



The Accuracy of Students' Novice, Apprentice, Proficient, and Distinguished Classifications for the 2014 Kentucky Performance Rating for Educational Progress (K-PREP) Tests

Prepared for: Kentucky Department of Education
Capital Plaza Tower, 18th Floor
500 Mero Street
Frankfort, KY 40601

Authors: Emily R. Dickinson
Arthur A. Thacker

Prepared under: 1200003303, FFP

Date: August 27, 2015

The Accuracy of Students’ Novice, Apprentice, Proficient, and Distinguished Classifications for the 2014 Kentucky Performance Rating for Educational Progress (K-PREP) Tests

Table of Contents

Background.....	1
Methods and Results	1
Classification Accuracy	1
Reading the Classification Accuracy Tables using Spring 2014 K-PREP Grade 4 Reading as an Example.....	2
Summary of the Results.....	3
Discussion and Conclusions	4
References.....	7
Appendix A. Classification Accuracy Tables for K-PREP Spring 2014	9
Appendix B. Technical Details	15

List of Tables

Table 1. Grade/Subject Combinations for the K-PREP test Analyzed for Classification Accuracy	1
Table 2. Grade 4 Reading 2014 Percentages of True Scores Being in Assigned Classification.....	3
Table 3. Total Percent Expected Correct Classifications and Average Category Distribution Error with Comparison Data from K-PREP 2012.....	4
Table 4. Comparison of K-PREP and Tests Total Percent Expected Correct Classifications with Three Other States.....	5

The Accuracy of Students’ Novice, Apprentice, Proficient, and Distinguished Classifications for the 2014 Kentucky Performance Rating for Educational Progress (K-PREP) Tests

Background

The Kentucky Performance Rating for Educational Progress (K-PREP) assessment has been administered in Kentucky since 2012. These tests were developed in 2012 and are aligned with the Kentucky Core Academic Standards (KCAS). For scoring and reporting, each grade/subject combination is treated as a separate test. Based on the results of these tests, each student is classified into one of four proficiency levels: Novice, Apprentice, Proficient, or Distinguished (NAPD). The purpose of this report is to present classification accuracy statistics for the Spring 2014 administration of K-PREP. Classification accuracy is a method for considering the reliability of a test.

Grade/subject combinations included in the classification accuracy calculations are presented in Table 1 below. Scoring is a two-step process described in the Kentucky Department of Education (KDE) technical manual (KDE, 2012). Students first receive a scale score derived from their responses to the items on the test. “Cut points” have been set to allow for the categorization of student scores into performance levels- Novice, Apprentice, Proficient, and Distinguished (NAPD). New performance level descriptors and cut score standards were developed to indicate mastery levels that are needed to be considered “on track” for college and career readiness.

Table 1. Grade/Subject Combinations for the K-PREP test Analyzed for Classification Accuracy

Subject	Grades					
	3	4	5	6	7	8
Reading	X	X	X	X	X	X
Mathematics	X	X	X	X	X	X
Science		X			X	
Social Studies			X			X

Tests are more useful when classification accuracy for NAPD levels is high—however, no test is perfect. This report examines the accuracy of the K-PREP NAPD assignments for Spring 2014. Results from prior years have been reported by Hoffman and Wise (2000b), Hoffman, Wise, and Thacker (2000), Hoffman (2002 and 2003), Hoffman and Dickinson (2004), Hoffman and Nemeth (2008), and Dickinson, Thacker, Levinson, and Hoffman (2013) and are available at <http://education.ky.gov> (search for “student classification accuracy”).

Methods and Results

Classification Accuracy

The methodology for this classification accuracy analysis was originally developed by Hoffman and Wise (1999) and presented to Kentucky’s National Technical Advisory Panel on Assessment and Accountability (NTAPAA) on two occasions (September 9-10, 1999 and

December 16-17, 1999). The method was approved by the NTAPAA during the September meeting. Preliminary results for the 1999 assessments were presented during the December meeting. The classification accuracy method was also presented to the National Council for Measurement in Education (NCME) at its annual meeting in April 2000. The NCME paper (Hoffman & Wise, 2000a) is available from the authors. While the present report conforms to the NTAPAA reporting specifications and uses the established methodology, we do note that student classification programs were rewritten in 2013 to adjust for changes in the item response model used by the new contractor (Generalized Partial Credit model to a Rasch model) and the removal of multiple-choice and open-response weighting procedures. More detailed methodological information is provided in Appendix B.

As mentioned above, no test is perfect, and the reliability of an assessment is of particular interest when scores are used to categorize students. Reliability of an observed test score is the product of two factors: true proficiency in the knowledge area being assessed and measurement error that comes from a variety of sources. Obtained scores are known, but true scores are unknown. Using test reliability-related statistics, however, it is possible to provide estimations that answer the following two questions:

1. For a given obtained score, what are the odds that true proficiency is in the same NAPD classification?
2. For that given obtained score, what are the odds that true proficiency falls into a different NAPD classification?

These two questions lead to 16 probability estimates: that is, for each of the four assigned NAPD proficiency levels, what are the odds of true proficiency at each level? These probability estimates are presented in classification accuracy tables for each grade/subject combination. Classification accuracy tables for the K-PREP are presented in Appendix A.

Reading the Classification Accuracy Tables using Spring 2014 K-PREP Grade 4 Reading as an Example

In Table 2 and similar tables in Appendix A, numbers represent percentages of all students, so that the sum of all of the italicized percentages is 100 (within rounding). The “Total % Assigned” row near the bottom of the table indicates the percent of students who were assigned each of the four NAPD classifications. For example, in Table 2, 27.05% of all Grade 4 students who took the K-PREP in Reading received test scores that placed them in the Novice category. Likewise, 31.36% of all students received test scores within the score range for the Apprentice category; 31.36% were Proficient; and 10.23% were Distinguished.

Since test scores are not perfect, some proportion of students are expected to have true achievement in categories matching their assigned categories, with the remaining students expected to have true achievement that falls in categories other than their assigned categories. The bold italicized numbers in Table 2 indicate proportions of accurate classifications. That is, 23.42% of all students are expected to be accurately classified as Novice, 21.99% of all students are expected to be accurately classified as Apprentice, 23.39% of all students are expected to be accurately classified as Proficient, and 5.78% of all students are expected to be accurately classified as Distinguished. The sum of these four percentages (74.59), labeled “Total % Expected Correct Assignments,” is the percent of all students expected to be classified accurately. That is, approximately 78% of all Grade 4 students would be assigned to the same category of reading proficiency if we actually knew their true achievement.

Table 2. Grade 4 Reading 2014 Percentages of True Scores Being in Assigned Classification

True Classification	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	23.42	<i>5.03</i>	<i>0.04</i>	<i>0.00</i>	<i>28.49</i>
Apprentice	<i>3.61</i>	21.99	<i>6.36</i>	<i>0.02</i>	<i>31.98</i>
Proficient	<i>0.01</i>	<i>4.34</i>	23.39	<i>4.43</i>	<i>32.17</i>
Distinguished	<i>0.00</i>	<i>0.00</i>	<i>1.57</i>	5.78	<i>7.36</i>
Total % Assigned	<i>27.05</i>	<i>31.36</i>	<i>31.36</i>	<i>10.23</i>	<i>100</i>
Total % Expected Correct Assignments	74.59%		Average Distribution Error:		1.44

The numbers in Table 2 in non-bold italics indicate the proportions of students expected to have true achievement classifications that are different from their assigned classification. For example, 3.61% of all students are expected to have obtained test scores that place them in the Novice range while their true achievement would place them one category higher in the Apprentice category. Conversely, 5.03% of all students are expected to have obtained test scores that place them in the Apprentice category, while their true achievement would place them one category lower in the Novice category. Another 4.34% of all students are expected to have obtained test scores that also place them in the Apprentice category, while their true achievement would place them one category higher in the Proficient category. In total, 25.41% (100-74.59%) of all students are expected to be misclassified in Grade 4 reading in 2014.

Student classification accuracy data also has important implications for school accountability scores. School accountability scores are a function of all students' classifications. Some of the inevitable classification error will be in one direction and some in the other. As seen in Table 2, some proportion of students are expected be classified higher than their true proficiency and some lower. The percentage of students actually assigned to a particular category versus our projections of the percent of students expected to have true achievement at that level shows that the misclassification errors tend to balance out for the total student population. For example, the last column in Table 2 shows that 31.98% of the students are expected to be Apprentice based on their unknowable true scores, while 31.36% of the students are assigned the Apprentice classification based on their test performance. The difference between these percentages is approximately 0.62. Differences between the expected and assignment distributions for the other three categories also tend to be small, with the average difference in category percentages being 1.44. Because this error refers to category total distributions, it is referred to as "Average Distribution Error" in the table. Results tables for all K-PREP grade/subjects, for 2014, are presented in the Appendix A.

Summary of the Results

Table 3 shows the total expected correct assignments for all K-PREP grade/subject combinations for 2014, with the results for K-PREP 2012 added for comparison. In 2014, student classification accuracy varies from approximately 72.2% (Grade 3 Reading) to approximately 81.3% (Grades 8 Social Studies). Overall, Reading has the lowest accuracy with results between 72.2%-76.0%. Mathematics, Science, and Social Studies have accuracies in the 78.0% to 81.3% range. Classification accuracy rates for 2014 K-PREP are very similar to the 2012 rates, with differences (i.e., absolute value(2014 accuracy rates - 2012 accuracy rates)) ranging from 0% to 2.8%.

Table 3. Total Percent Expected Correct Classifications and Average Category Distribution Error with Comparison Data from K-PREP 2012

Subject/Grade	Total Percent Expected Correct Classifications		Average Category Distribution Error	
	K-PREP 2012	K-PREP 2014	K-PREP 2012	K-PREP 2014
MA 03	78.9	78.0	2.1	2.2
MA 04	78.3	78.2	1.7	1.1
MA 05	79.9	79.5	1.1	1.5
MA 06	80.1	80.0	1.5	1.9
MA 07	80.2	80.0	2.5	2.4
MA 08	80.6	79.8	1.4	0.9
RD 03	73.4	72.2	0.4	2.1
RD 04	74.3	74.6	1.5	1.4
RD 05	74.8	72.8	1.0	1.7
RD 06	76.9	76.0	1.4	1.6
RD 07	74.5	74.5	1.5	2.2
RD 08	75.5	72.7	1.4	0.9
SC 04	78.1	78.9	2.1	1.3
SC 07	78.5	79.1	1.4	0.9
SS 05	80.5	80.9	1.7	1.5
SS 08	80.2	81.3	2.3	1.2

Turning from the individual level accuracy data to the distribution accuracy results, the two right-most columns in Table 3 summarize how well the distribution of assigned classifications matches the expected distribution of true classifications by showing the difference between expected and assigned total percentages averaged across the four achievement levels. As would be expected, this average error tends to be larger among the subject areas with lower level accuracy values.

Discussion and Conclusions

Test specialists are well aware of the need to study classification accuracy as well as more traditional measures of test reliability. Several methodological papers focusing on analytical variations of the accuracy theme have used operational data. For example, Rogosa (1994) examined 1993 California’s CLAS assessment, which uses six proficiency levels. He found that although the probability of classification within one category of true proficiency was nearly 95%, the probability of exact classification was only 51.72%.

In another example, Lee, Hanson, and Brennan (2000) used data from ACT’s Work Keys assessment. Their results confirm that the number of proficiency categories makes a difference – more categories mean more opportunities of classification error. For a Work Keys subtest with five categories, exact accuracy for several different forms was in the 70% range, while a subtest with six categories showed accuracy in the low- to mid-60% range. Lee et al. also looked at accuracy for classifying students simply above or below a single cutpoint, using each of the possible Work Keys cutpoints to look at these dichotomous classifications. Accuracy was in the upper 80% range to near 100% for classifying students into only one of two categories. The higher levels of accuracy occurred for classification of students into either extreme. When the

cutpoint was closer to the center, accuracy tended to be in the upper 80% range. Young and Yoon (1998) provide similar data from the New Standards assessments. Again, when making only a dichotomous (two-category) classification, they showed better accuracy (e.g., in the lower 90% range).

Table 4 presents the total percent of students that would be expected to be correctly classified in Kentucky, calculated for both four proficiency categories and two proficiency categories, along with percentages for three other states' past test administrations. When Kentucky students were classified as 'Proficient or Not' (by summing the four cells in the upper left and the four cells in the lower right of each table in the Appendix A), average classification accuracy across all tests jumped from 77.4% to 90.4%. California (Educational Testing Service, 2013) and Florida (HumRRO & Harcourt, 2007) use five proficiency categories, whereas Texas (Texas Education Agency, 2010) reports student classification accuracy as the percentage correctly classified as Meeting Standard or not. Table 4 shows that Kentucky demonstrates classification accuracy rates on par with other state assessments, and actually has higher accuracy rates when the two category system is compared to similar results reported by Texas. This is particularly pertinent given that the Kentucky uses a dichotomous classification system with the goal of increasing student performance to Proficient or greater as a part of the Unbridled Learning accountability system for students, schools, and districts.

Table 4. Comparison of K-PREP and Tests Total Percent Expected Correct Classifications with Three Other States

		Kentucky	Kentucky ¹	California ²	Florida ²	Texas ¹
Grade 3	Reading ³	72.2	89.3	76.0	71.8	80.1
	Math	78.0	90.6	81.0	70.1	80.9
Grade 4	Reading	74.6	89.2	81.0	66.6	81.6
	Math	78.2	91.1	82.0	67.4	81.6
	Science	78.9	90.7			
Grade 5	Reading	72.8	88.0	81.0	65.9	82.0
	Math	79.5	90.9	80.0	66.6	81.1
	Science			77.0		77.0
	Social Studies	80.9	90.7			
Grade 6	Reading	76.0	89.4	79.0	96.4	80.8
	Math	80.0	91.5	79.0	60.8	83.2
Grade 7	Reading	74.5	89.6	78.0	70.0	81.5
	Math	80.0	92.7	77.0	62.1	85.8
	Science	79.1	90.9			
Grade 8	Reading	72.7	87.6	79.0	64.2	80.3
	Math	79.8	91.5	74.0	61.8	84.7
	Science			77.0		82.1
	Social Studies	81.3	92.3	77.0		85.5

¹ Classification accuracy using two categories. ² Classification accuracy using five categories.

³ ELA is tested in California.

Given these examples, the K-PREP appears to have classification accuracy statistics that are similar to other educational proficiency assessments. We have also seen in this report that individual level inaccuracies tend to cancel out so that the distributions of students' scores on the aggregate level appear to be reasonably precise.

References

- Dickinson, E. R., Thacker, A. A., Levinson, H., & Hoffman, R. G. (2013). *The accuracy of students' novice, apprentice, proficient, and distinguished classifications for the 2012 Kentucky Performance Rating for Educational Progress (K-PREP) tests* (2013 No. 037). Alexandria, VA: Human Resources Research Organization.
- Educational Testing Service. (2013). California Standards Test Technical Report Spring 2012 Administration. Retrieved on June 17, 2013 from <http://www.cde.ca.gov/ta/tg/sr/documents/cst12techrpt.pdf#search=STARS%20student%20classification%20accuracy&view=FitH&page=ode=none>.
- Hoffman, R. G., Diaz, T. E., & Dickinson E. (2005). *Idaho Standards Achievement Test: Independent calculations of reliability estimates, standard errors of measurement, classification accuracy, and classification consistency*. (FR-04-87). Alexandria, VA: Human Resources Research Organization.
- Hoffman, R. G. (2003). *The accuracy of students' novice, apprentice, proficient, and distinguished classifications for the 2003 Kentucky Core Content Tests* (FR-03-73). Alexandria, VA: Human Resources Research Organization.
- Hoffman, R. G. (2002). *The accuracy of students' novice, apprentice, proficient, and distinguished classifications for the 2001 and 2002 Kentucky Core Content Tests* (FR-02-46). Alexandria, VA: Human Resources Research Organization.
- Hoffman, R. G., & Dickinson, E. R. (2004). *The accuracy of students' novice, apprentice, proficient, and distinguished classifications for the 2004 Kentucky Core Content Tests* (FR-04-77). Alexandria, VA: Human Resources Research Organization.
- Hoffman, R. G., & Nemeth, Y. M. (2008). *The accuracy of students' novice, apprentice, proficient, and distinguished classifications for the 2008 Kentucky Core Content Tests* (FR-08-108). Alexandria, VA: Human Resources Research Organization.
- Hoffman, R. G., & Wise, L. L. (1999). *Establishing the reliability of student level classifications: Analytic plan and demonstration* (FR-WATSD-99-34). Alexandria, VA: Human Resources Research Organization.
- Hoffman, R. G., & Wise, L. L. (2000a). *Establishing the reliability of student proficiency classifications: The accuracy of observed classifications*. Paper presented at the annual meeting of the National Council of Measurement in Education, New Orleans, April, 2000. Also available at www.Humrro.org.
- Hoffman, R. G., & Wise, L. L. (2000b). *The accuracy of students' novice, apprentice, proficient, and distinguished classification of the Kentucky Core Content Test for 1999*. (FR-WATSD-00-25). Alexandria, VA: Human Resources Research Organization.
- Hoffman, R. G., Wise, L. L., & Thacker, A. A. (2000). *The accuracy of students' novice, apprentice, proficient, and distinguished classification of the 2000 Kentucky Core Content Test for 1999*. (FR-00-41). Alexandria, VA: Human Resources Research Organization.

Human Resources Research Organization and Harcourt Educational Measurement (2004). *Florida Comprehensive Assessment Test for reading and mathematics: Technical report for test administrations of FCAT 2004*. San Antonio, TX: Harcourt Educational Measurement.

Kentucky Department of Education (2012). *Kentucky performance rating for educational progress 2011-12 Technical Report*. Frankfort, KY: Author.

Lee, W., Hanson, B. A., & Brennan, R. L. (2000). *Procedures for computing classification consistency and accuracy indices with multiple categories*. Paper presented at the annual meeting of the National Council of Measurement in Education, New Orleans, April, 2000.

Rogosa, D. (1994). Misclassification in student performance levels. In CTB/McGraw-Hill. (1994). *1994 CLAS Assessment Technical Report*. Monterrey, CA: Author.

Texas Education Agency. (2010). Technical digest for the academic year 2008-2009. Retrieved on June 17, 2013 from http://www.tea.state.tx.us/index3.aspx?id=2147484418&menu_id=793.

Young, M. J. & Yoon, B. (1998). *Estimating the consistency and accuracy of classifications in a standards-referenced assessment*. (CSE Technical Report 475). Los Angeles: Center for the Study of Evaluation, National Center for Research on Evaluation, Standards, and Student Testing, Graduate School of Education & Information Studies, University of California, Los Angeles.

Appendix A. Classification Accuracy Tables for K-PREP Spring 2014

Table A-1. Grade 3 Mathematics 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	18.66	3.41	0.00	0.00	22.07
Apprentice	3.24	28.72	5.37	0.00	37.33
Proficient	0.00	4.06	26.63	5.15	35.83
Distinguished	0.00	0.00	0.74	4.02	4.76
Total % Assigned	21.89	36.19	32.74	9.18	100
Total % Expected Correct Assignments: 78.03					Average Distribution Error: 2.21

Table A-2. Grade 4 Mathematics 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	17.64	3.97	0.00	0.00	21.62
Apprentice	4.00	25.17	5.46	0.00	34.63
Proficient	0.00	3.39	26.53	3.38	33.31
Distinguished	0.00	0.00	1.60	8.84	10.44
Total % Assigned	21.65	32.53	33.6	12.22	100
Total % Expected Correct Assignments: 78.19					Average Distribution Error: 1.05

Table A-3. Grade 5 Mathematics 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	15.17	3.03	0.00	0.00	18.19
Apprentice	4.11	31.10	5.55	0.00	40.77
Proficient	0.00	3.55	26.12	2.58	32.26
Distinguished	0.00	0.00	1.67	7.11	8.78
Total % Assigned	19.28	37.68	33.35	9.68	100
Total % Expected Correct Assignments: 79.50					Average Distribution Error: 1.54

Table A-4. Grade 6 Mathematics 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	18.89	2.98	0.02	0.00	21.90
Apprentice	4.22	15.10	3.39	0.00	22.70
Proficient	0.05	4.99	40.42	3.23	48.70
Distinguished	0.00	0.00	1.15	5.55	6.70
Total % Assigned	23.16	23.08	44.98	8.78	100
Total % Expected Correct Assignments: 79.97					Average Distribution Error: 1.86

Table A-5. Grade 7 Mathematics 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	21.79	2.88	0.00	0.00	24.67
Apprentice	6.41	28.52	3.98	0.00	38.91
Proficient	0.00	3.28	24.07	2.36	29.71
Distinguished	0.00	0.00	1.13	5.58	6.71
Total % Assigned	28.21	34.68	29.17	7.94	100
Total % Expected Correct Assignments: 79.96					Average Distribution Error: 2.39

Table A-6. Grade 8 Mathematics 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	18.08	3.86	0.00	0.00	21.94
Apprentice	4.24	30.88	4.87	0	39.99
Proficient	0.00	3.61	26.08	2.45	32.14
Distinguished	0.00	0.00	1.13	4.80	5.93
Total % Assigned	22.32	38.35	32.08	7.25	100
Total % Expected Correct Assignments: 79.83					Average Distribution Error: 0.85

Table A-7. Grade 3 Reading 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	23.38	4.74	0.15	0.00	28.27
Apprentice	3.19	14.61	5.79	0.08	23.66
Proficient	0.06	4.57	21.44	6.58	32.65
Distinguished	0.00	0.01	2.66	12.74	15.41
Total % Assigned	26.63	23.93	30.03	19.41	100
Total % Expected Correct Assignments: 72.18		Average Distribution Error: 2.13			

Table A-8. Grade 4 Reading 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	23.42	5.03	0.04	0.00	28.49
Apprentice	3.61	21.99	6.36	0.02	31.98
Proficient	0.01	4.34	23.39	4.43	32.17
Distinguished	0.00	0.00	1.57	5.78	7.36
Total % Assigned	27.05	31.36	31.36	10.23	100
Total % Expected Correct Assignments: 74.59		Average Distribution Error: 1.44			

Table A-9. Grade 5 Reading 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	19.26	5.37	0.18	0.00	24.81
Apprentice	3.40	16.55	7.48	0.01	27.43
Proficient	0.06	4.26	29.96	3.50	37.77
Distinguished	0.00	0.00	2.93	7.05	9.99
Total % Assigned	22.72	26.18	40.55	10.56	100
Total % Expected Correct Assignments: 72.82		Average Distribution Error: 1.67			

Table A-10. Grade 6 Reading 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	28.64	4.56	0.10	0.00	33.31
Apprentice	4.25	17.09	6.82	0.00	28.17
Proficient	0.04	3.63	26.09	2.41	32.16
Distinguished	0.00	0.00	2.16	4.20	6.36
Total % Assigned	32.93	25.29	35.17	6.61	100
Total % Expected Correct Assignments: 76.20		Average Distribution Error: 1.63			

Table A-11. Grade 7 Reading 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	25.72	5.60	0.05	0.00	31.37
Apprentice	3.94	20.70	4.72	0.01	29.37
Proficient	0.04	5.57	21.23	3.75	30.58
Distinguished	0.00	0.00	1.85	6.83	8.69
Total % Assigned	29.69	31.87	27.85	10.59	100
Total % Expected Correct Assignments: 74.47		Average Distribution Error: 2.20			

Table A-12. Grade 8 Reading 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	18.99	4.13	0.07	0.00	23.20
Apprentice	3.97	19.00	6.01	0.02	29.00
Proficient	0.06	6.22	28.38	4.16	38.83
Distinguished	0.00	0.01	2.68	6.30	8.98
Total % Assigned	23.03	29.36	37.14	10.47	100
Total % Expected Correct Assignments: 72.67		Average Distribution Error: 0.93			

Table A-13. Grade 4 Science 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	8.21	1.61	0.00	0.00	9.82
Apprentice	2.39	21.47	5.49	0.00	29.36
Proficient	0.00	3.78	38.81	4.14	46.72
Distinguished	0.00	0.00	3.68	10.42	14.10
Total % Assigned	10.6	26.86	47.98	14.56	100
Total % Expected Correct Assignments: 78.91		Average Distribution Error: 1.25			

Table A-14. Grade 7 Science 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	11.68	2.83	0.00	0.00	14.52
Apprentice	2.70	20.09	4.78	0.00	27.57
Proficient	0.00	4.35	37.83	4.01	46.20
Distinguished	0.00	0.00	2.27	9.45	11.72
Total % Assigned	14.38	27.27	44.89	13.46	100
Total % Expected Correct Assignments: 79.50		Average Distribution Error: 0.87			

Table A-15. Grade 5 Social Studies
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	12.47	2.47	0.00	0.00	14.93
Apprentice	3.23	26.45	5.78	0.00	35.46
Proficient	0.00	3.51	37.74	2.43	43.68
Distinguished	0.00	0.00	1.72	4.21	5.93
Total % Assigned	15.69	32.43	45.24	6.64	100
Total % Expected Correct Assignments:		80.86	Average Distribution Error: 1.52		

Table A-16. Grade 8 Social Studies 2014
 Percentages of True Scores Being in Assigned Classification

True Class	Assigned Classification				Total % Expected
	Novice	Apprentice	Proficient	Distinguished	
Novice	6.62	1.82	0.00	0.00	8.44
Apprentice	2.54	25.37	3.98	0.00	31.89
Proficient	0.00	3.75	39.51	4.16	47.43
Distinguished	0.00	0.00	2.45	9.80	12.25
Total % Assigned	9.16	30.94	45.94	13.96	100
Total % Expected Correct Assignments:		81.29	Average Distribution Error: 1.22		

Appendix B. Technical Details

Although student scoring for K-PREP will rest on IRT scaling, the end result for each student will be a classification of their test performance into one of four proficiency levels: Novice, Apprentice, Proficient, or Distinguished. Furthermore, school accountability indexes are calculated from these NAPD classification rather than from student IRT-based scale scores. Therefore, analysis of the potential for student classification errors is important at both the student and school levels of analysis.

Classical test theory is based on the concept of a “true” score for each examinee, defined as the expected or average score across an infinite number of repeated tests (see, for example, Lord and Novick, 1968). In most cases, we have only a score from a single administration of the test in question. The difference between this single “observed” score and the underlying “true” score is error. In the present case, we are concerned not just with the size of these “errors,” but with the impact of these errors on classifying students into NAPD categories. Our concern is estimating the rate of classification errors, defined as the probability that students could be classified in an NADP categories based on a single testing other than the one in which their “true” score falls. Classification accuracy is usually used to connote the likelihood that classification from observed scores agrees with classification expected by true scores. The converse, which we will call misclassification, is the likelihood that a student’s true score is in a proficiency classification different from his/her observed score. The phrasing of this last sentence is important because it expresses errors in a meaningful way for individual students and their parents. Parents know their students’ test scores; information should be provided which helps them with the uncertainty about their students’ true proficiency.

Test theory is concerned with the uncertainty of estimating students’ true achievement levels based on their fallible, observed test scores. Observed scores are assumed to vary in lawful ways around theoretical true scores or around domain scores (theoretical scores derived from all possible forms and conditions for measuring performance). This variation of observed scores around true or domain scores is calculated as the standard error of measurement (σ_e). In traditional reliability and generalizability theory, σ_e is a simple function of reliability (r_{tt}) and total test variability (σ_T):

$$\sigma_e = \sigma_T \sqrt{1 - r_{tt}} \tag{1}$$

Error bands around estimated scores often accompany reports of students’ test scores. These error bands are based on σ_e with the assumption that errors of measurement are normally distributed. Although σ_e reflects the variation of observed scores around true scores, because σ_e is constant across the scale of measurement, σ_e is also used to estimate the likelihood of true scores being within predictable intervals around observed scores. For students with a given observed score, the distribution of their true scores can be estimated by σ_e . For example, error bands may be constructed to show the interval around observed scores, which with 95 percent confidence, should contain students’ true scores.

If score intervals are divided into proficiency levels, such as NAPD, then σ_e may also be used to estimate misclassification. That is, for any given estimated score, the probably that a student’s true score lies in the score interval of an adjacent proficiency classification can be estimated from σ_e and normal distribution cumulative probabilities. An obvious implication of this is that the closer an student’s observed score is to a cut point between two proficiency levels, the greater

the chances that the student is misclassified. Indeed, the likelihood of misclassification approaches 50 percent for scores very near a cut point, regardless of σ_e .

IRT is one of several more elaborate methods that produces estimates of standard errors of measurement that vary along the true ability scale (Feldt and Brennan, 1989). That is, standard errors of measurement are conditioned on student ability. While estimates of conditional standard errors of measurement, i.e., $\sigma(x|\theta)$, increase precision in understanding the relationship between estimated scores and true scores, they create a complication for estimating classification accuracy. Misclassification is based on the distribution of true scores around observed scores, i.e., $\sigma(\theta|x)$, whereas $\sigma(x|\theta)$ represents the opposite – the distribution of observed scores round true scores. When $\sigma(x|\theta)$ varies across ability levels, it cannot confidently be used as a substitute for $\sigma(\theta|x)$.

For each content area, HumRRO produced an estimate of the expected raw score for each true ability (θ). Next, HumRRO calculated the probability of correct item responses (or the probability of each possible score point for constructed response) for each item at each possible θ . For multiple choice items, the formula applied was the following:

$$p1 = e^{th(i)-b_hum(j)} / 1 + e^{th(i)-b_hum(j)}, \quad (2)$$

where $p1$ is the probability of a correct item response, $th(i)$ is the ability of person I and b_hum is the difficulty of item j . For short answer items, the formula applied was the following:

$$p0 = e^{th(i)-D(0)} / e^{th(i)-D(0)} + e^{th(i)-D(0) + th(i)-D(1)} + e^{th(i)-D(0) + th(i)-D(1) + th(i)-D(2)} \quad (3)$$

$$p1 = e^{th(i)-D(0) + th(i)-D(1)} / e^{th(i)-D(0)} + e^{th(i)-D(0) + th(i)-D(1)} + e^{th(i)-D(0) + th(i)-D(1) + th(i)-D(2)}$$

$$p2 = e^{th(i)-D(0) + th(i)-D(1) + th(i)-D(2)} / e^{th(i)-D(0)} + e^{th(i)-D(0) + th(i)-D(1)} + e^{th(i)-D(0) + th(i)-D(1) + th(i)-D(2)},$$

where $p0$, $p1$, and $p2$ are the probabilities of each possible score point, respectively, and $D(0)$, $D(1)$, and $D(2)$ represent the step difficulty of each possible score point, which are essentially points on the x-axis at which the probability curves for each score point intersect. These step difficulty values are expected to increase for each subsequent score point. Similar equations were used for extended response options, but '4' being the maximum score possible.

Based on the item probabilities generated, probabilities for total test raw scores were next generated for each level of θ . Next, probabilities were generated for each possible θ on the 100-300 reporting scale. These two sets of probabilities were combined to produce a joint probability, or $P(\text{Obs}|\theta)$. By applying Bayes' Theorem, a probability density for true scale scores for a given observed score, $f(\theta|\text{Obs})$, can be estimated from $P(\text{Obs}|\theta)$ and some assumptions about the overall distribution of θ .

Bayes' Theorem, as applied to continuous variables, states that:

$$f(\theta|\text{Obs}_j) = \frac{f(\theta)P(\text{Obs}_j|\theta)}{P(\text{Obs}_j)} \quad (4)$$

where $P(\text{Obs}_j) = \int P(\text{Obs}_j|\theta) f(\theta) d\theta$ and $f(\theta)$ is the density function for the distribution of true scale scores.

To simplify both the computation and the presentation of the proposed approach, true ability is converted to a discrete variable with levels corresponding to the possible scale score estimates. Thus, $P(\theta_i) = P\{(ObS_{i-1} + ObS_i)/2 \leq \theta < (ObS_i + ObS_{i+1})/2\}$ where $-\infty$ and ∞ are substituted for ObS_0 and ObS_{k+1} respectively. With this change, the probability of different true scores, θ_i for any given observed score, ObS_j , can be rewritten as

$$P(\theta_i|ObS_j) = \frac{P(ObS_j|\theta_i)P(\theta_i)}{P(ObS_j|\theta_1)P(\theta_1)+P(ObS_j|\theta_2)P(\theta_2)+P(ObS_j|\theta_3)P(\theta_3)+\dots+P(ObS_j|\theta_k)P(\theta_k)} \quad (5)$$

Where ObS_j = observed scale score at level j , with k possible values in the scoring table, and θ_i = true ability at level i , with k levels of ability represented in the scoring table.

For each θ_i and its associated $\sigma(x|\theta_i)$, the probability of obtaining a given ObS_j , denoted $P(ObS_j|\theta_i)$, can be calculated directly from the IRT model and item parameter estimates or estimated using $\sigma(x|\theta_i)$ and assuming a normal distribution of errors.

For this report, we used the assumption of normally distribute errors. $P(ObS_j|\theta_i)$ was estimated for each of the k times k combinations of possible observed scores and true scores by assuming that each discrete scale score includes a hypothetical range of scale values from half of the distance to the next lower possible value to half the distance to the next higher value, i.e.,

$$\text{Score Range for } ObS_j = \frac{ObS_j + ObS_{(j-1)}}{2} \text{ to } \frac{ObS_j + ObS_{(j+1)}}{2} \quad (6)$$

For each θ_i and its associated $\sigma(x|\theta_i)$, the cumulative probability of scores within the score range for ObS_j can be calculated to estimate $P(ObS_j|\theta_i)$, assuming a normal distribution of errors. Raw-score-to-scale-score tables also limit the extremes of the distribution. For the lowest score level, $P(ObS_1|\theta_i)$ is calculated as the cumulative probability of $\frac{ObS_1 + ObS_2}{2}$, given θ_i and $\sigma(x|\theta_i)$. For the highest score level, $P(ObS_k|\theta_i)$ is calculated as the 1 - the cumulative probability of $\frac{ObS_{k-1} + ObS_k}{2}$, given θ_i and $\sigma(x|\theta_i)$.

The probability of misclassification is the probability that students with a given observed score, ObS_j , could have a true score in a proficiency (NAPD) level that is different than the level that contains that ObS_j . This probability can be calculated as

$$\begin{aligned} &P(\text{Proficiency Level for } \theta \neq \text{Proficiency Level from } ObS_j) \\ &= P(\text{Proficiency Level for } \theta < \text{Proficiency Level from } ObS_j) \\ &+ P(\text{Proficiency Level for } \theta > \text{Proficiency Level for } ObS_j) \\ &= \sum_{i=1}^m P(\theta_i|ObS_j) + \sum_{i=n}^k P(\theta_i|ObS_j), \end{aligned} \quad (7)$$

where m represents the highest level of θ in the next lower proficiency level from the observed score, n represents the lowest level of θ in the next higher proficiency level above the observed score, and k represents the highest θ represented in the scoring table.

Equation (7) provides the classification error information that will be the most meaningful for students interpreting their scores. If desirable, Equation (7) could be split into separate

estimates of the probability that a student's true classification is lower or is higher than indicated by his/her observed scale score, using only the first term or the second term of Equation (7), respectively. The later probability may be of more interest to students and parents.

The decision to place scores on transcripts, however, will not be made individually, but will be made on a system-wide level affecting all students. A system level estimate of the proportion of all students expected to be misclassified can be calculated by weighting the results of Equation (7) for each score level with the proportion of the sample who receive that score, and then summing overall score levels. That is,

Proportion of all student expected to be misclassified =

$$\sum_{j=1}^k \left\{ \left[\sum_{i=1}^m P(\theta_i | \text{Obs}_j) + \sum_{i=n}^k P(\theta_i | \text{Obs}_j) \right] * \frac{\text{Freq}_j}{\text{Total of All Students}} \right\}, \quad (8)$$

where m and n , defined as above, will vary with the proficiency level of Obs_j .