

# New York State Regents Examination in Algebra II (Common Core)

## 2016 Technical Report



Prepared for the New York State Education Department  
by Pearson

**March 2017**

# Copyright

---

Developed and published under contract with the New York State Education Department by Pearson.

Copyright © 2017 by the New York State Education Department.

## **Secure Materials.**

All rights reserved. No part of this document may be reproduced or transmitted by any means. Use of these materials is expressly limited to the New York State Education Department.

# Contents

---

<b>CHAPTER 1: INTRODUCTION AND HISTORY .....</b>	<b>1</b>
1.1 INTRODUCTION .....	1
1.2 HISTORY .....	1
1.3 PURPOSES OF THE EXAM .....	1
1.4 TARGET POPULATION (STANDARD 7.2) .....	2
<b>CHAPTER 2: CLASSICAL ITEM STATISTICS (STANDARD 4.10).....</b>	<b>4</b>
2.1 ITEM DIFFICULTY.....	4
2.2 ITEM DISCRIMINATION .....	5
2.3 DISCRIMINATION ON DIFFICULTY SCATTER PLOTS .....	6
2.4 OBSERVATIONS AND INTERPRETATIONS .....	7
<b>CHAPTER 3: IRT CALIBRATIONS, EQUATING, AND SCALING (STANDARDS 2, AND 4.10) .....</b>	<b>8</b>
3.1 DESCRIPTION OF THE RASCH MODEL.....	8
3.2 SOFTWARE AND ESTIMATION ALGORITHM .....	9
3.3 CHARACTERISTICS OF THE TESTING POPULATION.....	9
3.4. ITEM DIFFICULTY-STUDENT PERFORMANCE MAPS.....	9
3.5 CHECKING RASCH ASSUMPTIONS .....	10
<i>Unidimensionality</i> .....	10
<i>Local Independence</i> .....	12
<i>Item Fit</i> .....	14
3.6 SCALING OF OPERATIONAL TEST FORMS .....	15
<b>CHAPTER 4: RELIABILITY (STANDARD 2).....</b>	<b>18</b>
4.1 RELIABILITY INDICES (STANDARD 2.20).....	18
<i>Coefficient Alpha</i> .....	19
4.2 STANDARD ERROR OF MEASUREMENT (STANDARDS 2.13, 2.14, 2.15).....	19
<i>Traditional Standard Error of Measurement</i> .....	19
<i>Traditional Standard Error of Measurement Confidence Intervals</i> .....	20
<i>Conditional Standard Error of Measurement</i> .....	20
<i>Conditional Standard Error of Measurement Confidence Intervals</i> .....	21
<i>Conditional Standard Error of Measurement Characteristics</i> .....	21
<i>Results and Observations</i> .....	22
4.3 DECISION CONSISTENCY AND ACCURACY (STANDARD 2.16) .....	23
4.4 GROUP MEANS (STANDARD 2.17) .....	25
4.5 STATE PERCENTILE RANKINGS .....	26
<b>CHAPTER 5: VALIDITY (STANDARD 1).....</b>	<b>28</b>
5.1 EVIDENCE BASED ON TEST CONTENT .....	28
<i>Content Validity</i> .....	29
<i>Item Development Process</i> .....	29
<i>Item Review Criteria</i> .....	31
5.2 EVIDENCE BASED ON RESPONSE PROCESSES .....	32
<i>Administration and Scoring</i> .....	33
<i>Statistical Analysis</i> .....	35
5.3 EVIDENCE BASED ON INTERNAL STRUCTURE.....	35
<i>Item Difficulty</i> .....	36
<i>Item Discrimination</i> .....	36
<i>Differential Item Functioning</i> .....	36
<i>IRT Model Fit</i> .....	37
<i>Test Reliability</i> .....	37
<i>Classification Consistency and Accuracy</i> .....	37

<i>Dimensionality</i> .....	38
5.4 EVIDENCE BASED ON RELATIONS TO OTHER VARIABLES.....	38
5.5 EVIDENCE BASED ON TESTING CONSEQUENCES .....	39
<b>REFERENCES</b> .....	<b>40</b>
<b>APPENDIX A: OPERATIONAL TEST MAPS</b> .....	<b>44</b>
<b>APPENDIX B: RAW-TO-THETA-TO-SCALE SCORE CONVERSION TABLES</b> .....	<b>45</b>
<b>APPENDIX C: ITEM WRITING GUIDELINES</b> .....	<b>46</b>
GUIDELINES FOR WRITING CONSTRUCTED-RESPONSE MATH ITEMS.....	<b>ERROR! BOOKMARK NOT DEFINED.</b>

# List of Tables

---

TABLE 1 TOTAL EXAMINEE POPULATION: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	5
TABLE 2 MULTIPLE-CHOICE ITEM ANALYSIS SUMMARY: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	7
TABLE 3 CONSTRUCTED-RESPONSE ITEM ANALYSIS SUMMARY: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	8
TABLE 4 DESCRIPTIVE STATISTICS IN <i>p</i> -VALUE AND POINT BISERIAL CORRELATION: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	9
TABLE 5 SUMMARY OF ITEM RESIDUAL CORRELATIONS: ALGEBRA II (COMMON CORE).....	14
TABLE 6 SUMMARY OF INFIT MEAN SQUARE STATISTICS: ALGEBRA II (COMMON CORE) .....	15
TABLE 7 RELIABILITIES AND STANDARD ERRORS OF MEASUREMENT: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE).....	20
TABLE 8 DECISION CONSISTENCY AND ACCURACY RESULTS: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	24
TABLE 9 GROUP MEANS: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	25
TABLE 10 STATE PERCENTILE RANKING FOR RAW SCORE – REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	26
TABLE 11 TEST BLUEPRINT, REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	28

# List of Figures

---

FIGURE 1 SCATTERPLOT: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	9
FIGURE 2 STUDENT PERFORMANCE MAP: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	11
FIGURE 3 SCREE PLOTS: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	13
FIGURE 4 CONDITIONAL STANDARD ERROR PLOTS: REGENTS EXAMINATION IN ALGEBRA II (COMMON CORE) .....	22
FIGURE 5 PSEUDO-DECISION TABLE FOR TWO HYPOTHETICAL CATEGORIES .....	23
FIGURE 6 PSEUDO-DECISION TABLE FOR FOUR HYPOTHETICAL CATEGORIES .....	23
FIGURE 7 NEW YORK STATE EDUCATION DEPARTMENT TEST DEVELOPMENT PROCESS .....	29

# Chapter 1: Introduction and History

---

## 1.1 INTRODUCTION

This technical report for the Regents Examination in Algebra II (Common Core) will provide New York State with documentation on the purpose of the Regents Examination, scoring information, evidence of both reliability and validity of the exam, scaling information, and guidelines and reporting information for the June 2016 administration. As the *Standards for Education and Psychological Testing* discusses in Standard 7, “The objective of the documentation is to provide test users with the information needed to help them assess the nature and quality of the test, the resulting scores, and the interpretations based on the test scores” (American Educational Research Association [AERA], American Psychological Association [APA], & National Council on Measurement in Education [NCME], 2014, p.123).<sup>1</sup> Please note that a technical report, by design, addresses technical documentation of a testing program; other aspects of a testing program (content standards, scoring guides, guide to test interpretation, equating, etc.) are thoroughly addressed and referenced in supporting documents.

## 1.2 HISTORY

The Board of Regents adopted the Common Core State Standards (CCSS) for English Language Arts & Literacy and Mathematics at its July 2010 meeting and incorporated New York State-specific additions, creating the Common Core Learning Standards (CCLS), at its January 2011 meeting. Based on feedback from the field and to ensure adequate notice and time for students to be prepared to take the new Regents Exams measuring the CCLS, the Department provided an overlap in the administration of the Regents Exams measuring the 2005 Learning Standards with the Regents Exams measuring the CCLS and a phased-in sequence for mathematics.

Students who took the old Regents Exam in addition to the new Regents Exam were allowed to use the higher of the two scores for local transcript purposes, and, similarly, the higher of the two scores was used for institutional accountability for the 2015–2016 school year results. Such students were able to meet the mathematics exam requirement for graduation by passing either of these exams. The complete memo detailing transition to the Common Core examinations can be located at <http://www.p12.nysed.gov/assessment/commoncore/archive/transitionccregents1113rev-arc2.pdf>.

## 1.3 PURPOSES OF THE EXAM

The Regents Examination in Algebra II (Common Core) measures examinee achievement against the New York State (NYS) learning standards. The exam is prepared by teacher examination committees and New York State Education Department (NYSED) subject matter and testing specialists, and provides teachers and students with important information about student learning and performance against the established curriculum standards. Results of this

---

<sup>1</sup> References to specific *Standards* will be placed in parentheses throughout the technical report, to provide further context for each section.

exam may be used to identify student strengths and needs, in order to guide classroom teaching and learning. The exams also provide students, parents, counselors, administrators, and college admissions officers with objective and easily understood achievement information that may be used to inform empirically based educational and vocational decisions about students. As a state-provided objective benchmark, the Regents Examination in Algebra II (Common Core) is intended for use in satisfying state testing requirements for students who have finished a course in Algebra II. A passing score on the exam counts toward requirements for a high school diploma, as described in the New York State diploma requirements: <http://www.nysed.gov/common/nysed/files/programs/curriculum-instruction/currentdiplomarequirements2.pdf>. Results of the Regents Examination in Algebra II (Common Core) may also be used to satisfy various locally established requirements throughout the state.

### 1.4 TARGET POPULATION (STANDARD 7.2)

The examinee population for the Regents Examination in Algebra II (Common Core) is composed of students who have completed a course in Algebra II. Any student, regardless of grade level or cohort, who began their first commencement-level Algebra course in fall 2013 or later was provided with instruction aligned with the NYS P–12 Common Core Learning Standards for Algebra and, therefore, took or will take the Regents Examination in Algebra II (Common Core). More information about testing requirements can be found at <http://www.p12.nysed.gov/assessment/commoncore/transitionccregents1113rev.pdf>.

Table 1 provides a demographic breakdown of all students who took the June 2016 Regents Examination in Algebra II (Common Core). All analyses in this report are based on the population described in Table 1. Annual Regents Examination results in the New York State Report Cards are those reported in the Student Information Repository System (SIRS) as of the reporting deadline. The results include those exams administered in August, January, and June of the reporting year (see <http://data.nysed.gov/>). If a student takes the same exam multiple times in the year, only the highest score is included in these results. Item-level data used for the analyses in this report are reported by districts on a similar timeline, but through a different collection system. These data include all student results for each administration. Therefore, the n-sizes in this technical report will differ from publicly reported counts of student test-takers.

**Table 1 Total Examinee Population: Regents Examination in Algebra II (Common Core)**

Demographics	June Admin*	
	Number	Percent
All Students	91,478	100
<b>Race/Ethnicity</b>		
American Indian/Alaska Native	403	0.44
Asian/Native Hawaiian/Other Pacific Islander	13,394	14.64
Black/African American	9,176	10.03
Hispanic/Latino	13,296	14.54
Multiracial	1,300	1.42

	June Admin*	
Demographics	Number	Percent
White	53,902	58.93
<b>English Language Learner</b>		
No	90,428	98.85
Yes	1,050	1.15
<b>Economically Disadvantaged</b>		
No	60,560	66.20
Yes	30,918	33.80
<b>Gender</b>		
Female	49,025	53.60
Male	42,446	46.40
<b>Student with Disabilities</b>		
No	88,882	97.16
Yes	2,596	2.84

\*Note: Seven students were not reported in the Ethnicity and Gender group, but they are reflected in "All Students."



## Chapter 2: Classical Item Statistics (Standard 4.10)

---

This chapter provides an overview of the two most familiar item-level statistics obtained from classical item analysis: item difficulty and item discrimination. The following results pertain only to the operational Regents Examination in Algebra II (Common Core) items.

### 2.1 ITEM DIFFICULTY

At the most general level, an item's difficulty is indicated by its mean score in some specified group (e.g., grade level).

$$\bar{x} = \frac{1}{n} \cdot \sum_{i=1}^n x_i$$

In the mean score formula above, the individual item scores ( $x_i$ ) are summed and then divided by the total number of students ( $n$ ). For multiple-choice (MC) items, student scores are represented by 0s and 1s (0 = wrong, 1 = right). With 0–1 scoring, the equation above also represents the number of students correctly answering the item divided by the total number of students. Therefore, this is also the proportion correct for the item, or the  $p$ -value. In theory,  $p$ -values can range from 0.00 to 1.00 on the proportion-correct scale. For example, if an MC item has a  $p$ -value of 0.89, it means that 89 percent of the students answered the item correctly. Additionally, this value might also suggest that the item was relatively easy and/or that the students who attempted the item were relatively high achievers. For constructed-response (CR) items, mean scores can range from the minimum possible score (usually zero) to the maximum possible score (e.g., six points, in the case of some mathematics items). To facilitate average score comparability across MC and CR items, mean item performance for CR items is divided by the maximum score possible, so that the  $p$ -values for all items are reported as a ratio from 0.0 to 1.0.

Although the  $p$ -value statistic does not consider individual student ability in its computation, it provides a useful view of overall item difficulty and can provide an early and simple indication of items that are too difficult for the population of students taking the examination. Items with very high or very low  $p$ -values receive added scrutiny during all follow-up analyses, including item response theory analyses that factor student ability into estimates of item difficulty. Such items may be removed from the item pool during the test development process, as field testing typically reveals that they add very little measurement information. Items for the June 2016 Regents Examination in Algebra II (Common Core) show a range of  $p$ -values consistent with the targeted exam difficulty. Item  $p$ -values, presented in Table 2 and Table 3 for multiple-choice and constructed-response items, respectively, range from 0.11 to 0.76, with a mean of 0.49. Table 2 and Table 3 also show a standard deviation (SD) of item score and item mean (Table 3 only).

## 2.2 ITEM DISCRIMINATION

At the most general level, estimates of item discrimination indicate an item’s ability to differentiate between high and low performance on an item. It is expected that high-performing students (i.e., those who perform well on the Regents Examination in Algebra II [Common Core] overall) would be more likely to answer any given item correctly, while low-performing students (i.e., those who perform poorly on the exam overall) would be more likely to answer the same item incorrectly. Pearson’s product-moment correlation coefficient (also commonly referred to as a point-biserial correlation) between item scores and test scores is used to indicate discrimination (Pearson, 1896). The correlation coefficient can range from  $-1.0$  to  $+1.0$ . If high-scoring students tend to get the item right while low-scoring students do not, the correlation between the item score and the total test score will be both positive and noticeably large in its magnitude (i.e., above zero), meaning the item is likely discriminating well between high- and low-performing students. Point-biserials are computed for each answer option, including correct and incorrect options (commonly referred to as “distractors”). Finally, point-biserial values for each distractor are an important part of the analysis. The point-biserial values on the distractors are typically negative. Positive values can indicate that higher-performing students are selecting an incorrect answer or that the item key for the correct answer should be checked.

Refer to Table 2 and Table 3 for point-biserial values on the correct response and three distractors (Table 2 only). The values for correct answers are 0.20 or higher for all but one item (Item 5), indicating that the items are generally discriminating well between high- and low-performing examinees. Point-biserials for all distractors are negative, zero, or very close to zero, indicating that examinees are generally responding to the items as expected during item development.

**Table 2 Multiple-Choice Item Analysis Summary: Regents Examination in Algebra II (Common Core)**

Item	Number	$p$ -Value	SD	Point-Biserial	Point-Biserial Distractor 1	Point-Biserial Distractor 2	Point-Biserial Distractor 3
1	91,478	0.73	0.44	0.39	-0.22	-0.16	-0.22
2	91,478	0.65	0.48	0.37	-0.31	-0.13	-0.08
3	91,478	0.59	0.49	0.46	-0.23	-0.24	-0.22
4	91,478	0.79	0.41	0.37	-0.24	-0.16	-0.18
5	91,478	0.51	0.50	0.11	0.04	-0.20	-0.12
6	91,478	0.78	0.42	0.42	-0.23	-0.24	-0.19
7	91,478	0.43	0.50	0.20	-0.13	-0.21	-0.01
8	91,478	0.58	0.49	0.40	-0.22	-0.21	-0.16
9	91,478	0.67	0.47	0.41	-0.17	-0.31	-0.11
10	91,478	0.50	0.50	0.54	-0.19	-0.20	-0.32
11	91,478	0.54	0.50	0.26	-0.10	-0.12	-0.19
12	91,478	0.65	0.48	0.44	-0.21	-0.24	-0.22
13	91,478	0.61	0.49	0.41	-0.10	-0.34	-0.09

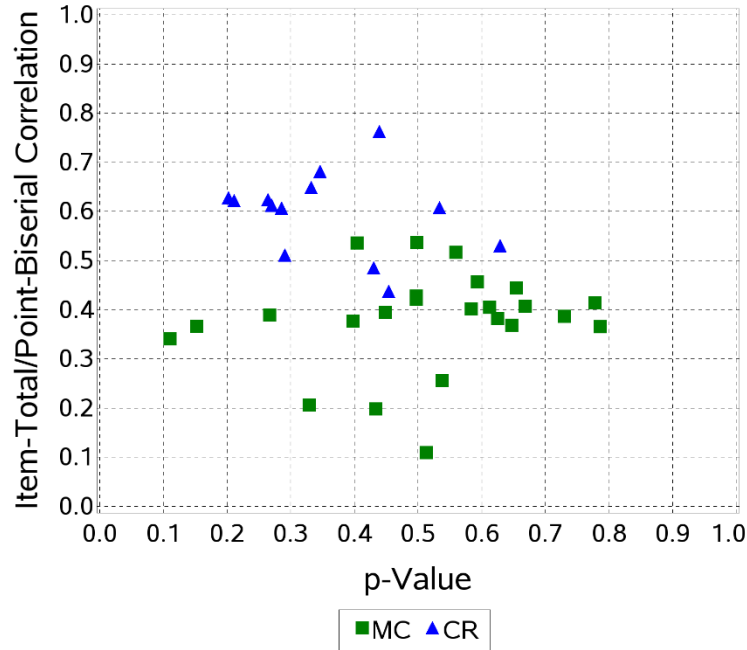
Item	Number	<i>p</i> -Value	SD	Point-Biserial	Point-Biserial Distractor 1	Point-Biserial Distractor 2	Point-Biserial Distractor 3
14	91,478	0.56	0.50	0.52	-0.17	-0.34	-0.20
15	91,478	0.63	0.48	0.38	-0.18	-0.15	-0.25
16	91,478	0.40	0.49	0.54	-0.11	-0.39	-0.14
17	91,478	0.45	0.50	0.39	-0.15	-0.14	-0.21
18	91,478	0.33	0.47	0.21	-0.21	-0.13	0.08
19	91,478	0.27	0.44	0.39	-0.34	-0.12	0.08
20	91,478	0.50	0.50	0.43	-0.17	-0.21	-0.21
21	91,478	0.15	0.36	0.37	-0.15	-0.16	0.03
22	91,478	0.50	0.50	0.42	-0.15	-0.18	-0.23
23	91,478	0.40	0.49	0.38	-0.10	-0.28	-0.13
24	91,478	0.11	0.31	0.34	-0.07	-0.20	0.01

**Table 3 Constructed-Response Item Analysis Summary: Regents Examination in Algebra II (Common Core)**

Item	Min. score	Max. score	Number of Students	Mean	SD	<i>p</i> -Value	Point-Biserial
25	0	2	91,478	1.26	0.87	0.63	0.53
26	0	2	91,478	0.91	0.75	0.45	0.44
27	0	2	91,478	1.07	0.91	0.53	0.61
28	0	2	91,478	0.57	0.74	0.29	0.61
29	0	2	91,478	0.86	0.86	0.43	0.49
30	0	2	91,478	0.42	0.76	0.21	0.62
31	0	2	91,478	0.58	0.78	0.29	0.51
32	0	2	91,478	0.53	0.79	0.26	0.62
33	0	4	91,478	1.08	1.49	0.27	0.61
34	0	4	91,478	1.33	1.33	0.33	0.65
35	0	4	91,478	0.81	1.28	0.20	0.63
36	0	4	91,478	1.39	1.68	0.35	0.68
37	0	6	91,478	2.63	2.26	0.44	0.76

### 2.3 DISCRIMINATION ON DIFFICULTY SCATTER PLOTS

Figure 1 shows a scatterplot of item discrimination values (*y*-axis) and item difficulty values (*x*-axis). The descriptive statistics of *p*-value and point-biserials, including mean, minimum, Q1, median, Q3, and maximum, are also presented in Table 4.



**Figure 1 Scatter Plot: Regents Examination in Algebra II (Common Core)**

**Table 4 Descriptive Statistics in  $p$ -value and Point-Biserial Correlation: Regents Examination in Algebra II (Common Core)**

Statistics	N	Mean	Min	Q1	Median	Q3	Max
$p$ -value	37	0.46	0.11	0.33	0.45	0.59	0.79
Point-Biserial	37	0.46	0.11	0.38	0.43	0.54	0.76

## 2.4 OBSERVATIONS AND INTERPRETATIONS

The  $p$ -values for the MC items ranged from about 0.10 to 0.80, while the mean proportion-correct values for the CR items (Table 3) ranged from about 0.20 to 0.60. From the difficulty distributions illustrated in the plot, a wide range of item difficulties appeared on each exam, which was one test development goal.

## Chapter 3: IRT Calibrations, Equating, and Scaling (Standards 2, and 4.10)

---

The item response theory (IRT) model used for the Regents Examination in Algebra II (Common Core) is based on the work of Georg Rasch (Rasch, 1960). The Rasch model has a long-standing presence in applied testing programs. IRT has several advantages over classical test theory, and has become the standard procedure for analyzing item response data in large-scale assessments. According to van der Linden and Hambleton (1997), “The central feature of IRT is the specification of a mathematical function relating the probability of an examinee’s response on a test item to an underlying ability.” Ability in this sense can be thought of as performance on the test and is defined as “the expected value of observed performance on the test of interest” (Hambleton, Swaminathan, and Roger, 1991). This performance value is often referred to as  $\theta$ . Performance and  $\theta$  will be used interchangeably throughout the remainder of this report.

A fundamental advantage of IRT is that it links examinee performance and item difficulty estimates and places them on the same scale, allowing for an evaluation of examinee performance that considers the difficulty of the test. This is particularly valuable for final test construction and test form equating, as it facilitates a fundamental attention to fairness for all examinees across items and test forms.

This chapter outlines the procedures used for calibrating the operational Regents Examination in Algebra II (Common Core) items. Generally, item calibration is the process of assigning a difficulty, or item “location,” estimate to each item on an assessment so that all items are placed onto a common scale. This chapter briefly introduces the Rasch model, reports the results from evaluations of the adequacy of the Rasch assumptions, and summarizes the Rasch item statistics.

### 3.1 DESCRIPTION OF THE RASCH MODEL

The Rasch model (Rasch, 1960) was used to calibrate multiple-choice items, and the partial credit model, or PCM (Wright and Masters, 1982), was used to calibrate constructed-response items. The PCM extends the Rasch model for dichotomous (0, 1) items so that it accommodates the polytomous CR item data. Under the PCM model, for a given item  $i$  with  $m_i$  score categories, the probability of person  $n$  scoring  $x$  ( $x = 0, 1, 2, \dots, m_i$ ) is given by

$$P_{ni}(X = x) = \frac{\exp \sum_{j=0}^x (\theta_n - D_{ij})}{\sum_{k=0}^{m_i} \exp \sum_{j=0}^k (\theta_n - D_{ij})}$$

where  $\theta_n$  represents examinee ability, and  $D_{ij}$  is the step difficulty of the  $j^{\text{th}}$  step on item  $i$ .  $D_{ij}$  can be expressed as  $D_{ij} = D_i - F_{ij}$ , where  $D_i$  is the difficulty for item  $i$  and  $F_{ij}$  is a step deviation value for the  $j^{\text{th}}$  step. For dichotomous MC items, the RPCM reduces to the standard Rasch model and the single step difficulty is referred to as the item’s difficulty. The Rasch model predicts the probability of person  $n$  getting item  $i$  correct as follows:

$$P_{ni}(X = 1) = \frac{\exp(\theta_n - D_{ij})}{1 + \exp(\theta_n - D_{ij})}$$

The Rasch model places both performance and item difficulty (estimated in terms of log-odds or logits) on the same continuum. When the model assumptions are met, the Rasch model provides estimates of examinee performance and item difficulty that are theoretically invariant across random samples of the same examinee population.

### **3.2 SOFTWARE AND ESTIMATION ALGORITHM**

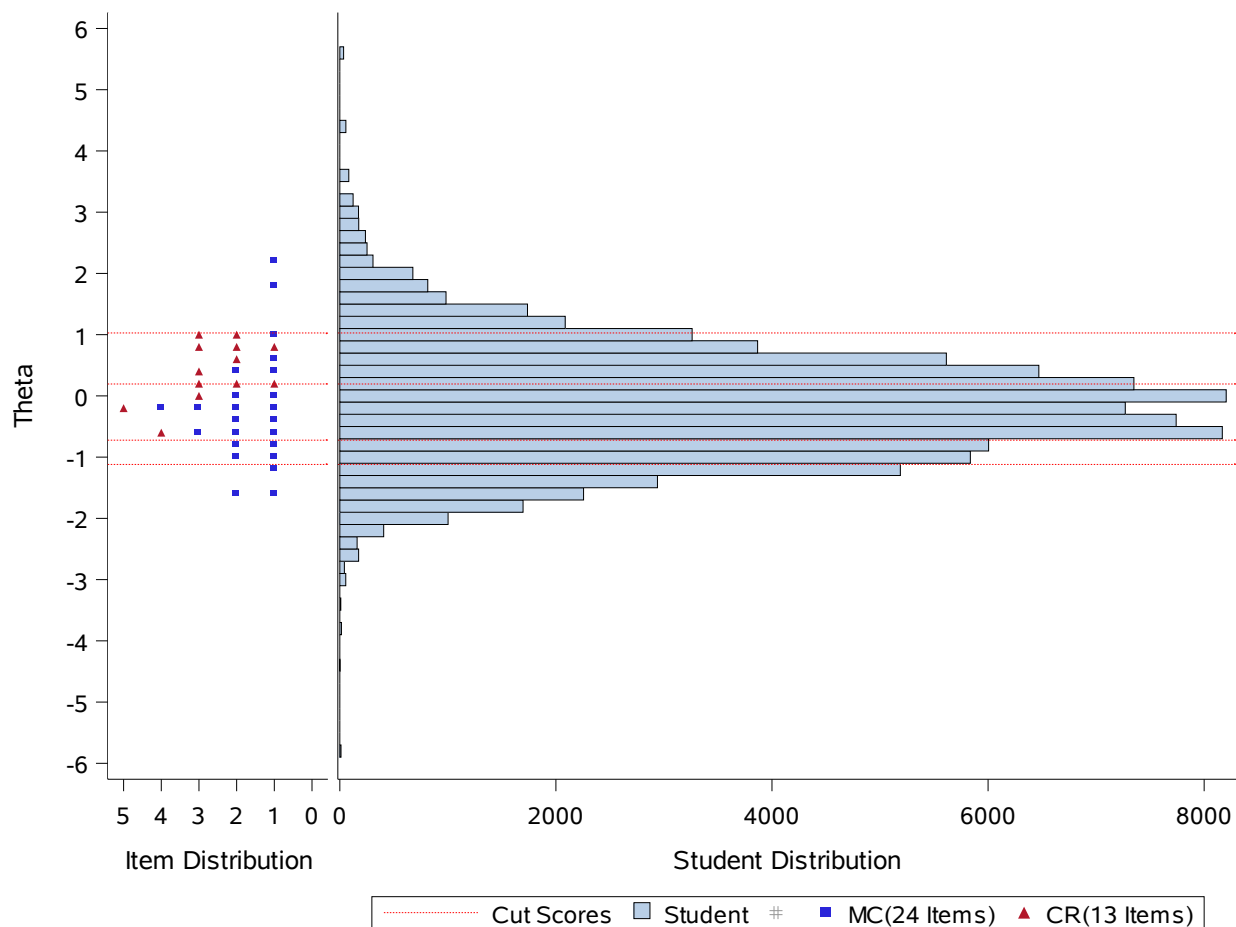
Item calibration was implemented via the WINSTEPS 3.60 computer program (Wright and Linacre, 2015), which employs unconditional (UCON), joint maximum likelihood estimation (JMLE).

### **3.3 CHARACTERISTICS OF THE TESTING POPULATION**

The data analyses reported here are based on all students who took the Regents Examination in Algebra II (Common Core) in the June 2016 administration. The characteristics of this population are provided in Table 1.

### **3.4. ITEM DIFFICULTY-STUDENT PERFORMANCE MAPS**

The distributions of the Rasch item logits (item difficulty estimates) and student performance are shown on the item difficulty-student performance map presented in Figure 2. This graphic illustrates the location of student performance and item difficulty on the same scale, along with their respective distributions and cut scores (indicated by the horizontal dotted lines). The figure shows more difficult items and higher examinee performance at the top and lower performance and easier items at the bottom.



**Figure 2 Student Performance Map: Regents Examination in Algebra II (Common Core)**

### 3.5 CHECKING RASCH ASSUMPTIONS

Since the Rasch model was the basis of all calibration, scoring, and scaling analyses associated with the Regents Examination in Algebra II (Common Core), the validity of the inferences from these results depends on the degree to which the assumptions of the model were met and how well the model fits the test data. Therefore, it is important to check these assumptions. This section evaluates the dimensionality of the data, local item independence, and item fit. It should be noted that only operational items were analyzed, since they are the basis of student scores.

#### *Unidimensionality*

Rasch models assume that one dominant dimension determines the differences between students' performances. Principal Components Analysis (PCA) can be used to assess the unidimensionality assumption. The purpose of the analysis is to verify if other dominant components exist among the items. If any other dimensions are found, the unidimensionality assumption would be violated.

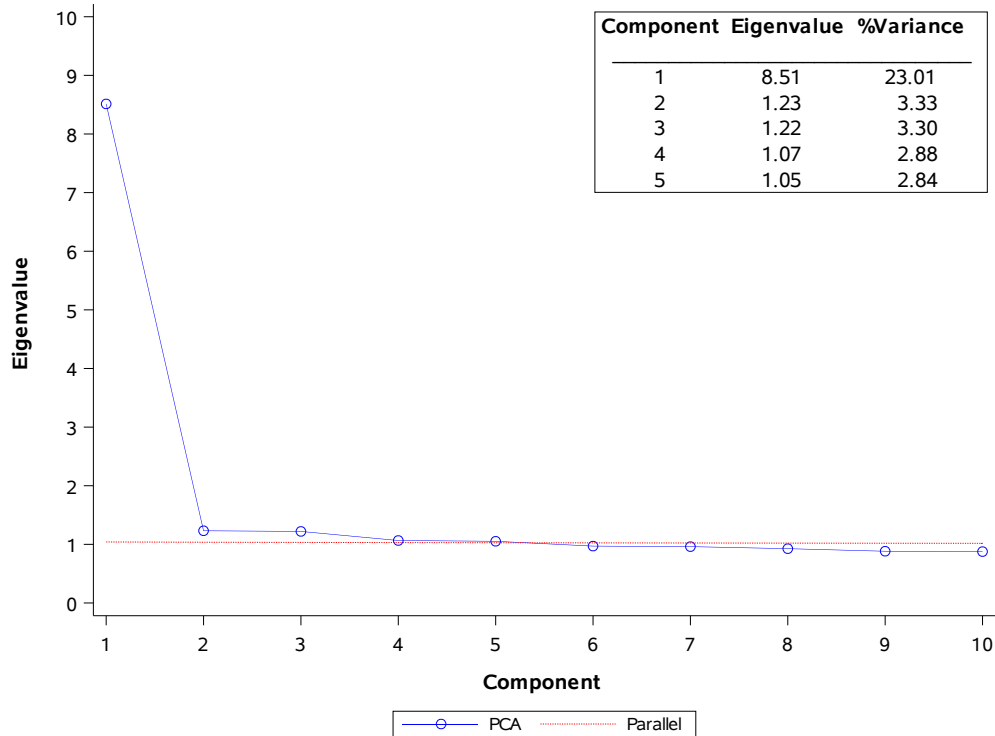
A parallel analysis (Horn, 1965) was conducted to help distinguish components that are real from components that are random. Parallel analysis is a technique to decide how many factors exist in principal components. For the parallel analysis, 100 random data sets of sizes equal to the original data were created. For each random data set, a PCA was performed and the resulting eigenvalues stored. Then, for each component, the upper 95th percentile value of the distribution of the 100 eigenvalues from the random data sets was plotted. Given the size of the data generated for the parallel analysis, the reference line is essentially equivalent to plotting a reference line for an eigenvalue of 1.

Figure 3 shows the PCA results for the Regents Examination in Algebra II (Common Core). The results include the eigenvalues and the percentage of variance explained for the first five components, as well as the scree plots. The scree plots show the eigenvalues plotted by component number and the results of a parallel analysis. Although the total number of components in PCA is same as the total number of items in a test, Figure 3 shows only the first 10 components. This view is sufficient for interpretation because components are listed in descending eigenvalue order. The fact that the eigenvalues for components 2 through 10 are much lower than the first component demonstrates that components beyond the first one are not, individually, adding much to the explanation of variance in the data.

As rule of thumb, Reckase (1979) proposed that the variance explained by the primary dimension should be greater than 20 percent, to indicate unidimensionality. However, as this rule is not absolute, it is helpful to consider three additional characteristics of the PCA and parallel analysis results: 1) whether or not the ratio of the first to the second eigenvalue is greater than 3, 2) whether the second value is not much larger than the third value, and 3) whether the second value is not significantly different from those from the parallel analysis.

As shown in Figure 3, the primary dimension explained 23.01 percent of the total variance for the Regents Examination in Algebra II (Common Core). The eigenvalue of the second dimension is less than one third of the first, at 1.23, and the second value is not significantly different from the parallel analysis. Overall, the PCA suggests that the test is reasonably unidimensional.





**Figure 3 Scree Plot: Regents Examination in Algebra II (Common Core)**

*Local Independence*

Local independence (LI) is a fundamental assumption of IRT. This means that, for statistical purposes, an examinee’s response to any one item should not depend on the examinee’s response to any other item on the test. In formal statistical terms, a test X that is comprised of items X1, X2, …Xn is locally independent with respect to the latent variable  $\theta$  if, for all  $x = (x_1, x_2, \dots, x_n)$  and  $\theta$ ,

$$P(\mathbf{X} = \mathbf{x} | \theta) = \prod_{i=1}^n P(X_i = x_i | \theta)$$

This formula essentially states that the probability of any pattern of responses across all items (x), after conditioning on the examinee’s true score ( $\theta$ ) as measured by the test, should be equal to the product of the conditional probabilities across each item (i.e., the multiplication rule for independent events where the joint probabilities are equal to the product of the associated marginal probabilities).

The equation above shows the condition after satisfying the strong form of local independence. A weak form of local independence (WLI) is proposed by McDonald (1979). The distinction is important because many indicators of local dependency are actually framed by WLI. For WLI, the conditional covariances of all pairs of item responses, conditioned on the abilities, are assumed to be equal to zero. When this assumption is met, the joint probability of responses to an item pair, conditioned on the abilities, is the product of the probabilities of

responses to these two items, as shown below. Based on the WLI, the following expression can be derived:

$$P(X_i = x_i, X_j = x_j | \theta) = P(X_i = x_i | \theta)P(X_j = x_j | \theta)$$

Marais and Andrich (2008) point out that local item dependence in the Rasch model can occur in two ways that may be difficult to distinguish. The first way occurs when the assumption of unidimensionality is violated. Here, other nuisance dimensions besides a dominant dimension determine student performance (this can be called “trait dependence”). The second way occurs when responses to an item depend on responses to another item. This is a violation of statistical independence and can be called response dependence. By distinguishing the two sources of local dependence, one can see that, while local independence can be related to unidimensionality, the two are different assumptions and therefore require different tests.

Residual item correlations provided in WINSTEPS for each item pair were used to assess the local dependence between the Regents Examination in Algebra II (Common Core) items. In general, these residuals are computed as follows. First, expected item performance based on the Rasch model is determined using  $(\theta)$  and item parameter estimates. Next, deviations (residuals) between the examinees’ expected and observed performance is determined for each item. Finally, for each item pair, a correlation between the respective deviations is computed.

Three types of residual correlations are available in WINSTEPS: raw, standardized, and logit. It is noted that the raw score residual correlation essentially corresponds to Yen’s Q3 index, a popular statistic used to assess local independence. The expected value for the Q3 statistic is approximately  $-1/(k - 1)$  when no local dependence exists, where  $k$  is test length (Yen, 1993). Thus, the expected Q3 values should be approximately  $-0.03$  for the items on the exam. Index values that are greater than 0.20 indicate a degree of local dependence that probably should be examined by test developers (Chen & Thissen, 1997).

Since the three residual correlations are very similar, the default “standardized residual correlation” in WINSTEPS was used for these analyses. Table 5 shows the summary statistics — mean, standard deviation, minimum, maximum, and several percentiles (P10, P25, P50, P75, P90) — for all the residual correlations for each test. The total number of item pairs (N) and the number of pairs with residual correlations greater than 0.20 are also reported in this table. There were two item pairs with residual correlations greater than 0.20. The mean residual correlations were slightly negative and the values were close to  $-0.02$ . The vast majority of the correlations were very small, suggesting that local item independence generally holds for the Regents Examination in Algebra II (Common Core).

**Table 5 Summary of Item Residual Correlations: Algebra II (Common Core)**

Statistic Type	Value
N	666
Mean	-0.02
SD	0.03

Minimum	-0.11
P <sub>10</sub>	-0.06
P <sub>25</sub>	-0.05
P <sub>50</sub>	-0.02
P <sub>75</sub>	0.00
P <sub>90</sub>	0.02
Maximum	0.13
> 0.20	0

---

### Item Fit

An important assumption of the Rasch model is that the data for each item fit the model. WINSTEPS provides two item fit statistics (INFIT and OUTFIT) for evaluating the degree to which the Rasch model predicts the observed item responses for a given set of test items. Each fit statistic can be expressed as a mean square (MnSq) statistic or on a standardized metric (Zstd with mean = 0 and variance = 1). MnSq values are more oriented toward practical significance, while Zstd values are more oriented toward statistical significance. INFIT MnSq values are the average of standardized residual variance (the difference between the observed score and the Rasch-estimated score divided by the square root of the Rasch-model variance). The INFIT statistic is weighted by the ( $\theta$ ) relative to item difficulty, and tends to be affected more by unexpected responses close to the person, item, or rating scale category measure (i.e., informative, on-target responses).

The expected MnSq value is 1.0 and can range from 0 to infinity. Deviation in excess of the expected value can be interpreted as noise, or lack of fit between the items and the model. Values lower than the expected value can be interpreted as item redundancy or overfitting items (too predictable, too much redundancy), and values greater than the expected value indicate underfitting items (too unpredictable, too much noise). Rules of thumb regarding “practically significant” MnSq values vary. Table 6 presents the summary statistics of INFIT mean square statistics for the Regents Examination in Algebra II (Common Core), including the mean, standard deviation, and minimum and maximum values.

The number of items within a targeted range of [0.7, 1.3] is also reported in Table 6. The mean INFIT value is 1.00, with 37 of the 37 items falling in a targeted range of [0.7, 1.3]. As the range of [0.7, 1.3] is used as a guide for ideal fit, fit values outside of the range are considered individually. Overall, these results indicate that, for all items, the Rasch model fits the Regents Examination in Algebra II (Common Core) item data well.

**Table 6 Summary of INFIT Mean Square Statistics: Algebra II (Common Core)**

	INFIT Mean Square					
	N	Mean	SD	Min	Max	[0.7, 1.3]
Algebra II (Common Core)	37	1.00	0.11	0.81	1.28	[37/37]

Items for the Regents Examination in Algebra II (Common Core) were field tested in 2015.

### 3.6 SCALING OF OPERATIONAL TEST FORMS

Operational test items were selected based on content coverage, content accuracy, and statistical quality. The sets of items on each operational test conformed to the coverage determined by content experts working from the learning standards established by the New York State Education Department and explicated in the test blueprint. Each item's classical and Rasch statistics were used to assess item quality. Items were selected to vary in difficulty to accurately measure students' abilities across the ability continuum. Appendix A contains the operational test map for the June 2016 administration. Note that statistics presented in the test map were generated based on the field test data.

All Regents examinations are pre-equated, meaning that the parameters used to derive the relationship between the raw and scale scores are estimated prior to the construction and administration of the operational form. These field tests are administered to as small a sample of students as possible to minimize the effect on student instructional time across the state. The small  $n$ -counts associated with such administrations are sufficient for reasonably accurate estimation of most items' parameters; however, for the six-point essay item, its parameters can be unstable when estimated across as small a sample as is typically used. Therefore, a set of constants is used for these items' parameters on operational examinations. These constants were set by the NYSED and are based on the values in the bank for all constructed response items. For Algebra II (Common Core) examination, there is only one six-point item with fixed constants as follows:  $D = -0.06$ ,  $F_0 = 0.00$ ,  $F_1 = -0.73$ ,  $F_2 = 0.59$ ,  $F_3 = -0.09$ ,  $F_4 = -0.15$ ,  $F_5 = 0.17$ , and  $F_6 = 0.21$ .

The New York State Regents Examination in Algebra II (Common Core) has four cut scores, which are set at the scale scores of 55, 65, 78 (floating), and 85. One of the primary considerations during test construction was to select items so as to minimize changes in the raw scores corresponding to these scale scores. Maintaining a consistent mean Rasch difficulty level from administration to administration facilitates this. For this assessment, the target value for the mean Rasch difficulty was set at 0.047. It should be noted that the raw scores corresponding to the scale score cut scores may still fluctuate, even if the mean Rasch difficulty level is maintained at the target value, due to differences in the distributions of the Rasch difficulty values among the items from administration to administration.

The relationship between raw and scale scores is explicated in the scoring tables for each administration. The table for the June 2016 administration can be found in Appendix B. This table is the end product of the following scaling procedure.

All Regents examinations are equated back to a base scale, which is held constant from year to year. Specifically, they are equated to the base scale through the use of a calibrated item pool. The Rasch difficulties from the items' initial administration in a previous year's field test are used to equate the scale for the current administration to the base administration. For this examination, the base administration was the June 2016 administration. Scale scores for the future administrations will be on the same scale, and can be directly compared to scale scores on all previous administrations back to the June 2016 administration.

When the base administration was concluded, the initial raw score to scale score relationship was established. Three raw scores were fixed at specific scale scores. Scale scores of 0 and 100 were fixed to correspond to the minimum and maximum possible raw scores. In addition, a standard setting had been held to determine the passing and passing with distinction cut scores in the raw score metric. The scale score points of 65, 78, and 85 were set to correspond to those raw score cuts. A fourth-degree polynomial is required to fit a line exactly to five arbitrary points (e.g., the raw scores corresponding to the five critical scale scores of 0, 65, 78, 85, and 100). The general form of this best-fitting line is:

$$SS = m4 * RS4 + m3 * RS3 + m2 * RS2 + m1 * RS1 + m0,$$

where SS is the scaled score, RS is the raw score, and m0 through m4 are the transformation constants that convert the raw score into the scale score (please note that m0 will always be equal to zero in this application, since a raw score of zero corresponds to a scale score of zero). A subscript for a person on both dependent and independent variables is not present for simplicity. The above relationship and the values of m1 to m4 specific to this subject were then used to determine the scale scores corresponding to the remainder of the raw scores on the examination. This initial relationship between the raw and scale scores became the base scale.

The Rasch difficulty parameters for the items on the base form were then used to derive a raw score-to-Rasch student ability (theta score) relationship. This allowed the relationship between the Rasch theta score and the scale score to be known, mediated through their common relationship with the raw scores.

In succeeding years, each test form was selected from the pool of items that had been tested in previous years' field tests, each of which had known Rasch item difficulty parameter(s). These known parameters were then used to construct the relationship between the raw and Rasch theta scores for that particular form. Because the Rasch difficulty parameters are all on a common scale, the Rasch theta scores were also on a common scale with previously administered forms. The remaining step in the scaling process was to find the scale score equivalent for the Rasch theta score corresponding to each raw score point on the new form, using the theta-to-scale score relationship established in the base year. This was done via linear interpolation.

This process results in a relationship between the raw scores on the form and the overall scale scores. The scale scores corresponding to each raw score are then rounded to the nearest integer for reporting on the conversion chart (posted at the close of each administration). The only exceptions are for the minimum and maximum raw scores and the raw scores that correspond to the scaled cut scores of 55, 65, 78, and 85.

The minimum (zero) and maximum possible raw scores are assigned scale scores of 0 and 100, respectively. In the event that there are raw scores less than the maximum with scale scores that round to 100, their scale scores are set equal to 99. A similar process is followed with the minimum score; if any raw scores other than zero have scale scores that round to zero, their scale scores are instead set equal to one.

With regard to the cuts, if two or more scale scores round to 55, 65, or 85, the lowest raw score's scale score is set equal to 55, 65, or 85, and the scale scores corresponding to the higher raw scores are set to 56, 66, or 86, as appropriate. This rule does not apply for the third cut at a scale score of 78. If no scale score rounds to these four critical cuts, then the raw score with the largest scale score that is less than the cut is set equal to the cut. The overarching principle when two raw scores both round to either scale score cut is that the lower of the raw scores is always assigned to be equal to the cut so that students are never penalized for this ambiguity.

## Chapter 4: Reliability (Standard 2)

---

Test reliability is a measure of the internal consistency of a test (Cronbach, 1951). It is a measure of the extent to which the items on a test provide consistent information about student mastery of a domain. Reliability should, ultimately, demonstrate that examinee score estimates maximize consistency and, therefore, minimize error, or, theoretically speaking, that examinees who take a test multiple times would get the same score each time.

According to the *Standards for Educational and Psychological Testing*, “A number of factors can have significant effects on reliability/precision, and in some cases, these factors can lead to misinterpretations of test scores, if not taken into account” (AERA et al., 2014, p. 38). First, test length and the variability of observed scores can both influence reliability estimates. Tests with fewer items or with a lack of heterogeneity in scores tend to produce lower reliability estimates. Second, reliability is specifically concerned with random sources of error. Accordingly, the degree of inconsistency due to random error sources is what determines reliability: less consistency is associated with lower reliability, and more consistency is associated with higher reliability. Of course, systematic error sources also exist.

The remainder of this chapter discusses reliability results for the Regents Examination in Algebra II (Common Core) and three additional statistical measures to address the multiple factors affecting an interpretation of the Exam’s reliability:

- standard errors of measurement
- decision consistency
- group means

### 4.1 RELIABILITY INDICES (STANDARD 2.20)

Classical test theory describes reliability as a measure of the internal consistency of test scores. The reliability ( $\rho_X^2$ ) is defined as the ratio of true score variance ( $\sigma_T^2$ ) to the observed score variance ( $\sigma_X^2$ ), as presented in the equation below. The total variance contains two components: 1) the variance in true scores and 2) the variance due to the imperfections in the measurement process ( $\sigma_E^2$ ). Put differently, total variance equals true score variance plus error variance.<sup>2</sup>

$$\rho_X^2 = \frac{\sigma_T^2}{\sigma_X^2} = \frac{\sigma_T^2}{\sigma_T^2 + \sigma_E^2}$$

Reliability coefficients indicate the degree to which differences in test scores reflect true differences in the attribute being tested rather than random fluctuations. Total test score variance (i.e., individual differences) is partly due to real differences in the construct (true variance) and partly due to random error in the measurement process (error variance).

Reliability coefficients range from 0.0 to 1.0. The index will be 0.0 if none of the test score variances is true. If all test score variances were true, the index would equal 1.0. Such scores

---

<sup>2</sup> A covariance term is not required, as true scores and error are assumed to be uncorrelated in classical test theory.

would be pure random noise (i.e., all measurement error). If the index achieved a value of 1.0, scores would be perfectly consistent (i.e., contain no measurement error). Although values of 1.0 are never achieved in practice, it is clear that larger coefficients are more desirable because they indicate that the test scores are less influenced by random error.

### *Coefficient Alpha*

Reliability is most often estimated using the formula for Coefficient Alpha, which provides a practical internal consistency index. Coefficient Alpha can be conceptualized as the extent to which an exchangeable set of items from the same domain would result in a similar rank ordering of students. Note that relative error is reflected in this index. Excessive variation in student performance from one sample of items to the next should be of particular concern for any achievement test user.

A general computational formula for Coefficient Alpha is as follows:

$$\alpha = \frac{N}{N-1} \left( 1 - \frac{\sum_{i=1}^N \sigma_{Yi}^2}{\sigma_X^2} \right),$$

where  $N$  is the number of parts (items),  $\sigma_X^2$  is the variance of the observed total test scores, and  $\sigma_{Yi}^2$  is the variance of part  $i$ .

## **4.2 STANDARD ERROR OF MEASUREMENT (STANDARDS 2.13, 2.14, 2.15)**

Reliability coefficients best reflect the extent to which measurement inconsistencies may be present or absent. The standard error of measurement (SEM) is another indicator of test score precision that is better suited for determining the effect of measurement inconsistencies for the scores obtained by individual examinees. This is particularly so for conditional SEMs (CSEMs), discussed further below.

### *Traditional Standard Error of Measurement*

The standard error of measurement is defined as the standard deviation of the distribution of observed scores for students with identical true scores. Because the SEM is an index of the random variability in test scores in test score units, it represents important information for test score users.

The SEM formula is provided below.

$$SEM = SD\sqrt{1 - \alpha}$$

This formula indicates that the value of the SEM depends on both the reliability coefficient (the Coefficient Alpha, as detailed previously) and the standard deviation of test scores. If the reliability were equal to 0.00 (the lowest possible value), the SEM would be equal to the standard deviation of the test scores. If test reliability were equal to 1.00 (the highest possible value), the SEM would be 0.0. In other words, a perfectly reliable test has no measurement error (Harvill, 1991). Additionally, the value of the SEM takes the group variation (i.e., score standard deviation) into account. Consider that a SEM of 3 on a 10-point test would be very different from a SEM of 3 on a 100-point test.



### *Traditional Standard Error of Measurement Confidence Intervals*

The SEM is an index of the random variability in test scores reported in actual score units, which is why it has such great utility for test score users. SEMs allow statements regarding the precision of individual test scores. SEMs help place “reasonable limits” (Gulliksen, 1950) around observed scores, through construction of an approximate score band. Often referred to as confidence intervals, these bands are constructed by taking the observed scores,  $X$ , and adding and subtracting a multiplicative factor of the SEM. As an example, students with a given true score will have observed scores that fall between  $\pm 1$  SEM about two-thirds of the time.<sup>3</sup> For  $\pm 2$  SEM confidence intervals, this increases to about 95 percent.

The Coefficient Alpha and associated SEM for the Regents Examination in Algebra II (Common Core) are provided in Table 7.

**Table 7 Reliabilities and Standard Errors of Measurement: Regents Examination in Algebra II (Common Core)**

Subject	Coefficient Alpha	SEM
Algebra II (Common Core)	0.89	5.83

Assuming normally distributed scores, one would expect about two-thirds of the observations to be within one standard deviation of the mean. An estimate of the standard deviation of the true scores can be computed as

$$\hat{\sigma}_T = \sqrt{\hat{\sigma}_x^2 - \hat{\sigma}_x^2(1 - \hat{\rho}_{xx})}$$

### *Conditional Standard Error of Measurement*

Every time that an assessment is administered, the score that the student receives contains some error. If the same exam were administered an infinite number of times to the same student, the mean of the distribution of the student’s raw scores would be equal to their true score ( $\theta$ , the score obtained with no error), and the standard deviation of the distribution of their raw scores would be the conditional standard error. Since there is a one-to-one correspondence between the raw score and  $\theta$  in the Rasch model, we can apply this concept more generally to all students who obtained a particular raw score and calculate the probability of obtaining each possible raw score, given the student’s estimated  $\theta$ . The standard deviation of this conditional distribution is defined as the conditional standard error of measurement (CSEM). The computer program POLYCSEM (Kolen, 2004) was used to carry out the mechanics of this computation.

The relationship between  $\theta$  and the scale score is not expressible in a simple mathematical form because it is a blend of the third-degree polynomial relationship between the raw and scale scores and the nonlinear relationship between the expected raw and  $\theta$  scores. In addition,

<sup>3</sup> Some prefer the following interpretation: if a student were tested an infinite number of times, the  $\pm 1$  SEM confidence intervals constructed for each score would capture the student’s true score 68 percent of the time.

as the exam is equated from year to year, the relationship between the raw and scale scores moves away from the original fourth-degree polynomial relationship to one that is also no longer expressible in simple mathematical form. In the absence of a simple mathematical relationship between  $\theta$  and the scale scores, the CSEMs that are available for each  $\theta$  score via Rasch IRT cannot be converted directly to the scale score metric.

The use of Rasch IRT to scale and equate the Regents Exams does, however, make it possible to calculate CSEMs by using the procedures described by Kolen, Zeng, and Hanson (1996) for dichotomously scored items and extended by Wang, Kolen, and Harris (2000) to polytomously scored items. For tests such as the Regents Examination in Algebra II (Common Core) that do not have a one-to-one relationship between raw and scale scores, the CSEM for each achievable scale score can be calculated using the compound multinomial distribution to represent the conditional distribution of raw scores for each level of  $\theta$ .

Consider an examinee with a certain performance level. If it were possible to measure this examinee's performance perfectly, without any error, this measure could be called the examinee's "true score," as discussed earlier. This score is equal to the expected raw score. However, whenever an examinee takes a test, their observed test score always includes some level of measurement error. Sometimes, this error is positive, and the examinee achieves a higher score than would be expected given their level of  $\theta$ ; other times, it is negative, and the examinee achieves a lower-than-expected score. If we could give an examinee the same test multiple times and record their observed test scores, the resulting distribution would be the conditional distribution of raw scores for that examinee's level of  $\theta$  with a mean value equal to the examinee's expected raw (true) score. The CSEM for that level of  $\theta$  in the raw score metric is the square root of the variance of this conditional distribution.

The conditional distribution of raw scores for any level of  $\theta$  is the compound multinomial distribution (Wang et al., 2000). An algorithm to compute this can be found in Hanson (1994) and Thissen, Pommerich, Billeaud, and Williams (1995) and is also implemented in the computer program POLYCSEM (Kolen, 2004). The compound multinomial distribution yields the probabilities that an examinee with a given level of  $\theta$  has of achieving each achievable raw (and accompanying scale) score. The point values associated with each achievable raw or scale score point can be used to calculate the mean and variance of this distribution in the raw or scale score metric, respectively; the square root of the variance is the CSEM of the raw or scale score point associated with the current level of  $\theta$ .

#### *Conditional Standard Error of Measurement Confidence Intervals*

CSEMs allow statements regarding the precision of individual tests scores. Like SEMs, they help place reasonable limits around observed scaled scores through the construction of an approximate score band. The confidence intervals are constructed by adding or subtracting a multiplicative factor of the CSEM.

#### *Conditional Standard Error of Measurement Characteristics*

The relationship between the scale score CSEM and  $\theta$  depends both on the nature of the raw-to-scale score transformation (Kolen and Brennan, 2005; Kolen and Lee, 2011) and on whether the CSEM is derived from the raw scores or from  $\theta$  (Lord, 1980). The pattern of CSEMs for raw scores and linear transformations of the raw score tend to have a characteristic

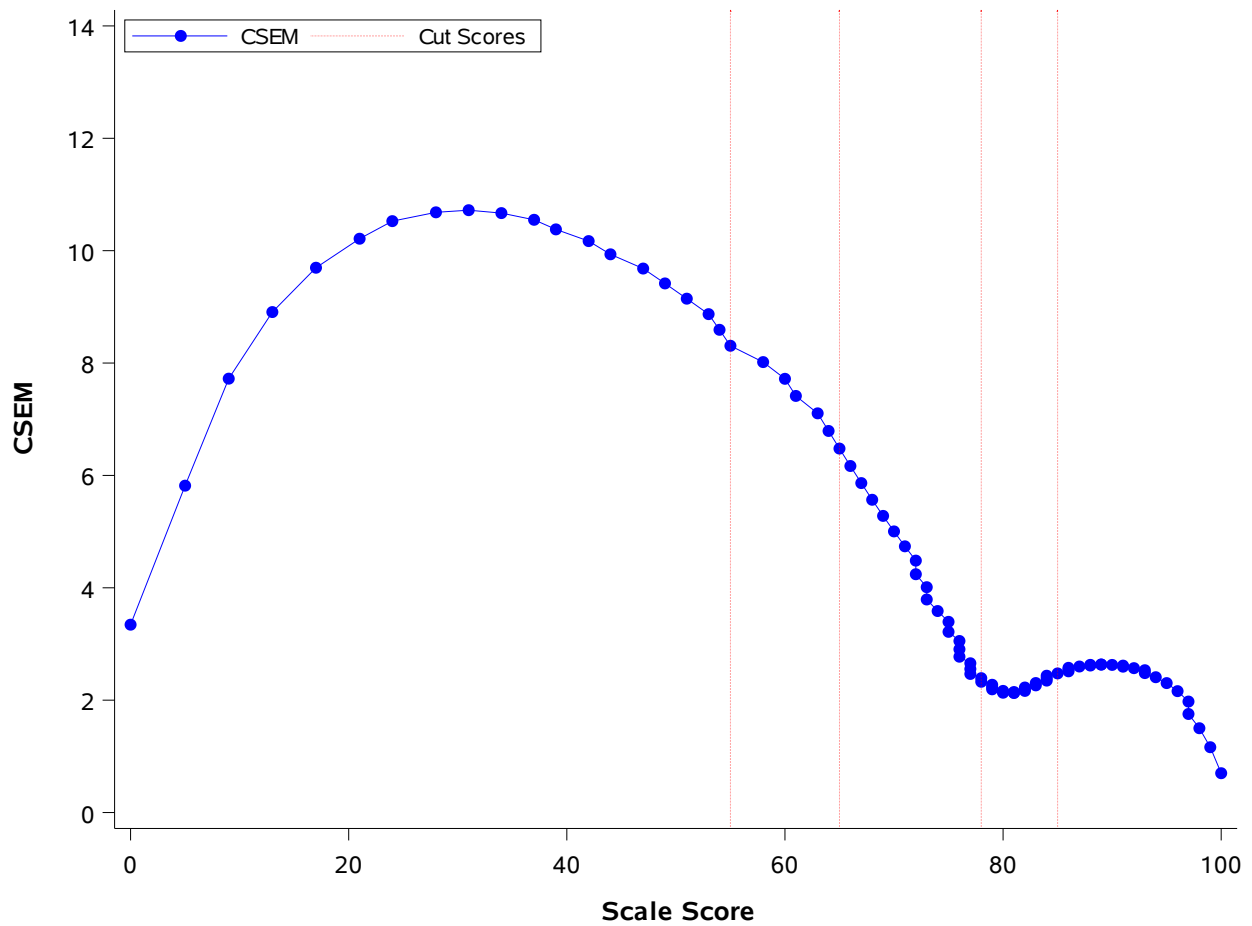
“inverted-U” shape, with smaller CSEMs at the ends of the score continuum and larger CSEMs toward the middle of the distribution.

Achievable raw score points for these distributions are spaced equally across the score range. Kolen and Brennan (2005, p. 357) state, “When, relative to raw scores, the transformation compresses the scale in the middle and stretches it at the ends, the pattern of the conditional standard errors of measurement will be concave up (U-shaped), even though the pattern for the raw scores was concave down (inverted-U shape).”

### *Results and Observations*

The relationship between raw and scale scores for the Regents Exams tends to be roughly linear from scale scores of 0 to 65 and then concave down from about 65 to 100. In other words, the scale scores track linearly with the raw scores for the lower two-thirds of the scale score range and then are compressed relative to the raw scores for the remaining one-third of the range, though there are variations. The CSEMs for the Regents Exams can be expected to have inverted-U shaped patterns, with some variations.

Figure 4 shows this type of CSEM variation for the Regents Examination in Algebra II (Common Core) where the compression of raw score to scale scores around the cut score of 65 changes the shape of the curve very noticeably. This type of expansion and compression can be seen in Figure 4 by looking at the changing density of raw score points along the scale score range on the horizontal axis. Specifically, the largest compression can be seen between about 65 to 90 scale score points.



**Figure 4 Conditional Standard Error Plot: Regents Examination in Algebra II (Common Core)**

**4.3 DECISION CONSISTENCY AND ACCURACY (STANDARD 2.16)**

In a standards-based testing program there is interest in knowing how accurately students are classified into performance categories. In contrast to the Coefficient Alpha, which is concerned with the relative rank-ordering of students, it is the absolute values of student scores that are important in decision consistency and accuracy.

Classification consistency refers to the degree to which the achievement level for each student can be replicated upon retesting using an equivalent form (Huynh, 1976). Decision consistency answers the following question: What is the agreement in classifications between the two non-overlapping, equally difficult forms of the test? If two parallel forms of the test were given to the same students, the consistency of the measure would be reflected by the extent to which the classification decisions based on the first set of test scores matched the decisions based on the second set of test scores. Consider the tables below.

		TEST ONE		
		LEVEL I	LEVEL II	MARGINAL
TEST TWO	LEVEL I	$\phi_{11}$	$\phi_{12}$	$\phi_{1\bullet}$
	LEVEL II	$\phi_{21}$	$\phi_{22}$	$\phi_{2\bullet}$
	MARGINAL	$\phi_{\bullet 1}$	$\phi_{\bullet 2}$	1

**Figure 5 Pseudo-Decision Table for Two Hypothetical Categories**

		TEST ONE				
		LEVEL I	LEVEL II	LEVEL III	LEVEL IV	MARGINAL
TEST TWO	LEVEL I	$\phi_{11}$	$\phi_{12}$	$\phi_{13}$	$\phi_{14}$	$\phi_{1\bullet}$
	LEVEL II	$\phi_{21}$		$\phi_{23}$	$\phi_{24}$	$\phi_{2\bullet}$
	LEVEL III	$\phi_{31}$	$\phi_{32}$		$\phi_{34}$	$\phi_{3\bullet}$
	LEVEL IV	$\phi_{41}$	$\phi_{42}$	$\phi_{43}$		$\phi_{4\bullet}$
	MARGINAL	$\phi_{\bullet 1}$	$\phi_{\bullet 2}$	$\phi_{\bullet 3}$	$\phi_{\bullet 4}$	1

**Figure 6 Pseudo-Decision Table for Four Hypothetical Categories**

If a student is classified as being in one category based on Test One’s score, how probable would it be that the student would be reclassified as being in the same category if he or she took Test Two (a non-overlapping, equally difficult form of the test)? This proportion is a measure of decision consistency.

The proportions of correct decisions,  $\phi$ , for two and four categories are computed by the following two formulas, respectively:

$$\phi = \phi_{11} + \phi_{22}$$

$$\phi = \phi_{11} + \phi_{22} + \phi_{33} + \phi_{44}$$

The sum of the diagonal entries — that is, the proportion of students classified by the two forms into exactly the same achievement level — signifies the overall consistency.

Classification accuracy refers to the agreement of the observed classifications of students with the classifications made on the basis of their true scores. As discussed above, an observed score contains measurement error while a true score is theoretically free of measurement error. A student’s observed score can be formulated by the sum of his or her true score plus measurement error, or *Observed = True + Error*. Decision accuracy is an index to determine the extent to which measurement error causes a classification different from the one expected from the true score.

Since true scores are unobserved and decision consistency is computed based on a single administration of the Regents Examination in Algebra II (Common Core), a statistical model using solely data from the available administration is used to estimate the true scores and to project the consistency and accuracy of classifications (Hambleton & Novick, 1973). Although a number of procedures are available, a well-known method developed by Livingston and Lewis (1995) that utilizes a specific true score model is used.

Several factors might affect decision consistency and accuracy. One important factor is the reliability of the scores. All other things being equal, more reliable test scores tend to result in more similar reclassifications and less measurement error. Another factor is the location of the cut score in the score distribution. More consistent and accurate classifications are observed when the cut scores are located away from the mass of the score distribution. The number of performance levels is also a consideration. Consistency and accuracy indices based on four performance levels should be lower than those based on two performance levels. This is not surprising, since classification and accuracy using four performance levels would allow more opportunity to change achievement levels. Hence, there would be more classification errors and less accuracy with four performance levels, resulting in lower consistency indices.

**Results and Observations** The results for the dichotomies created by the four corresponding cut scores are presented in Table 8. The tabled values are derived with the program *BB-Class* (Brennan, 2004), using the Livingston and Lewis method. The decision consistency ranged from 0.84 to 0.94, and the decision accuracy ranged from 0.88 to 0.95. For the Regents Examination in Algebra II (Common Core), both decision consistency and accuracy values are high, indicating very good consistency and accuracy of examinee classifications, as shown in Table 8.

**Table 8 Decision Consistency and Accuracy Results: Regents Examination in Algebra II (Common Core)**

Statistic	1/2	2/3	3/4	4/5
Consistency	0.94	0.89	0.84	0.88
Accuracy	0.95	0.92	0.88	0.91

#### 4.4 GROUP MEANS (STANDARD 2.17)

Mean scale scores were computed based on reported gender, race/ethnicity, English Language Learner status, economically disadvantaged status, and student with disability status. The results are reported in Table 9.

**Table 9 Group Means: Regents Examination in Algebra II (Common Core)**

Demographics	Number	Mean Scale Score	SD Scale Score
All Students	91,478	70.96	13.09
<b>Ethnicity</b>			
American Indian/Alaska Native	403	66.76	12.74
Asian/Native Hawaiian/Other Pacific Islander	13,394	74.15	12.66
Black/African American	9,176	62.07	14.07
Hispanic/Latino	13,296	63.89	13.77
Multiracial	1,300	71.63	13.21
White	53,902	73.43	11.42
<b>English Language Learner</b>			
No	90,428	71.06	13.02
Yes	1,050	62.47	15.83
<b>Economically Disadvantaged</b>			
No	60,560	73.22	12.07
Yes	30,918	66.53	13.85
<b>Gender</b>			
Female	49,025	70.45	13.00
Male	42,446	71.54	13.16
<b>Student with Disabilities</b>			
No	88,882	71.18	12.96
Yes	2,596	63.28	14.90

\*Note: Seven students were not reported in the Ethnicity and Gender group, but they are reflected in “All Students.”

#### 4.5 STATE PERCENTILE RANKINGS

State percentile rankings based on raw score distributions are noted in Table 10. The percentiles are based on the distribution of all students taking the Regents Examination in Algebra II (Common Core) for the June 2016 administration. Note that the scale scores for the Regent Examination in Algebra II (Common Core) range from 0 to 100, but some scale scores may not be obtainable, depending on the raw score-to-scale score relationship for a specific administration. The percentile ranks are computed in the following manner:

- A student’s assigned “state percentile rank” will be the cumulative percentage of students scoring at the immediate lower score plus half of the percentage of students obtaining the given score.
- Students who obtain the highest possible score will receive a percentile rank of 99.

**Table 10 State Percentile Ranking for Raw Score – Regents Examination in Algebra II (Common Core)**

Scale Score	Percentile Rank	Scale Score	Percentile Rank	Scale Score	Percentile Rank	Scale Score	Percentile Rank
0	1	26	1	52	10	78	67
1	1	27	1	53	10	79	71
2	1	28	1	54	12	80	76
3	1	29	1	55	14	81	80
4	1	30	1	56	15	82	83
5	1	31	1	57	15	83	86
6	1	32	1	58	16	84	88
7	1	33	1	59	17	85	90
8	1	34	1	60	18	86	92
9	1	35	1	61	21	87	93
10	1	36	1	62	22	88	94
11	1	37	2	63	23	89	96
12	1	38	2	64	25	90	97
13	1	39	2	65	27	91	97
14	1	40	3	66	29	92	98
15	1	41	3	67	32	93	98
16	1	42	3	68	34	94	99
17	1	43	4	69	36	95	99
18	1	44	4	70	38	96	99
19	1	45	5	71	40	97	99
20	1	46	5	72	44	98	99
21	1	47	6	73	48	99	99
22	1	48	6	74	51	100	99
23	1	49	7	75	54		
24	1	50	8	76	58		
25	1	51	9	77	63		



## Chapter 5: Validity (Standard 1)

---

Restating the purpose and uses of the Regents Examination in Algebra II (Common Core), this exam measures examinee achievement against the New York State learning standards. The exam is prepared by teacher examination committees and New York State Education Department subject matter and testing specialists, and it provides teachers and students with important information about student learning and performance against the established curriculum standards. Results of this exam may be used to identify student strengths and needs, in order to guide classroom teaching and learning. The exams also provide students, parents, counselors, administrators, and college admissions officers with objective and easily understood achievement information that may be used to inform empirically based educational and vocational decisions about students. As a state-provided objective benchmark, the Regents Examination in Algebra II (Common Core) is intended for use in satisfying state testing requirements for students who have finished a course in Algebra II. A passing score on the exam counts toward requirements for a high school diploma, as described in the New York State diploma requirements: <http://www.nysed.gov/common/nysed/files/programs/curriculum-instruction/currentdiplomarequirements2.pdf>. Results of the Regents Examination in Algebra II (Common Core) may also be used to satisfy various locally established requirements throughout the state.

The validity of score interpretations for the Regents Examination in Algebra II (Common Core) is supported by multiple sources of evidence. Chapter 1 of the *Standards for Educational Psychological Testing* (AERA et al., 2014) specifies five sources of validity evidence that are important to gather and document in order to support validity claims for an assessment:

- test content
- response processes
- internal test structure
- relation to other variables
- consequences of testing

It is important to note that these categories are not mutually exclusive. One source of validity evidence often falls into more than one category, as discussed in more detail in this chapter. Nevertheless, these classifications provide a useful framework within the *Standards* (AERA et al., 2014) for the discussion and documentation of validity evidence, so they are used here. The process of gathering evidence of the validity of score interpretations is best characterized as ongoing throughout test development, administration, scoring, reporting, and beyond.

### 5.1 EVIDENCE BASED ON TEST CONTENT

The validity of test content is fundamental to arguments that test scores are valid for their intended purpose. It demands that a test developer provide evidence that test content is well-aligned with the framework and standards used in curriculum and instruction. Accordingly, detailed attention was given to this correspondence between standards and test content during test design and construction.

The Regents Examination in Algebra II (Common Core) measures student achievement on the NYS P–12 Common Core Learning Standards for Mathematics, consistent with the Model

Content Frameworks for Mathematics provided by the Partnership for the Assessment of Readiness for College and Career (PARCC, 2014). The model content frameworks are located at <http://www.parcconline.org/resources/educator-resources/model-content-frameworks/mathematics-model-content-framework>. The standards for mathematics are located at <http://www.engageny.org/resource/new-york-state-p-1-12-common-core-learning-standards-for-mathematics>. Clarifications for Algebra II (Common Core) standards are located at <http://www.engageny.org/resource/regents-exams-mathematics-algebra-i-standards-clarifications>.

### *Content Validity*

Content validity is necessarily concerned with the proper definition of the construct and evidence that the test provides an accurate measure of examinee performance within the defined construct. The test blueprint for the Regents Examination in Algebra II (Common Core) is essentially the design document for constructing the exam. It provides explicit definition of the construct domain that is to be represented on the exam. The test development process (discussed in the next section) is in place to ensure, to the extent possible, that the blueprint is met in all operational forms of the exam.

Table 11 displays domain titles along with their cluster, standard, and targeted proportions of conceptual categories on the exam.

**Table 11 Test Blueprint, Regents Examination in Algebra II (Common Core)**

Conceptual Category	Percent of Test by Credits	Domains in Algebra II
Number & Quantity	5–12%	The Real Number System (N-RN) Quantities (N-Q)
Algebra	35–44%	The complex Number System (N-CN) Seeing Structure in Expressions (A-SSE) Arithmetic with Polynomials and Rational Expressions (A-APR) Creating Equations (A-CED) Reasoning with Equations and Inequalities (A-REI) Expressing Geometric Properties with Equations (G-GPE)*
Functions	30–40%	Interpreting Functions (F-IF) Building Functions (F-BF) Linear, Quadratic, and Exponential Models (F-LE) Trigonometric Function (F-TF)
Statistics & Probability	14–21%	Interpreting categorical and quantitative data (S-ID) Making Inferences and Justifying Conclusions (S-IC) Conditional Probability and the Rules of Probability (S-CP)

\*Although the organization of the CCLS places one standard from the G-GPE domain into the Geometry Conceptual Category, the content within this domain will be assessed as part of the Algebra Conceptual Category for the Regents Examination in Algebra II (Common Core).

### *Item Development Process*

Test development for the Regents Examination in Algebra II (Common Core) is a detailed, step-by-step process of development and review cycles. An important element of this process

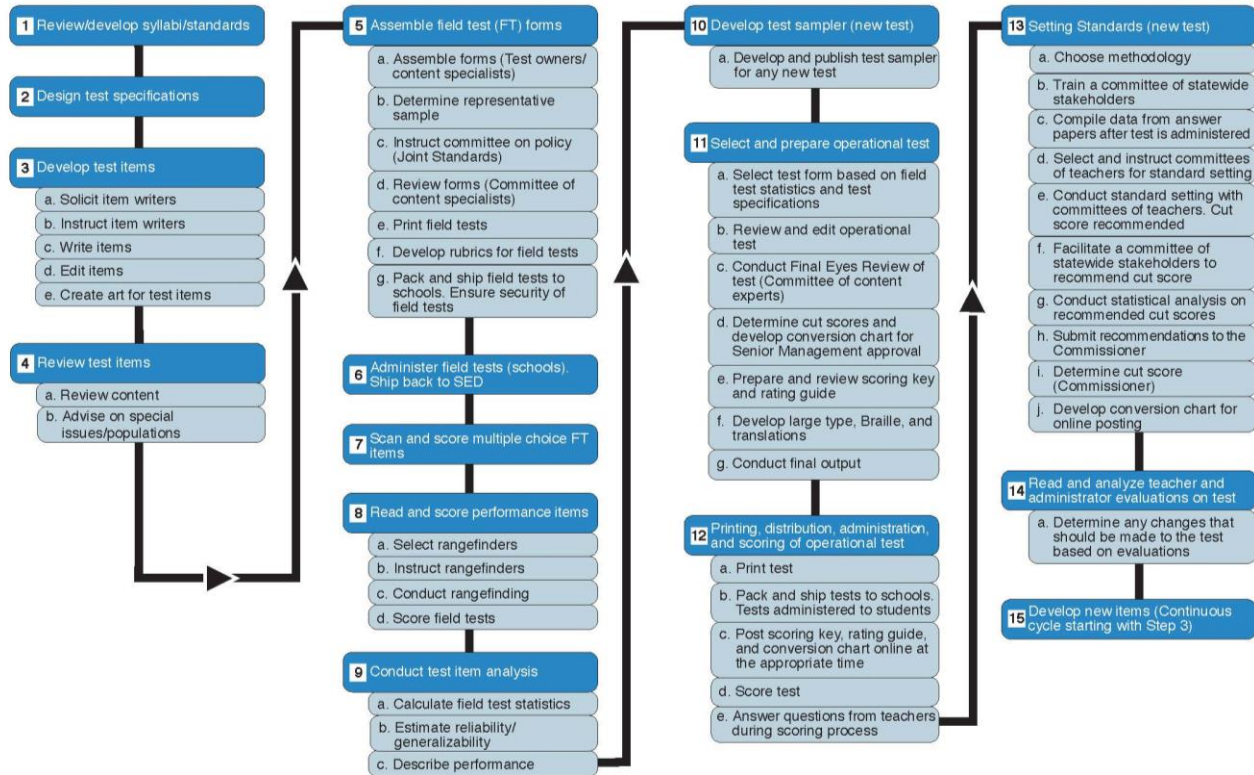
is that all test items are developed by New York State educators in a process facilitated by state subject matter and testing experts. Bringing experienced classroom teachers into this central item development role serves to draw a strong connection between classroom and test content.

Only New York State-certified educators may participate in this process. The New York State Education Department asks for nominations from districts, and all recruiting is done with diversity of participants in mind, including diversity in gender, ethnicity, geographic region, and teaching experience. Educators with item-writing skills from around the state are retained to write all items for the Regents Examination in Algebra II (Common Core), under strict guidelines that leverage best practices (see Appendix C). State educators also conduct all item quality and bias reviews, to ensure that item content is appropriate to the construct being measured and fair for all students. Finally, educators use the defined standards, test blueprint targets, and statistical information generated during field testing to select the highest quality items for use in the operational test.

Figure 7 summarizes the full test development process, with steps 3 and 4 addressing initial item development and review. This figure also demonstrates the ongoing nature of ensuring the content validity of items through field test trials, and final item selection for operational testing.

Initial item development is conducted under the criteria and guidance provided by multiple documents, including the blueprint, item writing criteria, and a content verification checklist. Both multiple-choice and constructed-response items are included in the Regents Examination in Algebra II (Common Core), to ensure appropriate coverage of the construct domain. The *Guidelines for Writing Multiple-Choice Math Items* and the *Guidelines for Writing Constructed-Response Math Items* provide detailed information about how items are developed for the Regents examinations. The guidelines are included in Appendix C.

## NEW YORK STATE EDUCATION DEPARTMENT TEST DEVELOPMENT PROCESS



**Figure 7 New York State Education Department Test Development Process**

### *Item Review Criteria*

Item Review Criteria assist in the consistent application of rigorous item reviews intended to assess the quality of the items developed and identify items that require edits or removal from the pool of items to be field tested. The criteria that follow help ensure that high-quality items are continually developed in a manner that is consistent with the test blueprint. All reviewers participate in rigorous training designed to assist in a consistent interpretation of the standards throughout the item review process. This is a critical step in item development because consistency between the standards and what the items are asking examinees is a fundamental form of evidence of the validity of the intended score interpretations. Another integral component of this item review process is to review the scoring rules, or “rubrics,” for their clarity and consistency in what the examinee is being asked to demonstrate by responding to each item. Each of these elements of the review process is in place, ultimately, to target fairness for all students by targeting consistency in examinee scores and providing evidence of the validity of their interpretations.

Specifically, the item review criteria articulate the four major item characteristics that the New York State Education Department looks for in developing quality items:

1. language and graphical appropriateness
2. sensitivity/bias

3. fidelity of measurement to CCLS
4. conformity to the expectations for the specific item types and formats (e.g., multiple-choice questions, 2-point constructed-response questions, 4-point constructed-response questions, and 6-point constructed-response questions).

Each section of the criteria includes pertinent questions that help reviewers determine whether or not an item is of sufficient quality. Within the first two categories, the headings Language Appropriateness, Sensitivity/Bias, and Math Art identify the basic components of quality assessment items. The criteria for language appropriateness are used to help ensure that students understand what is asked in each question and that the language in the question does not adversely affect a student's ability to perform the required task. Similarly, the sensitivity/bias criteria are used to evaluate whether questions are unbiased, non-offensive, and not disadvantageous to any given subgroup(s). The math art criteria assess the appropriateness and clarity, when graphics are used within questions.

The third category of the item review criteria framework, Item Alignment, addresses how each item measures a given mathematics standard. This criterion asks the reviewer to comment on key aspects of how the item addresses and calls for the skills demanded by the standards. Additionally, these criteria prompt reviewers to comment on how more than one standard is addressed by a given question.

The fourth category of the item review criteria framework addresses the specific demands for different item types and formats. Reviewers evaluate each item, to ensure that it conforms to the given requirements. For example, multiple-choice items must have, among other characteristics, one unambiguously correct answer and several plausible, but incorrect, answer choices.

Refer to the following link for more detail on the item review criteria: <https://www.engageny.org/resource/regents-exams-mathematics-item-criteria-checklist>.

Following these reviews, only items that are approved by an assigned educator panel move forward for field testing.

Ongoing attention is also given to the relevance of the standards used to guide curriculum and assessment. Consistent with a desire to assess this relevance, the New York State Education Department (NYSED) is committed to ongoing standards review over time and periodically solicits thoughtful, specific responses from stakeholders about individual standards within the NYS P–12 Standards.

## **5.2 EVIDENCE BASED ON RESPONSE PROCESSES**

The second source of validity evidence is based on examinee response processes. This standard requires evidence that examinees are responding in the manner intended by the test items and rubrics and that raters are scoring those responses in a manner that is consistent with the rubrics. Accordingly, it is important to control and monitor whether or not construct-irrelevant variance in response patterns has been introduced at any point in the test development, administration, or scoring processes.

The controls and monitoring in place for the Regents Examination in Algebra II (Common Core) include the item development process, with attention paid to mitigating the introduction of construct-irrelevant variance. The development process described in the previous sections details the process and attention given to reducing the potential for construct irrelevance in response processes by attending to the quality and alignment of test content to the test blueprint and to the item development guidelines (Appendix C). Further evidence is documented in the test administration and scoring procedures, as well as the results of statistical analyses, which are covered in the following two sections.

### *Administration and Scoring*

Adherence to standardized administration procedures is fundamental to the validity of test scores and their interpretation, as such procedures allow for adequate and consistently applied conditions for scoring the work of every student who takes the examination. For this reason, guidelines, which are contained in the *School Administrator's Manual, Secondary Level Examinations* (<http://www.p12.nysed.gov/assessment/sam/secondary/hssam-update.html>), have been developed and implemented for the New York State Regents testing program. All secondary-level Regents examinations are administered under these standard conditions, in order to support valid inferences for all students. These standard procedures also cover testing students with disabilities who are provided testing accommodations consistent with their Individualized Education Programs (IEPs) or Section 504 Accommodation Plans (504 Plans). Full test administration procedures are available at <http://www.p12.nysed.gov/assessment/hsgen/>.

The implementation of rigorous scoring procedures directly supports the validity of the scores. Regents test-scoring practices therefore focus on producing high-quality scores. Multiple-choice items are scored via local scanning at testing centers, and trained educators score constructed-response items. There are many studies that focus on various elements of producing valid and reliable scores for constructed-response items, but generally, attention to the following all contribute to valid and reliable scores for constructed-response items:

1. Quality training (Hoyt & Kerns, 1999; Lumley & McNamara, 1995; Wang, Wong, and Kwong, 2010; Gorman & Rentsch, 2009; Schleicher, Day, Bronston, Mayes, and Riggo, 2002; Woehr & Huffcutt, 1994; Henry et al., 2010; Johnson, Penny, and Gordon, 2008; Weigle, 1998)
2. Detection and correction of rating bias (McQueen & Congdon, 1997; Congdon & McQueen, 2000; Myford, & Wolfe, 2009; Barkaoui, 2011; Patz, Junker, Johnson, and Mariano, 2002)
3. Consistency or reliability of ratings (Congdon & McQueen, 2000; Harik, Clauser, Grabovsky, Nungester, Swanson, & Nandakumar, 2009; McQueen & Congdon, 1997; Myford & Wolfe, 2009; Mero & Motowidlo, 1995; Weinrott & Jones, 1984)
4. Rubric designs that facilitate consistency of ratings (Pecheone & Chung, 2006; Wolfe & Gitomer, 2000; Cronbach, Linn, Brennan, & Haertel, 1995; Cook & Beckman, 2009; Penny, Johnson, & Gordon, 2000; Smith, 1993; Leacock, Gonzalez, and Conarroe, 2014)

The distinct steps for operational test scoring include close attention to each of these elements and begin before the operational test is even selected. After the field test process, during which many more items than appear on the operational test are administered to a representative sample of students, a set of “anchor” papers representing student responses across the range of possible responses for constructed-response items is selected. The objective of these “range-finding” efforts is to create a training set for scorer training and execution, the scores from which are used to generate important statistical information about the item. Training scorers to produce reliable and valid scores is the basis for creating rating guides and scoring ancillaries to be used during operational scoring.

To review and select these anchor papers, NYS educators serve as table leaders during the range-finding session. In the range-finding process, committees of educators receive a set of student papers for each field-tested question. Committee members familiarize themselves with each item type and score a number of responses that are representative of each of the different score points. After the independent scoring is completed, the committee reviews and discusses their results and determines consensus scores for the student responses. During this process, atypical responses are important to identify and annotate for use in training and live scoring. The range-finding results are then used to build training materials for the vendor’s scorers, who then score the rest of the field test responses to constructed-response items. The final model response set for the June 2016 administration of the Regents Examination in Algebra II (Common Core) is located at <http://www.nysedregents.org/algebratwo/616/algtwo62016-mrs.pdf>.

During the range-finding and field test scoring processes, it is important to be aware of and control for sources of variation in scoring. One possible source of variation in constructed-response scores is unintended rater bias associated with items and examinee responses. Because the rater is often unaware of such bias, this type of variation may be the most challenging source of variation in scoring to control and measure. Rater biases can appear as severity or leniency in applying the scoring rubric. Bias also includes phenomena such as the halo effect, which occurs when good or poor performance on one element of the rubric encourages inaccurate scoring of other elements. These types of rater bias can be effectively controlled by training practices with a strict focus on rubric requirements.

The training process for operational scoring by state educators begins with a review and discussion of actual student work on constructed-response test items. This helps raters understand the range and characteristics typical of examinee responses, as well as the kinds of mistakes that students commonly make. This information is used to train raters on how to consistently apply key elements of the scoring rubric across the domain of student responses.

Raters then receive training consistent with the guidelines and ancillaries produced after field testing, and are allowed to practice scoring prior to the start of live scoring. Throughout the scoring process, there are important procedures for correcting inconsistent scoring or the misapplication of scoring rubrics for constructed-response items. When monitoring and correction do not occur during scoring, construct-irrelevant variation may be introduced. Accordingly, a scoring lead may be assigned to review the consistency of scoring for their assigned staff against model responses and to be available for consultation throughout the scoring process.

Attention to the rubric design also fundamentally contributes to the validity of examinee response processes. The rubric specifies what the examinee needs to provide as evidence of learning based on the question asked. The more explicit the rubric (and the item), the more clear the response expectations are for examinees. To facilitate the development of constructed-response scoring rubrics, the NYSED training for writing items includes specific attention to rubric development, as follows:

- The rubric should clearly specify the criteria for awarding each credit.
- The rubric should be aligned to what is asked for in the item and correspond to the knowledge or skill being assessed.
- Whenever possible, the rubric should be written to allow for alternate approaches and other legitimate methods.

In support of the goal of valid score interpretations for each examinee, then, such scoring training procedures are implemented for the Regents Examination in Algebra II (Common Core). Operational raters are selected based on expertise in the exam subject and are assigned a specific set of items to score. No more than one-third of the items on the test are assigned to any one rater. This has the effect of increasing the consistency of scoring across examinee responses by allowing each rater to focus on a subset of items. It also assures that no one rater is allowed to score the entire test for any one student. This practice reduces the effect of any potential bias of a single rater on individual examinees. Additionally, no rater is allowed to score the responses of his or her own students.

### *Statistical Analysis*

One statistic that is useful for evaluating the response processes for multiple-choice items is an item's point-biserial correlation on the distractors. A high point-biserial on a distractor may indicate that students are not able to identify the correct response for a reason other than the difficulty of the item. A finding of poor model fit for an item may also support a finding that examinees are not responding the way that the item developer intended them to. As documented in Table 2, the point-biserial statistics for distractors in the multiple-choice items all appear to be very low, indicating that, for the most part, examinees are not being drawn to an unintended construct.

## **5.3 EVIDENCE BASED ON INTERNAL STRUCTURE**

The third source of validity evidence comes from the internal structure of the test. This requires that test developers evaluate the test structure to ensure that the test is functioning as intended. Such an evaluation may include attention to item interactions, tests of dimensionality, or indications of test bias for or against one or more subgroups of examinees detected by differential item functioning (DIF) analysis. Evaluation of internal test structure also includes a review of the results of classical item analyses, test reliability, and the IRT scaling and equating.



The following analyses were conducted for the Regents Examination in Algebra II (Common Core):

- item difficulty
- item discrimination
- differential item functioning
- IRT model fit
- test reliability
- classification consistency
- test dimensionality

### *Item Difficulty*

Multiple analyses allow an evaluation of item difficulty. For this exam,  $p$ -values and Rasch difficulty (item location) estimates were computed for MC and CR items. Items for the June 2016 Regents Examination in Algebra II (Common Core) show a range of  $p$ -values consistent with the targeted exam difficulty. Item  $p$ -values range from 0.11 to 0.76, with a mean of 0.49.

### *Item Discrimination*

How well the items on a test discriminate between high- and low-performing examinees is an important measure of the structure of a test. Items that do not discriminate well generally provide less reliable information about student performance. Tables 2 and 3 provide point-biserial values on the correct responses, and Table 2 also provides point-biserial values on the three distractors. The values for correct answers are 0.20 or higher for all but one item, indicating that most items are discriminating well between high- and low-performing examinees. Point-biserials for all distractors are negative or very close to zero, indicating that examinees are responding to the items as expected during item development. Refer to section 2 of this report for additional details.

### *Differential Item Functioning*

Differential item functioning (DIF) for gender was conducted following field testing of the items in 2015. Sample sizes for subgroups based on ethnicity and English language learner status were, unfortunately, too small to reliably compute DIF statistics, so only gender DIF analyses were conducted. The Mantel-Haenszel  $\chi^2$  and standardized mean difference were used to detect items that may function differently for any of these subgroups. The Mantel-Haenszel  $\chi^2$  is a conditional mean comparison of the ordered response categories for reference and focal groups combined over values of the matching variable score. “Ordered” means that a response earning a score of “1” on an item is better than a response earning a score of “0,” a “2” is better than “1,” and so on. “Conditional,” on the other hand, refers to the comparison of members from the two groups who received the same score on the matching variable — the total test score in our analysis.

Two operational items on the June 2016 administration had DIF flags from the field test. One item (#6) had a moderate DIF favoring female students while the other item (#13) had a moderate DIF favoring male students. The items were subsequently reviewed by content specialists. They were unable to identify content-based reasons why the items might be functioning differently between male students and female students and did not see any issue with using them for the operational exam.

Full differential item functioning results are reported in Appendix E of the field test reports for 2015.

#### *IRT Model Fit*

Model fit for the Rasch method used to estimate location (difficulty) parameters for the items on the Regents Examination in Algebra II (Common Core) provide important evidence that the internal structure of the test is of high technical quality. The number of items within a targeted range of [0.7, 1.3] is reported in Table 5. The mean INFIT value is 1.00, with 37 of the 37 items falling in a targeted range of [0.7, 1.3]. As the range of [0.7, 1.3] is used as a guide for ideal fit, fit values outside of the range are considered individually. Overall, these results indicate that, for all items, the Rasch model fits the Regents Examination in Algebra II (Common Core) item data well.

#### *Test Reliability*

As discussed, test reliability is a measure of the internal consistency of a test (Cronbach, 1951). It is a measure of the extent to which the items on a test provide consistent information about student mastery of the domain. Reliability should, ultimately, demonstrate that examinee score estimates maximize consistency and, therefore, minimize error or, theoretically speaking, that examinees who take a test multiple times would get the same score each time. The reliability estimate for the Regents Examination in Algebra II (Common Core) is .89, showing high reliability of examinee scores. Refer to section 4 of this report for additional details.

#### *Classification Consistency and Accuracy*

A decision consistency analysis measures the agreement between the classifications based on two non-overlapping, equally difficult forms of the test. If two parallel forms of the test were given to the same students, the consistency of the measure would be reflected by the extent to which the classification decisions based on the first set of test scores matched the decisions based on the second set of test scores. Decision accuracy is an index to determine the extent to which measurement error causes a classification different from that expected from the true score. High decision consistency and accuracy provides strong evidence that the internal structure of a test is sound.

For the Regents Examination in Algebra II (Common Core), both decision consistency and accuracy values are high, indicating very good consistency and accuracy of examinee classifications. The results for the overall consistency across all five performance levels, as well as for the dichotomies created by the four corresponding cut scores, are presented in Table 7. The tabled values are derived with the program BB-Class (Brennan, 2004), using the Livingston and Lewis method. The decision consistency ranged from 0.84 to 0.94, and the decision accuracy ranged from 0.88 to 0.95.

## *Dimensionality*

In addition to model fit, a strong assumption of the Rasch model is that the construct measured by a test is unidimensional. Violation of this assumption might suggest that the test is measuring something other than the intended content and indicate that the quality of the test structure is compromised. A principal components analysis was conducted to test the assumption of unidimensionality, and the results provide strong evidence that a single dimension in the Regents Examination in Algebra II (Common Core) is explaining a large portion of the variance in student response data. This analysis does not characterize or explain the dimension, but a reasonable assumption can be made that the test is largely unidimensional and that the dimension most present is the targeted construct. Refer to section 3 for details of this analysis.

Considering this collection of detailed analyses on the internal structure of the Regents Examination in Algebra II (Common Core), strong evidence exists that the exam is functioning as intended and is providing valid and reliable information about examinee performance.

## **5.4 EVIDENCE BASED ON RELATIONS TO OTHER VARIABLES**

Another source of validity evidence is based on the relation of the test to other variables. This source commonly encompasses two validity categories prevalent in the literature and practice — concurrent and predictive validity. To make claims about the validity of a test that is to be used for high stakes purposes, such as the Regents Examination in Algebra II (Common Core), these claims could be supported by providing evidence that performance on the Algebra II (Common Core) test correlates well with other tests that measure the same or similar constructs. Although not absolute in its ability to offer evidence that concurrent test score validity exists, such correlations can be helpful for supporting a claim of concurrent validity, if the correlation is high. To conduct such studies, matched examinee score data for other tests measuring the same content as the Regents Examination in Algebra II (Common Core) is ideal, but the systematic acquisition of such data is complex and costly.

Importantly, a strong connection between classroom curriculum and test content may be inferred by the fact that New York State educators, deeply familiar with both the curriculum standards and their enactment in the classroom, develop all content for the Regents Examination in Algebra II (Common Core).

In terms of predictive validity, time is a fundamental constraint on gathering evidence. The gold standard for supporting the validity of predictive statements about test scores requires empirical evidence of the relationship between test scores and future performance on a defined characteristic. To the extent that the objective of the CCLS is to prepare students for college and career, it will be important to gather evidence of this empirical relationship over time.

Currently, the predictive validity is supported by expert judgments gathered during the standard-setting process for the Regents Examination in Algebra II (Common Core). During this process, subject matter experts described the performance of examinees across five levels and made recommendations on the cut scores to be used in distinguishing such performance. The process reflected best psychometric practice as articulated in the Standards for Educational and Psychological Measurement (AERA et al., 2014) and proceeded according to the plans reviewed by the New York State Technical Advisory Committee and an independent

consultant. This effort inherently represents further expert review of the test content and its alignment with the objectives of the CCLS. Participating subject matter experts made explicit judgments about what each item was asking of examinees and what successful performance on the items means for progress toward college and career readiness as defined by the standards.

After careful consideration of the nature of the new examinations, including their goal of providing evidence to support readiness claims, the rigor of the new curricula, the transitional and aspirational aspects of the state policy directives, and the role of the assessment in student learning throughout high school and beyond, the standard setting committees made recommendations on the cut scores to the New York State Commissioner of Education. The Commissioner accepted the recommendations of the standard setting panelists. More information is available in the Standard Setting technical report at <http://www.p12.nysed.gov/assessment/reports/>.

## **5.5 EVIDENCE BASED ON TESTING CONSEQUENCES**

There are two general approaches in the literature to evaluating consequential validity. Messick (1995) points out that adverse social consequences invalidate test use mainly if they are due to flaws in the test. In this sense, the sources of evidence documented in this report (based on the construct, internal test structure, response processes, and relation to other variables) serve as a consequential validity argument, as well. This evidence supports conclusions, based on test scores, that social consequences are not likely to be traced to characteristics or qualities of the test itself.

Cronbach (1988), on the other hand, argues that negative consequences could invalidate test use. From this perspective, the test user is obligated to make the case for test use and to ensure appropriate and supported uses. Regardless of perspective on the nature of consequential validity, it is important to caution against uses that are not supported by the validity claims documented for this test. For example, use of this test to predict examinee scores on other tests is not directly supported by either the stated purposes or by the development process and research conducted on examinee data. A brief survey of websites for New York State universities and colleges finds that, beyond the explicitly defined use as a testing requirement toward graduation for students who have completed a course in Algebra II, the exam is most commonly used to inform admissions and course placement decisions. Such uses can be considered reasonable, assuming that the competencies demonstrated in the Regents Examination in Algebra II (Common Core) are consistent with those required in the courses for which a student is seeking enrollment or placement. Educational institutions using the exam for placement purposes are advised to examine the scoring rules for the Regents Examination in Algebra II (Common Core) and to assess their appropriateness for the inferences being made about course placement.

As stated, the nature of validity arguments is not absolute, but it is supported through ongoing processes and studies designed to accumulate support for validity claims. The evidence provided in this report documents the evidence to date that supports the use of the Regents Examination in Algebra II (Common Core) scores for the purposes described.

## References

---

- American Educational Research Association, American Psychological Association, & National Council on Measurement in Education. (2014). *Standards for educational and psychological testing*. Washington, DC: American Educational Research Association.
- Barkaoui, Khaled. (2011). Effects of marking method and rater experience on ESL essay scores and rater performance. *Assessment in Education: Principles, Policy & Practice*, 18:3.
- Congdon, P. J., & McQueen, J. (2000). The stability of rater severity in large-scale assessment programs. *Journal of Educational Measurement*, 37(2), 163–178.
- Cook, D. A., & Beckman, T. J. (2009). Does scale length matter? A comparison of nine- versus five-point rating scales for mini-CEX. *Advances in Health Sciences Education*, 14, 655–684.
- Cronbach, L. J., Linn, R. L., Brennan, R. T., & Haertel, E. (1995, Summer). Generalizability analysis for educational assessments. Los Angeles, CA: University of California, Center for the Study of Evaluation/National Center for Research on Evaluation, Standards, and Student Testing. Retrieved February 17, 2016, from [www.cse.ucla.edu/products/evaluation/cresst\\_ec1995\\_3.pdf](http://www.cse.ucla.edu/products/evaluation/cresst_ec1995_3.pdf).
- Cronbach, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, 16, 297–334.
- Cronbach, L. J. (1988). Five Perspectives on validity argument. In H. Wainer & H. Braun (Eds.), *Test Validity* (pp. 3–17) Hillsdale, NJ: Lawrence Erlbaum.
- National Governors Association Center for Best Practices, & Council of Chief State School Officers. (2010). *Common Core State Standards*. Washington, DC: Authors.
- Gorman, C. A., & Rentsch, J. R. (2009). Evaluating frame-of-reference rater training effectiveness using performance schema accuracy. *Journal of Applied Psychology*, 94(5), 1336–1344.
- Gulliksen, H. (1950). *Theory of mental tests*. New York: Wiley.
- Hambleton, R. K., & Novak, M. R. (1973). Toward an integration of theory and methods for criterion-referenced tests. *Journal of Educational Measurement*, 10, 159–170.
- Hambleton, R. K., Swaminathan, H., & Rogers, H. J. (1991). *Item response theory*. Newbury Park, CA: Sage Publications.
- Hanson, B. A. (1994). Extension of Lord-Wingersky algorithm to computing test scores for polytomous items. Retrieved February 17, 2016 from <http://www.b-a-h.com/papers/note9401.pdf>.

- Harik, P., Clauser, B. E., Grabovsky, I., Nungester, R. J., Swanson, D., & Nandakumar, R. (2009, Spring). An examination of rater drift within a generalizability theory framework. *Journal of Educational Measurement*, 46(1), 43–58.
- Harvill, L. M. (1991). Standard error of measurement. *Educational Measurement: Issues and Practices*, 10(2), 33–41.
- Horn, J. L. (1965). A rationale and test for the number of factors in factor analysis. *Psychometrika*, 32, 179–185.
- Hoyt, W. T., & Kerns, M. (1999). Magnitude and moderators of bias in observer ratings: A meta-analysis. *Psychological Methods*, 4, 403–424.
- Huynh, H. (1976). Statistical consideration of mastery scores. *Psychometrika*, 41, 65–78.
- Johnson, R. L., Penny, J. A., & Gordon, B. (2008). *Assessing performance*. New York, NY: The Guilford Press.
- Kolen, M. J. (2004). POLYCSEM [Computer program]. University of Iowa. Retrieved August 1, 2012, from [http://www.education.uiowa.edu/casma/computer\\_programs.htm](http://www.education.uiowa.edu/casma/computer_programs.htm).
- Kolen, M. J., & Brennan, R. L. (2005). *Test equating, scaling, and linking: Methods and practices* (2nd ed.). New York: Springer-Verlag.
- Kolen, M. J. & Lee, W. (2011). Psychometric Properties of Raw and Scale Scores on Mixed-Format Tests. *Educational Measurement: Issues and Practice* 30(2), 15–24.
- Kolen, M. J., Zeng, L., & Hanson, B. A. (1996). Conditional standard errors of measurement for scale scores using IRT. *Journal of Educational Measurement*, 33(2), 129–140.
- Leacock, Claudia, Gonzalez, Erin, Conarro, Mike. (2014). *Developing effective scoring rubrics for AI short answer scoring*. CTB McGraw-Hill Innovative Research and Development Grant.
- Le Breton, J. M., & Sentor, J. L. (2008). Answers to 20 questions about interrater reliability and interrater agreement. *Organizational Research Methods*, 11(4), 815–852.
- Livingston, S. A., & Lewis, C. (1995). Estimating the consistency and accuracy of classifications based on test scores. *Journal of Educational Measurement*, 32, 179–197.
- Lumley, T., & McNamara, T. F. (1995). Rater characteristics and rater bias: Implications for training. *Language Testing*, 12, 54–72.
- Lord, F. M. (1980). *Applications of item response theory to practical testing problems*. Hillsdale, NJ: Lawrence Erlbaum.
- Messick, S. (1995). Standards of Validity and the validity of and standards in performance assessment. *Educational Measurement: Issues and Practice*, 14(4), 5–8.

- McDonald, R.P. (1979). The structural analysis of multivariate data: A sketch of general theory. *Multivariate Behavioral Research*, 14, 21–38.
- McQueen, J., & Congdon, P. J. (1997, March). *Rater severity in large-scale assessment: Is it invariant?* Paper presented at the Annual Meeting of the American Educational Research Association, Chicago, IL.
- Mero, N. P., & Motowidlo, S. J. (1995). Effects of rater accountability on the accuracy and the favorability of performance ratings. *Journal of Applied Psychology*, 80(4), 517–524.
- Myford, C. M., & Wolfe, E. W. (2009). Monitoring rater performance over time: A framework for detecting differential accuracy and differential scale use. *Journal of Educational Measurement*, 46(4), 371–389.
- Nunnally, J. C., & Bernstein. I. H. (1994). *Psychometric theory* (3rd ed.). New York: McGraw-Hill.
- Patz, R. J., Junker, B. W., Johnson, M. S., and Mariano, L. T. (2002). The Hierarchical Rater Model for Rated Test Items and its Application to Large-Scale Educational Assessment Data. *Journal of Educational and Behavioral Statistics*, 27: 341.
- Pearson, K. (1896). Mathematical Contributions to the Theory of Evolution. III. Regression, Heredity and Panmixia. *Philosophical Transactions of the Royal Society of London*, 187, 253–318.
- Pecheone, R. L., & Chung Wei, R. R. (2007). Performance assessment for California teachers: Summary of validity and reliability studies for the 2003–04 pilot year. Palo Alto, CA: Stanford University PACT Consortium.
- Penny, J., Johnson, R. L., & Gordon, B. (2000). Using rating augmentation to expand the scale of an analytic rubric. *The Journal of Experimental Education*, 68(3), 269–287.
- Schleicher, D. J., Day, D. V., Bronston, T., Mayes, B. T., & Riggo, R. E. (2002). A new frame for frame-of-reference training: Enhancing the construct validity of assessment centers. *Journal of Applied Psychology*, 87(4), 735–746.
- Shrout, P. E., & Fleiss, J. L. (1979). Intraclass correlations: Uses in assessing rater reliability. *Psychological Bulletin*, 86(2), 420–428.
- Thissen, D., Pommerich, M., Billeaud, K., & Williams, V. S. L. (1995). Item response theory for scores on tests including polytomous items with ordered responses. *Applied Psychological Measurement*, 19, 39–49.
- van der Linden, W. J., & Hambleton, R. K. (Eds.). (1997). *Handbook of modern item response theory*. New York, NY: Springer-Verlag.

- Wang, T., Kolen, M. J., & Harris, D. J. (2000). Psychometric properties of scale scores and performance levels for performance assessments using polytomous IRT. *Journal of Educational Measurement*, 37(2), 141–162.
- Wang, X. M., Wong, K. F. E., & Kwong, J. Y. Y. (2010). The roles of rater goals and rater performance levels in the distortion of performance ratings. *Journal of Applied Psychology*, 95(3), 546–561.
- Weigle, S.C. (1998). Using FACETS to model rater training effects. *Language Testing* 15, 263–287.
- Weinrott, L., & Jones, B. (1984). Overt versus covert assessment of observer reliability. *Child Development*, 55, 1125–1137.
- Wilson, Mark and Hoskens, Machteld. (2001). The Rater Bundle Model. *Journal of Educational and Behavioral Statistics*, 26: 283
- Woehr, D. J., & Huffcutt, A. I. (1994). Rater training for performance appraisal: A quantitative review. *Journal of Occupational and Organizational Psychology*, 67, 189–205.
- Wolfe, E. W., & Gitomer, D. H. (2000). *The influence of changes in assessment design on the psychometric quality of scores*. Princeton, NJ: Educational Testing Service.



## Appendix A: Operational Test Maps

Table A.1 Test Map for June 2016 Administration

Position	Item Type	Max Points	Weight	Cluster	Mean	Point - Biserial	Rasch Difficulty	INFIT
1	MC	1	2	N-RN.A	0.71	0.34	-1.2887	0.99
2	MC	1	2	A-CED.A	0.64	0.33	-0.9186	1.02
3	MC	1	2	N-CN.A	0.57	0.42	-0.5852	0.95
4	MC	1	2	F-IF.C	0.76	0.34	-1.5766	0.98
5	MC	1	2	A-REI.A	0.50	0.10	-0.2571	1.24
6	MC	1	2	A-APR.B	0.77	0.38	-1.6477	0.92
7	MC	1	2	S-IC.A	0.43	0.17	0.0711	1.18
8	MC	1	2	F-BF.A	0.56	0.38	-0.5700	0.99
9	MC	1	2	S-ID.A	0.63	0.37	-0.8996	0.99
10	MC	1	2	F-BF.A	0.48	0.50	-0.1875	0.89
11	MC	1	2	S-CP.A	0.53	0.24	-0.3893	1.12
12	MC	1	2	N-CN.C	0.63	0.42	-0.8580	0.94
13	MC	1	2	F-IF.B	0.58	0.37	-0.6324	1.00
14	MC	1	2	A-APR.D	0.52	0.48	-0.3547	0.91
15	MC	1	2	F-IF.C	0.60	0.37	-0.7536	0.99
16	MC	1	2	F-BF.B	0.38	0.49	0.2978	0.90
17	MC	1	2	F-TF.A	0.42	0.35	0.0817	1.02
18	MC	1	2	F-IF.C	0.32	0.16	0.5814	1.18
19	MC	1	2	A-SSE.A	0.26	0.37	0.9477	0.99
20	MC	1	2	F-IF.B	0.46	0.37	-0.1062	1.01
21	MC	1	2	A-SSE.B	0.13	0.33	1.8488	0.97
22	MC	1	2	A-REI.D	0.48	0.40	-0.1690	0.98
23	MC	1	2	F-BF.A	0.38	0.33	0.3011	1.04
24	MC	1	2	F-TF.B	0.10	0.31	2.2594	0.96
25	CR	2	1	A-REI.A	1.23	0.51	-0.6895	1.00
26	CR	2	1	S-IC.B	0.83	0.46	0.1348	1.05
27	CR	2	1	A-APR.B	0.91	0.63	-0.1248	0.87
28	CR	2	1	F-IF.C	0.52	0.61	0.7824	0.86
29	CR	2	1	S-CP.B	0.79	0.49	0.1139	1.09
30	CR	2	1	G-GPE.A	0.34	0.62	1.0152	0.81
31	CR	2	1	A-SSE.A	0.50	0.52	0.7543	1.01
32	CR	2	1	F-LE.A	0.47	0.61	0.7375	0.88
33	CR	4	1	A-REI.C	1.01	0.63	0.5297	1.26
34	CR	4	1	A-SSE.B	1.25	0.65	0.4003	1.07
35	CR	4	1	S-IC.B	0.69	0.64	0.9591	1.05
36	CR	4	1	F-IF.C	1.27	0.70	0.2492	1.08
37	CR	6	1	A-REI.D	2.50	0.78	-0.0567	0.97

# Appendix B: Raw-to-Theta-to-Scale Score Conversion Tables

**Table B.1 Score Table for June 2016 Administration**

Raw Score	Ability	Scale Score
0	-5.7169	0.000
1	-4.4952	4.484
2	-3.7767	8.767
3	-3.3460	12.855
4	-3.0333	16.754
5	-2.7853	20.470
6	-2.5784	24.010
7	-2.3997	27.379
8	-2.2419	30.583
9	-2.1000	33.628
10	-1.9708	36.519
11	-1.8517	39.262
12	-1.7412	41.863
13	-1.6379	44.326
14	-1.5407	46.657
15	-1.4488	48.861
16	-1.3616	50.944
17	-1.2786	52.909
18	-1.1994	54.763
19	-1.1234	56.510
20	-1.0506	58.154
21	-0.9806	59.701
22	-0.9133	61.154
23	-0.8483	62.519
24	-0.7857	63.800
25	-0.7252	65.000
26	-0.6668	66.125
27	-0.6103	67.178
28	-0.5555	68.163
29	-0.5026	69.084
30	-0.4512	69.945
31	-0.4014	70.750
32	-0.3529	71.503
33	-0.3058	72.206
34	-0.2598	72.864
35	-0.2149	73.480
36	-0.1710	74.057
37	-0.1279	74.598
38	-0.0856	75.107
39	-0.0441	75.586
40	-0.0030	76.039

Raw Score	Ability	Scale Score
41	0.0376	76.468
42	0.0777	76.876
43	0.1177	77.265
44	0.1573	77.639
45	0.1969	78.000
46	0.2364	78.350
47	0.2760	78.691
48	0.3157	79.026
49	0.3556	79.357
50	0.3957	79.686
51	0.4361	80.015
52	0.4770	80.345
53	0.5185	80.679
54	0.5604	81.018
55	0.6030	81.364
56	0.6464	81.718
57	0.6907	82.082
58	0.7358	82.457
59	0.7820	82.844
60	0.8294	83.244
61	0.8781	83.659
62	0.9283	84.090
63	0.9800	84.536
64	1.0335	85.000
65	1.0890	85.482
66	1.1467	85.982
67	1.2070	86.501
68	1.2701	87.039
69	1.3364	87.597
70	1.4063	88.175
71	1.4804	88.773
72	1.5593	89.391
73	1.6437	90.030
74	1.7346	90.689
75	1.8332	91.367
76	1.9410	92.065
77	2.0597	92.783
78	2.1922	93.519
79	2.3419	94.273
80	2.5139	95.046
81	2.7161	95.834

Raw Score	Ability	Scale Score
82	2.9615	96.639
83	3.2741	97.459
84	3.7075	98.294
85	4.4317	99.141
86	5.6595	100.000

## Appendix C: Item Writing Guidelines

---

### Guidelines for Writing Multiple-Choice Math Items

- 1. The item measures the knowledge, skills, and proficiencies characterized by the standards within the identified cluster.**
- 2. The focus of the problem or topic should be stated clearly and concisely.**  
The stem should be meaningful and convey the central problem. A multiple-choice item functions most effectively when a student is required to compare specific alternatives related to the stem. It should not be necessary for the student to read all of the alternatives to understand an item. (*Hint: Cover the alternatives and read the stem on its own. Then ask yourself if the question includes the essential elements or if the essential elements are lost somewhere in the alternatives.*)
- 3. Include problems that come from a real-world context or problems that make use of multiple representations.**  
When using real-world problems, use formulas and equations that are real-world (e.g., the kinetic energy of an object with mass,  $m$ , and velocity,  $V$ , is  $k = \frac{1}{2} mv^2$ ). Use real-world statistics whenever possible.
- 4. The item should be written in clear and simple language, with vocabulary and sentence structure kept as simple as possible.**  
Each multiple-choice item should be specific and clear. The important elements should generally appear early in the stem of an item, with qualifications and explanations following. Difficult and technical vocabulary should be avoided, unless essential for the purpose of the question.
- 5. The stem should be written as a direct question or an incomplete statement**  
Direct questions are often more straightforward. However, an incomplete statement may be used to achieve simplicity, clarity, and effectiveness. Use whichever format seems more appropriate to present the item effectively.
- 6. The stem should not contain irrelevant or unnecessary detail.**  
Be sure that sufficient information is provided to answer the question, but avoid excessive detail or “window dressing.”
- 7. The phrase *which of the following* should not be used to refer to the alternatives; instead, use *which* followed by a noun.**  
In the stem, *which of the following* requires the student to read all of the alternatives before knowing what is being asked and assessed. Expressions such as *which statement*, *which expression*, *which equation*, and/or *which graph* are acceptable.

- 8. The stem should include any words that must otherwise be repeated in each alternative.**  
In general, the stem should contain everything the alternatives have in common or as much as possible of their common content. This practice makes an item concise. Exceptions include alternatives containing units and alternatives stated as complete sentences.
- 9. The item should have one and only one correct answer.**  
Items should not have two or more correct alternatives. *All of the above* and *none of the above* are not acceptable alternatives.
- 10. The distractors should be plausible and attractive to students who lack the knowledge, understanding, or ability assessed by the item.**  
Distractors should be designed to reflect common errors or misconceptions of students.
- 11. The alternatives should be grammatically consistent with the stem.**  
Use similar terminology, phrasing or sentence structure in the alternatives. Alternatives must use consistent language, including verb tense, nouns, singular/plurals, and declarative statements. Place a period at the end of an alternative *only* if the alternative by itself is a complete sentence.
- 12. The alternatives should be parallel with one another in form.**  
The length, complexity and specificity of the alternatives should be similar. For example, if the stem refers to a process, then all the alternatives must be processes. Avoid the use of absolutes such as *always* and *never* in phrasing alternatives.
- 13. The alternatives should be arranged in logical order, when possible.**  
When the alternatives consist of numbers and letters, they should ordinarily be arranged in ascending or descending order. An exception would be when the number of an alternative and the value of that alternative are the same. For example: (1) 1 (2) 2 (3) 0 (4) 4.
- 14. The alternatives should be independent and mutually exclusive.**  
Alternatives that are synonymous or overlap in meaning often assist the student in eliminating distractors.
- 15. The item should not contain extraneous clues to the correct answer.**  
Any aspect of the item that provides an unintended clue that can be used to select or eliminate an alternative should be avoided. For example, any term that appears in the stem should not appear in only one of the alternatives.
- 16. Notation and symbols as presented on Common Core examinations should be used consistently.**  
For example,  $AB$  means the length of line segment  $AB$ ,  $\overline{AB}$  means line segment  $AB$ ,  $m\angle A$  means the number of degrees in the measure of angle  $A$ , etc.

## REVIEW CRITERIA CHECKLIST FOR POTENTIAL MATH ITEMS

The following list of criteria will be used to train item writers and then to review items for possible inclusion on test forms.

<b>Language Appropriateness</b>	<b>Yes</b>	<b>No</b>	<b>n/a</b>	<b>Explain or Describe</b>
1. Item: Uses grade-level vocabulary. Uses the simplest terms possible to convey information. Avoids technical terms unrelated to content.				
2. Sentence complexity well within grade expectations.				
3. Avoids ambiguous or double-meaning words.				
4. Pronouns have clear referents.				
5. Item avoids irregularly spelled words. <i>Use most common spelling of words.</i>				
6. Item can be put into Braille. Item can be translated appropriately according to the specific accommodations as outlined in universal design guidelines.				

Sensitivity/Bias	Yes	No	n/a	Explain or Describe
1. The item is free of content that might be deemed offensive to groups of students, based upon culture, religion, race, ethnicity, gender, geographic location, ability, socioeconomic status, etc.				
2. The item is free of content that contains stereotyping.				
3. The item is free of content that might unfairly advantage or disadvantage subgroups of students (ethnicity, gender, geographic location, ability, socioeconomic status, etc.) by containing unfamiliar contexts or examples, unusual names of people or places, or references to local events or issues.				

Math Art	Yes	No	n/a	Explain or Describe
<p>1. The artwork clearly relates to the item and is important as an aspect of the problem-solving experience.</p>				
<p>2. The details in the artwork accurately and appropriately portray numbers/concepts contained in text or in lieu of text.</p> <p><i>Items should be drawn to scale as much as possible. By default, we do not include the text “Not drawn to scale” on every item; however, if a figure is drawn and there is a distortion in the figure, it should be indicated under the art that the figure is “not drawn to scale.” The degree of distortion should not be actively misleading.</i></p>				
<p>3. Graphics are clear (symbols are highly distinguished, free from clutter, at a reasonable scale, etc.).</p>				
<p>4. Visual load requirements are reasonable (interpreting graphic does not confuse underlying construct) and as simple as possible to present the prompt.</p> <p><i>“Visual load” refers to the amount of visual/graphic material included within a contained space. When graphics become overly busy, they break the cognitive process for different people or trip people up.</i></p>				

Item Alignment	Yes	No	n/a	Explain or Describe
<p>1. Is the item aligned to the standard to which it is written?</p> <p><i>List the primary standard to which the item is aligned and explain the degree to which there is alignment/lack of alignment.</i></p>				
<p>2. Is the item aligned to the correct secondary/tertiary standard(s)?</p>				
<p>3. The stem is reflective of the concept embedded within the standard and is representative of the goal of the standard.</p>				
<p>4. The item requires students to show understanding of key aspects of the standard.</p> <p><i>If "No," which aspects are not attended to?</i></p> <p><i>For constructed response items, it is important that the item be solved through an understanding of the key point of the standard. For example, if the language of the standard calls for "prove" or "show," items should actually involve proof to be aligned, not simply the ability to solve a related problem or perform a related manipulation.</i></p>				
<p>5. Does the question lend itself to being answered using a below-grade-level standard rather than the skills/concepts references in the on-grade-level standard?</p>				
<p>6. The item requires the student to use skills referenced in the primary standard and any additional standards listed.</p>				



<p>7. The item includes grade/course-appropriate standard numbers/variables (e.g., students are asked to solve questions using numbers/variables that are grade-appropriate).</p> <p><i>Note: This includes the parameters outlined in the PARCC Pathways document for guidance on how some standards are split across A1 and A2.</i></p>				
<p>8. The item is aligned to the correct primary Multiple Representations(s). <i>If "No," indicate the correct MR code(s).</i></p>				
<p>9. The item expects students to use a formula that is:</p> <ul style="list-style-type: none"> <li>- from a standard for an earlier grade level (i.e., prior knowledge);</li> <li>- part of the current mathematics curriculum;</li> <li>- not from another content area (e.g., physics).</li> </ul> <p>If "No," the formula should be in the item stem.</p> <p><i>For example, the formula for kinetic energy from physics should be included in the item stem.</i></p>				

Application/Modeling Items	Yes	No	n/a	Explain or Describe
<p>1. The item is aligned to a standard that requires modeling/application.</p> <p><i>Note: See starred items in CCSS for high school math. These items are identified as lending themselves to modeling.</i></p>				
<p>2. Does the language of the item obscure the math concept being assessed?</p> <p><i>Students should not stumble over irrelevant information.</i></p>				
<p>3. Modeling/application scenario is realistic and appropriate to the grade level (the situation is one that a reasonable person would encounter in everyday life—no stretching velvet ropes or weighing kittens in milligrams).</p> <p><i>If “No,” explain why it’s not.</i></p>				
<p>4. Standard does not call for modeling/application, but there is a reason for it to be represented as such.</p> <p><i>Even non-starred standards can and should involve appropriate applications where possible.</i></p>				

<p>5. Figures/numbers/concepts used in modeling/application as well as in the response are realistic (e.g., downloads cost 99 cents, the side of a house isn't 2x-32 long).</p>				
<p>6. Modeling scenario is presented in the most realistic and simple manner possible.</p>				
<p>7. Modeling/application scenario does not assume outside knowledge (e.g., approximate weight of paper, definition of a micron).</p>				
<p>8. Modeling/application scenario provides all necessary information for student to apply math concepts.</p>				
<p>9. Item does not clue students to which math strategy is needed to solve, but rather allows the student to choose a strategy to solve the item correctly.</p> <p><i>For example, we should not tell students to use Pythagorean theorem, but rather allow them to decide which approach to solving is appropriate.</i></p>				

<b>Mathematic Correctness</b>	<b>Yes</b>	<b>No</b>	<b>n/a</b>	<b>Explain or Describe</b>
1. The stem addresses a central math concept, either implicitly or explicitly.				
2. The math presented in stem is clear, accurate, and conceptually plausible.				
3. At least one strategy exists that is on grade level to solve the problem.				
4. If there is more than one strategy, regardless of the strategy employed, the same correct answer will be achieved.				
5. There is a rationale for the correct response that is aligned to the language of the Standards and that demonstrates knowledge and/or application of the Standards.				
6. For MCQs: Is answer Choice 1 plausible or the correct answer?  <i>If not, why?</i>				
7. For MCQs: Is answer Choice 2 plausible or the correct answer?  <i>If not, why?</i>				
8. For MCQs: Is answer Choice 3 plausible or the correct answer?  <i>If not, why?</i>				
9. For MCQs: Is answer Choice 4 plausible or the correct answer?  <i>If not, why?</i>				

Constructed Response and All Regents	Yes	No	n/a	Explain or Describe
1. The item involves a multi-step process.				
2. The item requires students to show work.  <i>Work referenced in item should not be trivial (e.g., if work was not shown, it would be likely that mistakes would be made).</i>				
3. The item assesses more than computation.				
4. The item asks student to explain a concept or procedure used to solve the problem.  <i>Note: Not always applicable.</i>				
5. If students are asked to describe what they did, clear direction is given as to what they should describe (the theory, the rationale for the answer, the reason a strategy is wrong, etc.).				
6. The item explicitly describes what we're trying to elicit from the student.				
7. The item is presented in a manner consistent with the Application MRs.				

<b>Overarching Comments</b>	<b>Yes</b>	<b>No</b>	<b>n/a</b>	<b>Explain or Describe</b>
1. The item is aligned to standard.				
2. The item is rigorous.  <i>The math should be sound, tight, challenging, and at the appropriate level of difficulty.</i>				
3. The item is fair.				
4. The item is mathematically correct.				
5. The item is coded correctly for MR.				

<b>Final Recommendation</b>	<b>Yes</b>	<b>No</b>	<b>n/a</b>	<b>Explain or Describe</b>
1. Accept.				
2. Accept with Edits.  <i>Are suggested edits minor (won't impact stats)?</i>  <i>Note: Does not apply if at final typesetting phase.</i>				
3. Reject.				

## Guidelines for Writing Constructed-Response Math Items

- 1. The item measures the knowledge, skills, and proficiencies characterized by the standards within the identified cluster.**
- 2. The focus of the problem or topic should be stated clearly and concisely.**  
The item should be meaningful, address important knowledge and skills, and focus on key concepts.
- 3. Include problems that come from a real-world context or problems that make use of multiple representations.**  
When using real-world problems, use formulas and equations that are real-world (e.g., *the kinetic energy of an object with mass,  $m$ , and velocity,  $V$  is  $k = \frac{1}{2} mv^2$* ). Use real-world statistics whenever possible.
- 4. The item should be written with terminology, vocabulary and sentence structure kept as simple as possible. The item should be free of irrelevant or unnecessary detail.**  
The important elements should generally appear early in the item, with qualifications and explanations following. Present only the information needed to make the context/scenario clear.
- 5. The item should not contain extraneous clues to the correct answer.**  
The item should not provide unintended clues that allow a student to obtain credit without the appropriate knowledge or skill.
- 6. The item should require students to demonstrate depth of understanding and higher-order thinking skills through written expression, numerical evidence, and/or diagrams.**  
An open-ended item should require more than an either/or answer or any variation such as yes/no, decrease/increase, and faster/slower. Often either/or items can be improved by asking for an explanation.
- 7. The item should require work rather than just recall.**  
Students need to show their mathematical thinking in symbols or words.
- 8. The stimulus should provide information/data that is mathematically accurate.**  
Examples of stimuli include, but are not limited to, art, data tables, and diagrams. It is best to use actual data whenever possible. Hypothetical data, if used, should be plausible and clearly identified as hypothetical.

- 9. The item should be written so that the student does not have to identify units of measurement in the answer, unless the question is testing dimensional analysis.**  
For example, consider the question: “A circle has a radius of length 4 centimeters. Find the number of centimeters in the length of the arc intercepted by a central angle measuring 2 radians.” Students would receive credit for an answer of “8” and would not be penalized for writing “8 cm.”
- 10. The item should be written to require a specific form of answer.**  
Phrases like “in terms of  $\pi$ ,” “to the nearest tenth,” and “in simplest radical form” may simplify the writing of the rubric for these types of items.
- 11. Items that require students to explain in words are encouraged.**  
One of the emphases of the Common Core standards is to foster student ability to communicate mathematical thinking. An example is to have students construct viable arguments such as to make conjectures, analyze situations or justify conclusions. These items would require students to demonstrate precision of knowledge in their responses.
- 12. Items may be broken into multiple parts that may be labeled a, b, c, etc.**  
Clear division of the parts of the problems may simplify the writing of the rubric for these types of items.
- 13. Notation and symbols as presented on Common Core examinations should be used consistently.**  
For example,  $AB$  means the length of line segment  $AB$ ,  $\overline{AB}$  means line segment  $AB$ ,  $m\angle A$  means the number of degrees in the measure of angle A, etc.