

New York State Regents Comprehensive Examination in English

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

CHAPTER 1: INTRODUCTION	1
1.1 INTRODUCTION	1
1.2 PURPOSES OF THE EXAM	1
1.3 TARGET POPULATION (STANDARD 7.2)	2
CHAPTER 2: CLASSICAL ITEM STATISTICS (STANDARD 4.10)	4
2.1 ITEM DIFFICULTY.....	4
2.2 ITEM DISCRIMINATION	4
2.3 DISCRIMINATION ON DIFFICULTY SCATTER PLOTS	6
2.4 OBSERVATIONS AND INTERPRETATIONS	7
CHAPTER 3: IRT CALIBRATIONS, EQUATING, AND SCALING (STANDARDS 2, AND 4.10)	8
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
CHAPTER 4: RELIABILITY (STANDARD 2)	18
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>	22
<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
CHAPTER 5: VALIDITY (STANDARD 1)	28
5.1 EVIDENCE BASED ON TEST CONTENT	28
<i>Content Validity</i>	29
<i>Item Development Process</i>	29
<i>Item Review Process</i>	30
5.2 EVIDENCE BASED ON RESPONSE PROCESSES	31
<i>Administration and Scoring</i>	32
<i>Statistical Analysis</i>	34
5.3 EVIDENCE BASED ON INTERNAL STRUCTURE.....	34
<i>Item Difficulty</i>	35
<i>Item Discrimination</i>	35
<i>Differential Item Functioning</i>	35
<i>IRT Model Fit</i>	36
<i>Test Reliability</i>	36
<i>Classification Consistency and Accuracy</i>	36
<i>Dimensionality</i>	36

5.4 EVIDENCE BASED ON RELATIONS TO OTHER VARIABLES.....	37
5.5 EVIDENCE BASED ON TESTING CONSEQUENCES	37
REFERENCES.....	39
APPENDIX A: OPERATIONAL TEST MAPS	43
APPENDIX B: RAW-TO-THETA-TO-SCALE SCORE CONVERSION TABLES.....	46
APPENDIX C: ITEM WRITING GUIDELINES.....	49
GENERAL RULES FOR WRITING MULTIPLE-CHOICE ITEMS.....	49
APPENDIX D: TABLES AND FIGURES FOR AUGUST 2015 ADMINISTRATION	52
APPENDIX E: TABLES AND FIGURES FOR JANUARY 2016 ADMINISTRATION	59

List of Tables

TABLE 1 TOTAL EXAMINEE POPULATION: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	2
TABLE 2 MULTIPLE-CHOICE ITEM ANALYSIS SUMMARY: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	4
TABLE 3 CONSTRUCTED-RESPONSE ITEM ANALYSIS SUMMARY: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	5
TABLE 4 DESCRIPTIVE STATISTICS IN <i>p</i> -VALUE AND POINT BISERIAL CORRELATION: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	6
TABLE 5 SUMMARY OF ITEM RESIDUAL CORRELATIONS: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	11
TABLE 6 SUMMARY OF INFIT MEAN SQUARE STATISTICS: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH.....	12
TABLE 7 RELIABILITIES AND STANDARD ERRORS OF MEASUREMENT: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH.....	17
TABLE 8 DECISION CONSISTENCY AND ACCURACY RESULTS: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	21
TABLE 9 GROUP MEANS: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	22
TABLE 10 STATE PERCENTILE RANKING FOR RAW SCORE – REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	23
TABLE 11 TEST BLUEPRINT, REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH.....	25

List of Figures

FIGURE 1 SCATTERPLOT: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH.....	5
FIGURE 2 STUDENT PERFORMANCE MAP: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH.....	8
FIGURE 3 SCREE PLOTS: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH	10
FIGURE 4 CONDITIONAL STANDARD ERROR PLOTS: REGENTS COMPREHENSIVE EXAMINATION IN ENGLISH.....	19
FIGURE 5 PSEUDO-DECISION TABLE FOR TWO HYPOTHETICAL CATEGORIES	20
FIGURE 6 PSEUDO-DECISION TABLE FOR FOUR HYPOTHETICAL CATEGORIES	20
FIGURE 7 NEW YORK STATE EDUCATION DEPARTMENT TEST DEVELOPMENT PROCESS	26

Chapter 1: Introduction

1.1 INTRODUCTION

This technical report for the Regents Comprehensive Examination in English will provide New York State with documentation on the purpose of the Regents Examination, scoring information, evidence of both reliability and validity of the exams, scaling information, and guidelines and reporting information for the August 2015, January 2016, and June 2016 administrations. Chapters 1–5 detail results for the June 2016 administration. Results for the August 2015 and January 2016 administrations are provided in Appendices D and E, respectively. 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).¹ 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.

The Regents Comprehensive Examination in English is given in June, August, and January to students enrolled in New York State schools. The examination is based on the English Core Curriculum which is based the New York State Learning Standards for English Language Arts.

1.2 PURPOSES OF THE EXAM

The Regents Comprehensive Examination in English 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 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 Comprehensive Examination in English is intended for use in satisfying state testing requirements for students who have finished a course in English. 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 Comprehensive

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

Examination in English may also be used to satisfy various locally established requirements throughout the state.

1.3 TARGET POPULATION (STANDARD 7.2)

The examinee population for the Regents Comprehensive Examination in English is composed of students who have completed a course in English.

Table 1 provides a demographic breakdown of all students who took the August 2015, January 2016, and June 2016 administrations of the Regents Comprehensive Examination in English. 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 2015, January 2016, and June 2016 (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 publically reported counts of student test-takers.

Table 1 Total Examinee Population: Regents Comprehensive Examination in English

Demographics	August Admin*		January Admin**		June Admin***	
	Number	Percent	Number	Percent	Number	Percent
All Students	9,997	100	16,535	100	10,398	100
Race/Ethnicity						
American Indian/Alaska Native	79	0.80	135	0.82	78	0.75
Asian/Native Hawaiian/Other Pacific Islander	1,029	10.37	1,756	10.63	1,046	10.06
Black/African American	2,889	29.12	4,765	28.85	2,474	23.80
Hispanic/Latino	3,475	35.02	6,260	37.91	3,547	34.12
Multiracial	71	0.72	125	0.76	52	0.50
White	2,379	23.98	3,473	21.03	3,200	30.78
English Language Learner						
No	7,786	77.88	12,637	76.43	7,640	73.48
Yes	2,211	22.12	3,898	23.57	2,758	26.52
Economically Disadvantaged						
No	3,637	36.38	5,054	30.57	4,413	42.44
Yes	6,360	63.62	11,481	69.43	5,985	57.56
Gender						
Female	4,257	42.90	6,785	41.09	4,842	46.57
Male	5,665	57.10	9,729	58.91	5,555	53.43
Student with Disabilities						
No	7,176	71.78	11,800	71.36	7,809	75.10

Yes	2,821	28.22	4,735	28.64	2,589	24.90
-----	-------	-------	-------	-------	-------	-------

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

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

***Note: One student was not reported in the Ethnicity and Gender group, but that student is 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 Comprehensive Examination in English 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.² 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 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. 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 Regents Comprehensive Examination in English 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 0.41 to 0.83, with a mean of 0.68.

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 students who perform well on the Regents Comprehensive Examination in English would be more likely to answer any given item correctly, while low-performing students (i.e., those who perform

² For MC items with four response options, pure random guessing would lead to an expected p -value of 0.25.

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 that 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 0.25 or higher for all items, indicating that the items are discriminating well between high- and low-performing examinees. Point-biserials for all distractors are negative, indicating that examinees are responding to the items as expected during item and rubric development.

Table 2 Multiple-Choice Item Analysis Summary: Regents Comprehensive Examination in English

Item	Number of Students	p-Value	SD	Point-Biserial	Point - Biserial Distractor 1	Point - Biserial Distractor 2	Point-Biserial Distractor 3
1	10,398	0.77	0.42	0.44	-0.24	-0.22	-0.22
2	10,398	0.81	0.39	0.45	-0.17	-0.36	-0.15
3	10,398	0.41	0.49	0.34	-0.23	0.08	-0.22
4	10,398	0.73	0.44	0.43	-0.13	-0.33	-0.18
5	10,398	0.50	0.50	0.32	-0.20	-0.03	-0.20
6	10,398	0.67	0.47	0.34	-0.15	-0.12	-0.26
7	10,398	0.77	0.42	0.31	-0.16	-0.18	-0.16
8	10,398	0.69	0.46	0.44	-0.19	-0.19	-0.28
9	10,398	0.74	0.44	0.43	-0.21	-0.22	-0.23
10	10,398	0.57	0.50	0.25	-0.08	-0.19	-0.12
11	10,398	0.70	0.46	0.42	-0.19	-0.22	-0.24
12	10,398	0.76	0.43	0.36	-0.24	-0.14	-0.19
13	10,398	0.74	0.44	0.38	-0.21	-0.22	-0.18
14	10,398	0.43	0.50	0.34	-0.15	-0.20	-0.12
15	10,398	0.61	0.49	0.43	-0.24	-0.17	-0.22
16	10,398	0.77	0.42	0.34	-0.14	-0.18	-0.24
17	10,398	0.83	0.37	0.38	-0.21	-0.22	-0.19
18	10,398	0.70	0.46	0.39	-0.27	-0.18	-0.14
19	10,398	0.54	0.50	0.46	-0.20	-0.23	-0.22

Item	Number of Students	<i>p</i> -Value	SD	Point-Biserial	Point - Biserial Distractor 1	Point - Biserial Distractor 2	Point-Biserial Distractor 3
20	10,398	0.83	0.37	0.39	-0.22	-0.19	-0.22
21	10,398	0.73	0.45	0.34	-0.20	-0.21	-0.13
22	10,398	0.71	0.46	0.48	-0.21	-0.27	-0.24
23	10,398	0.55	0.50	0.41	-0.25	-0.26	-0.13
24	10,398	0.67	0.47	0.36	-0.19	-0.10	-0.25
25	10,398	0.80	0.40	0.38	-0.16	-0.22	-0.22

Table 3 Constructed-Response Item Analysis Summary: Regents Comprehensive Examination in English

Item	Min. score	Max. score	Number of Students	Mean	SD	<i>p</i> -Value	Point-Biserial
26	0	2	10,398	1.42	0.65	0.71	0.68
27	0	2	10,398	1.39	0.67	0.70	0.70
28	0	6	10,398	3.12	1.36	0.52	0.83

2.3 DISCRIMINATION ON DIFFICULTY SCATTER PLOTS

Figure 1 shows a scatter plot 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.

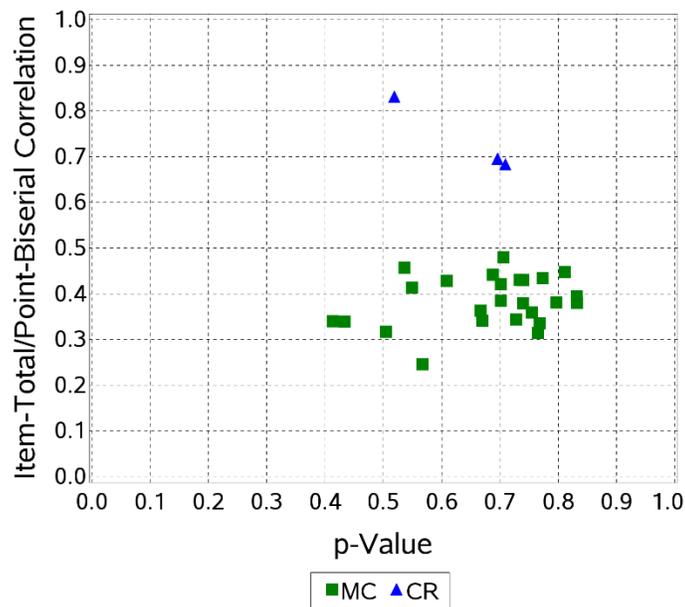


Figure 1 Scatter Plot: Regents Comprehensive Examination in English

Table 4 Descriptive Statistics in p -value and Point-Biserial Correlation: Regents Comprehensive Examination in English

Statistics	N	Mean	Min	Q1	Median	Q3	Max
p -value	28	0.68	0.41	0.59	0.70	0.76	0.83
Point-Biserial	28	0.42	0.25	0.34	0.39	0.44	0.83

2.4 OBSERVATIONS AND INTERPRETATIONS

The p -values for the MC items ranged from about 0.40 to 0.80, while the proportion-correct values for the CR items (Table 3) ranged from about 0.50 to 0.70. From the difficulty distributions illustrated in the plot, the range of item difficulties reflects a relatively high-performing cohort for this administration.

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

The Item Response Theory (IRT) model used for the Regents Comprehensive Examination in English 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 θ . Performance and θ 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 Comprehensive Examination in English 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 θ_n represents examinee ability, and D_{ij} is the step difficulty of the j^{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^{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 Comprehensive English in the June 2016 administration. The characteristics of this population are provided in Table 1 Regents Comprehensive Examination in English.

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 illustrates the high-performing nature of this cohort.

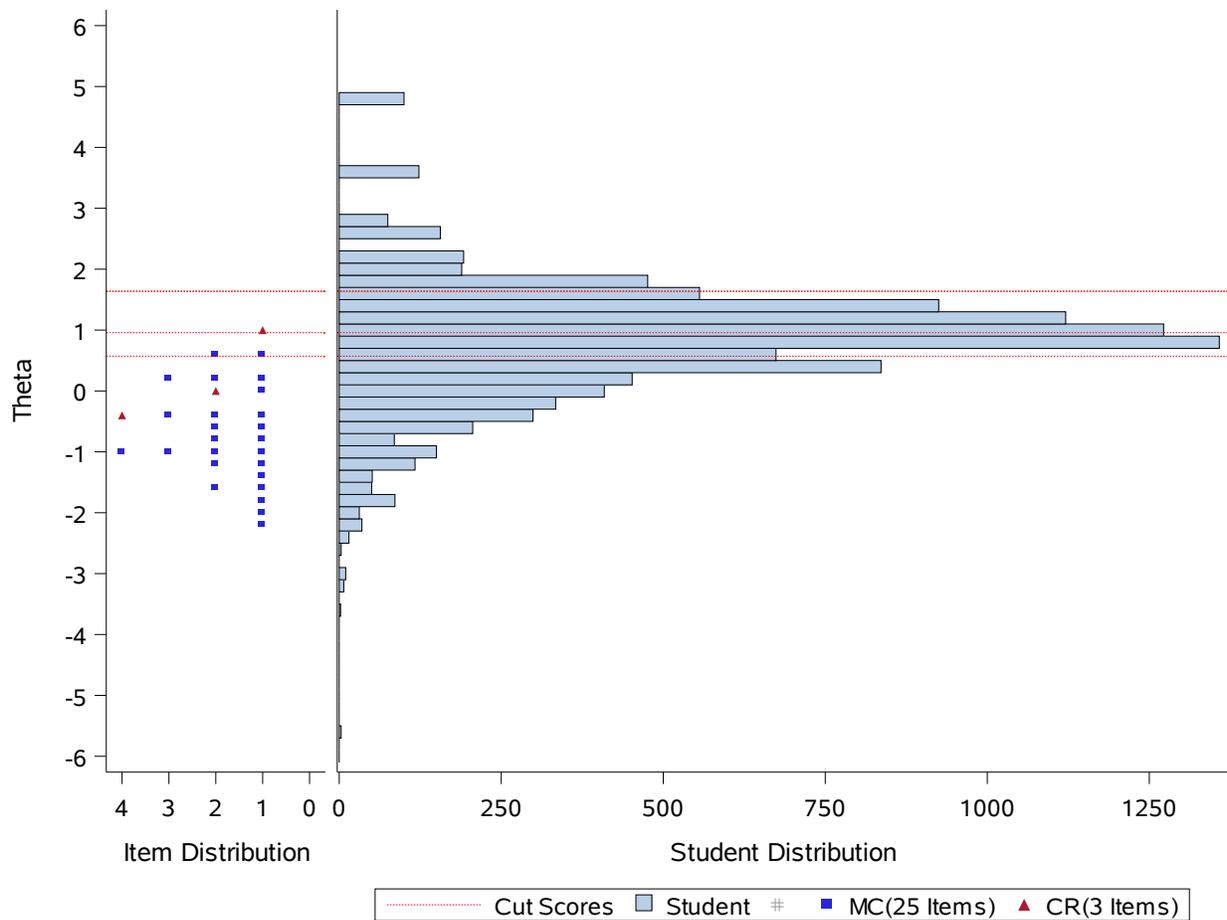


Figure 2 Student Performance Map: Regents Comprehensive Examination in English

3.5 CHECKING RASCH ASSUMPTIONS

Since the Rasch model was the basis of all calibration, scoring, and scaling analyses associated with the Regents Comprehensive Examination in English, 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 any 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 and parallel analysis results for the Regents Comprehensive Examination in English. 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 there is only one dominant component, showing an evidence of unidimensionality.

As a rule of thumb, Reckase (1979) proposed that the variance explained by the primary dimension should be greater than 20 percent, in order 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 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 22.23 percent of the total variance for the Regents Comprehensive Examination in English. The eigenvalue of the second dimension is less than one third of the first, at 1.47, and the second value is not significantly different from the parallel analysis. Overall, the PCA suggests that the test is reasonably unidimensional.

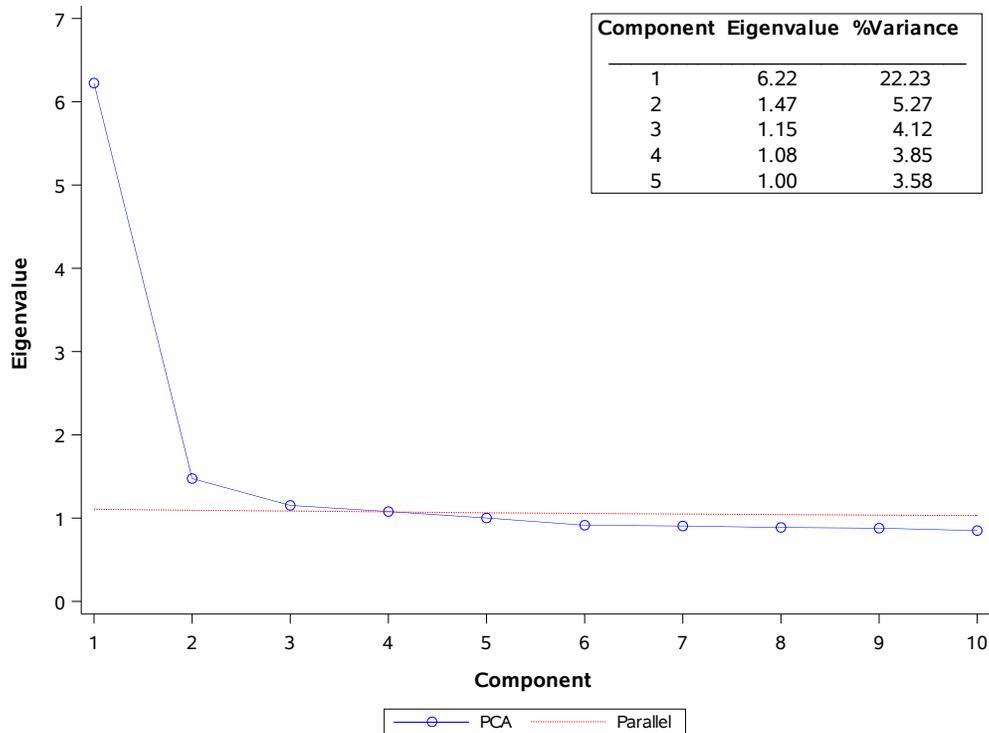


Figure 3 Scree Plot: Regents Comprehensive Examination in English

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 X_1, X_2, \dots, X_n is locally independent with respect to the latent variable θ if, for all $x = (x_1, x_2, \dots, x_n)$ and θ ,

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

This formula essentially states that the probability of any pattern of responses across all items (\mathbf{x}), after conditioning on the examinee’s true score (θ) 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 Comprehensive Examination in English items. In general, these residuals are computed as follows. First, expected item performance based on the Rasch model is determined using (θ) 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 Q_3 index, a popular statistic used to assess local independence. The expected value for the Q_3 statistic is approximately $-1/(k - 1)$ when no local dependence exists, where k is test length (Yen, 1993). Thus, the expected Q_3 values should be approximately -0.02 for the items on the exam. Absolute 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 (P_{10} , P_{25} , P_{50} , P_{75} , P_{90}) — for all the residual correlations for each test. The total number of item pairs (N) and the number of pairs with the absolute residual correlations greater than 0.20 are also reported in this table. There were three item pairs with absolute residual correlations greater than 0.20. The mean residual correlations were very slightly negative, at -0.03 . All residual correlations were small, with a maximum absolute value of 0.35, suggesting that local item independence generally holds for the Regents Comprehensive Examination in English.

Table 5 Summary of Item Residual Correlations: Regents Comprehensive Examination in English

Statistic Type	Value
N	378
Mean	-0.03
SD	0.05
Minimum	-0.18
P ₁₀	-0.10
P ₂₅	-0.05
P ₅₀	-0.03
P ₇₅	-0.01
P ₉₀	0.01
Maximum	0.35
> 0.20	3

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 (θ) relative to item difficulty.

The expected MnSq value is 1.0 and can range from 0.0 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 Comprehensive English, 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 all 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. These results indicate that the Rasch model fits the Regents Comprehensive Examination in English item data well.

Table 6 Summary of INFIT Mean Square Statistics: Regents Comprehensive Examination in English

	INFIT Mean Square					
	N	Mean	SD	Min	Max	[0.7, 1.3]
English	28	1	0.09	0.87	1.21	[28/28]

Items for the Regents Comprehensive Examination in English were field tested in 2014.

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, in order to accurately measure students’ abilities across the ability continuum. Appendix A contains the operational test maps for the August 2015, January 2016, and June 2016 administrations. Note that statistics presented in the test maps 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 essay items. For the Regents Comprehensive Examination in English, there is only one six-point item with fixed constants.

The New York State Regents Comprehensive Examination in English has three cut scores, which are set at the scale scores of 55, 65, 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 0.451. 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. These tables for the August 2015, January 2016, and June 2016 administrations can be found in Appendix B. These tables are 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 2004 administration. Scale scores from the August 2015, January 2016, and June 2016 administrations are on the same scale, and can be directly compared to scale scores on all previous administrations, back to the June 2004 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 and 85 were set to correspond to those raw score cuts. A third-degree polynomial is required to fit a line exactly to four arbitrary points (e.g., the raw scores corresponding to the four critical scale scores of 0, 65, 85, and 100). The general form of this best-fitting line is:

$$SS = m3 * RS^3 + m2 * RS^2 + m1 * RS1 + m0,$$

where SS is the scaled score, RS is the raw score, and m0 through m3 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 m3 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, 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. If no scale score rounds to these 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 Regents Comprehensive Examination in English 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 (ρ_X^2) is defined as the ratio of true score variance (σ_T^2) to the observed score variance (σ_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 (σ_E^2). Put differently, total variance equals true score variance plus error variance.³

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

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).

³ A covariance term is not required, as true scores and error are assumed to be uncorrelated in classical test theory.

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 of the test score variances were true, the index would equal 1.0. Such scores 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. It 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:

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

where N is the number of parts (items), σ_X^2 is the variance of the observed total test scores, and $\sigma_{Y_i}^2$ 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 ± 1 SEM about two-thirds of the time.⁴ For ± 2 SEM confidence intervals, this increases to about 95 percent.

The Coefficient Alpha and associated SEM for the Regents Comprehensive Examination in English are provided in Table 7. The reliability of 0.86 reflects the relatively short test length, as well as the presence of items with high score point ranges.

Table 7 Reliabilities and Standard Errors of Measurement: Regents Comprehensive Examination in English

Subject	Coefficient Alpha	SEM
English	0.86	4.02

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 (θ , 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 θ 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 θ . 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.

⁴ Some prefer the following interpretation: If a student were tested an infinite number of times, the ± 1 SEM confidence intervals constructed for each score would capture the student’s true score 68 percent of the time.

The relationship between θ 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 θ scores. In addition, as the exam is equated from year to year, the relationship between the raw and scale scores moves away from the original third-degree polynomial relationship to one that is also no longer expressible in a simple mathematical form. In the absence of a simple mathematical relationship between θ and the scale scores, the CSEMs that are available for each θ 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, 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 Comprehensive Examination in English that have a one-to-one relationship between raw (θ) 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 θ .

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 θ ; 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 θ with a mean value equal to the examinee's expected raw (true) score. The CSEM for that level of θ 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 θ 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 θ 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 θ .

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 construction of an approximate score band. The confidence intervals are constructed by adding and subtracting a multiplicative factor of the CSEM.

Conditional Standard Error of Measurement Characteristics

The relationship between the scale score CSEM and θ 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 θ (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 20 and then concave down from about 20 to 100. In other words, the scale scores track linearly with the raw scores for the first quarter of the scale score range and then are compressed relative to the raw scores for the remaining three quarters of the range, though there are slight 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 Comprehensive Examination in English, in which the compression of raw score to scale scores around the cut score of 55 changes the shape of the curve slightly. 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, at the lower end of the scale, scale scores 0 through 34 span raw scores 0 through 19 (20 raw score points for 35 scale score points). Over the range of scale scores from 36 to 78, the raw score range is 20 to 43 (24 raw score points for 43 scale score points). Finally, scale scores over the range of 80 to 100 span raw scores of 44 to 55 (12 raw score points for 21 scale score points).

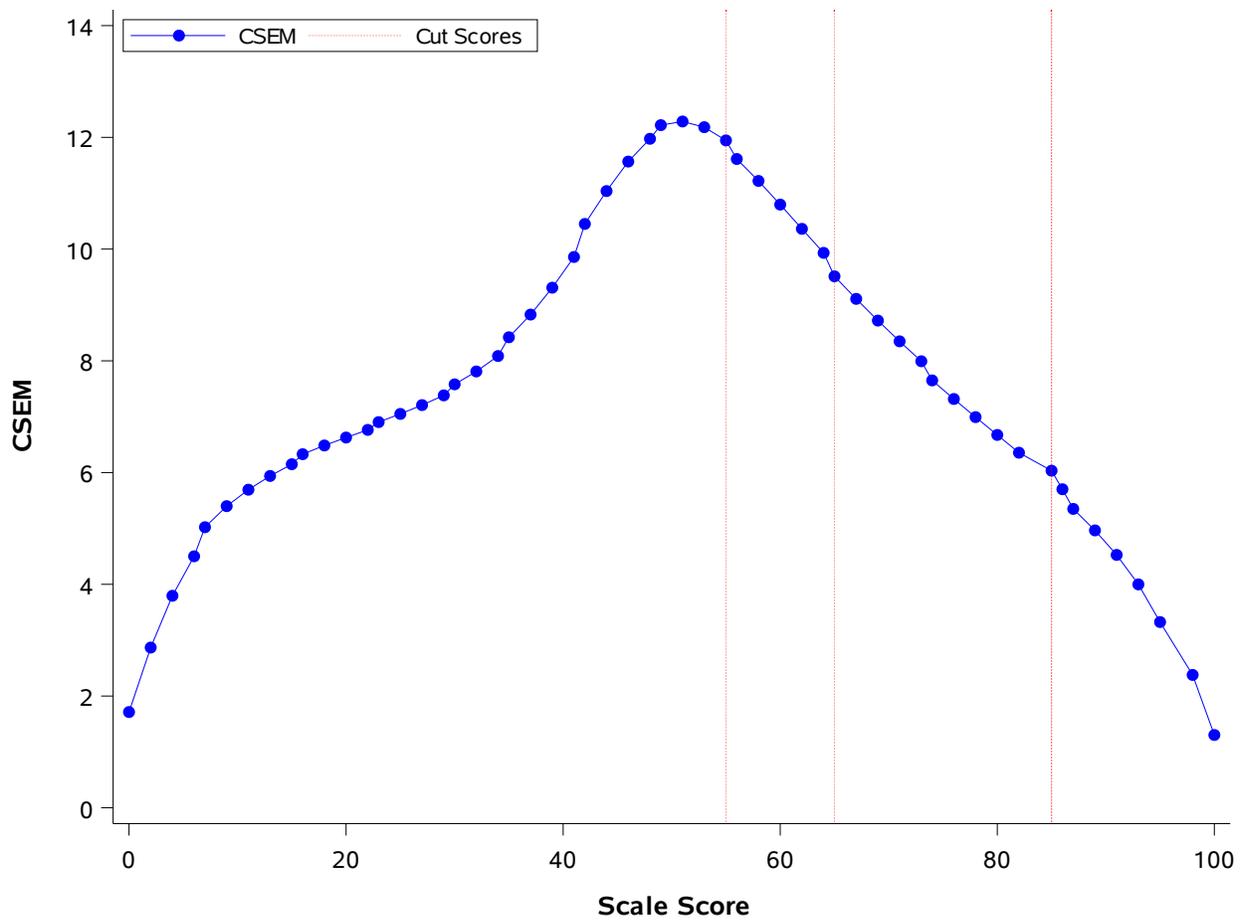


Figure 4 Conditional Standard Error Plot: Regents Comprehensive Examination in English

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	ϕ_{11}	ϕ_{12}	$\phi_{1\bullet}$
	LEVEL II	ϕ_{21}	ϕ_{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	ϕ_{11}	ϕ_{12}	ϕ_{13}	ϕ_{14}	$\phi_{1\bullet}$
	LEVEL II	ϕ_{21}	ϕ_{22}	ϕ_{23}	ϕ_{24}	$\phi_{2\bullet}$
	LEVEL III	ϕ_{31}	ϕ_{32}	ϕ_{33}	ϕ_{34}	$\phi_{3\bullet}$
	LEVEL IV	ϕ_{41}	ϕ_{42}	ϕ_{43}	ϕ_{44}	$\phi_{4\bullet}$
	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, ϕ , 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 Comprehensive Examination in English, a statistical model using data solely 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 performance 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 three 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. Decision consistency ranged from 0.84 to 0.87, and the decision accuracy ranged from 0.88 to 0.90. Both decision consistency and accuracy values based on individual cut points indicate generally good consistency and accuracy of examinee classifications. The consistency associated with the cut point for the highest performance level is not as high as desired; however, this highest cut point is not critical to high stakes decisions for examinees. Refer to Table 8 for details.

Table 8 Decision Consistency and Accuracy Results: Regents Comprehensive Examination in English

Statistic	1/2	2/3	3/4
Consistency	0.86	0.84	0.87
Accuracy	0.90	0.88	0.90

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 Comprehensive Examination in English

Demographics	Number	Mean Scale Score	SD Scale Score
All Students*	10,398	62.13	19.25
Ethnicity			
American Indian/Alaska Native	78	58.03	16.42
Asian/Native Hawaiian/Other Pacific Islander	1,046	58.15	17.67
Black/African American	2,474	57.74	17.12
Hispanic/Latino	3,547	56.33	17.56
Multiracial	52	65.44	16.78
White	3,200	73.31	18.45
English Language Learner			
No	7,640	66.48	18.18
Yes	2,758	50.09	16.86
Economically Disadvantaged			
No	4,413	69.94	19.38
Yes	5,985	56.38	17.00
Gender			
Female	4,842	65.71	19.24
Male	5,555	59.02	18.71
Student with Disabilities			
No	7,809	65.41	18.75
Yes	2,589	52.26	17.26

*Note: One student was not reported in the Ethnicity and Gender group, but that student is 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 Comprehensive Examination in English for the June 2016 administration. Note that the scale score for the Regents Examination ranges 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 Comprehensive Examination in English

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

Chapter 5: Validity (Standard 1)

Restating the purpose and uses of the Regents Comprehensive Examination in English, 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 exam also provides 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 Comprehensive Examination in English is intended for use in satisfying state testing requirements for students who have finished a course in English. 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 Comprehensive Examination in English may also be used to satisfy various locally established requirements throughout the state.

The validity of score interpretations for the Regents Comprehensive Examination in English 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 the 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 Comprehensive Examination in English measures student achievement on the New York State Learning Standards for English Language Arts. The standards can be found at: <http://www.nysed.gov/curriculum-instruction/english-language-arts-ela-literacy/>.

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 Comprehensive Examination in English is essentially the design document for constructing the exam. It provides explicit definition of the content 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 the targeted proportions of core performance indicators and standards on the exam.

Table 11 Test Blueprint, Regents Comprehensive Examination in English

	Core Performance Indicators	Standard 1	Standard 2	Standard 3
Listening 14%	0 – 2%	2 – 5%	2 – 5%	4 – 7%
Reading 31%	4 – 9%	4 – 9%	4 – 9%	13 – 20%
Writing 55%	30 – 40%	6 – 10%	7 – 11%	3 – 7%

Item Development Process

Test development for the Regents Comprehensive Examination in English 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 Comprehensive Examination in English, under strict guidelines that leverage best practices (see Appendix C). State educators also conduct all item quality and bias reviews, in order 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 the Department. Both multiple-choice and constructed-response items are included in the Regents Comprehensive Examination in English, in order to ensure appropriate coverage of the construct domain.

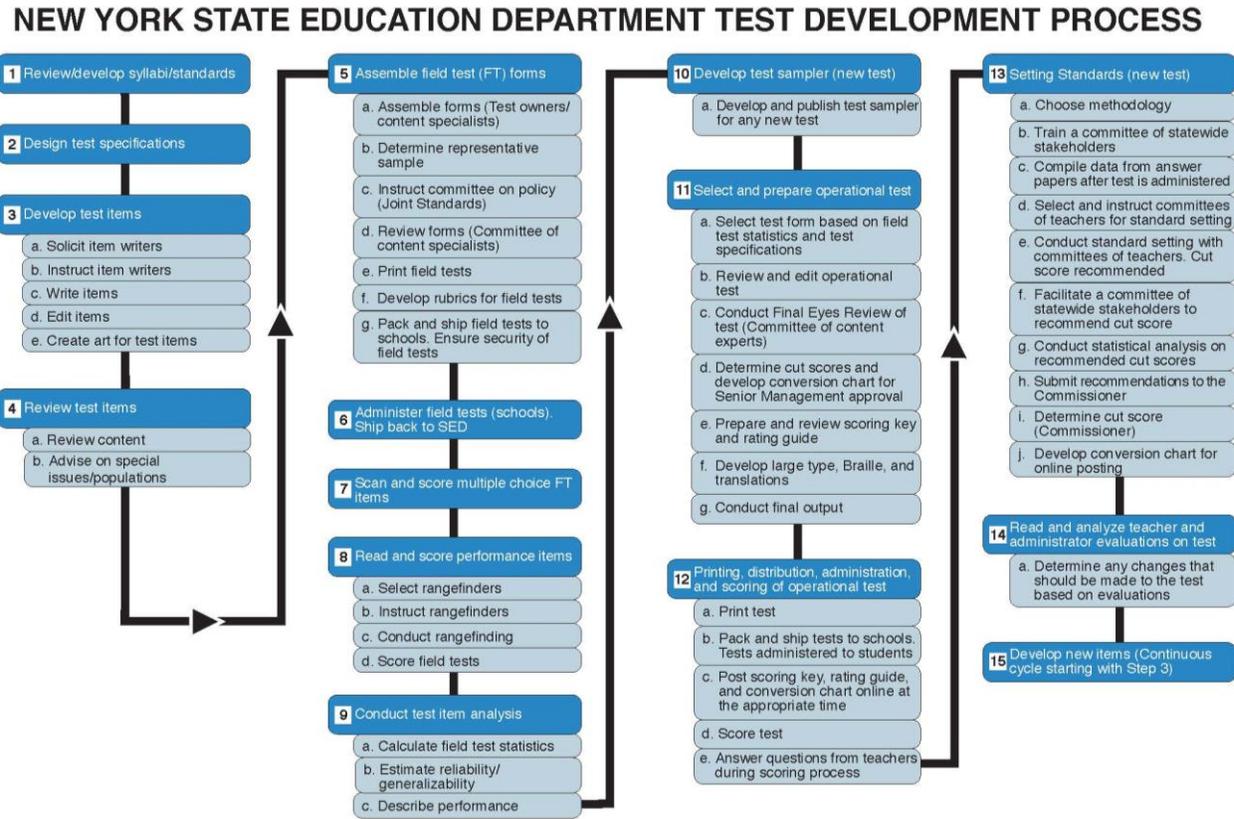


Figure 7 New York State Education Department Test Development Process

Item Review Process

The item review process helps to ensure 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. This process allows high-quality items to be 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 process articulates 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 standards
4. conformity to the expectations for the specific item types and formats

Each section of the criteria includes pertinent questions that help reviewers determine whether an item is of sufficient quality. Within the first two categories, 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, sensitivity/bias criteria are used to evaluate whether questions are unbiased, non-offensive, and not disadvantageous to any given subgroup(s).

The third category of item review, alignment, addresses how each item measures a given standard. This category asks the reviewer to comment on key aspects of how the item addresses and calls for the skills demanded by the standards.

The fourth category addresses the specific demands for different item types and formats. Reviewers evaluate each item, in order 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. 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 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 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 Comprehensive Examination in English 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; 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 sets for the 2015–2016 administrations of the Regents Comprehensive Examination in English are located at <http://www.nysedregents.org/ComprehensiveEnglish/home.html>.

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 alternative 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 Comprehensive Examination in English. Operational raters are selected based on expertise in the exam subject and are assigned a specific set of items to score. No more than approximately one-half 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.

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, in order 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 Comprehensive Examination in English:

- 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 Regents Comprehensive Examination in English show a range of p -values consistent with the targeted exam difficulty. Item p -values range from 0.41 to 0.83, with a mean of 0.68, showing that the items were relatively easy for this group of examinees.

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. Table 2 and Table 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 0.25 or higher for all items, and all distractors are negative, indicating that examinees are responding to the items as expected during item and rubric development.

Differential Item Functioning

Differential item functioning (DIF) for gender was conducted following field testing of the items in 2014. 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 χ^2 and standardized mean difference were used to detect items that may function differently for any of these subgroups. The Mantel-Haenszel χ^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.

Three operational items for the June 2016 administration had DIF flags from the field test. Two of the items (#s 16, 24) showed a moderate DIF favoring female students, and the other one (#6) had a strong DIF favoring female students. These 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 report.

IRT Model Fit

Model fit for the Rasch method used to estimate location (difficulty) parameters for the items on the Regents Comprehensive Examination in English provides 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 all items falling in a targeted range of [0.7, 1.3]. As the range of [0.7, 1.3] is used as guide for ideal fit, fit values outside of the range are considered individually. These results indicate that the Rasch model fits the Regents Comprehensive Examination in English 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 Comprehensive Examination in English is 0.86. 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 provide strong evidence that the internal structure of a test is sound.

For the Regents Comprehensive Examination in English, both decision consistency and accuracy values generally indicate good consistency and accuracy of examinee classifications. Decision consistency ranged from 0.84 to 0.87, and the decision accuracy ranged from 0.88 to 0.90. The consistency associated with the cut point for the highest performance level is not as high as desired; however, this highest cut point is not critical to high stakes decisions for examinees. Refer to Table 7 for details.

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 Comprehensive Examination in English 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 Comprehensive Examination in English, 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 Comprehensive Examination in English, these claims could be supported by providing evidence that performance on this 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 Comprehensive Examination in English are 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 Comprehensive Examination in English.

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 standards is to prepare students for meeting graduation requirements, it will be important to gather evidence of this empirical relationship over time.

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 of 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 English, 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 Comprehensive Examination in English 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 Comprehensive Examination in English 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 Comprehensive Examination in English 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.
- 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.
- 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.
- 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, Conarroe, Mike. (2014). *Developing effective scoring rubrics for AI short answer scoring*. CTB McGraw-Hill Education Innovative Research and Development Grant.
- 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.
- National Governors Association Center for Best Practices, & Council of Chief State School Officers. (2010). *Common Core State Standards*. Washington, DC: Authors.
- 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.
- 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.
- 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 August 2015 Administration

Position	Item Type	Max Points	Weight	Strand/ Standard	Mean	Point- Biserial	Rasch Difficulty	INFIT
1	MC	1	1	S1	0.90	0.32	-2.3069	1.10
2	MC	1	1	S2	0.82	0.53	-1.3863	0.92
3	MC	1	1	S2	0.89	0.42	-2.0954	0.99
4	MC	1	1	S2	0.88	0.51	-1.9612	0.87
5	MC	1	1	S3	0.95	0.30	-3.0776	1.02
6	MC	1	1	S3	0.77	0.50	-0.9972	0.99
7	MC	1	1	S3	0.80	0.51	-1.1716	0.96
8	MC	1	1	S3	0.62	0.45	-0.0001	1.12
9	MC	1	1	CPI	0.80	0.53	-0.9782	0.88
10	MC	1	1	CPI	0.78	0.50	-0.8107	0.92
11	MC	1	1	S2	0.75	0.58	-0.6080	0.85
12	MC	1	1	S3	0.75	0.48	-0.5780	0.96
13	MC	1	1	S3	0.61	0.40	0.2597	1.12
14	MC	1	1	S1	0.73	0.50	-0.4755	0.96
15	MC	1	1	CPI	0.74	0.60	-0.5337	0.83
16	MC	1	1	S1	0.59	0.59	0.3449	0.88
17	MC	1	1	S3	0.60	0.56	0.2841	0.91
18	MC	1	1	S3	0.53	0.59	0.6915	0.85
19	MC	1	1	S1	0.79	0.59	-0.8598	0.81
20	MC	1	1	CPI	0.64	0.59	0.0608	0.87
21	MC	1	1	S2	0.78	0.54	-1.7744	0.96
22	MC	1	1	S3	0.71	0.46	-1.2327	1.06
23	MC	1	1	S1	0.84	0.49	-2.2897	0.91
24	MC	1	1	S2	0.79	0.49	-1.8266	0.96
25	MC	1	1	S3	0.72	0.53	-1.2983	0.99
26	CR	2	3	CPI,S1,S2, S3	1.01	0.50	-0.0157	1.25
27	CR	2	3	CPI,S1,S2	0.93	0.56	0.2475	1.16
28	CR	6	3	CPI,S1,S2, S3	2.32	0.66	1.4051	1.45

Table A.2 Test Map for January 2016 Administration

Position	Item Type	Max Points	Weight	Strand/ Standard	Mean	Point- Biserial	Rasch Difficulty	INFIT
1	MC	1	1	S1	0.88	0.39	-1.7327	0.98
2	MC	1	1	S2	0.80	0.32	-1.0075	1.13
3	MC	1	1	S1	0.89	0.38	-1.9100	0.97
4	MC	1	1	S2	0.80	0.47	-1.0075	0.95
5	MC	1	1	S2	0.89	0.37	-1.8867	1.00
6	MC	1	1	S3	0.76	0.52	-0.6921	0.93
7	MC	1	1	S3	0.76	0.53	-0.7181	0.90
8	MC	1	1	S3	0.82	0.42	-1.1274	1.00
9	MC	1	1	S3	0.79	0.53	-0.9543	0.90
10	MC	1	1	S3	0.68	0.53	-0.2550	0.91
11	MC	1	1	CPI	0.74	0.54	-0.6423	0.88
12	MC	1	1	CPI	0.54	0.45	0.4836	1.01
13	MC	1	1	S3	0.70	0.53	-0.4154	0.90
14	MC	1	1	S3	0.65	0.42	-0.1138	1.05
15	MC	1	1	S1	0.75	0.65	-0.7018	0.77
16	MC	1	1	S1	0.60	0.60	0.1562	0.83
17	MC	1	1	CPI	0.58	0.47	0.2747	0.99
18	MC	1	1	S3	0.58	0.51	0.2629	0.95
19	MC	1	1	S1	0.75	0.65	-0.6719	0.77
20	MC	1	1	S3	0.65	0.54	-0.1011	0.90
21	MC	1	1	CPI	0.81	0.39	-1.3818	0.99
22	MC	1	1	S2	0.69	0.36	-0.5873	1.09
23	MC	1	1	S2	0.66	0.41	-0.4009	1.04
24	MC	1	1	S3	0.63	0.35	-0.2049	1.12
25	MC	1	1	S3	0.60	0.42	-0.0474	1.04
26	CR	2	3	CPI,S1,S2, S3	1.34	0.57	-0.4993	0.99
27	CR	2	3	CPI,S1,S2	1.24	0.61	-0.1632	0.96
28	CR	6	3	CPI,S1,S2, S3	2.46	0.65	1.1391	1.69

Table A.3 Test Map for June 2016 Administration

Position	Item Type	Max Points	Weight	Strand/ Standard	Mean	Point- Biserial	Rasch Difficulty	INFIT
1	MC	1	1	S3	0.94	0.36	-3.0123	0.92
2	MC	1	1	S1	0.96	0.24	-3.4488	0.99
3	MC	1	1	S2	0.61	0.40	0.0048	1.22
4	MC	1	1	S3	0.90	0.32	-2.2746	1.06
5	MC	1	1	CPI	0.71	0.46	-0.6310	1.10
6	MC	1	1	CPI	0.80	0.43	-1.2409	1.08
7	MC	1	1	S2	0.84	0.34	-1.6119	1.16
8	MC	1	1	S1	0.88	0.44	-2.0126	0.95
9	MC	1	1	S3	0.77	0.61	-1.0107	0.82
10	MC	1	1	S3	0.57	0.54	0.2464	1.01
11	MC	1	1	S2	0.58	0.59	0.1575	0.94
12	MC	1	1	S3	0.77	0.55	-1.0435	0.90
13	MC	1	1	S3	0.66	0.65	-0.3298	0.81
14	MC	1	1	S3	0.50	0.54	0.6219	1.00
15	MC	1	1	S1	0.59	0.66	0.1447	0.81
16	MC	1	1	S1	0.70	0.61	-0.5397	0.87
17	MC	1	1	CPI	0.73	0.68	-0.7300	0.74
18	MC	1	1	S2	0.68	0.66	-0.3987	0.80
19	MC	1	1	S3	0.50	0.60	0.6094	0.87
20	MC	1	1	S1	0.73	0.66	-0.7300	0.78
21	MC	1	1	S3	0.84	0.40	-1.4785	0.99
22	MC	1	1	CPI	0.85	0.44	-1.5367	0.94
23	MC	1	1	S1	0.69	0.51	-0.3468	0.97
24	MC	1	1	S3	0.79	0.39	-1.0177	1.05
25	MC	1	1	S3	0.87	0.39	-1.7229	0.99
26	CR	2	3	CPI,S1,S2,S3	1.36	0.60	-0.3272	1.02
27	CR	2	3	CPI,S1,S2	1.24	0.59	0.0302	1.06
28	CR	6	3	CPI,S1,S2,S3	2.66	0.57	0.9347	1.82

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

Table B.1 Score Table for August 2015 Administration

Raw Score	Ability	Scale Score
0	-6.0069	0.000
1	-4.7582	1.363
2	-4.0023	2.931
3	-3.5351	4.340
4	-3.1865	5.830
5	-2.9028	7.397
6	-2.6600	9.032
7	-2.4453	10.692
8	-2.2510	12.379
9	-2.0719	14.094
10	-1.9047	15.835
11	-1.7468	17.601
12	-1.5964	19.390
13	-1.4519	21.200
14	-1.3124	23.030
15	-1.1768	24.878
16	-1.0446	26.743
17	-0.9150	28.623
18	-0.7878	30.515
19	-0.6627	32.419
20	-0.5394	34.332
21	-0.4180	36.251
22	-0.2987	38.173
23	-0.1817	40.097
24	-0.0676	42.022
25	0.0429	43.942
26	0.1491	45.854
27	0.2504	47.760
28	0.3461	49.654
29	0.4359	51.536
30	0.5199	53.410
31	0.5984	55.275
32	0.6721	57.131
33	0.7414	58.980
34	0.8072	60.826
35	0.8702	62.667
36	0.9312	64.509
37	0.9909	66.349
38	1.0498	68.188
39	1.1088	70.030
40	1.1687	71.874

Raw Score	Ability	Scale Score
41	1.2301	73.722
42	1.2941	75.571
43	1.3617	77.425
44	1.4345	79.286
45	1.5140	81.150
46	1.6028	83.022
47	1.7040	84.899
48	1.8220	86.781
49	1.9634	88.669
50	2.1371	90.563
51	2.3573	92.454
52	2.6491	94.344
53	3.0657	96.230
54	3.7768	98.119
55	4.9967	100.000

Table B.2 Score Table for January 2016 Administration

Raw Score	Ability	Scale Score
0	-5.5300	0.000
1	-4.2975	2.186
2	-3.5635	4.224
3	-3.1170	6.214
4	-2.7878	8.163
5	-2.5229	10.075
6	-2.2985	11.952
7	-2.1018	13.797
8	-1.9252	15.614
9	-1.7637	17.405
10	-1.6141	19.173
11	-1.4738	20.918
12	-1.3411	22.645
13	-1.2147	24.353
14	-1.0933	26.046
15	-0.9761	27.727
16	-0.8624	29.396
17	-0.7518	31.057
18	-0.6436	32.712
19	-0.5375	34.362
20	-0.4334	36.007
21	-0.3310	37.652
22	-0.2302	39.298
23	-0.1312	40.947
24	-0.0340	42.602
25	0.0611	44.265
26	0.1536	45.937
27	0.2431	47.618
28	0.3293	49.311
29	0.4119	51.016
30	0.4904	52.732
31	0.5651	54.460
32	0.6359	56.197
33	0.7034	57.946
34	0.7678	59.702
35	0.8297	61.469
36	0.8897	63.244
37	0.9483	65.027
38	1.0060	66.818
39	1.0635	68.617
40	1.1214	70.423

Raw Score	Ability	Scale Score
41	1.1805	72.237
42	1.2414	74.059
43	1.3053	75.891
44	1.3732	77.731
45	1.4466	79.586
46	1.5276	81.451
47	1.6188	83.332
48	1.7242	85.232
49	1.8494	87.156
50	2.0032	89.109
51	2.1998	91.104
52	2.4645	93.152
53	2.8519	95.271
54	3.5325	97.480
55	4.7310	100.000

Table B.3 Score Table for June 2016 Administration

Raw Score	Ability	Scale Score
0	-6.0719	0.000
1	-4.8031	1.295
2	-4.0227	2.880
3	-3.5357	4.338
4	-3.1713	5.914
5	-2.8755	7.568
6	-2.6236	9.279
7	-2.4024	11.033
8	-2.2037	12.804
9	-2.0221	14.589
10	-1.8541	16.381
11	-1.6966	18.180
12	-1.5479	19.982
13	-1.4064	21.785
14	-1.2707	23.587
15	-1.1399	25.390
16	-1.0133	27.190
17	-0.8901	28.988
18	-0.7699	30.785
19	-0.6523	32.579
20	-0.5370	34.371
21	-0.4237	36.162
22	-0.3125	37.951
23	-0.2034	39.740
24	-0.0966	41.529
25	0.0075	43.318
26	0.1084	45.109
27	0.2056	46.901
28	0.2986	48.696
29	0.3870	50.491
30	0.4706	52.288
31	0.5494	54.090
32	0.6237	55.891
33	0.6939	57.696
34	0.7606	59.502
35	0.8243	61.313
36	0.8858	63.127
37	0.9455	64.944
38	1.0043	66.766
39	1.0628	68.593
40	1.1215	70.426

Raw Score	Ability	Scale Score
41	1.1815	72.266
42	1.2432	74.111
43	1.3079	75.961
44	1.3767	77.821
45	1.4512	79.692
46	1.5334	81.573
47	1.6262	83.469
48	1.7335	85.381
49	1.8612	87.313
50	2.0182	89.272
51	2.2189	91.268
52	2.4885	93.309
53	2.8817	95.407
54	3.5687	97.577
55	4.7719	100.000

Appendix C: Item Writing Guidelines

GENERAL RULES FOR WRITING MULTIPLE-CHOICE ITEMS

1. Use either a direct question or an incomplete statement as the item stem, whichever seems more appropriate to effective presentation of the item.
Some item ideas can be expressed more simply and clearly in the incomplete statement style of question. On the other hand, some items seem to require direct question stems for the most effective expression. Teachers should use the item style that seems most appropriate.
2. Items should be written in clear and simple language, with vocabulary kept as simple as possible.
Like any other item, the multiple-choice item should be perfectly clear. Difficult and technical vocabulary should be avoided unless essential for the purpose of the question. The important elements should generally appear early in the statement of the item, with qualifications and explanations following.
3. Each item should have one and only one correct answer.
While this requirement is obvious, it is not always fulfilled. Sometimes writers produce items involving issues so controversial and debatable that even experts are unable to agree on one correct answer. More often the trouble is failure to consider the full implications of each response.
4. Base each item on a single central problem.
A multiple-choice item functions most effectively when the student is required to compare directly the relative merits of a number of specific responses to a definite problem. An item consisting merely of a series of unrelated true-false statements, all of which happen to begin with the same phrase, is unacceptable.
5. State the central problem of the item clearly and completely in the stem.
The stem should be meaningful by itself. It should be clear and should convey the central problem of the item. It should not be necessary for the student to read and reread all the responses before he/she can understand the basis upon which he/she is to make a choice.
6. In general, include in the stem any words that must otherwise be repeated in each response.
The stem should contain everything the answers have in common or as much as possible of their common content. This practice serves to make the item shorter, so that it can be read and grasped more quickly.
7. Avoid negative statements.
Negative statements in multiple-choice items lead to unnecessary difficulties and confusion. Special care must be exercised against the double negative.

8. Avoid excessive “window dressing.”
The item should contain only material relevant to its solution, unless selection of what is relevant is part of the problem.
9. Make the responses grammatically consistent with the stem and parallel with one another in form.
10. Make all responses plausible and attractive to students who lack the information or ability tested by the item.
The incorrect responses should be plausible answers. So far as possible, each response should be designed specifically to attract students who have certain misconceptions or who tend to make certain common errors.
11. Arrange the responses in logical order, if one exists.
Where the responses consist of numbers or letters, they should ordinarily be arranged in ascending order. Events should be listed in the order in which they occurred, from earliest to most recent, except when this order would clue the answer. This practice helps insure the student will mark the answer correctly.
12. Make the responses independent and mutually exclusive.
Responses should not be interrelated in meaning. Responses that are not mutually exclusive, aid the student in eliminating wrong answers and reduce the reliability of the item by decreasing the number of effective, functioning responses.
13. Avoid extraneous clues.
Since the student is required to associate one of several alternative responses with the stem, any aspect of the question that provides an extraneous basis for correctly associating the right answer or for eliminating a wrong response constitutes an undesirable clue.
14. Avoid using “all of the above” and “none of the above” as alternatives.
15. Avoid using the phrase “of the following” in the stem.

CHECKLIST OF TEST CONSTRUCTION PRINCIPLES
(Multiple-Choice Items)

	YES	NO
1. Is the item significant?		
2. Does the item have curricular validity?		
3. Is the item presented in clear and simple language, with vocabulary kept as simple as possible?		
4. Does the item have one and only one correct answer?		
5. Does the item state one single central problem completely in the stem? (See Helpful Hint below.)		
6. Does the stem include any extraneous material (“window dressing”)?		
7. Are all responses grammatically consistent with the stem and parallel with one another in form?		
8. Are all responses plausible (attractive to students who lack the information tested by the item)?		
9. Are all responses independent and mutually exclusive?		
10. Are there any extraneous clues due to grammatical inconsistencies, verbal associations, length of response, etc.?		
11. Were the principles of Universal Design used in constructing the item?		

HELPFUL HINT

To determine if the stem is complete (meaningful all by itself):

1. Cover up the responses and read just the stem.
2. Try to turn the stem into a short-answer question by drawing a line after the last word. (If it would not be a good-short answer item you may have a problem with the stem.)
3. The stem must consist of a statement that contains a verb.

Appendix D: Tables and Figures for August 2015 Administration

Table D.1 Multiple-Choice Item Analysis Summary: Regents Comprehensive Examination in English

Item	Number of Students	<i>p</i> -Value	SD	Point-Biserial	Point-Biserial Distractor 1	Point-Biserial Distractor 2	Point-Biserial Distractor 3
1	9,997	0.78	0.41	0.30	-0.17	-0.18	-0.14
2	9,997	0.55	0.50	0.34	-0.18	-0.18	-0.10
3	9,997	0.80	0.40	0.31	-0.15	-0.16	-0.17
4	9,997	0.82	0.39	0.36	-0.21	-0.22	-0.14
5	9,997	0.90	0.30	0.30	-0.20	-0.14	-0.15
6	9,997	0.58	0.49	0.37	-0.14	-0.19	-0.19
7	9,997	0.56	0.50	0.31	-0.13	-0.19	-0.13
8	9,997	0.45	0.50	0.35	-0.19	-0.13	-0.18
9	9,997	0.83	0.38	0.32	-0.15	-0.24	-0.10
10	9,997	0.72	0.45	0.29	-0.19	-0.17	-0.10
11	9,997	0.61	0.49	0.36	-0.16	-0.17	-0.19
12	9,997	0.63	0.48	0.25	-0.14	-0.11	-0.11
13	9,997	0.45	0.50	0.17	-0.22	-0.16	0.08
14	9,997	0.72	0.45	0.16	-0.07	-0.14	-0.05
15	9,997	0.77	0.42	0.30	-0.13	-0.21	-0.11
16	9,997	0.63	0.48	0.31	-0.20	-0.12	-0.13
17	9,997	0.58	0.49	0.27	-0.16	-0.13	-0.10
18	9,997	0.37	0.48	0.28	-0.08	-0.11	-0.13
19	9,997	0.86	0.35	0.37	-0.24	-0.17	-0.19
20	9,997	0.50	0.50	0.34	-0.19	-0.11	-0.17
21	9,997	0.63	0.48	0.36	-0.19	-0.22	-0.11
22	9,997	0.53	0.50	0.34	-0.18	-0.15	-0.16
23	9,997	0.86	0.34	0.27	-0.13	-0.16	-0.15
24	9,997	0.70	0.46	0.31	-0.15	-0.12	-0.18
25	9,997	0.61	0.49	0.33	-0.14	-0.18	-0.16

Table D.2 Constructed-Response Item Analysis Summary: Regents Comprehensive Examination in English

Item	Min. score	Max. score	Number of Students	Mean	SD	<i>p</i> -Value	Point-Biserial
26	0	2	9,997	1.39	0.63	0.70	0.68
27	0	2	9,997	1.31	0.67	0.66	0.69
28	0	6	9,997	2.76	1.29	0.46	0.80

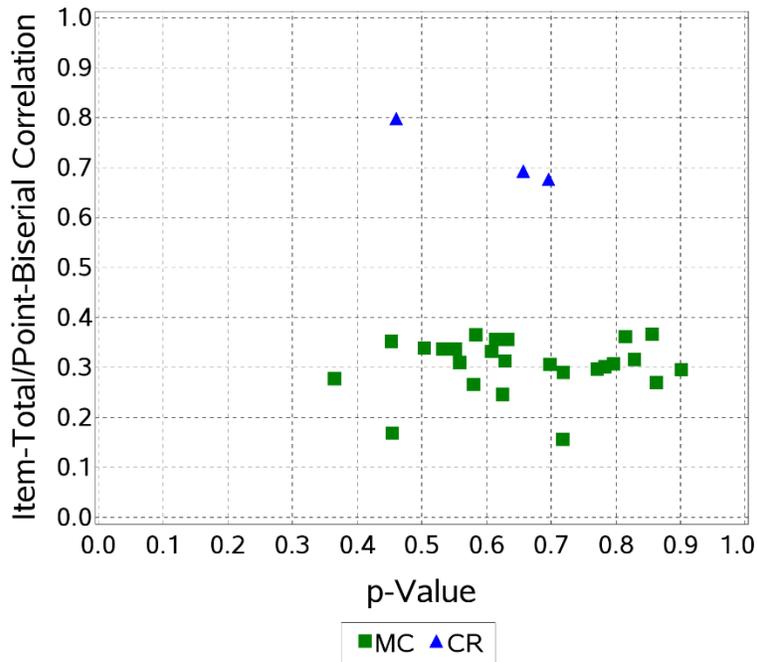


Figure D.1 Scatter Plot: Regents Comprehensive Examination in English

Table D.3 Descriptive Statistics in *p*-value and Point-Biserial Correlation: Regents Comprehensive Examination in English

Statistics	N	Mean	Min	Q1	Median	Q3	Max
<i>p</i> -value	28	0.65	0.37	0.56	0.63	0.78	0.90
Point-Biserial	28	0.35	0.16	0.29	0.32	0.36	0.80

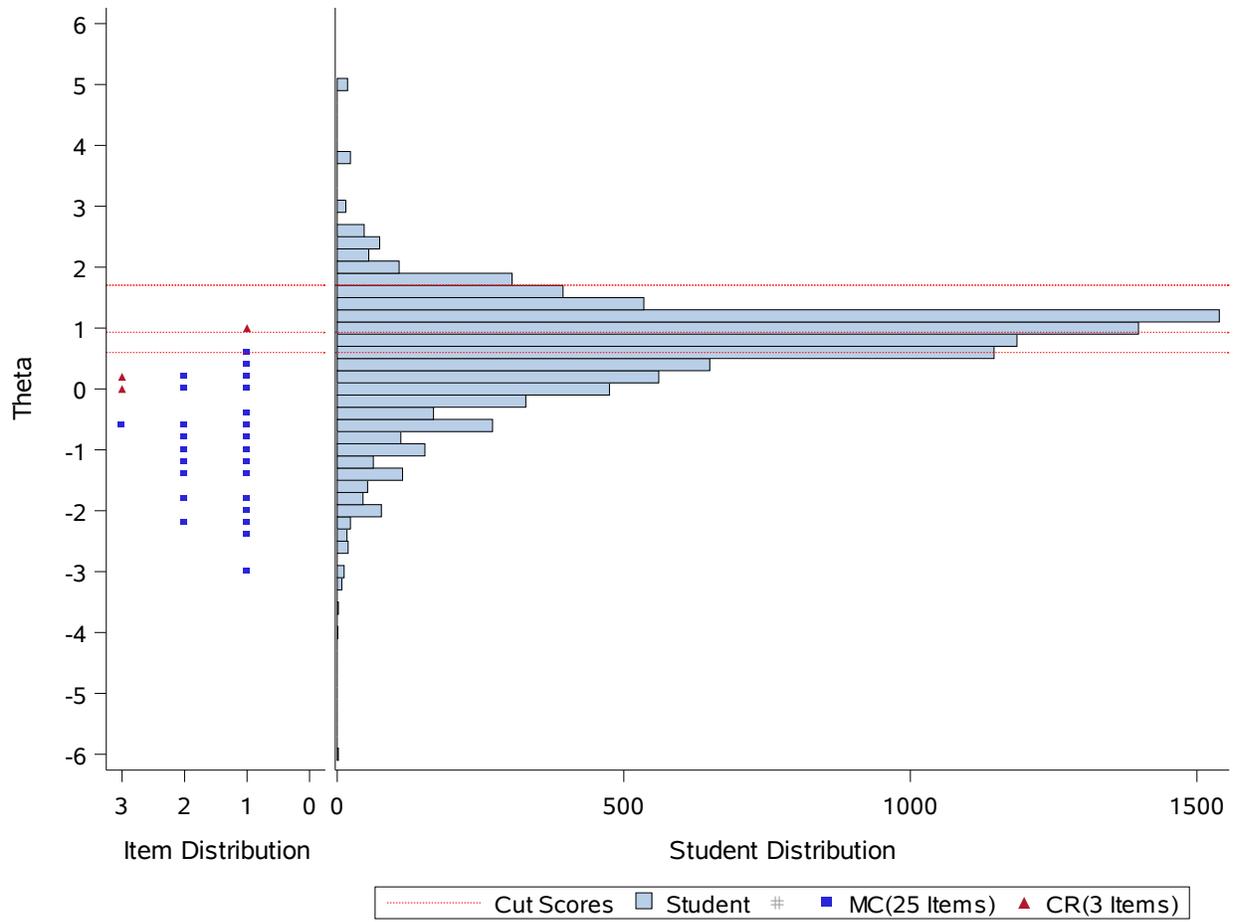


Figure D.2 Student Performance Map: Regents Comprehensive Examination in English

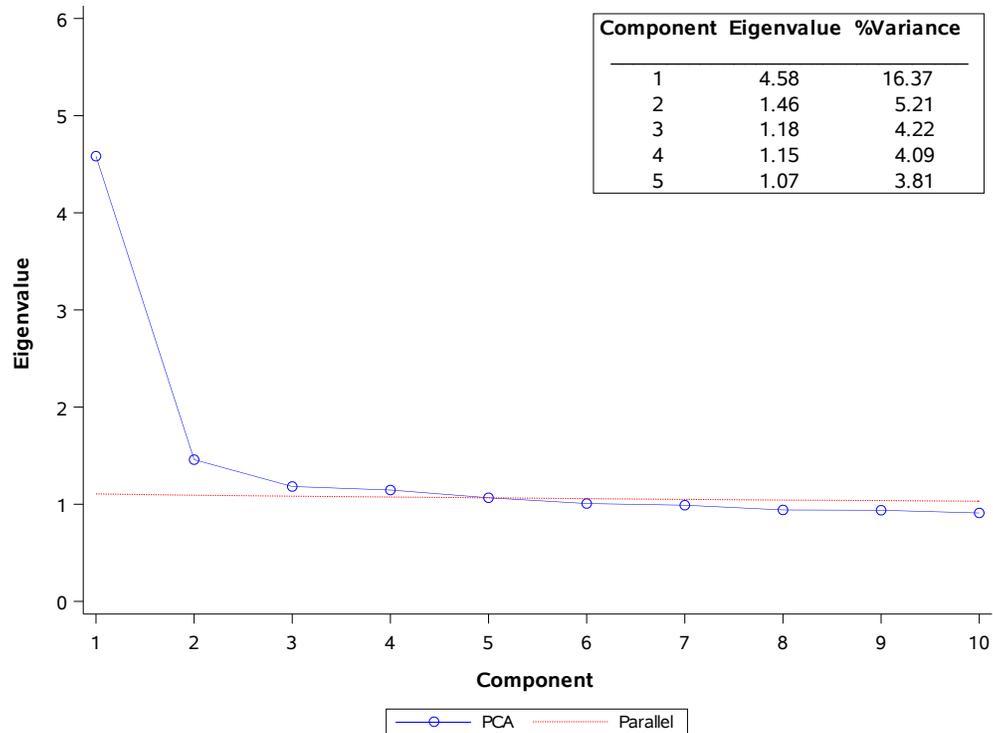


Figure D.3 Scree Plot: Regents Comprehensive Examination in English

Table D.4 Summary of Item Residual Correlations: Regents Comprehensive Examination in English

Statistic Type	Value
N	378
Mean	-0.03
SD	0.05
Minimum	-0.17
P ₁₀	-0.09
P ₂₅	-0.05
P ₅₀	-0.03
P ₇₅	-0.01
P ₉₀	0.01
Maximum	0.31
> 0.20	1

Table D.5 Summary of INFIT Mean Square Statistics: Regents Comprehensive Examination in English

	INFIT Mean Square					
	N	Mean	SD	Min	Max	[0.7, 1.3]
English	28	1	0.07	0.89	1.18	[28/28]

Table D.6 Reliabilities and Standard Errors of Measurement: Regents Comprehensive Examination in English

Subject	Coefficient Alpha	SEM
English	0.79	4.24

Table D.7 Decision Consistency and Accuracy Results: Regents Comprehensive Examination in English

Statistic	1/2	2/3	3/4
Consistency	0.82	0.79	0.92
Accuracy	0.87	0.85	0.94

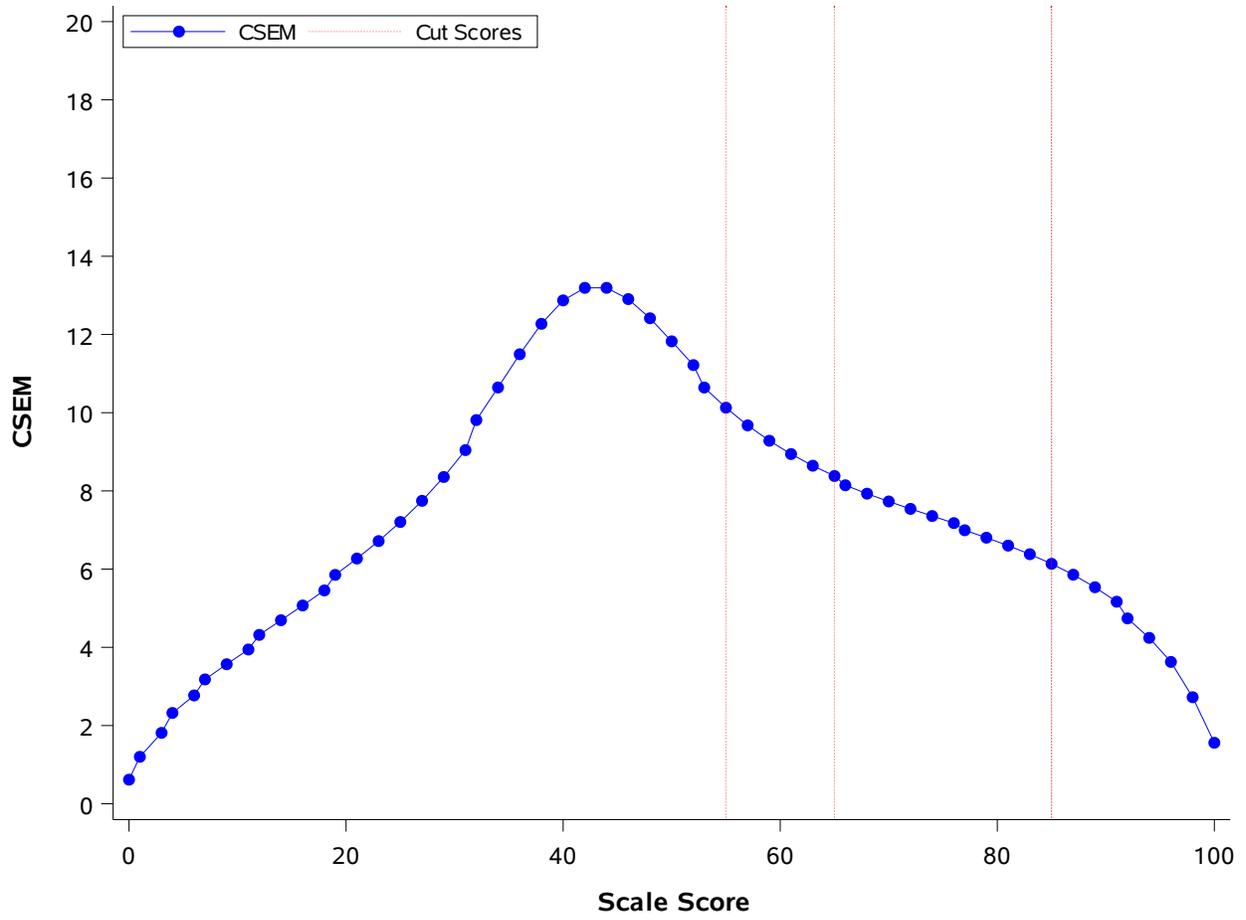


Figure D.4 Conditional Standard Error Plot: Regents Comprehensive Examination in English

Table D.8 Group Means: Regents Comprehensive Examination in English

Demographics	Number	Mean Scale Score	SD Scale Score
All Students*	9,997	59.53	17.38
Ethnicity			
American Indian/Alaska Native	79	57.96	16.28
Asian/Native Hawaiian/Other Pacific Islander	1,029	60.46	18.31
Black/African American	2,889	58.27	16.61
Hispanic/Latino	3,475	57.74	17.63
Multiracial	71	65.89	15.37
White	2,379	63.12	17.04
English Language Learner			
No	7,786	62.13	16.48
Yes	2,211	50.37	17.40
Economically Disadvantaged			
No	3,637	61.51	17.40
Yes	6,360	58.39	17.27
Gender			
Female	4,257	61.38	17.12
Male	5,665	58.13	17.45
Student with Disabilities			
No	7,176	61.95	17.19
Yes	2,821	53.36	16.32

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

Appendix E: Tables and Figures for January 2016 Administration

Table E.1 Multiple-Choice Item Analysis Summary: Regents Comprehensive Examination in English

Item	Number of Students	<i>p</i> -Value	SD	Point-Biserial	Point-Biserial Distractor 1	Point-Biserial Distractor 2	Point-Biserial Distractor 3
1	16,535	0.70	0.46	0.38	-0.16	-0.21	-0.22
2	16,535	0.66	0.47	0.32	-0.13	-0.14	-0.22
3	16,535	0.78	0.41	0.34	-0.16	-0.16	-0.21
4	16,535	0.70	0.46	0.29	-0.13	-0.18	-0.14
5	16,535	0.76	0.43	0.31	-0.14	-0.17	-0.17
6	16,535	0.52	0.50	0.35	-0.21	-0.16	-0.10
7	16,535	0.57	0.50	0.36	-0.16	-0.23	-0.12
8	16,535	0.68	0.47	0.32	-0.17	-0.18	-0.14
9	16,535	0.76	0.43	0.37	-0.21	-0.22	-0.13
10	16,535	0.60	0.49	0.30	-0.10	-0.20	-0.11
11	16,535	0.70	0.46	0.36	-0.16	-0.20	-0.19
12	16,535	0.41	0.49	0.21	-0.01	-0.15	-0.14
13	16,535	0.66	0.47	0.32	-0.16	-0.11	-0.23
14	16,535	0.76	0.43	0.21	-0.12	-0.09	-0.10
15	16,535	0.83	0.37	0.35	-0.16	-0.24	-0.15
16	16,535	0.59	0.49	0.36	-0.13	-0.17	-0.21
17	16,535	0.61	0.49	0.22	-0.15	-0.13	-0.05
18	16,535	0.68	0.47	0.19	0.00	-0.20	-0.17
19	16,535	0.82	0.38	0.34	-0.22	-0.19	-0.12
20	16,535	0.62	0.48	0.29	-0.20	-0.08	-0.17
21	16,535	0.67	0.47	0.31	-0.19	-0.19	-0.14
22	16,535	0.41	0.49	0.29	-0.11	-0.16	-0.09
23	16,535	0.52	0.50	0.26	-0.07	-0.16	-0.14
24	16,535	0.40	0.49	0.23	-0.05	-0.10	-0.13
25	16,535	0.47	0.50	0.24	-0.12	-0.15	-0.06

Table E.2 Constructed-Response Item Analysis Summary: Regents Comprehensive Examination in English

Item	Min. score	Max. score	Number of Students	Mean	SD	<i>p</i> -Value	Point-Biserial
26	0	2	16,535	1.31	0.66	0.66	0.69
27	0	2	16,535	1.27	0.67	0.63	0.70
28	0	6	16,535	2.75	1.27	0.46	0.79

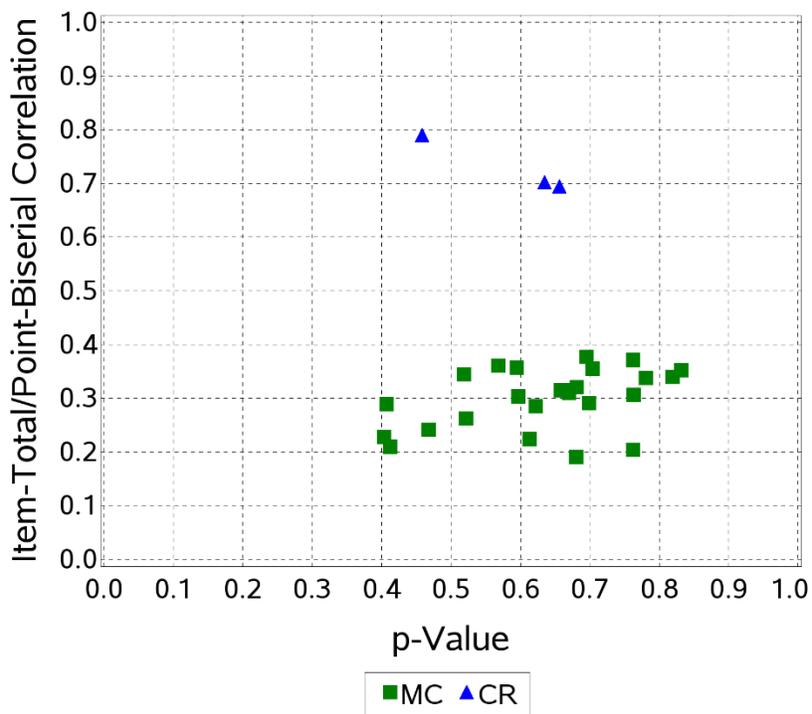


Figure E.1 Scatter Plot: Regents Comprehensive Examination in English

Table E.3 Descriptive Statistics in *p*-value and Point-Biserial Correlation: Regents Comprehensive Examination in English

Statistics	N	Mean	Min	Q1	Median	Q3	Max
<i>p</i> -value	28	0.63	0.40	0.55	0.66	0.70	0.83
Point-Biserial	28	0.35	0.19	0.27	0.32	0.36	0.79

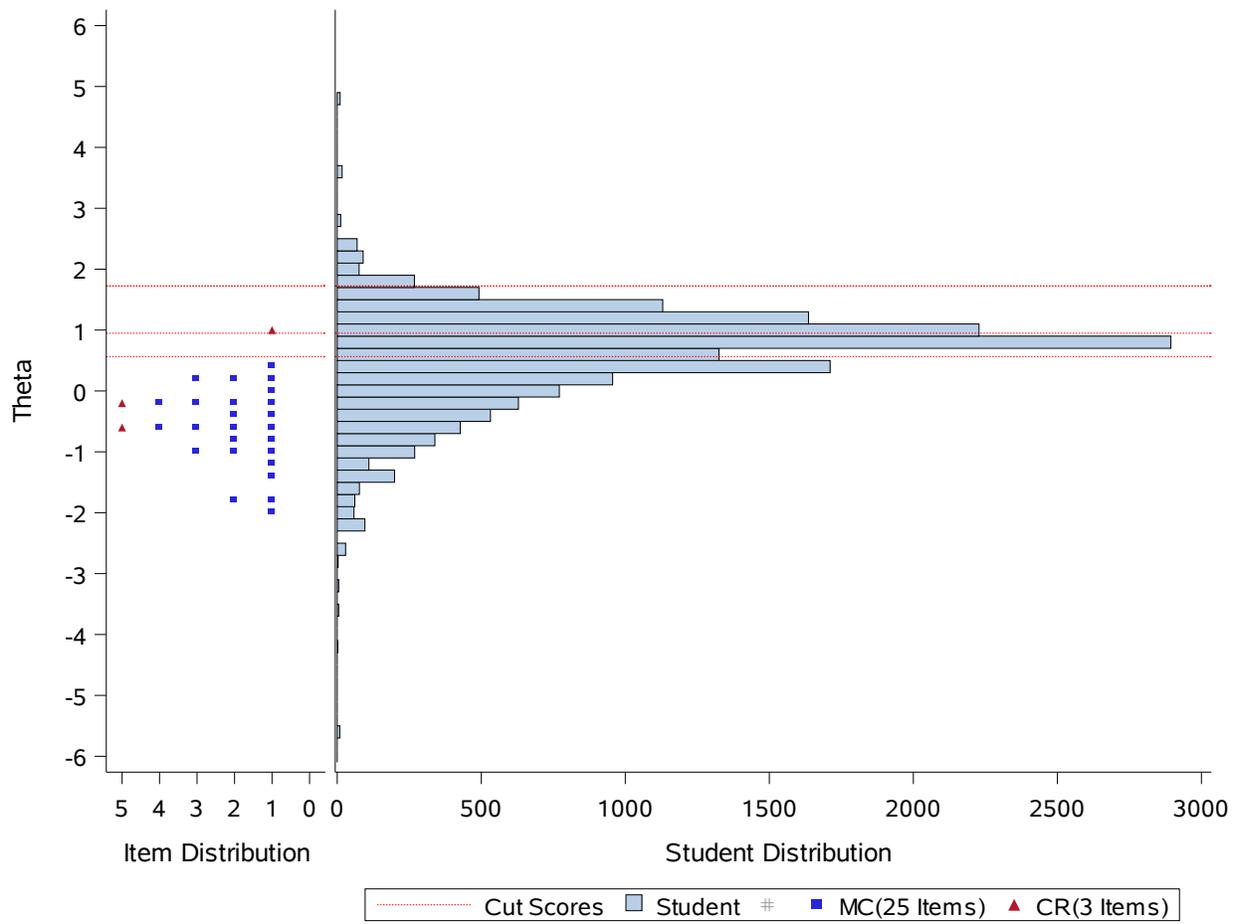


Figure E.2 Student Performance Map: Regents Comprehensive Examination in English

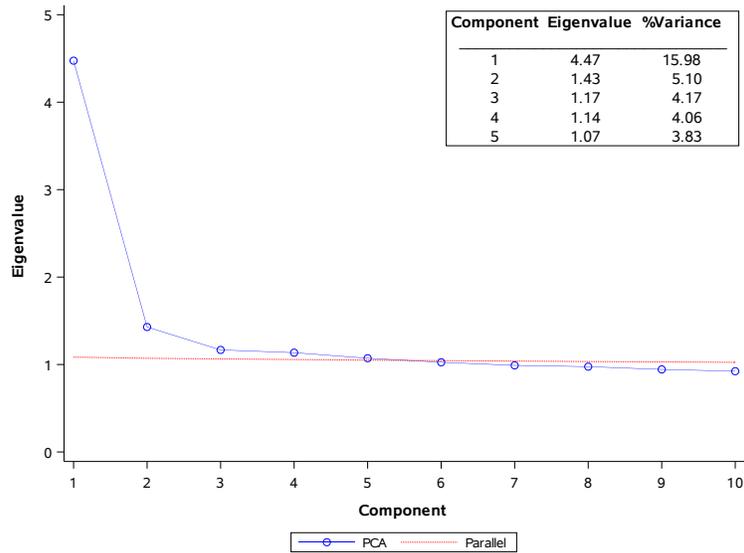


Figure E.3 Scree Plot: Regents Comprehensive Examination in English

Table E.4 Summary of Item Residual Correlations: Regents Comprehensive Examination in English

Statistic Type	Value
N	378
Mean	-0.03
SD	0.05
Minimum	-0.16
P ₁₀	-0.09
P ₂₅	-0.05
P ₅₀	-0.03
P ₇₅	-0.01
P ₉₀	0.01
Maximum	0.34
> 0.20	1

Table E.5 Summary of INFIT Mean Square Statistics: Regents Comprehensive Examination in English

	INFIT Mean Square					
	N	Mean	SD	Min	Max	[0.7, 1.3]
English	28	1	0.07	0.89	1.13	[28/28]

Table E.6 Reliabilities and Standard Errors of Measurement: Regents Comprehensive Examination in English

Subject	Coefficient Alpha	SEM
English	0.79	4.33

Table E.7 Decision Consistency and Accuracy Results: Regents Comprehensive Examination in English

Statistic	1/2	2/3	3/4
Consistency	0.80	0.79	0.96
Accuracy	0.86	0.85	0.97

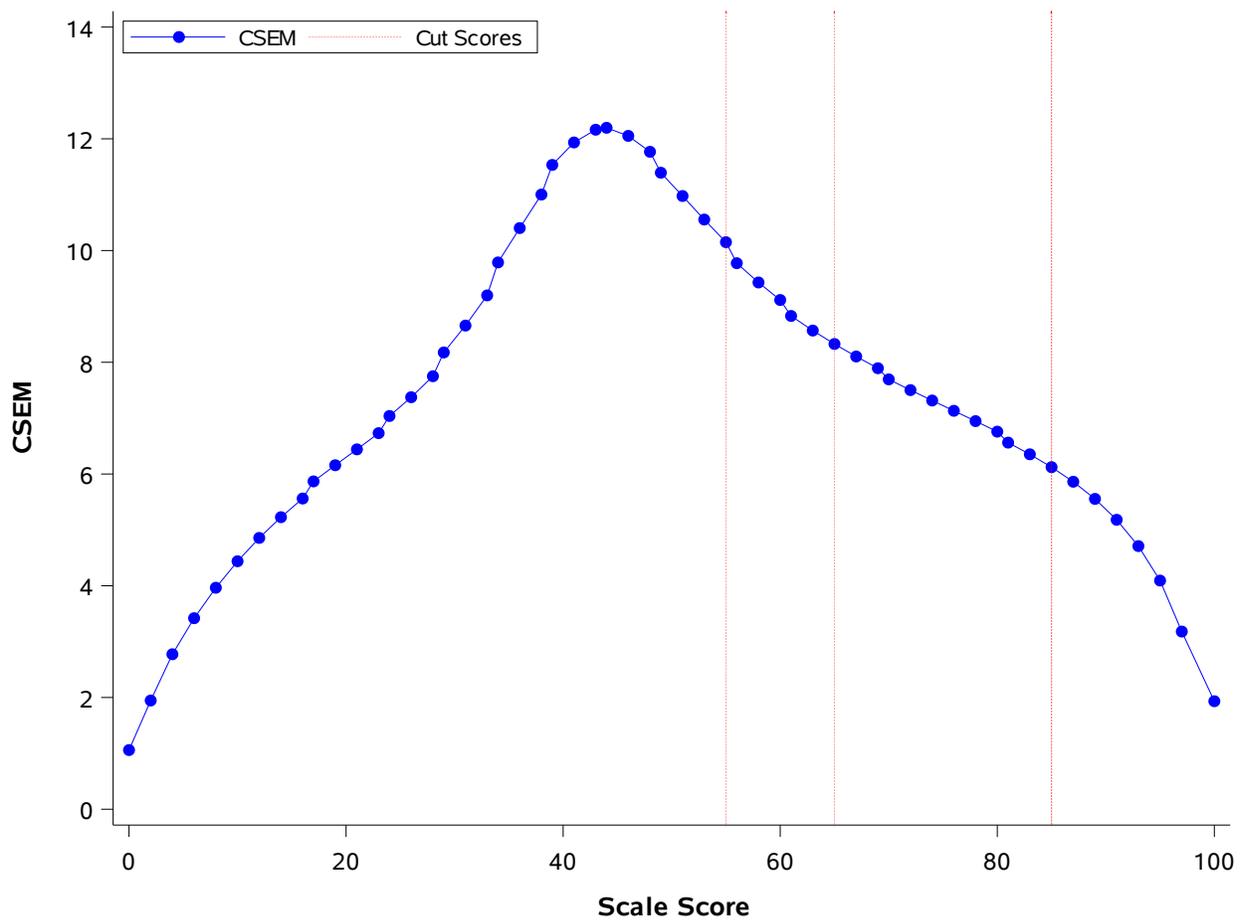


Figure E.4 Conditional Standard Error Plot: Regents Comprehensive Examination in English

Table E.8 Group Means: Regents Comprehensive Examination in English

Demographics	Number	Mean Scale Score	SD Scale Score
All Students*	16,535	57.22	16.31
Ethnicity			
American Indian/Alaska Native	135	55.72	14.80
Asian/Native Hawaiian/Other Pacific Islander	1,756	55.51	16.50
Black/African American	4,765	57.05	15.46
Hispanic/Latino	6,260	55.40	16.37
Multiracial	125	62.43	14.29
White	3,473	61.56	16.53
English Language Learner			
No	12,637	60.51	15.10
Yes	3,898	46.58	15.53
Economically Disadvantaged			
No	5,054	59.81	17.02
Yes	11,481	56.09	15.86
Gender			
Female	6,785	58.79	15.79
Male	9,729	56.16	16.58
Student with Disabilities			
No	11,800	59.42	16.26
Yes	4,735	51.74	15.11

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