

New York State Testing Program

Mathematics Grade 4

Technical Report 2003



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 MATHEMATICS	4
TEST CONFIGURATION	4
STUDENT PARTICIPATION AND TESTING ACCOMMODATIONS	ERROR! BOOKMARK NOT DEFINED.
<i>Students to be Tested</i>	<i>Error! Bookmark not defined.</i>
<i>Testing Accommodations</i>	<i>Error! Bookmark not defined.</i>
ITEM DEVELOPMENT	7
ITEM REVIEW PROCESS	7
<i>Documenting Content</i>	7
<i>Minimizing Bias</i>	7
<i>Minimizing Speededness</i>	8
TEST CONSTRUCTION AND PRE-EQUATING	8
<i>Calibration Samples</i>	8
<i>Answer Choice Information</i>	8
<i>Item Response Theory Models</i>	8
<i>Equating Method</i>	10
<i>Item Selection Criteria and Process</i>	10
<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	17
DIFFERENTIAL ITEM FUNCTIONING ANALYSIS OF OPERATIONAL DATA	17
PART 3: SCORING AND RELIABILITY	19
RAW SCORE TO SCALE SCORE CONVERSION	19
RELIABILITY	19
ESTIMATED CONDITIONAL STANDARD ERRORS OF SCALE SCORES	21
LOWEST AND HIGHEST OBTAINABLE SCALE SCORES	22
INTER-RATER AGREEMENT	22
EXPECTED SPI SCORES ON THE STANDARDS AT THE DECISION POINTS	24
PART 4: DESCRIPTIVE STATISTICS	25
SCALE-SCORE FREQUENCY DISTRIBUTIONS FOR THE STATE AND SUBGROUPS	25
G4 MA SCALE SCORE MEANS AND STANDARD DEVIATIONS	26
REFERENCES	27

List of Tables

Table 1 Key Ideas for Grade 4 Mathematics	4
Table 2 G4 MA Test Design	4
Table 3 Condition Codes for the MA CR Items.....	5
Table 4 Steps Involved in Data Clean-up for Analysis Preparation.....	13
Table 5 G4 MA Item Level Statistics	15
Table 6 Number of Students in each Gender or Ethnic Group	17
Table 7 The Numbers of Items Flagged for DIF in G4 MA Summary	18
Table 8 Raw Score to Scale Score with SE for G4 MA 2003	20
Table 9 G4 MA 2003 Inter-Rater Agreement	22
Table 10 Percentages of Inter-Rater Score Differences	23
Table 11 Reliability Indices of Hand Scoring	23
Table 12 G4 MA 2003 Standard Performance Index Information.....	24
Table 13 G4 MA 2003 Summary of Scale Score Information.....	25
Table 14 G4 MA Statewide Scale Score Information.....	26

Part 1: Test Design

The New York State Learning Standards for Mathematics

The test measures progress toward the seven Key Ideas described in Standard 3 of the *Learning Standards for Mathematics, Science, and Technology* at:

<http://www.emsc.nysed.gov/ciai/pub/standards.pdf>. The Grade 4 Mathematics test is written to test students in all seven Key Ideas and for each Key Idea students had the opportunity to demonstrate their knowledge both by selecting and generating responses. The seven Key Ideas are listed in Table 1 below with the approximate percent emphasis that is placed on each for Grade 4 Mathematics.

Table 1 Key Ideas for Grade 4 Mathematics

Key Ideas	Emphasis for Grade 4
Mathematical Reasoning	10–15%
Number and Numeration	15–25%
Operations	20–25%
Modeling/ Multiple Representation	5–10%
Measurement	15–20%
Uncertainty	5–10%
Patterns / Functions	10–15%

Test Configuration

Table 2 provides the test design for Grade 4 Mathematics, including the number of questions, question types, number of points, and time allotted for each testing session. Table 3 indicates the conditions codes used in scoring the responses to the CR items.

Table 2 G4 MA Test Design

	Number of Questions	Number of Points	Time in Minutes
Session 1	30 MC	30	40
Session 2	7 SR 2 ER	14 6	50
Session 3	7 SR 2 ER	14 6	50

Table 3 Condition Codes for the MA CR Items

Condition Code	Meaning
A	Blank
B	Refusal
C	Illegible
D	Other language

Student Participation and Testing Accommodations

Students to be Tested

The New York State Testing Program (NYSTP) Grade 4 Mathematics test must be administered to all public school students in Grade 4 and all ungraded students who are age-equivalent to students in Grade 4. This includes students who have been retained in Grade 4. Nonpublic schools are strongly encouraged to administer the tests. The exceptions noted below apply to students in public and nonpublic schools participating in the NYSTP.

Testing Accommodations

Accommodations were used in the NYSTP operational tests to provide equal access to assessments for students with disabilities. These accommodations are used to increase the validity of test scores by offsetting behavioral constraints due to 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

The Committee on Special Education (CSE) must decide for each student on a case-by-case basis (and document on the student's Individualized Education Program) whether the student will participate in the general State assessment, in a locally selected assessment, or in the New York State Alternate Assessment for Students with Severe Disabilities (NYSAA). The criteria that the CSE must use to determine eligibility for a locally selected assessment is available at <http://www.emsc.nysed.gov/deputy/Documents/disabilities-assess.htm>. The criteria to determine eligibility for the NYSAA is available on <http://www.vesid.nysed.gov/specialed/alterassessment/alterassess.htm>.

It is the responsibility of the principal to ensure that testing accommodations specified in the IEP or 504 Plan are provided to students with disabilities as long as they do not alter a construct being measured by the test. 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 that alter the construct being measured are not permitted on elementary- and intermediate-level State assessments. For more information, see <http://web.nysed.gov/vesid/sped/policy/changeaccomm.htm>.

Principals may modify testing procedures for General Education students who incur an injury (for example, a broken arm) or experience the onset of a short- or long-term disability (for example, epilepsy) sustained or diagnosed within 30 days prior to the administration of State tests. In such cases, when sufficient time is not available for the development of an Individualized Education Program (IEP) or a Section 504 Accommodation Plan (504 Plan), principals may authorize certain accommodations that will not significantly change the skills being tested.

Eligibility for such accommodations is based on the principal's professional discretion, but the principal may confer with members of the Committee for Special Education (CSE) or with other school personnel in making such a determination. Pursuant to Section 100.3 of the Regulations of the Commissioner of Education, building principals are responsible for administering State assessments and for maintaining the integrity of test content and programs in accordance with directions and procedures established by the Commissioner of Education.

Limited English Proficient (LEP) Students

The provisions of the NCLB Act do not permit any exemption of LEP students from the State's Grades 4 and 8 Mathematics tests. All LEP students in these grades must take the Grade 4 or 8 Mathematics test. These tests are available in Chinese, Haitian Creole, Korean, Russian, and Spanish. They can be translated orally into other languages for those LEP students whose first language is one for which a written translation is not available from the Department. Schools are permitted to offer LEP students specific testing accommodations when taking State examinations to ensure valid and reliable test results.

Additional information concerning the inclusion of LEP students in State examinations in English Language Arts and Mathematics will be provided on the Department's website <http://www.emsc.nysed.gov/osa>.

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.

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 Key Ideas for Mathematics. 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. Tables 5 shows the omit rates for items on the G4 MA 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 2002. Efforts were made to ensure that the sample of students was representative of the state tested population. The 2002 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 Mathematics, grade 8 ELA, and grade 8 Mathematics).

Since these items are calibrated and on a common scale, the pool of available grade 4 Mathematics items can be used to construct a test form and to produce a raw-score-to-scale-score table for that form. The 2003 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 (mean item score divided by the total number of points possible) are examined.

Item Response Theory Models

Although useful, the differences in proportion of points received (p-values) limit the degree to which one can compare important characteristics of the test items. Item-response theory (IRT) allows one to make better comparisons among items, even those from different test forms, 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 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 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 performance 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 (3PL) model (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 (2PPC) model (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 (Fitzpatrick, 1990; Fitzpatrick, 1994; Fitzpatrick and Julian, 1996).

Equating Method

After the item calibration, all of the Grade 4 Mathematics field test items were placed on the NYS G4 MA scale using the 2002 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 the production of scoring tables (see Part 3) ahead of the 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 the standard error of ability estimate (a function 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. Examples of possible faults that may occur are when the test is too easy or difficult, contains items demonstrating differential item functioning or DIF (see below), does not meet the requirements to match a parallel form, or does not adequately measure 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.

Procedures for Eliminating Bias and Minimizing Differential Item Functioning

Statistical differential item functioning (DIF) analyses were conducted for gender groups and such ethnic groups as African-American, Hispanic-American and Asian-American in the sample.

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. These 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, New York State educational community professionals who represent various ethnic groups reviewed all tryout materials. These professionals 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 the cities of Rochester, Syracuse, or Yonkers, and
- Includes at least two of the cities of Mount Vernon, Albany, Binghamton, Schenectady, or New Rochelle.

Initially, the state data set contained 249,861 cases. Table 4 below shows the data cleaning steps and the resulting size of the 85% sample used for conducting item analyses.

Table 4 Steps Involved in Data Clean-up for Analysis Preparation

Steps Taken	# Cases Deleted	Ending N
Original Data		249,861
Duplicate Records	1	249,860
Grade Not Equal to 4	225	249,635
LEP3 Data	2,786	246,849
Non-LEP3 Data	0	246,849
Non-LEP3 Data after Exclusion Rules	4,157	242,692

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. These students are not required to take the Mathematics test, unless a version of the test in their native language is available, and were consequently removed from consideration,
- All students who took the test in Chinese, Haitian-Creole, Russian, Korean, or Spanish, plus all students who took the English version AND whose LEP status is not 3, and
- Students who did not have a valid attempt in each of three sections as determined by the application of CTB's Invalidation / Omission / Suppression rules (approved by NYSED).

Item Analysis

Table 5 presents the results of item analyses conducted using the scaling sample for the G4 MA test. The labels for the variables denote the following:

ITEM	Item number.
OMIT	Proportion 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 through the fourth answer option. For CR items, it is the percentage of students who received a score of 0 through the maximum number of points possible. 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 proportion 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 G4 MA Item Level Statistics

Raw Score Data			Test Administration Data					Reliability Feldt-Raju			P-Value Mean	
Mean		SD	Number of Items		Number of Students							
47.77		13.02	48		242,621			0.932			0.7038	
Item	Omit	Pctsel0	Pctsel1	Pctsel2	Pctsel3	Pctsel4	Ptbis1	Ptbis2	Ptbis3	Ptbis4	Key	p-value ⁺
1	0.0002		2.46%	2.79%	2.31%	*92.41%	-0.14	-0.20	-0.11	0.27	4	0.9241
2	0.0004		*88.65%	1.43%	7.19%	2.67%	0.38	-0.12	-0.29	-0.19	1	0.8865
3	0.0014		*88.26%	6.57%	2.96%	2.05%	0.40	-0.20	-0.25	-0.23	1	0.8826
4	0.0004		1.94%	*92.82%	1.81%	3.38%	-0.20	0.32	-0.22	-0.14	2	0.9282
5	0.0008		4.88%	3.71%	14.61%	*76.66%	-0.34	-0.22	-0.25	0.48	4	0.7666
6	0.0010		*60.31%	28.01%	2.06%	9.47%	0.31	-0.24	-0.21	-0.04	1	0.6031
7	0.0006		*88.44%	6.27%	1.88%	3.34%	0.43	-0.24	-0.19	-0.29	1	0.8844
8	0.0036		6.47%	*67.67%	16.52%	8.95%	-0.21	0.42	-0.22	-0.21	2	0.6767
9	0.0009		7.68%	6.68%	*78.25%	7.27%	-0.20	-0.18	0.43	-0.29	3	0.7825
10	0.0005		0.54%	3.82%	*94.93%	0.65%	-0.13	-0.27	0.32	-0.11	3	0.9493
11	0.0026		4.90%	*69.85%	13.73%	11.23%	-0.13	0.38	-0.21	-0.22	2	0.6985
12	0.0010		5.51%	18.68%	*73.30%	2.38%	-0.26	-0.25	0.43	-0.22	3	0.7330
13	0.0014		15.03%	5.08%	3.81%	*75.88%	-0.28	-0.25	-0.19	0.45	4	0.7588
14	0.0033		9.77%	22.88%	*53.09%	13.91%	-0.19	-0.10	0.42	-0.31	3	0.5309
15	0.0009		4.20%	2.35%	*88.45%	4.88%	-0.30	-0.21	0.45	-0.23	3	0.8845
16	0.0011		13.93%	6.01%	3.63%	*76.31%	-0.37	-0.19	-0.22	0.51	4	0.7631
17	0.0020		11.09%	12.86%	*67.16%	8.66%	-0.23	-0.07	0.28	-0.11	3	0.6716
18	0.0011		2.93%	2.91%	*89.69%	4.34%	-0.25	-0.24	0.42	-0.21	3	0.8969
19	0.0021		6.91%	*82.23%	6.88%	3.75%	-0.26	0.43	-0.21	-0.22	2	0.8223
20	0.0016		7.33%	*74.81%	14.25%	3.43%	-0.12	0.37	-0.28	-0.16	2	0.7481
21	0.0025		30.60%	9.05%	8.01%	*51.99%	-0.35	-0.10	-0.12	0.46	4	0.5199
22	0.0027		*55.67%	32.56%	8.24%	3.21%	0.40	-0.24	-0.20	-0.17	1	0.5567
23	0.0087		19.88%	*61.41%	8.89%	8.89%	-0.26	0.42	-0.16	-0.16	2	0.6141
24	0.0035		4.80%	2.54%	8.64%	*83.60%	-0.21	-0.18	-0.11	0.29	4	0.8360
25	0.0039		8.49%	6.64%	*64.51%	19.91%	-0.31	-0.26	0.41	-0.10	3	0.6451
26	0.0073		*52.35%	18.02%	16.30%	12.56%	0.38	-0.22	-0.16	-0.12	1	0.5235
27	0.0067		3.80%	*83.23%	4.18%	8.08%	-0.11	0.34	-0.23	-0.19	2	0.8323
28	0.0103		*48.09%	16.14%	17.18%	17.52%	0.33	-0.18	-0.11	-0.13	1	0.4809
29	0.0109		24.60%	20.76%	*48.49%	5.04%	-0.15	-0.12	0.31	-0.16	3	0.4849
30	0.0120		2.29%	5.01%	20.78%	*70.69%	-0.18	-0.20	-0.16	0.31	4	0.7069
31	0.0037	16.19%	28.88%	54.56%							CR	0.6900
32	0.0039	2.53%	14.53%	82.55%							CR	0.8982
*33	1.0000	0%	0%	0%							CR	0.0000
34	0.0035	6.60%	13.41%	79.64%							CR	0.8634
35	0.0041	16.85%	37.19%	45.55%							CR	0.6415
36	0.0051	19.10%	26.04%	54.36%							CR	0.6737
37	0.0050	3.09%	11.71%	35.73%	48.97%						CR	0.7669
38	0.0037	9.18%	37.30%	53.16%							CR	0.7181
39	0.0057	24.00%	16.12%	16.83%	42.48%						CR	0.5907
40	0.0036	7.05%	23.34%	69.25%							CR	0.8092

(Table 5 continues)

Table 5 G4 MA Item Level Statistics (continued)

Item	Omit	Pctsel0	Pctsel1	Pctsel2	Pctsel3	Pctsel4	Ptbis1	Ptbis2	Ptbis3	Ptbis4	Key	p-value ⁺
41	0.0043	14.51%	19.79%	65.27%							CR	0.7516
42	0.0046	18.75%	13.94%	66.84%							CR	0.7382
43	0.0040	2.47%	14.17%	82.95%							CR	0.9004
44	0.0045	25.23%	16.90%	57.42%							CR	0.6587
45	0.0063	27.28%	24.10%	47.99%							CR	0.6005
46	0.0101	25.74%	18.14%	55.12%							CR	0.6419
47	0.0061	22.68%	12.16%	27.56%	36.99%						CR	0.5942
48	0.0092	46.23%	36.01%	10.58%	6.26%						CR	0.2532
* Item 33 was omitted due to a misprinted manipulative.												

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 4, students bubble in this information. These ethnic groups were chosen for DIF analyses because these populations are the largest 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 for any group of students whose count in the sample was less than 500 in an attempt to enhance the reliability of the DIF analyses. The numbers of cases for the groups are reported in Table 6 below.

Table 6 Number of Students in each Gender or Ethnic Group

Test	Female	Male	African-American	Asian-American	Hispanic-American
Grade 4 Mathematics	118,958	123,639	48,338	14,310	45,136

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.

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

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 k th 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 G4 Mathematics are summarized below.

Table 7 The Numbers of Items Flagged for DIF in G4 MA Summary

Focal Group	Direction of DIF	Number of Items
Female	In favor of	1 ¹
	Against	1 ²
African-American	In favor of	0
	Against	0
Asian-American	In favor of	0
	Against	1 ³
Hispanic-American	In favor of	0
	Against	1 ⁴

- 1 Item #37 (D = .12)
- 2 Item #42 (D = -.13)
- 3 Item #45 (D = -.11)
- 4 Item #36 (D = -.12)

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 2003. 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 a *conditional* error – the standard error (SE) for each scale score for G4 MA 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. Cronbach's alpha is a measure of internal consistency (i.e., how homogeneous test items are) that is appropriate for a test containing only MC items. Since the G4 MA 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 G4 MA test. It was 0.932, a value comparable to that for 2002.

Table 8 Raw Score to Scale Score with SE for G4 MA 2003

No. Correct (RS)	2003 G4 Mathematics			
	Non-Braille		Braille	
	Scale Score ¹	SE ¹	Scale Score ²	SE ²
0	448	126	448	127
1	448	126	448	127
2	448	126	448	127
3	448	126	448	127
4	448	126	448	127
5	448	126	448	127
6	448	126	448	127
7	500	73	502	73
8	527	47	529	47
9	541	32	543	32
10	552	25	553	25
11	559	21	561	21
12	566	18	568	18
13	571	16	573	16
14	576	15	578	15
15	581	14	582	14
16	584	13	586	13
17	588	12	590	12
18	591	12	593	12
19	594	11	596	11
20	597	11	599	11
21	600	10	602	10
22	603	10	604	10
23	605	10	607	10
24	608	9	610	9
25	610	9	612	9
26	613	9	614	9
27	615	9	616	9
28	617	9	619	9
29	619	8	621	8
30	621	8	623	8
31	623	8	625	8
32	625	8	627	8
33	627	8	629	8
34	629	8	631	8
35	631	8	633	8
36	633	8	635	8
37	635	8	637	8
38	637	8	639	8
39	639	8	641	8
40	641	8	643	8
41	642	8	645	8
42	644	8	647	8
43	646	8	649	8
44	648	8	651	8
45	650	8	653	8
46	652	8	655	8

(Table 8 continues)

Table 8 Raw Score to Scale Score with SEM for G4 MA 2003 (continued)

No. Correct (RS)	2003 G4 Mathematics			
	Non-Braille		Braille	
	Scale Score ¹	SE ¹	Scale Score ²	SE ²
47	654	8	657	8
48	656	8	659	8
49	658	8	661	8
50	661	8	664	8
51	663	8	666	8
52	665	8	669	9
53	668	8	671	9
54	670	9	674	9
55	673	9	677	10
56	676	9	681	10
57	679	10	684	10
58	682	10	688	11
59	686	11	693	12
60	690	11	698	13
61	694	12	704	14
62	699	13	711	15
63	705	14	719	18
64	712	15	732	22
65	721	18	753	32
66	734	22	810	84
67	755	32		
68	810	81		

1 Scale scores for students who took the non-Braille versions of the Grade 4 MA exam were computed using all items but item 33 due to a misprinted manipulative.

2 Scale scores for students who took the Braille version of the Grade 4 MA exam were computed using all but items 33 and 38 due to a misprinted manipulative and an item that could not be Brailled, respectively.

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 contains some measurement error. The classical SEM presumes the amount of measurement error is constant throughout the range of student ability. However, this is not 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 a standard error for each scale score.

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 smallest in the middle. A SS plus and minus one SE constitutes a 68% confidence interval. For example, for a student who took the non-Braille version and whose grade 4 MA SS is 641, we are 68% confident that his or her true score lies within the range 641 plus or minus 8, that is, between 633 and 649.

Lowest and Highest Obtainable Scale Scores

A maximum likelihood procedure cannot produce scale score estimates for students with zero or perfect scores. Scale score estimates below the level expected by guessing are unreliable and subsequently not reported. Also, while maximum likelihood estimates may be available for students with extreme scores other than a perfect score, occasionally these estimates have standard errors 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 used for either number-correct or item-pattern scoring. For the New York State G4 MA test, LOSS and HOSS values were set at 448 and 810.

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 G4 MA 2003 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
Item 31	2	88.6	11.1	99.7	1.5	1.5	0.69	0.75
Item 32	2	96.5	3.2	99.8	1.8	1.8	0.40	0.41
Item 33	
Item 34	2	96.8	2.8	99.7	1.8	1.8	0.48	0.47
Item 35	2	98.2	1.6	99.8	1.3	1.3	0.70	0.70
Item 36	2	59.9	14.1	74.0	1.4	0.8	0.75	0.90
Item 37	3	78.1	20.7	98.8	2.4	2.4	0.73	0.76
Item 38	2	95.8	4.0	99.8	1.5	1.5	0.61	0.62
Item 39	3	61.8	31.9	93.7	1.9	1.7	1.17	1.15
Item 40	2	90.3	9.5	99.8	1.7	1.7	0.55	0.57
Item 41	2	91.2	8.6	99.8	1.6	1.6	0.65	0.68
Item 42	2	96.7	2.1	98.9	1.6	1.6	0.73	0.74
Item 43	2	92.2	7.4	99.6	1.9	1.9	0.39	0.39
Item 44	2	88.2	11.2	99.5	1.4	1.4	0.81	0.84
Item 45	2	80.7	18.4	99.1	1.3	1.2	0.81	0.85
Item 46	2	89.1	10.7	99.8	1.4	1.4	0.78	0.85
Item 47	3	69.0	29.8	98.8	1.9	1.8	1.13	1.10
Item 48	3	81.4	18.2	99.6	0.8	0.7	0.88	0.84
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	-3	-2	-1	0	1	2	3
Item 31	0.13	3.44	88.62	7.66	0.15		
Item 32	0.04	1.72	96.52	1.52	0.20		
Item 33							
Item 34	0.18	1.88	96.84	0.95	0.15		
Item 35	0.04	0.68	98.17	0.97	0.15		
Item 36	0.11	1.28	59.85	12.86	25.90		
Item 37	0.40	8.80	78.05	11.91	0.71	0.13	
Item 38		1.55	95.79	2.47	0.18		
Item 39	1.85	10.55	61.84	21.35	4.06	0.29	0.05
Item 40	0.04	3.75	90.27	5.78	0.16		
Item 41	0.04	3.37	91.17	5.25	0.18		
Item 42	0.53	0.77	96.73	1.37	0.60		
Item 43	0.15	4.39	92.15	3.02	0.29		
Item 44	0.24	3.68	88.24	7.54	0.31		
Item 45	0.16	3.75	80.68	14.69	0.71		
Item 46	0.07	1.98	89.13	8.71	0.11		
Item 47	0.20	9.51	69.03	20.27	0.86	0.09	0.04
Item 48	0.09	5.80	81.42	12.37	0.26	0.07	

Table 11 Reliability Indices of Hand Scoring

CR Item	Intra-Class Correlation ¹	Weighted Kappa ²
Item 31	0.94	0.88
Item 32	0.94	0.87
Item 33	.	.
Item 34	0.95	0.91
Item 35	0.99	0.98
Item 36	0.63	0.34
Item 37	0.88	0.76
Item 38	0.97	0.94
Item 39	0.89	0.78
Item 40	0.92	0.84
Item 41	0.95	0.89
Item 42	0.97	0.94
Item 43	0.85	0.7
Item 44	0.95	0.9
Item 45	0.92	0.84
Item 46	0.96	0.91
Item 47	0.93	0.86
Item 48	0.93	0.87

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 for a student – for a given Key Idea – is an estimate of the percent of maximum raw score that the student would get if he or she took a large sample of items in that Key Idea. 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 one key idea and a low SPI on another key idea 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 one key idea tend to be easy, while items measuring another key idea tend to be hard.

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, if a student's SPI on Key Idea 1 is 50 and the expected SPI for the Level 3 Student is 47, the student's 50, although seemingly low compared with the perfect 100, is still higher than what is expected for the Level 3 Student on the key idea. Expected SPIs for the 2003 Grade 8 Mathematics exam are listed in Table 12.

Table 12 G4 MA 2003 Standard Performance Index Information

Key Idea	Expected Percent of the Max. Raw Score at each of the Cut Points for the Non-Braille version of the test.					Expected Percent of the Max. Raw Score at each of the Cut Points for the Braille version of the test.				
	# Items	Max Pts.	Level 2	Level 3	Level 4	# Items	Max Pts.	Level 2	Level 3	Level 4
			At SS=602	At SS=637	At SS=678			At SS=602	At SS=637	At SS=678
1	5	8	25	47	82	5	8	25	47	82
2	11	14	46	73	93	11	14	46	73	93
3	11	16	30	59	88	11	16	30	59	88
4	3	5	42	68	87	*2	3	50	79	95
5	6	10	28	44	66	6	10	28	44	66
6	5	6	25	47	81	5	6	25	47	81
7	6	9	23	46	80	6	9	23	46	80

* Item 38 was not translatable into Braille so Key Idea 4 has one fewer scored item in it worth 2 points.

Part 4: Descriptive Statistics

Scale-Score Frequency Distributions for the State and Subgroups

Table 13 summarizes the scale-score frequency distributions for the state and the following groups of students:

- public schools,
- 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.

A summary table of the scale score frequency distributions containing the SSs at the 10th, 25th, 50th, 75th, and 90th percentiles is provided below. No interpolation was employed in computing the percentiles. As an example, in the row of Statewide Inclusive at the 25th percentile the number 641 represents the highest scale score achieved by the lowest 25 percent of the population.

Table 13 G4 MA 2003 Summary of Scale Score Information

Sub Groups - Percentages	10th	25th	50th	75th	90th
Statewide Inclusive	619	641	661	682	699
LEP = 2	613	631	652	676	690
LEP = 3	600	619	641	665	686
Public	617	641	661	682	705
Non-Public	623	642	661	679	669
Disabled	581	610	635	658	676
Visually Impaired	588	627	650	669	690
Non-Disabled	625	644	665	686	705

G4 MA Scale Score Means and Standard Deviations

The total number of students and the percent of students in each performance level in the statewide final general research file are shown in the table below. Statistics for the three previous years are also included.

Table 14 G4 MA Statewide Scale Score Information

Year	Population Sub Grouping	Number of Students (N)	PCT in PL1	PCT in PL2	PCT in PL3	PCT in PL4
2003	All Students	245,579	4.6	16.93	48.19	30.27
2002	All Students	245,022	6.88	25.20	45.49	22.44
2001	All Students	249,119	8.16	22.57	43.14	26.14
2000	All Students	249,797	8.83	26.03	46.68	18.45
1999	All Students	245,358	9.61	23.59	43.15	23.64

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.
- Fitzpatrick, A. R. (1990) *Status Report on the results of Preliminary Analysis of Dichotomous and Multi-Level Items Using the PARMATE Program*. Unpublished manuscript
- Fitzpatrick, A. R. (1994) *Two studies comparing parameter estimates produced by PARDUX and BIGSTEPS*. Unpublished manuscript.
- Fitzpatrick, A. R., & Julian, M. W. (1996) *Two studies comparing the parameter estimates produced by PARDUX and PARSCALE*. Monterey, CA:CTB/McGraw-Hill.
- 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*. Hillsdale, NJ: Lawrence Erlbaum.
- Lord, F. M., & Novick, M. R. (1968). *Statistical theories of mental test scores*. Menlo Park, CA: 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>.

New York State Department of Education. (2003). *Learning Standards for Mathematics, Science, and Technology*. Available at <http://www.emsc.nysed.gov/ciai/pub/standards.pdf>.

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.

Zwick, R., Donoghue, J. R., & Grima, A. (1993). Assessment of differential item functioning for performance tasks. *Journal of Educational Measurement*, 36, 233-25.