

New York State Testing Program English Language Arts Grade 8

Technical Report 2002



Developed and published under contract with New York State Department of Education by CTB/McGraw-Hill LLC, a subsidiary of The McGraw-Hill Companies, Inc., 20 Ryan Ranch Road, Monterey, California 93940-5703. Copyright © 2003 by New York State Department of Education. Only State of New York educators and citizens may copy, download, and/or print the document located online at <http://www.emsc.nysed.gov/ciai/testing/pubs.html>. Any other use or reproduction of this document, in whole or in part, requires written permission of New York State Department of Education.

Foreword

The technical information herein is 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 described in *Standards for educational and psychological testing* (American Educational Research Association, American Psychological Association, & National Council on Measurement in Education, 1999).

Table of Contents

FOREWORD	1
TABLE OF CONTENTS	2
LIST OF TABLES	3
PART 1: TEST DESIGN	4
THE NEW YORK STATE LEARNING STANDARDS FOR ENGLISH LANGUAGE ARTS	4
TEST CONFIGURATION.....	4
<i>Session 1</i>	4
<i>Session 2</i>	4
<i>Writing Mechanics Score</i>	5
STUDENT PARTICIPATION AND TESTING ACCOMMODATIONS	6
<i>Students to be Tested</i>	6
<i>Testing Accommodations</i>	7
ITEM DEVELOPMENT	7
ITEM REVIEW PROCESS	8
<i>Documenting Content</i>	8
<i>Minimizing Bias</i>	8
<i>Minimizing Speededness</i>	8
TEST CONSTRUCTION AND PRE-EQUATING.....	9
<i>Calibration Samples</i>	9
<i>Answer Choice Information</i>	9
<i>Item Response Theory Models</i>	9
<i>Equating Method</i>	11
<i>Item Selection Criteria and Process</i>	11
DIFFERENTIAL ITEM FUNCTIONING.....	12
<i>Procedures for Eliminating Bias and Minimizing Differential Item Functioning</i>	12
PART 2: ITEM STATISTICS FOR THE OPERATIONAL DATA	13
DATA CLEANING	13
ITEM ANALYSIS.....	14
DIFFERENTIAL ITEM FUNCTIONING ANALYSIS OF OPERATIONAL DATA.....	16
PART 3: SCORING AND RELIABILITY	18
RAW SCORE TO SCALE SCORE CONVERSION	18
RELIABILITY	18
ESTIMATED CONDITIONAL STANDARD ERRORS OF SCALE SCORES	19
LOWEST AND HIGHEST OBTAINABLE SCALE SCORES	20
INTER-RATER AGREEMENT	20
EXPECTED SPI SCORES ON THE STANDARDS AT THE DECISION POINTS	21
PART 4: DESCRIPTIVE STATISTICS	24
SCALE-SCORE FREQUENCY DISTRIBUTIONS FOR THE STATE AND SUBGROUPS.....	24
G8 ELA SCALE SCORE MEANS AND STANDARD DEVIATIONS	25
REFERENCES	26

List of Tables

Table 1 New York State Learning Standards for English Language Arts	4
Table 2 Points per item type for GD 8 ELA scores.....	5
Table 3 Condition Codes for the ELA CR items.....	5
Table 4 Steps involved in data clean-up for analysis preparation	13
Table 5 G8 ELA item level statistics.....	15
Table 6 Number of Students in each Gender or Ethnic Group.....	16
Table 7 The numbers of items flagged for DIF in G8 ELA summary	17
Table 8 Raw Score to Scale Score with SEM for G8 ELA 2002	19
Table 9 G8 ELA Inter-rater agreement	21
Table 10 Percentages of inter-rater score differences	21
Table 11 Reliability indices of hand scoring	21
Table 12 G8 ELA standard performance index information	22
Table 13 G8 ELA summary of scale score information	24
Table 14 G8 ELA statewide scale score information	25

Part 1: Test Design

The New York State Learning Standards for English Language Arts

The New York State *Learning Standards for English Language Arts* is available from the New York State Education Department web site, at <http://www.emsc.nysed.gov/ciai/ela/pub/elalearn.pdf>. The four learning standards are listed in Table 1 below. The G8 ELA is written to test students in Standards 1, 2, and 3.

Table 1 New York State Learning Standards for English Language Arts

*Standard 1	Students will read, write, listen, and speak for information and understanding.
*Standard 2	Students will read, write, listen, and speak for literary response and expression.
*Standard 3	Students will read, write, listen, and speak for critical analysis and evaluation.
Standard 4	Students will read, write, listen, and speak for social interaction.

Test Configuration

Similar to the 1999, 2000, and 2001 forms, the 2002 G8 ELA test has the following configuration. The test is divided into two sessions. There are 25 multiple choice (MC) items worth a total of 25 points; there are nine constructed response (CR) items worth a total of 18 points. The CR items may be short response (SR) or extended response (ER) items. The total number of items on the test is 34, and the maximum raw score total is 43 points.

Session 1

Session 1 is comprised of 25 MC items, together with 3 SR items and 1 ER item. Each item in session 1 addresses one of the three tested New York State Learning Standards for English Language Arts. The four CR items in session 1 follow a listening passage, and these items make up the listening cluster. These items are scored together to derive a listening score, which can range from zero to six points.

Session 2

Session 2 contains linked information stimuli, accompanied by three SR items and one ER item, which are scored together to derive a reading cluster score (zero to six points), which addresses Standard 3. Session 2 also contains an independent writing prompt addressing Standard 1. The prompt is followed by an ER item, which is scored to derive an independent writing score (zero to three points).

Writing Mechanics Score

As part of the ELA test, the three ER responses across sessions 1 and 2 are scored together to derive a writing mechanics cluster (zero to three points). Although the writing mechanics is not linked to any of the New York State Learning Standards for English Language Arts, it contributes to the overall ELA score. Table 2 shows the numbers of score points by the item type or cluster, and the total numbers of items and clusters, for the grade 8 ELA test.

Table 2 shows the numbers of score points by the item type or cluster, and the total numbers of items and clusters, for the grade 8 ELA test.

Table 2 Points per item type for GD 8 ELA scores

Item type or cluster	Grade 8
	ELA
Multiple choice (MC)	25 pts
Listening cluster	6 pts
Reading cluster	6 pts
Independent writing item	3 pts
Writing mechanics cluster	3 pts
Total points	43 pts
Total items or clusters	29 items

Table 3 Condition Codes for the ELA CR items

Condition code	Meaning
A	Blank
B	Refusal
C	Insufficient to score
D	Illegible
E	Other language

In scaling and scoring, each of the clusters is treated as a constructed response (CR) item. The following condition codes were used in scoring the responses to the CR items:

Student Participation and Testing Accommodations

Students to be Tested

It is the policy of the New York State Department of Education (NYSED) that the NYSTP grade 8 ELA test be administered to all public school students. Nonpublic schools are strongly encouraged to administer the tests. The exceptions noted below, which represent the policy of the NYSED, apply to students in public and nonpublic schools participating in the NYSTP.

Students with Disabilities

All students with disabilities must be provided full access to the tests, to the extent that such testing is consistent with their individual needs. The Committee on Special Education (CSE) must decide for each student on a case-by-case basis and document, on the individual student's individualized education program (IEP), whether the student will participate in the NYSTP or in the New York State Alternate Assessment for Students with Severe Disabilities. Criteria that the CSE must use to determine if a student should participate in the Alternate Assessment are given in the NYSTP School Administrator's Manual (SAM).

Students in Ungraded Classes

Both students with disabilities and general education students who are in ungraded classes are required to take the NYSTP tests, unless they meet the criteria for LEP exemption or are eligible for the Alternate Assessment. The chronological ages of students in ungraded classes should be used to determine who must be tested. Ungraded students should be tested on the grade 8 assessments no later than the school year (July 1 - June 30) in which they reach their 15th birthday.

Limited English Proficient (LEP) Students

LEP students scoring at or above the 30th percentile on a norm-referenced English reading test or the publisher's recommended score on an approved measure of English as a Second Language (ESL) in reading are required to participate fully in the ELA test. Those students must take the ELA test in English.

LEP students scoring below the 30th percentile on a norm-reference English reading test or the publisher's recommended score on an approved measure of ESL in reading may be exempted from taking the English Language Arts examination.

Other Considerations

When determining who will participate in the NYSTP and who will participate in the Alternate Assessment, school administrators must consider those students who attend programs operated by the Board of Cooperative Educational Services (BOCES), or who are in approved private school placements, as well as in any other programs located outside the school district. Students who are absent during the testing administrations should be tested during the designated makeup period.

Testing Accommodations

Accommodations were used in the NYSTP operational tests to provide equal access to assessments for students with disabilities. Such accommodations are used to increase the validity of test scores by offsetting distortions introduced by the disability and retaining the essential features of the assessment. The following represents the policy of the NYSED for the use of testing accommodations.

Students with Disabilities

It is the responsibility of the principal to ensure that the testing accommodations specified in the IEP or Section 504 Accommodation Plan are provided to students with disabilities. Students who have been declassified may continue to be provided testing accommodations if recommended by the local CSE at the time of declassification and in the student's declassification IEP.

Testing accommodations for students with disabilities are discussed in detail in the NYSED's 1995 publication titled *Test Access and Modification for Individuals with Disabilities*.

School administrators are to indicate in writing on the test book whether the student received "deletion of spelling, punctuation, and/or paragraphing requirements" or "use of scribe, tape recorder, word processor, or typewriter." If a student uses a typewriter or word processor, administrators are to staple the printed pages to the test book. For student receiving "deletion of spelling, paragraphing, and/or punctuation," teachers are to cross out and correct misspelled words and/or provide correct paragraphing and/or punctuation. For students using scribes, tape recorder or large type or Braille editions, teachers are to transcribe the students' text onto regular test answer documents and test books exactly as dictated or recorded.

Students Who Incur Disabilities Shortly Before Test Administration

School principals may modify testing procedures for general education students who incur an injury (e.g., a broken arm) or experience the onset of a short- or long-term disability (e.g., epilepsy) sustained or diagnosed within 30 days prior to the administration of State tests. More details are available in the SAM.

Limited English Proficiency

Limited English Proficient (LEP) students are allowed extended test time, and tests may be administered to LEP students individually or in small groups in a separate location. LEP students may use bilingual dictionaries or glossaries when taking State examinations in ELA and mathematics. The bilingual dictionaries and glossaries must not provide definitions or other explanations, only word-for-word translations. In addition, the Listening passage may be read a third time to LEP students taking the ELA tests.

Item Development

A staff of professional item writers researched, collected, and wrote the field test material. All assessment materials were carefully reviewed for content and editorial accuracy. Artists and designers worked with the writers during development to ensure graphic and textual consistency. With assistance from the New York State Department of Education, all test

items were developed to align with the content and measure the State Learning Standards for English Language Arts. Each of the three tested Learning Standards is a reporting category for the English Language Arts tests. Standards Performance Index (SPI) scores are assigned to students on each of these reporting categories.

Item Review Process

Documenting Content

An integral part of the development process was documentation of content using New York State's Learning Standards. All items used on the New York State tests are reviewed for content by both CTB Development staff and by New York State Department of Education staff and New York State teachers. This procedure ensures that items would be sound in content and format, and targeted appropriately to the courses in which the associated concepts are typically taught.

Minimizing Bias

The developers of the NYSTP tests gave careful attention to questions of ethnic, racial, gender, regional, and age bias. All materials were written and reviewed to conform to the company's editorial policies and guidelines for equitable assessment, as well as NYSED's guidelines for item development.

In addition, educators and other stakeholders from different parts of the state reviewed the items from their perspective as members of various ethnic groups. They identified assessment materials that might reflect possible bias in language, subject matter, or representation of people. Their comments and suggestions were considered carefully during the revision and selection of items for the calibration tests. All materials were written to SED specifications and carefully checked by groups of trained New York community participants.

Minimizing Speededness

Test developers also considered speededness in the development of the NYSTP tests. CTB believes that achievement tests should not be speeded; little or no useful instructional information can be obtained from the fact that a student did not finish a test, while a great deal can be learned from student responses to questions. For that reason, sufficient administration time limits were set for the NYSTP tests.

The Research Department at CTB routinely conducts additional speededness analyses based on actual test data. Table 5 shows the omit rates for items on the G8 ELA test. These results provide little evidence of speededness on these tests.

Test Construction and Pre-equating

Calibration Samples

Three field test forms for the NYSTP tests were administered to students in public and private schools across the State in 2001. Each year of field testing, efforts are made to ensure that the sample of students was representative of the state tested population. The 2001 field test items were calibrated and equated to the existing scale. Thus, parameters for these items were already on the appropriate New York State scale (one each for grade 4 ELA, grade 4 Math, grade 8 ELA, and grade 8 Math).

Since these items are calibrated, the pool of available grade 8 ELA items can be used to construct a test form and to produce a raw-score-to-scale-score table for that form. The 2002 operational NYSTP tests were constructed using the items in the pool. What follows is an overview of the analysis of field test data which results in the calibration of items.

Answer Choice Information

Statistical information about student performance is produced for each multiple choice item. Specifically, three statistics are examined for each item: (1) the proportion of students choosing each answer, (2) the point-biserial correlation between the answer choice and the number-correct score on the rest of the test, and (3) omit rates. For each constructed response item, the proportion of students at each score level, omit rates, and p-values are examined. (The p-value for a constructed response item is the mean item score divided by the total number of points obtainable.)

Item Response Theory Models

Although useful, the differences in proportions of correct responses (p-values) limit the degree to which one can compare important characteristics about the test items. Item-response theory (IRT) allows one to make better comparisons among items, even those from different test forms (there were three for each subject area and grade of the NYSTP), by using a common scale for all items (i.e., as if there were a hypothetical test that contained items from all forms). The three-parameter logistic (3PL) model (Lord & Novick, 1968; Lord, 1980) was used to analyze item responses on the multiple choice items. For analysis of the constructed response items, the two-parameter partial credit model (2PPC) (Muraki, 1992; Yen, 1993) was used.

Item response theory is a statistical procedure that takes into account the fact that not all test items are alike and that all items do not provide the same amount of information in determining how much a student knows or can do. Computer programs that implement IRT models use actual students' data to estimate the characteristics of the items on a test -- called "parameters." The parameter estimation process is also called "item calibration."

IRT models typically vary according to the number of parameters estimated. For the New York State tests, three parameters are estimated: The discrimination parameter, the difficulty parameter(s), and, for multiple choice items, the guessing parameter. The discrimination

parameter is an index of how well an item differentiates between high-performing and low-performing students. An item that low-performing students cannot answer correctly, but high-performing students can, will have a high discrimination value. The difficulty parameter is an index of how easy or difficult an item is. The higher the difficulty parameter, the harder the item. The guessing parameter is the probability that a student with very low ability will answer the item correctly.

The estimated parameters are then used to determine weights for the items in computing student scale scores. The scale score (SS) is the basic score for the New York State tests. It is used to derive other scores that describe test performance, such as the four performance levels and the standard-based index scores (SPIs). Scale scores can be obtained by one of two scoring methods: IRT item-pattern scoring, or number-correct scoring. Starting in 2002, scores on the New York State tests are determined using number-correct scoring.

Because the characteristics of MC and CR items are different, two IRT models were used in item calibration. The three-parameter logistic model (3PL) (Lord & Novick, 1968; Lord, 1980) was used in the analysis of MC items. In this model, the probability that a student with ability θ responds correctly to item i is

$$P_i(\theta) = c_i + \frac{1 - c_i}{1 + \exp[-1.7a_i(\theta - b_i)]},$$

where a_i is the item discrimination, b_i is the item difficulty, and c_i is the probability of a correct response by a very low-scoring student.

For analysis of the constructed response items, the two-parameter partial credit model (2PPC) (Muraki, 1992; Yen, 1993) was used. The 2PPC model is a special case of Bock's (1972) nominal model. Bock's model states that the probability of an examinee with ability θ having a score ($k-1$) at the k -th level of the j -th item is

$$P_{jk}(\theta) = P(x_j = k - 1 | \theta) = \frac{\exp Z_{jk}}{\sum_{i=1}^{m_j} \exp Z_{ji}}, \quad k = 1 \dots m_j,$$

where

$$Z_{jk} = A_{jk}\theta + C_{jk}.$$

The m_j denotes the number of score levels for the j -th item, and typically, the highest score level is assigned ($m_j - 1$) score points. For the special case of the 2PPC model used here, the following constraints were used:

$$A_{jk} = \alpha_j(k - 1),$$

and

$$C_{jk} = -\sum_{i=0}^{k-1} \gamma_{ji}, \text{ where } \gamma_{j0} = 0,$$

where α_j and γ_{ji} are the free parameters to be estimated from the data. Each item has $(m_j - 1)$ independent γ_{ji} parameters and one α_j parameter; a total of m_j parameters are estimated for each item.

The IRT model parameters were estimated using CTB's PARDUX software (Burket, 1991). PARDUX estimates parameters simultaneously for MC and CR items using marginal maximum likelihood procedures implemented via the EM (expectation-maximization) algorithm (Bock & Aitkin, 1981; Thissen, 1982).

Simulation studies have compared PARDUX with MULTILOG (Thissen, 1991), PARSCALE (Muraki & Bock, 1991), and BIGSTEPS (Wright & Linacre, 1992). PARSCALE, MULTILOG, and BIGSTEPS are among the most widely known and used IRT programs. PARDUX was found to perform at least as well as these other programs.

Equating Method

After the item calibration, all of the G8 ELA items, when they were field-tested, were placed on the NYS G8 ELA scale using the 2001 operational items as anchors. This was possible, because the operational items were taken by the same students who took the field test items within the same testing window. The equating was performed using the test characteristic curve method (Stocking & Lord, 1983) implemented by PARDUX. In previous years, operational data were used to re-calibrate items and re-equate them. NYSED, however, made a decision in 2002 to use the pre-equating model, which is similar to what is done for the New York State Regents program. This allows for the production of scoring tables (see Part 3) ahead of actual operational administration, once the operational form is selected.

Item Selection Criteria and Process

Item selection for the NYSTP tests was based on the classical and IRT statistics of the test items. Selection was conducted by content experts from CTB and NYSED and reviewed by psychometricians at CTB. Final approval of the items selected was given by NYSED. Two criteria governed the item selection process. The first of these was to meet the content specifications specified by the New York State Department of Education. Within the limits set by these requirements, developers selected from the pool of field test items, those with the best psychometric characteristics. Developers selected items that minimized measurement error throughout the range of expected achievement, as indicated by the reciprocal of the square root of the IRT information function (Lord, 1980, p. 71). Developers aimed to create forms with the content and psychometric properties of previous operational forms.

Item selection for the calibration tests was facilitated using the Windows version of the program ITEMSYS (Burket, 1988). ITEMSYS creates an interactive connection between the developer selecting the test items and the item database. This program monitors the impact of each decision made during the item selection process and offers a variety of options for

grouping, classifying, sorting, and ranking items to highlight key information as it is needed (see Green, Yen, & Burket, 1989).

ITEMSYS has three parts. The first part selects a working item pool of manageable size from the larger tryout pool. The second part of the program uses this selected item pool to perform the final test selection. In the third part of the program, a table shows both expected number correct and standard error of measurement as functions of scale score, as well as statistical and graphic summaries on bias, fit, and the standard error of the final test. Any fault in the final selection becomes immediately apparent as the final statistics are generated: Whether the test is too easy or too difficult, suggests differential item functioning or DIF (see below), does not meet the requirements to match a parallel form, or does not adequately cover part of the range of performance. A developer detecting any such problems can then return to the second stage of the program and revise the selection. The flexibility and utility of the program encourages multiple attempts at fine-tuning the item selection.

Differential Item Functioning

Procedures for Eliminating Bias and Minimizing Differential Item Functioning

Three procedures were used to eliminate bias and minimize differential item functioning (DIF) in the New York State Tests.

The first was based on the premise that careful editorial attention to validity is an essential step in keeping bias to a minimum. Bias can occur only if the test is differentially valid for a given group of test takers. If the test entails irrelevant skills or knowledge (however common), the possibility of DIF is increased. Thus, preserving content validity is essential.

The second step was to follow the item writing guidelines established by NYSED. Developers reviewed NYSTP materials with these guidelines in mind. Such internal editorial reviews were done by at least four separate people: the content editor, who directly supervises the item writers; the project director; a style editor; and a proofreader. The final test built from the tryout materials was reviewed by at least these same people.

In the third procedure, educational community professionals who represent various ethnic groups reviewed all tryout materials. These groups included representatives from New York State. They were asked to consider and comment on the appropriateness of language, subject matter, and representation of people.

It is believed that these three procedures improved the quality of the New York State Tests and reduced bias. However, current evidence suggests that expertise in this area is no substitute for data; reviewers are often wrong about which items work to the disadvantage of a group, apparently because some of their ideas about how students will react to items may be faulty (Sandoval & Mille, 1979; Jensen, 1980). Thus, an empirical approach is desirable.

Part 2: Item Statistics for the Operational Data

Data Cleaning

Item analyses were conducted once CTB received data that met the following requirements established by NYSED:

- Comprises at least 85% of the estimated number of students in the State,
- Includes New York City and Buffalo,
- Includes at least one of Rochester, Syracuse, or Yonkers, and
- Includes at least two of Mount Vernon, Albany, Binghamton, Schenectady, or New Rochelle.

Initially, the state data set contained 220,297 cases. Table 4 below shows the data cleaning steps and the resulting size of the 85% data for conducting item analyses.

Table 4 Steps involved in data clean-up for analysis preparation

Steps Taken	# cases deleted	Ending N
original data		220,297
duplicate completely identical	424	219,873
duplicate identical personal info	578	219,295
grade not equal to 8	0	219,295
lep3 data	11,086	208,209
non-lep3 data	0	208,209
non-lep3 data after exclusion	6,830	201,379
Bedscodes = (parameter A)	68	201,311

Students whose LEP status = 3 are not required to take the test.

As Table 4 shows, the following records were eliminated, in the order listed:

- Duplicated records,
- Students whose limited English proficient (LEP) status was "3," indicating that they scored below the thirtieth percentile on a norm-referenced English reading test or the publisher's recommended score on an approved measure of English as a second language in reading,

- Students who are invalidated. Invalidated students are those who do not have a valid attempt in each test section as defined by CTB's Technology and Scoring Departments and who will not receive the ELA score or be categorized into a performance level, and
- Students from one bedscore who were removed due to very similar response patterns.

Item Analysis

Table 5 shows the results of item analyses conducted using the scaling sample for the G8 ELA test. The labels for the variables denote the following:

ITEM	Item number.
OMIT	Percentage of students who had blank response or double marks on MC items, or condition codes on the CR items.
PCTSEL*	For MC items, this is the percentage of students who chose the first (or second, etc.) answer option. For CR items, it is the percentage of students who received score 1 (or 2, etc.). Asterisked numbers indicate values for the correct response option.
PTBIS*	Point-biserial correlations for each response option. Asterisked numbers indicate values for the correct response option.
P_VAL	Item difficulty after omitted responses are converted to 0s (wrong). For MC items, p-value is the percentage of students responding correctly; for a CR item, p-value is the mean raw score divided by the maximum number of score points for an item.

Table 5 G8 ELA item level statistics

Raw Score Data		Test Administration Data						Reliability		P-Value	
Mean	SD	Number of Items			Number of Students			Feldt-Raju		Mean	
26.89	6.77	29			201,311			0.884		0.625	
Item	Omit	Pctsel1	Pctsel2	Pctsel3	Pctsel4	Ptbis1	Ptbis2	Ptbis3	Ptbis4	KEY	P-Value
1	0.06%	7.80%	*83.95%	4.55%	3.65%	-0.273	*0.397	-0.259	-0.272	2	0.84
2	0.15%	5.58%	*61.57%	21.83%	10.88%	-0.199	*0.393	-0.335	-0.229	2	0.62
3	0.18%	*57.38%	13.73%	17.27%	11.45%	*0.292	-0.194	-0.160	-0.285	1	0.57
4	0.12%	15.74%	35.07%	2.65%	*46.42%	-0.121	-0.402	-0.156	*0.378	4	0.46
5	0.10%	8.15%	*68.72%	16.09%	6.94%	-0.270	*0.320	-0.203	-0.230	2	0.69
6	0.19%	5.72%	5.58%	*64.79%	23.72%	-0.239	-0.212	*0.450	-0.391	3	0.65
7	0.15%	13.54%	16.37%	*67.56%	2.37%	-0.287	-0.230	*0.307	-0.126	3	0.68
8	0.10%	14.02%	*56.89%	4.95%	24.05%	-0.287	*0.244	-0.198	-0.122	2	0.57
9	0.13%	*60.50%	8.16%	4.57%	26.64%	*0.309	-0.213	-0.058	-0.330	1	0.60
10	0.14%	16.13%	6.34%	21.28%	*56.12%	-0.235	-0.237	-0.164	*0.276	4	0.56
11	0.11%	2.71%	*82.23%	7.98%	6.97%	-0.143	*0.277	-0.326	-0.121	2	0.82
12	0.18%	8.76%	*65.00%	10.72%	15.33%	-0.291	*0.341	-0.244	-0.186	2	0.65
13	0.13%	10.89%	1.20%	2.97%	*84.82%	-0.383	-0.173	-0.229	*0.416	4	0.85
14	0.18%	13.72%	24.69%	12.68%	*48.73%	-0.222	-0.309	-0.112	*0.338	4	0.49
15	0.14%	3.51%	6.37%	*82.91%	7.06%	-0.123	-0.133	*0.153	-0.152	3	0.83
16	0.19%	37.64%	*59.94%	1.41%	0.83%	-0.210	*0.153	-0.177	-0.136	2	0.60
17	0.34%	3.55%	2.79%	12.63%	*80.69%	-0.239	-0.245	-0.225	*0.315	4	0.81
18	0.56%	*72.88%	3.32%	7.30%	15.93%	*0.443	-0.204	-0.255	-0.369	1	0.73
19	1.13%	8.60%	13.35%	*62.36%	14.56%	-0.224	-0.362	*0.371	-0.133	3	0.62
20	0.98%	6.50%	11.76%	*73.78%	6.98%	-0.222	-0.285	*0.440	-0.324	3	0.74
21	1.31%	12.40%	*48.70%	8.09%	29.50%	-0.207	*0.375	-0.210	-0.271	2	0.49
22	1.32%	*51.60%	15.46%	17.05%	14.57%	*0.265	-0.233	-0.239	-0.066	1	0.52
23	1.59%	10.59%	4.45%	*68.52%	14.85%	-0.372	-0.238	*0.432	-0.208	3	0.69
24	1.61%	*69.024%	3.70%	23.06%	2.61%	*0.442	-0.208	-0.407	-0.158	1	0.69
25	1.75%	8.95%	*54.42%	8.51%	26.36%	-0.273	*0.394	-0.297	-0.200	2	0.54
		Pctsel1	Pctsel2	Pctsel3	Pctsel4	Pctsel5	Pctsel6				
26	0.61%	3.99%	11.64%	23.56%	30.71%	21.69%	7.79%	Listening		T	0.63
27	0.86%	5.72%	14.22%	23.06%	26.83%	19.82%	9.49%	Ind. Wr.		T	0.61
28	3.45%	21.34%	46.25%	28.97%				Wr. Mec.		T	0.67
29	1.25%	13.34%	49.66%	35.75%				Reading		T	0.73

Differential Item Functioning Analysis of Operational Data

To assess DIF for the New York State tests, students were identified as African-American, White, Hispanic, or Asian-American. For grade 8, students bubble in this information. These ethnic subgroups were chosen for DIF analyses because these subgroups contained the largest proportions of students in the State. Gender analyses were also conducted.

Developers strive to produce tests that minimize DIF. The DIF results reported here are those obtained when scoring students on the operational test using the pre-equated field test parameters. Thus, they may differ from DIF results obtained at the time of the field test administration.

Using demographic information, statistical DIF analyses were conducted for various ethnic groups and for males and females. A random sample was drawn from the final state GRT. Next, the sample was augmented by randomly selecting additional cases from any group of students whose count in the sample was less than 500, to supplement for reliable DIF analyses. The numbers of cases for the groups are reported Table 6 below.

Table 6 Number of Students in each Gender or Ethnic Group

Test	Female	Male	African-American	Asian-American	Hispanic-American
G8 ELA	3,650	3,640	1,357	500	1,194

The standardized mean difference (SMD) statistic (Zwick, Donoghue, & Grima, 1993) was used to examine DIF on the operational data. The SMD statistics can provide DIF information for both multiple choice and constructed response items. The SMD takes into account the natural ordering of the response levels of the items and has the desirable property of being based on those ability levels where members of the focal group are present. The standardized mean difference output results in a single statistic for each item.

$$\text{SMD} = \sum p_{Fk} m_{Fk} - \sum p_{Rk} m_{Rk},$$

where p_{Fk} is the proportion of focal group members who are at the k th level of the matching variable,

m_{Fk} is the mean item score for the focal group at the k th level, and

m_{Rk} is the analogous value for the reference group.

The matching variable is raw score and the kth level refers to the each successive raw score point.

A moderate amount of practically significant DIF, for or against the focal group, is represented by an SMD with an absolute value between .10 and .19, inclusive. A large amount of practically significant DIF is represented by an SMD with an absolute value of .20 or greater. SMD DIF results using operational data for G8 ELA are summarized below. There were no items with practically significant DIF.

Table 7 The numbers of items flagged for DIF in G8 ELA summary

Focal Group	Direction of DIF	G8 ELA
Female	In favor of	3 ¹
	Against	1 ²
African-American	In favor of	2 ³
	Against	1 ⁴
Asian-American	In favor of	2 ⁵
	Against	2 ⁶
Hispanic-American	In favor of	1 ⁷
	Against	1 ⁸

¹ Items #26, 27, and 29 (D = .18, .15, and .12).

² Item #3 (D = -.11)

³ Items #26, and 27 (D = .15, .18)

⁴ Item #4 (D = -.11)

⁵ Items #26, and 27 (D = .19, and .22)

⁶ Items #22, and 24 (D = -.11, and -.11)

⁷ Item #27, (D = .12)

⁸ Item #24 (D = -.11)

Part 3: Scoring and Reliability

Raw Score to Scale Score Conversion

To facilitate ease of interpretation and implementation, number-correct scoring was used on the New York State Tests in 2002. In number-correct scoring, a student's scale score is derived directly from his or her raw, or number-correct, score. The relationship between raw scores and their corresponding scale scores is expressed in a raw-score-to-scale-score (RS-SS) table.

In IRT, all the item characteristic curves for the items on a test can be added together to yield a function - the test characteristic curve (TCC) - that shows the expected raw score for each given scale score. By inverting the TCC, an expected scale score can be computed for each raw score. This new function - the inverse of the TCC - can be summarized in an RS-SS table. An advantage of RS-SS tables is that they make scoring relatively straightforward: With number-correct scoring, it is sufficient to know how many raw score points a student obtained on the test, to determine a student's scale score. The RS-SS conversion tables for both content areas appear in Table 8.

Reliability

The reliability of measurement refers to the reproducibility or consistency of an individual's tests scores. The two most frequently reported indices of reliability are the standard error of measurement and the reliability coefficient.

The standard error of measurement is a measure of the extent to which an individual's scores vary over numerous parallel tests. We computed *conditional* SEMs - SEMs for each scale score for G8 ELA, and these are reported below in Table 8. See also the section on estimated conditional standard errors of scale scores, below.

The reliability coefficient is the correlation coefficient between scores on parallel tests and is an index of how well scores on one parallel test predict scores from another parallel test. Among several ways to estimate the reliability of a test, Cronbach's alpha (Cronbach, Schönemann, & McKie, 1965) probably is the most frequently used. It is a measure of internal consistency (i.e., how homogeneous test items are) appropriate for a test containing only MC items. Since the G8 ELA test contains MC and CR items, Cronbach's alpha would underestimate reliability because of the effect of variance attributable to item types. A more appropriate index of internal consistency, the Feldt-Raju index, was used to estimate the reliability of the G8 ELA test. It was 0.884, and comparable to that for 2001.

Table 8 Raw Score to Scale Score with SEM for G8 ELA 2002

No. Correct (RS)	2003 G8 ELA	
	Scale Score	SEM
0	527	115
1	527	115
2	527	115
3	527	115
4	527	115
5	527	115
6	591	51
7	610	32
8	621	22
9	629	18
10	635	15
11	641	14
12	645	12
13	649	11
14	653	10
15	656	10
16	660	9
17	663	9
18	666	9
19	669	9
20	671	8
21	674	8
22	677	8
23	680	8
24	683	8
25	685	8
26	688	8
27	691	8
28	694	8
29	697	8
30	701	9
31	704	9
32	708	9
33	712	9
34	716	10
35	721	10
36	726	11
37	731	12
38	738	13
39	745	14
40	755	16
41	769	21
42	794	30
43	830	50

Estimated Conditional Standard Errors of Scale Scores

Each student's scale score is based on a sample of student performance at a given time and inherently has some measurement error. This understanding has led to notions like reliability and standard error of measurement. Reliability (i.e., classical SEM) presumes that the amount of measurement error is constant throughout the range of student ability. However, this is not

very realistic. Measurement error is less, and reliability greater, where more items exist and items are more informative. Item response theory lends itself to the calculation of measurement error for each scale score (i.e., conditional SEM).

Table 8 lists standard errors for selected scale scores. These standard errors are "constrained" so that the upper and lower limits of one standard error band around a scale score are below the upper and lower limits of the band for the next higher scale score. Typically, only standard errors on extreme ends are constrained. Because more items exist in the middle range of scale scores, the standard error is typically the lowest in the middle. A SS plus and minus one SEM constitutes a 68% confidence interval. For example, for a student whose grade 8 ELA SS is 680, we are 68% confident that his or her true score lies within the range 680 plus or minus 8, that is, between 672 and 688.

Lowest and Highest Obtainable Scale Scores

A maximum likelihood procedure cannot produce scale score estimates for students with perfect scores or scores below the level expected by guessing. Also, while maximum likelihood estimates are available for students with extreme scores other than zero or perfect, occasionally these estimates have standard errors of measurement that are very large, and differences between these extreme values have little meaning. Therefore, scores are established for these students based on a rational but necessarily non-maximum likelihood procedure. These values are called the lowest obtainable scale score (LOSS) and the highest obtainable scale score (HOSS). The same LOSS and HOSS values are also used for either number-correct or item-pattern scoring. For the New York State G8 ELA test, LOSS and HOSS values were set at 527 and 830.

Inter-Rater Agreement

In order to monitor the reliability of scoring among the teachers who scored the student responses, approximately 10% of the student papers were submitted to a second group of raters provided by Measurement Incorporated. Note that the teachers were trained by Measurement Incorporated. The results of the inter-rater agreement analyses for public schools and outside of New York City are provided in Tables 9-11. Additional results for public schools in New York City and non-public school will be reported as they become available.

Table 9 G8 ELA Inter-rater agreement

Inter-rater agreement (Read 1 : Non-NYC public school teachers; Read 2 : MI readers)								
CR item	Score Points	Agreement (%)			RS Mean		RS SD	
		Exact	Approx.	TOTAL	Read 1	Read 2	Read 1	Read 2
Listening	6	42.13	47.8	90.0	3.8	3.8	1.24	1.09
Ind. Writing	3	60.23	38.0	98.2	2.0	2.0	0.77	0.76
Wr. Mechanics	3	61.64	37.3	99.0	2.2	2.2	0.70	0.67
Reading	6	42.18	48.4	90.5	3.7	3.7	1.34	1.15

Approximate agreement (%) is the percent of pairs of reads that differ by one score point.
Total agreement (%) is the sum of exact and approximate agreement percents

Table 10 Percentages of inter-rater score differences

Reader 1 = Non-NYC public school teachers minus Reader 2 = MI readers											
CR item	-6	-5	-4	-3	-2	-1	0	1	2	3	4
Listening			.01	.27	3.96	22.78	42.13	25.06	5.38	.35	.06
Ind. Writing				.01	.90	18.44	60.23	19.55	.83	.03	
Wr. Mechanics					.47	16.74	61.64	20.58	.55	.02	
Reading	0.01		0.02	.36	4.73	25.05	42.18	23.31	4.10	.24	

Table 11 Reliability indices of hand scoring

CR item	Intra-Class Correlation ¹	Weighted Kappa ²
Listening	.83	.67
Ind. Writing	.81	.61
Wr. Mechanics	.78	.56
Reading	.86	.71

1 Agresti, A. (1990). Categorical data analysis (pp.366-367). New York: Wiley. Intra-class correlation is the percent of overall score variance accounted for by the variance of mean response scores.
2 Weighted kappa is a measure of association in contingency tables, and is 1 when agreement is perfect and 0 when agreement is what would be expected by chance.

Expected SPI Scores on the Standards at the Decision Points

The current New York State Grades 4 and 8 Score Reports for students report a Standard Performance Index (SPI) score for each of the standards or key ideas. The SPI is a diagnostic tool in the sense that it provides a profile of the student's relative strengths and weaknesses in terms of the content standards. However, just because a student has a high SPI on Standard 1 and a low SPI on Standard 2 does not necessary mean that she or he is strong on the former

standard and weak on the latter. This can occur if items measuring Standard 1 tend to be easy, while items measuring Standard 2 tend to be hard.

Table 12 G8 ELA standard performance index information

Book Item #	Max. Pts.	Standard	Level 2		Level 3		Level 4	
			At SS=660		At SS=699		At SS=738	
			Exp'd Diff.	Diff.*Max Pts	Exp'd Diff.	Diff.*Max Pts	Exp'd Diff.	Diff.*Max Pts
1	1	Std 1	0.57	0.57	0.95	0.95	1.00	1.00
2	1	Std 1	0.33	0.33	0.59	0.59	0.91	0.91
3	1	Std 1	0.54	0.54	0.75	0.75	0.89	0.89
4	1	Std 1	0.28	0.28	0.62	0.62	0.96	0.96
11	1	Std 1	0.64	0.64	0.90	0.90	0.98	0.98
13	1	Std 1	0.60	0.60	0.96	0.96	1.00	1.00
14	1	Std 1	0.28	0.28	0.57	0.57	0.90	0.90
17	1	Std 1	0.63	0.63	0.88	0.88	0.97	0.97
18	1	Std 1	0.56	0.56	0.86	0.86	0.97	0.97
22	1	Std 1	0.33	0.33	0.64	0.64	0.87	0.87
26 (Listen)	6	Std 1	0.34	2.04	0.55	3.30	0.73	4.38
28 (Ind. Wr)	3	Std 1	0.32	0.96	0.65	1.95	0.85	2.55
SUM	19	Total	5.42	7.76	8.92	12.97	11.03	16.38
			Exp'd Prop. NC	.41		.68		.86
6	1	Std 2	0.40	0.40	0.86	0.86	0.99	0.99
7	1	Std 2	0.60	0.60	0.84	0.84	0.95	0.95
8	1	Std 2	0.42	0.42	0.66	0.66	0.86	0.86
9	1	Std 2	0.42	0.42	0.71	0.71	0.91	0.91
10	1	Std 2	0.42	0.42	0.67	0.67	0.87	0.87
19	1	Std 2	0.36	0.36	0.79	0.79	0.97	0.97
20	1	Std 2	0.34	0.34	0.87	0.87	0.99	0.99
23	1	Std 2	0.39	0.39	0.90	0.90	0.99	0.99
24	1	Std 2	0.28	0.28	0.88	0.88	1.0	1.00
25	1	Std 2	0.29	0.29	0.65	0.65	0.95	0.95
SUM	10	Total	3.92	3.92	7.83	7.83	9.48	9.48
			Exp'd Prop. NC	.39		.78		.95
Table Continued								

Table 12 continued

Book Item #	Max. Pts.	Standard	Level 2		Level 3		Level 4	
			At SS=660		At SS=699		At SS=738	
			Exp'd Diff.	Diff.*Max Pts	Exp'd Diff.	Diff.*Max Pts	Exp'd Diff.	Diff.*Max Pts
5	1	Std 3	0.43	0.43	0.70	0.70	0.90	0.90
12	1	Std 3	0.46	0.46	0.76	0.76	0.94	0.94
15	1	Std 3	0.67	0.67	0.86	0.86	0.95	0.95
16	1	Std 3	0.46	0.46	0.66	0.66	0.83	0.83
21	1	Std 3	0.23	0.23	0.50	0.50	0.91	0.91
27 (Read)	6	Std 3	0.20	1.20	0.51	3.06	0.83	4.98
31	3		Does not contribute to any of the Standards					
Sum	11	Total	2.45	3.45	3.99	6.54	5.36	9.51
			Exp'd Prop. NC	.31		.59		.86

What teachers and students seem to need in order to better understand the SPIs are the SPIs expected of students who are just at each of the New York State decision points. These expected SPIs at the decision points can be used as "reference points" against which each student's SPIs are compared. For example, even if a student's SPI on Standard 3 is 67, if the expected SPI for the Level 3 Student is 59, the student's 67, although seemingly low compared with the perfect 100, is still higher than what is expected for the Level 3 Student on the standard.

Such expected SPIs for the 2002 Grade 8 English Language Arts exam are listed in Table 12.

Part 4: Descriptive Statistics

Scale-Score Frequency Distributions for the State and Subgroups

Tables 13 summarizes the scale-score frequency distributions for the state and the subgroups of students in public schools, students in non-public schools, two groups of limited-English-proficient (LEP) students, non-disabled students, and students with disabilities.

The public vs. non-public distinction was identified by the 9th character of the BEDs LEA code for each school. The non-disabled vs. disabled distinction was identified in the final state dataset. Additionally, two groups of LEP students are defined as those who have either "2" or "3" in the appropriate column of the final state dataset. The "LEP2" group is identified as having limited English proficiency and scored at or above either the 30th percentile on a norm-referenced English reading test or the publisher's recommended score on an approved measure of English as a Second Language (ESL) in reading. Similarly, the "LEP3" group is identified as having limited English proficiency and scored below either the 30th percentile on a norm-referenced English reading test or the publisher's recommended score on an approved measure of English as a Second Language (ESL) in reading.

As a summary table of the scale score frequency distributions, the SSs at the 10th, 25th, 50th, 75th, and 90th percentiles are listed below. No interpolation was employed in computing the percentiles. As an example, in the row of Statewide Inclusive at the 25th percentile the number 634 represents the highest scale score achieved by the lower 25 percent of the population.

Table 13 G8 ELA summary of scale score information

Sub Groups - Percentages	10 th	25 th	50 th	75 th	90 th
Statewide Inclusive	666	680	697	721	745
LEP = 2	653	666	680	697	716
LEP = 3	641	653	666	677	688
Public	666	680	697	721	745
Non-Public	674	688	704	726	745
Disabled	649	660	674	691	708
Non-Disabled	671	685	704	721	745

G8 ELA Scale Score Means and Standard Deviations

The scale score mean, standard deviation, and the total number of students in the statewide final general research file are shown in the table below.

Table 14 G8 ELA statewide scale score information

Population Sub Grouping	Number of Students (N)	Scale Score Mean	Scale Score Standard Deviation
All Students	228,619	697.85	30.00

References

- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (1999). *Standards for educational and psychological testing*. Washington, DC: American Psychological Association, Inc.
- Bock, R. D. (1972). Estimating item parameters and latent ability when responses are scored in two or more nominal categories. *Psychometrika*, 37, 29-51.
- Bock, R. D., & Aitkin, M. (1981). Marginal maximum likelihood estimation of item parameters: An application of an EM algorithm. *Psychometrika*, 46, 443-459.
- Burket, G. R. (1988). *ITEMSYS* [Computer program]. Unpublished.
- Burket, G. R. (1991). *PARDUX* [Computer program]. Unpublished.
- Cronbach, L. J., Schönemann, P., & McKie, D. (1965). Alpha coefficients for stratified-parallel tests. *Educational and Psychological Measurement*, 25, 291-312.
- Green, D. R., Yen, W. M., & Burket, G. R. (1989). Experiences in the application of item response theory in test construction. *Applied Measurement in Education*, 2, 297-312.
- Jensen, A. R. (1980). *Bias in mental testing*. New York: Free Press.
- Lord, F. M. (1980). *Applications of item response theory to practical testing problems* (71, 179-181). Hillsdale, NJ: Lawrence Erlbaum.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theories of mental test scores*. Reading, MA: Addison-Wesley.
- Muraki, E., (1992). A generalized partial credit model: Application of an EM algorithm. *Applied Psychological Measurement*, 16, 159-176.
- Muraki, E., & Bock, R. D. (1991). *PARSCALE: Parameter Scaling of Rating Data* [Computer program]. Chicago, IL: Scientific Software, Inc.
- New York State Department of Education. (1995). *Test Access and Modification for Individuals with Disabilities*. Available at <ftp://unix2.nysed.gov/pub/education.dept.pubs/vesid/oses/test.access.mod/testacce.txt>.
- Sandoval, J. H., & Mille, M. P. (1979, August). *Accuracy of judgments of WISC-R item difficulty for minority groups*. Paper presented at the annual meeting of the American Psychological Association, New York.
- Stocking, M. L., Lord, F. M., (1983). Developing a common metric in item response theory. *Applied Psychological Measurement*, 7, 201-210.

Thissen, D. (1982). Marginal maximum-likelihood estimation for the one-parameter logistic model. *Psychometrika*, 47, 175–186.

Thissen, D. (1991). MULTILOG: [Computer program]. Chicago, IL: Scientific Software, Inc.

Wright, B. D., & Linacre, J. M. (1992). *BIGSTEPS Rasch Analysis* [Computer program]. Chicago, IL: MESA Press.

Yen, W. M., (1993). Scaling performance assessments: Strategies for managing local item dependence. *Journal of Educational Measurement*, 30, 187-213

Yen, W. M., Burket, G. R., & Sykes, R. C. (1988). *Non-unique solutions to the likelihood equation for the three-parameter logistic model*. Paper presented at the meeting of the Psychometric Society, Los Angeles.

Zwick, R., Donoghue, J.R., & Grima, A. (1993). Assessment of differential item functioning for performance tasks. *Journal of Educational Measurement*, 36, 233-251.