥49370	0.98	94.37
792	*****	≥530	≥49900	1.02	95.39
795	*****	≥420	≥50320	0.81	96.20
798	*****	≥380	≥50700	0.73	96.93
801	*****	≥320	≥51030	0.63	97.56
804	*****	≥270	≥51300	0.52	98.08
808	*****	≥250	≥51560	0.49	98.57
812	****	≥200	≥51760	0.38	98.95
816	***	≥150	≥51910	0.29	99.24
821	***	≥140	≥52060	0.27	99.52
826	**	≥80	≥52140	0.17	99.68
831	*	≥50	≥52200	0.11	99.79
838	*	≥40	≥52250	0.09	99.88
847		≥10	≥52260	0.04	99.91
850		≥40	≥52310	0.09	100.00

-----+-----+-----+-----+-----+-----+-----  
 400 800 1200 1600 2000 2400

Frequency

Table D.1.11 Scale Score Descriptive Statistics and Plots for Spring 2021 Social Studies Grade 8

DESCRIPTIVE STATISTICS - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 08

N	≥51830	Median	733.00
Mean	730.88	Variance	1576.07
Std deviation	39.70	Kurtosis	-0.5051
Skewness	-0.1235	Std Error Mean	0.1744
Mode	650.00	Interquartile Range	57.00
Range	200.00		

Quantile	Estimate
100% Max	850
99%	815
95%	793
90%	779
75% Q3	760
50% Median	733
25% Q1	703
10%	676
5%	662
1%	650
0% Min	650

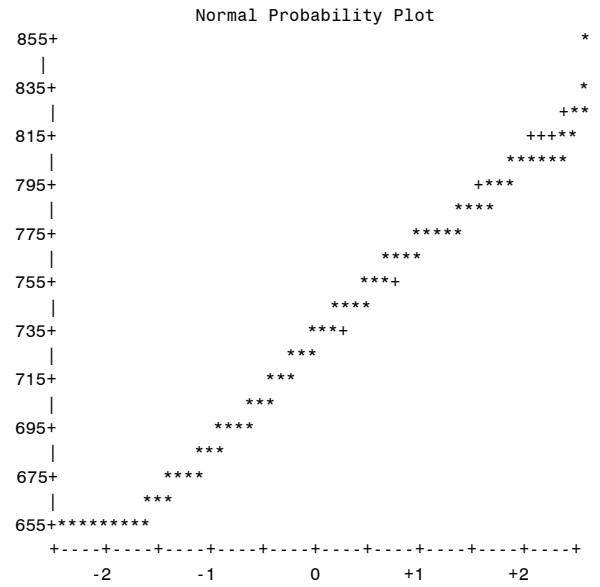
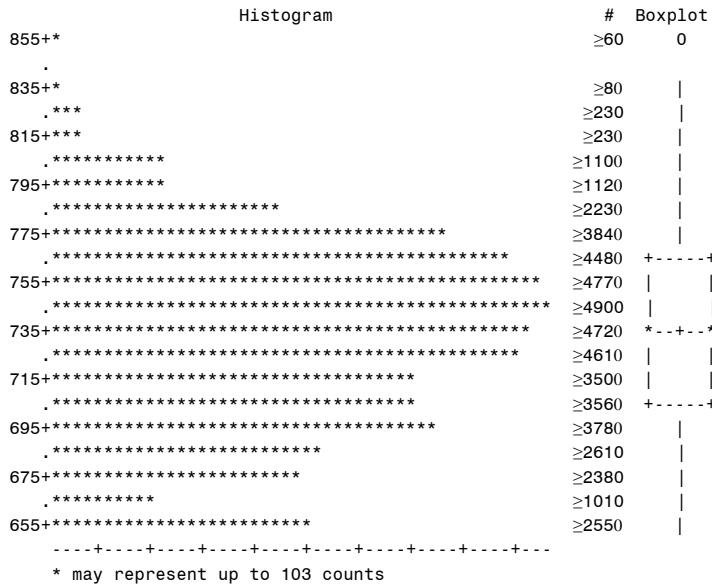


Table D.1.12 Frequency Distribution of Scale Scores for Spring 2021 Social Studies Grade 8

FREQUENCY DISTRIBUTION - SCALE SCORES  
 Social Studies  
 ALL STUDENTS  
 GRADE 08

SCALE_SCORE		Freq	Cum. Freq	Percent	Cum. Percent
650	*****	≥1690	≥1690	3.26	3.26
653	*****	≥860	≥2550	1.67	4.93
662	*****	≥1010	≥3560	1.95	6.88
670	*****	≥1110	≥4680	2.16	9.03
676	*****	≥1260	≥5940	2.43	11.47
682	*****	≥1330	≥7270	2.57	14.04
687	*****	≥1270	≥8550	2.46	16.50
691	*****	≥1290	≥9840	2.49	19.00
695	*****	≥1230	≥11080	2.39	21.38
699	*****	≥1250	≥12340	2.43	23.81
703	*****	≥1190	≥13530	2.30	26.11
706	*****	≥1180	≥14710	2.28	28.39
709	*****	≥1190	≥15900	2.30	30.69
712	*****	≥1190	≥17090	2.30	32.98
715	*****	≥1100	≥18200	2.13	35.12
718	*****	≥1210	≥19410	2.34	37.45
721	*****	≥1130	≥20550	2.20	39.65
723	*****	≥1150	≥21700	2.22	41.87
726	*****	≥1150	≥22860	2.23	44.10
728	*****	≥1160	≥24020	2.25	46.35
731	*****	≥1170	≥25200	2.26	48.61
733	*****	≥1140	≥26350	2.22	50.83
736	*****	≥1210	≥27560	2.35	53.18
738	*****	≥1180	≥28750	2.28	55.46
740	*****	≥1210	≥29960	2.34	57.80
743	*****	≥1270	≥31230	2.46	60.26
745	*****	≥1180	≥32420	2.29	62.55
748	*****	≥1220	≥33650	2.36	64.91
750	*****	≥1190	≥34840	2.31	67.22
753	*****	≥1250	≥36100	2.42	69.65
755	*****	≥1160	≥37270	2.25	71.90
758	*****	≥1150	≥38420	2.23	74.13
760	*****	≥1140	≥39570	2.20	76.33
763	*****	≥1130	≥40700	2.19	78.53
765	*****	≥1140	≥41840	2.20	80.73
768	*****	≥1060	≥42910	2.05	82.78
771	*****	≥1000	≥43920	1.95	84.73
774	*****	≥1000	≥44930	1.95	86.68
776	*****	≥960	≥45890	1.87	88.54
779	*****	≥860	≥46750	1.66	90.20
783	*****	≥850	≥47610	1.66	91.86
786	*****	≥740	≥48360	1.44	93.30
789	*****	≥620	≥48990	1.21	94.51
793	*****	≥620	≥49610	1.20	95.71
796	*****	≥490	≥50110	0.96	96.68
800	*****	≥440	≥50550	0.86	97.53
805	*****	≥360	≥50910	0.69	98.23
809	*****	≥300	≥51220	0.58	98.81
815	*****	≥230	≥51450	0.45	99.26
821	*****	≥130	≥51590	0.26	99.52
828	****	≥90	≥51680	0.18	99.70
838	***	≥80	≥51770	0.16	99.87
850	***	≥60	≥51830	0.13	100.00

# Appendix E: Reliability and Classification Accuracy

## *Reliability and Classification Accuracy Reports*

### *Social Studies*

Contents
Table E.1.1 Reliability for Overall and Subgroups: Spring 2021 Operational Social Studies
Table E.2.1 Cronbach's Alpha Reliability: Spring 2021 Operational Social Studies
Table E.3.1 Classification Accuracy and Decision Consistency: Spring 2021 Operational Social Studies

- Because the spring 2021 test was administered during the COVID-19 pandemic, great caution should be applied when any statistical inference is drawn.

Table E.1.1

*Reliability for Overall and Subgroups: Spring 2021 Operational Social Studies*

<b>Subgroup</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>
All Students	0.840	0.860	0.883	0.906	0.910	0.926
Female	0.832	0.846	0.877	0.900	0.904	0.920
Male	0.848	0.872	0.889	0.912	0.917	0.931
African American	0.781	0.805	0.828	0.872	0.875	0.907
American Indian or Alaska Native	0.841	0.816	0.857	0.890	0.905	0.920
Asian	0.859	0.873	0.915	0.915	0.929	0.929
Hispanic/Latino	0.820	0.843	0.873	0.907	0.907	0.927
Two or More Races	0.836	0.855	0.883	0.901	0.909	0.924
Native Hawaiian or Other Pacific Islander	0.692	0.836	N/A	0.908	0.928	0.911
White	0.846	0.860	0.889	0.903	0.914	0.921
Economically Disadvantaged: No	0.850	0.859	0.890	0.902	0.913	0.915
Economically Disadvantaged: Yes	0.805	0.831	0.853	0.890	0.891	0.915
EL: No	0.841	0.861	0.883	0.906	0.911	0.925
EL: Yes	0.711	0.728	0.743	0.834	0.816	0.877
Regular Education	0.841	0.859	0.882	0.904	0.909	0.922
Special Education	0.810	0.823	0.827	0.852	0.851	0.887
Section 504 Status: No	0.841	0.862	0.885	0.907	0.911	0.926
Section 504 Status: Yes	0.815	0.824	0.838	0.877	0.878	0.911

N/A means that estimate was not calculated since the *n* count is smaller than 30.

Table E.2.1

*Cronbach's Alpha Reliability: Spring 2021 Operational Social Studies*

<b>Grade</b>	<b>Cronbach's Alpha</b>
3	0.840
4	0.860
5	0.883
6	0.906
7	0.910
8	0.926

**Table E.3.1*****Classification Accuracy and Decision Consistency: Spring 2021 Operational Social Studies***

Table E.3.1.1

*Estimates of Accuracy and Consistency of Achievement-Level Classification for Each Grade*

<b>Grade</b>	<b>Accuracy</b>	<b>Consistency</b>	<b>PChance</b>	<b>Kappa</b>
3	0.625	0.530	0.245	0.377
4	0.656	0.558	0.246	0.413
5	0.662	0.565	0.241	0.427
6	0.697	0.602	0.234	0.481
7	0.686	0.591	0.224	0.473
8	0.716	0.618	0.218	0.512

Table E.3.1.2

*Accuracy of Classification at Each Achievement Level for Each Grade*

<b>Grade</b>	<b>Unsatisfactory (1)</b>	<b>Approaching Basic (2)</b>	<b>Basic (3)</b>	<b>Mastery (4)</b>	<b>Advanced (5)</b>
3	0.837	0.640	0.451	0.463	N/A
4	0.853	0.635	0.544	0.538	N/A
5	0.840	0.617	0.590	0.567	N/A
6	0.885	0.654	0.648	0.517	0.647
7	0.874	0.532	0.597	0.598	0.731
8	0.884	0.626	0.651	0.688	0.701

N/A indicates that it was inestimable due to restricted sample size.



Table E.3.1.3

*Accuracy of Dichotomous Categorizations for Each Grade (PAC Metric)*

<b>Grade</b>	<b>1 / 2+3+4+5</b>	<b>1+2 / 3+4+5</b>	<b>1+2+3 / 4+5</b>	<b>1+2+3+4 / 5</b>
3	0.905	0.867	0.882	0.947
4	0.917	0.881	0.882	0.963
5	0.917	0.885	0.900	0.951
6	0.930	0.906	0.915	0.942
7	0.925	0.905	0.911	0.934
8	0.944	0.919	0.914	0.936

Table E.3.1.4

*Consistency of Dichotomous Categorizations for Each Grade (PAC Metric)*

<b>Grade</b>	<b>1 / 2+3+4+5</b>	<b>1+2 / 3+4+5</b>	<b>1+2+3 / 4+5</b>	<b>1+2+3+4 / 5</b>
3	0.863	0.819	0.834	0.924
4	0.880	0.836	0.835	0.949
5	0.881	0.844	0.857	0.938
6	0.899	0.869	0.879	0.922
7	0.893	0.869	0.875	0.907
8	0.919	0.886	0.880	0.912

Table E.3.1.5

*Kappa of Dichotomous Categorizations for Each Grade (PAC Metric)*

<b>Grade</b>	<b>1 / 2+3+4+5</b>	<b>1+2 / 3+4+5</b>	<b>1+2+3 / 4+5</b>	<b>1+2+3+4 / 5</b>
3	0.691	0.625	0.512	0.144
4	0.721	0.669	0.526	-0.108
5	0.719	0.688	0.619	0.042
6	0.766	0.736	0.658	0.367
7	0.754	0.738	0.712	0.547
8	0.783	0.768	0.740	0.499

Table E.3.1.6

*Accuracy of Dichotomous Categorizations: False Positive Rates for Each Grade (PAC Metric)*

<b>Grade</b>	<b>1 / 2+3+4+5</b>	<b>1+2 / 3+4+5</b>	<b>1+2+3 / 4+5</b>	<b>1+2+3+4 / 5</b>
3	0.050	0.057	0.054	0.054
4	0.043	0.052	0.049	0.037
5	0.046	0.049	0.042	0.049
6	0.035	0.046	0.038	0.043
7	0.039	0.045	0.041	0.038
8	0.028	0.040	0.038	0.038

Table E.3.1.7

*Accuracy of Dichotomous Categorizations: False Negative Rates for Each Grade (PAC Metric)*

<b>Grade</b>	<b>1 / 2+3+4+5</b>	<b>1+2 / 3+4+5</b>	<b>1+2+3 / 4+5</b>	<b>1+2+3+4 / 5</b>
3	0.045	0.075	0.064	N/A
4	0.040	0.066	0.068	N/A
5	0.037	0.066	0.058	N/A
6	0.036	0.048	0.047	0.015
7	0.036	0.050	0.047	0.028
8	0.029	0.041	0.048	0.026

N/A indicates that it was inestimable due to restricted sample size.