

CHAPTER 1: INTRODUCTION	3
Background	3
Statement of the Problem	5
Purpose and Research Questions	6
Significance of the Study	7
CHAPTER 2: LITERATURE REVIEW	10
Definition and Methods of Manipulation	10
Prevalence of Test Score Manipulations	19
Prevalence Estimates Based on Teacher Surveys	23
Prevalence Estimates Based on State Surveys	28
Prevalence Estimates Based on Direct Observation	30
Prevalence Estimates Based on Statistical Detection	32
Prevalence Estimates Based on Targeted Research	35
Why Do Educators Manipulate Test Scores?	40
Educators Are Former Students	40
Pressure From State Accountability Systems	42
Educators Unaware of Manipulations & Their Impact	45
Lack of Oversight and Policies	52
How to Prevent Manipulations: Evaluation of State Policies	55
Test Security Policy Content: Evaluative Framework	58
Relationship Between Test Security Policies and Score Trend Discrepancies	67
Single-Year Comparisons of State and NAEP Results	70
Trend Comparisons of State and NAEP Results	76
Scale-Invariant Trend Comparison Methods	81
State and NAEP Trend Discrepancies: Plausible Rival Hypotheses	82
Summary	85
METHODOLOGY	86
Independent Variable: Test Security Policy Quality	87
Data Collection and Verification	88
Sampling	88
Analysis	92
Technical Quality of Policy Evaluation Data	94
Assumptions and Limits	95
Dependent Variable: Scale-Invariant State Test and NAEP Score Trend Discrepancies	96
Data Collection and Verification	97
Sampling	101
NAEP Data Collection	102
Analysis	102
Reporting	115
Assumptions and Limitations	116
Analysis	118
Scope	120
Assumptions, Limitations, and Confounding Variables	120
Summary	121
APPENDIX A: PUBLISHED NEWS SUMMARIES	124
APPENDIX B: STATISTICAL DETECTION INDICES	137
Early Developments	137

From Empirical to Chance Models.....	138
Incorporating More Information.....	140
Controlling for False Positives	143
Incorporating Item Response Theory.....	144
Person-Fit and Aberrant Response Indices	146
Aberrant Response Indices to Detect Examinee Cheating.....	148
Adjacent Seating Methods	151
Methods to Detect Educator Cheating.....	152
Index #1: Unusual Test Score Fluctuations.....	154
Index #2: Unexpected Patterns in Student Answers.....	156
Combining Indices to Detect Cheating Classrooms	160
APPENDIX C: NATIONAL TESTING CODES AND STANDARDS	161
REFERENCES.....	164

CHAPTER 1: INTRODUCTION

Preventing cheating by those who give tests is a particularly underresearched topic. It is ironic that much attention has been given to preventing cheating by individual students – behavior that can cause a single score to be of questionable value – and so little attention has been paid to cheating by those who give tests, which can invalidate the scores of entire groups of students.

Gregory J. Cizek, *Cheating on Tests: How to Do It, Detect It, and Prevent It*

Background

In an effort to increase the academic achievement of all students and confront the “soft bigotry of low expectations” (Bush, 2000), the *No Child Left Behind* (NCLB) Act was passed into law on January 8, 2002. Like the *Improving America’s Schools Act* (IASA) signed in 1994, NCLB required all states to develop content and performance standards; implement assessment systems to track student performance against those standards; and create adequate yearly progress (AYP) goals to ensure all students reach a proficient level of achievement (IASA, 1994; NCLB, 2001). Believing the IASA was ineffective in improving student achievement due to its status as “an undertaking without consequences,” (Rotherham, 1999) NCLB granted the federal government the authority to impose sanctions upon schools, school districts, and states failing to meet AYP goals. The sanctions were intended to provide an incentive for educators to improve the quality of education provided to students and, ultimately, to improve student achievement.

Incentive theory predicts that by tying the threat of sanctions to assessment results, NCLB would motivate educators to increase test scores by implementing more effective instructional programs and, as a result, student achievement will increase (Laffont & Martimort, 2001; Jacob, 2007). Researchers have found evidence of this effect, finding that students in states with accountability systems achieve significantly higher gains on the NAEP (National Assessment of Educational Progress) math test than students in states having no sanctions tied to assessment results (Carnoy & Loeb, 2002; Hanushek & Raymond, 2004ab, 2006). These researchers conclude that “the introduction

of consequential accountability systems has a clearly beneficial impact on overall performance” (Hanushek & Raymond, 2004a, p. 32) for students of all racial groups even after controlling for test participation rates and state characteristics (Carnoy & Loeb, 2002, p. 318). Other researchers have also found this beneficial impact of NCLB on student achievement, concluding that the strength of a state’s accountability system “is indeed an important predictor of student performance at all points on the distribution curve, and especially so for students at the basic level” (Loeb & Strunk, 2005, p. 23) and that “students perform better than expected when their test score is particularly important for their schools’ accountability rating” (Reback, 2007, p. 1).

Although these studies found positive effects of accountability systems on student achievement, other evidence suggests educators are “gaming the system” to increase test scores without a corresponding increase in student achievement. Some educators game the system by manipulating the teaching process or their teaching philosophies. They do this by narrowing the curriculum to primarily teach content found on the test, using actual test items as practice for the test, focusing instructional resources only on the students most likely to improve the school’s test scores, spending an inordinate amount of time on test preparation, or by bribing students for higher test scores (Neal & Schazenbach, 2007; Jacob, 2005; Nichols & Berliner, 2004). Other educators have been found to manipulate the test administration by giving students hints on test questions, changing student answers, reading questions aloud to students, or by providing students extra time to complete the test (Cullen & Reback, 2006; Figlio & Getzler, 2002; Jacob & Levitt, 2003; Nichols & Berliner, 2005). Others attempt to game the system by excluding students from testing, inappropriately classifying examinees as disabled, or by using other methods to manipulate the examinee pool (Cullen & Reback, 2006; Figlio, 2005; Figlio & Getzler, 2002). Still others game the system by manipulating student test scores or by lowering proficiency standards (King, 2007; Nichols & Berliner, 2005).

Statement of the Problem

Reports of these manipulations cast doubt on educators and the inferences made from test scores. Tests questions are intended to represent a sample of a larger domain of interest and test scores are intended to represent examinees' performance in this domain (Haladyna & Downing, 2004). Manipulations that increase test scores without correspondingly increasing examinee performance in the larger domain destroy the validity of inferences made from the test scores. Therefore, these attempts to game the system negatively impact any decisions made on the basis of test scores, including the evaluation of instructional programs and the allocation of educational resources.

To protect the validity of inferences made from test scores, several methods have been used to deter educators from gaming the system. Professional organizations have developed ethical codes, guidelines, and standards to inform educators of the negative impact of test manipulations (AERA, APA, & NCME, 1999; JCTP, 2004; NEA, 1990; Schmeiser et al., 1995), but research suggests that educators are still unaware of what behaviors are appropriate or inappropriate (Kher-Durlabhji & Lacina-Gifford, 1992; Lai & Waltman, 2007; Moore, 1994). Several state departments of education and test publishers employ statistical methods in an attempt to detect educators who game the system, but these methods can only detect the most blatant manipulations, and research has shown their limited statistical power (Chason & Maller, 1996; Iwamoto, Nungester, & Luecht, 1996; Impara, Kingsbury, Maynes, & Fitzgerald, 2005). Some states have tried to deter educators by outlining harsh sanctions for anyone caught cheating, but research has shown that test manipulations are not frequently reported (Gay, 1990) and that many states do not follow through with the sanctions (Mehrens, Phillips, & Schram, 1993; Sorensen, 2006).

One promising way in which state officials have attempted to deter educators from gaming the system is through the development, implementation, and dissemination of comprehensive test security policies to both discourage test manipulations and

encourage ethical behavior. While little research has been conducted to determine the effectiveness of these policies on deterring educators from gaming the system (Cizek, 1999), similar policies have been shown to be effective in reducing student cheating on tests at the postsecondary level (McCabe & Trevino, 1993, 2002). If found to be effective in deterring educators from manipulating test scores, states could develop and implement test security policies to ensure inferences made from test scores are valid.

Purpose and Research Questions

This study will first document and classify methods educators use to manipulate test scores. News reports, surveys, observational studies, and statistical detection studies will be used to develop a taxonomy of manipulations and estimate the prevalence of various manipulation methods. This study will then explore possible reasons why educators manipulate test scores.

With information about how and why educators manipulate test scores, this study will then document and classify policies and procedures used by states in an attempt to deter educators from manipulating test scores. A framework will then be developed to evaluate the quality of these state test security policies. This framework will be based on the work of Cizek (1999), McCabe and Trevino (1993, 2002) in designing honor codes and test security policies to deter student cheating on tests.

Finally, this study will attempt to determine if a relationship exists between the strength of a state's test security policy and the discrepancy in score trends between the state tests and an audit test. The logic is that if manipulations increase test scores without a corresponding increase in student achievement, then discrepancies between score trends on the state test and another (audit) test of the same domain could possibly provide evidence of test score manipulations. A scale-invariant framework will be used to estimate trends on high-stakes reading and mathematics tests used in state accountability

systems and trends on the relatively low-stakes NAEP reading and mathematics tests for grades 4 and 8.

While the discrepancies between state test and NAEP score trends could possibly provide evidence of test score manipulations, it is important to note that any discrepancies could be explained by any combination of other plausible rival hypotheses. This study will discuss other possible explanations for the discrepancies between state and NAEP score trends, including possible differences in test content and administration; examinee pool and examinee motivation; and the strength of a state's testing program.

In summary, this study will attempt to address the following research questions:

1. What kinds of manipulations do educators use to increase test scores? Why do educators manipulate test scores? What is the estimated prevalence of each type of manipulation?
2. What test security policies and practices do states implement in an attempt to deter educators from manipulating test scores? What is the quality of each state's test security policy?
3. What is the relationship between the quality of a state's test security policy and any discrepancies between score trends on state and NAEP tests? Which aspects of a state's test security policy seem to have the strongest relationship with score trend discrepancies? What are some potential explanations for the discrepancies between state test and NAEP score trends?

Significance of the Study

This study will synthesize research on inappropriate testing practices to provide a taxonomy of methods used by educators to game the system. Results from surveys (Gay, 1990; Hall & Kleine, 1992; Lai & Waltman, 2007; Mehrens, Phillips, & Schram, 1993; Nolen, Haladyna, & Haas, 1992; Pedulla, et al., 2003; Shepard & Dougherty, 1991; Sorenson, 2006), direct observations (Horne & Gary, 1981; White, Taylor, Carcelli, &

Eldred, 1981; Wodtke, Harper, Schommer, & Brunelli, 1989), test preparation research (Moore, 1994; Popham, 1991), and statistical analyses (Perlman, 1985; Jacob & Levitt, 2004; Wesolowsky, 2000) will be combined with published news reports (Nichols & Berliner, 2004; Thiessen, 2007) to estimate the prevalence of each method. This extends the work of Haladyna, Nolen, and Haas (1991) to determine sources of test score pollution and the work of Jacob and Levitt (2004), Cizek (1999), and Wesolowsky (2000) in determining the prevalence of educator cheating on achievement tests. This study will also synthesize the above research in an attempt to explain why educators manipulate test scores.

This study will extend the work of McCabe and Trevino (1993, 2002) and Cizek (1999) in examining test security practices and honor codes in educational organizations. This study will synthesize professional codes and state policies designed to deter educators from manipulating test scores and develop a framework to evaluate test security policies.

Instead of analyzing examinee- and classroom-level data to estimate the prevalence of test score manipulations (Jacob & Levitt, 2004; Wesolowsky, 2000), this study will attempt to determine if a relationship exists at the state-level between the strength of a state's test security policy and the discrepancy in score trends on two tests of the same domain. This will extend the research of Klein, Hamilton, McCaffrey, and Stecher (2000); Linn, Baker, and Betebenner (2002); Peterson & Hess (2005, 2006); Koretz (2005), and Wei, Shen, Lukoff, Ho, & Haertel (2006) into the discrepancy between state and NAEP test scores.

Rather than using scale-dependent methods of comparing proficiency rates (Education Week, 2006; Lee, 2006; Thomas B. Fordham Foundation, 2005), or mapping state cut scores onto the NAEP score scale (Braun & Qian, 2007; Jacob, 2007; Linn, 2005; McLaughlin et al., 2002), this study will promote the use of scale-invariant methods to compare discrepancies in score trends between state and NAEP tests. This

will extend the work of Ho (2005, 2007) and help support or refute conclusions regarding the discrepancies between state and NAEP trends. By discussing possible explanations for these discrepancies, this study will extend the work of Hill (1998), Koretz (1999, 2005), and Koretz, McCaffrey, and Hamilton (2001).

This study will contribute to an understanding of test score manipulations and test security policies. Results of this research could be used by states in developing, improving, or auditing their current test security policies. It could also provide information that could be used in professional development to train teachers in appropriate test preparation and administration activities. This study could also contribute to the debate over the effectiveness of accountability systems and sanctions in improving student achievement.

CHAPTER 2: LITERATURE REVIEW

The review of literature is divided into five sections. The first section defines test score manipulations and provides examples of the four main ways in which educators manipulate test scores. The second section synthesizes estimates of the prevalence of each manipulation method to demonstrate the problem of test score manipulation. The third section then summarizes research into why educators would manipulate test scores to provide a better understanding of the problem. The fourth section then examines methods used to prevent test score manipulations and provides a framework for evaluating the quality of state test security policies and practices. The fifth section then summarizes research into discrepancies between state test and NAEP results, including a discussion of possible explanations for discrepancies.

Definition and Methods of Manipulation

Inferences made about a student's performance on an underlying construct, such as reading comprehension or mathematics problem solving, are made on the basis of observed test scores. It is normally assumed that an increase in test scores reflects an increase in student performance on the underlying construct. This is not always the case, however, as educators can implement practices to artificially increase test scores. The term *manipulation* will be used to describe any practice used by educators to increase student test scores without an equal, corresponding increase in student performance on the underlying construct.

The definition of *manipulation* is influenced by related concepts in the literature. Messick (1984) used the more general term *construct-irrelevant variance* to refer to the influence on test scores of any factor unrelated to the underlying construct (p. 216). Haladyna, Haas, & Nolen (1990) defined a similar concept of *test score pollution* to refer to situations in which test scores are distorted by factors unrelated to the construct being

tested (p. 9). The term *manipulation* is more specific than these other terms in that it refers only to test score distortions caused by educators' practices.

The term manipulation is also defined to be as general and value-neutral as possible. Research into educator test preparation practices use the terms *inappropriate* (Moore, 1994; Popham, 1991) or *unethical* (Lai & Waltman, 2007) to refer to practices that may distort test scores. Likewise, research into institutional cheating defines *cheating* as "a deception used to misrepresent student achievement" (Haladyna & Downing, 2004, p. 25). These terms imply malice on the part of educators. A manipulation is defined as any practice that distorts test scores, whether the practice was implemented maliciously or with the best of intentions. Also, while cheating and inappropriate or unethical test preparation practices are specific types of manipulations, the term manipulation refers to a broader collection of practices used by educators to inflate test scores.

In an attempt to develop a comprehensive list of manipulations used by educators, published news reports of alleged educator misconduct were collected using the LexisNexis® database and subscriptions to the Google® news alerts system. These reports, summarized in Appendix A, were combined with reports summarized by Nichols and Berliner (2004) in their critique of accountability systems in public education, *The Inevitable Corruption of Indicators and Educators Through High-Stakes Testing*, to yield a total of 186 published news reports of incidents from 1994 through 2007 in which educators in American public schools manipulated test scores. These reports of incidents were then combined with results from related research (discussed in the next section) to develop a list of 36 manipulations used by educators.

In order to better understand the methods educators use to manipulate test scores, a taxonomy was developed to categorize the 35 manipulations. Table 2.1 displays this taxonomy along with the number of published incidents for each manipulation. Under this taxonomy, manipulations are classified into one of four categories: manipulations of

the teaching process or philosophy, manipulations of the examinee pool, manipulations of test administration, or manipulations of score reports or scoring standards.

The first, and broadest, category of methods used by educators to manipulate test scores is through manipulations of the process or philosophy of teaching. These manipulations take place before the test is administered to students. These methods include questionable test preparation practices such as practicing with items identical or similar to those on the actual test, practicing with items from previous years' tests, purchasing commercial test preparation packages, and teaching test-taking skills. Practicing with items identical or similar to those on the test fits the definition of manipulation, because students who practice with these items will earn test scores that may not accurately represent their achievement in the domain of interest ([Moore, 1994](#); [Popham, 1991](#)). The use of commercial test preparation packages and the teaching of test-taking skills also fit within the definition of manipulation, because they serve to increase test scores on a specific test or test format without necessarily increasing student achievement in the underlying domain ([Lai & Waltman, 2007](#)). Finally, educators who focus instructional resources on specific students who have the best chance at improving test scores at a classroom, school, district, or state level at the expense of other students are also manipulating the teaching process to increase test scores without a corresponding increase in overall student achievement ([Neal & Schanzenbach, 2007](#)).

Table 2.1 Taxonomy of Manipulations in Published News Reports (1994 - 2007)

	News Reports 1994-2000	New Reports 2001-2007
Manipulate Teaching Philosophy or Process (before test administration)	18	67
Examining the test or making copies prior to test administration (piracy)	8	35
Practicing with items identical or similar to the test	7	24
Practice with last year's (alternate form) test items	---	1
Practice with items of the same format as the test	---	---
Use commercial test preparation packages	---	---
Teaching test-taking skills; test-wiseness	---	---
Teaching content from specific test items	3	5
Focusing resources on students who are closest to proficiency	---	2
Primarily teaching content found on the test	---	---
Manipulate Examinee Pool (before or during test administration)	6	18
Excluding students from testing (encouraging drop outs; suspending students)	5	13
Bribing or paying students to increase test scores	1	1
Having high-scoring students take the test multiple times	---	2
Providing inappropriate special education placement	---	1
Increasing the caloric content of school meals to increase scores	---	1
Manipulate Test Administration (during test administration)	25	127
Altering a student's answer sheet (changing student answers)	8	33
Sanitizing answer sheets (cleaning answer sheets before scoring)	---	1
Not following test administration procedures exactly	---	---
Giving students answers	7	16
Checking student answers and/or pointing out incorrect answers	5	25
Giving students hints on test items (verbal or nonverbal)	2	19
Rephrasing test items for students	---	2
Allowing students to work together during testing	---	1
Ignoring students who are cheating	---	2
Giving students additional examples	1	7
Providing students extra time	1	7
Reading items that are supposed to be read by students	---	5
Answering questions about test content	---	1
Providing students with reference materials or tools during testing	1	7
Instructing students to fill-in specifics for unanswered items	---	1
Providing inappropriate accommodations to students	---	---
Review skills that will be on tomorrow's test	---	---
Manipulate Score Reports or Standards (after test administration)		
Removing or changing student test scores on official records	1	10
Moving or providing students with false IDs so scores won't count	1	1
Misrepresenting data	---	---
Changing the criteria for proficiency or making the test easier	---	---

The following three news report summaries show recent incidents in which educators have been caught allegedly manipulating the teaching process:

The Dallas Morning News, July 13, 2007 (Benton, 2007b).

A state investigation finds that David Tamez, an elementary school teacher in Amarillo, Texas, leaked the fourth-grade writing test prompt on the spring TAKS writing test to colleagues before the test administration. Tamez reportedly leaked the test information because he believed educators in other districts were doing it as well. The teacher obtained the test information by volunteering to serve on the committee that selects questions for the final form of the TAKS. He alleges that committee members “regularly smuggle out secret TAKS information to share in the home districts.” Another teacher interviewed by investigators signed a statement indicating that Tamez “bragged that the source of his insider test information was... a person he had sex with who works for a company that helps build the TAKS.” The Amarillo Independent School District concluded that the teacher obtained the information from an unidentified employee at Pearson Educational Measurement. Tamez resigned from his position, but will retain his teaching certificate if he cooperates with the investigation.

Dayton Daily News, February 4, 2007 (Elliott, 2007)

A newspaper investigation found that students at City Day Elementary School in Dayton, Ohio were given 44 practice questions that were identical or “substantially the same” as questions from the actual state exam. In some questions on the practice test, only names or small details were changed from the real test questions. The investigation was launched due to the suspiciously large amount of improvement shown by the school. In 2005, no sixth grade student in the school passed the math subtest of the Ohio Achievement Test. One year later, 100% of these students (now in 7th grade) passed the math test.

The Columbus Dispatch, October 22, 2006 (Richards, 2006a).

Of the 28 Ohio school districts analyzed by The Columbus Dispatch, 15 had instances of educators cheating on standardized tests. Barbara Oaks, a teacher in the Coventry district, looked through the test and wrote out a geometry problem she thought her students would have trouble with. Winifred Shima, a teacher from the Parma district, used a copy of the test to create a study guide for students that included 45 of the 46 actual test questions. Brian Wirick (East Knox) and Heather Buchanan (Wapakoneta) both used the test to create study guides for students. Judy Wray, a veteran teacher in Marietta, made copies of the actual state test to help students prepare. Wray is reported to have said that teachers cheat more than administrators know.

A second way in which educators manipulate test scores is by manipulating the examinee pool. Rather than increasing the achievement of all students, educators using this method attempt to exclude low-ability students from taking the test or convince high-ability students to take the test multiple times. To exclude low-ability students from testing, educators have resorted to suspending students during the test administration period (Figlio, 2005) or inappropriately classifying students as being disabled (Cullen & Reback, 2006; Figlio & Getzler, 2002). These actions would increase test scores at a classroom, school, district, or state level without actually increasing the achievement of all students, so they fit the definition of manipulations.

The following three news report summaries show recent incidents in which educators have been caught allegedly manipulating the examinee pool:

San Francisco Chronicle, July 16, 2007 (Asimov, 2007ab).

The California Department of Education concludes that for the second consecutive year, educators at University Preparatory Charter High School in San Francisco interfered with state-mandated testing. State investigators seized illegal copies of the 2005 form of the test that was used to prepare students for the exams. Eight former teachers at the school assert the existence of a culture of cheating at the school. According to those former teachers, student grades are frequently falsified and low-scoring students are excluded from state-mandated testing. Last year, the state found that hundreds of answers on the ninth-grade English and math tests had been changed from wrong to right. A counselor from Oakland's Skyline High school reports that a student earning D's and F's transferred to University Preparatory Charter High School and received A's and B's while taking 16 classes in a single semester. When the student returned to Skyline High, he once again earned D's and F's. Last year, investigators concluded that educators at the school changed hundreds of test answers before they were sent for scoring. Former testing coordinator Mike Schwartz is suing school founder and director Isaac Haqq for breach of contract, claiming Haqq was responsible for the altered answer sheets.

Brevard School District, June 30, 2006 (Brevard SD, 2006)

Lori Backus, principal of Cocoa High School in Brevard, FL is accused of moving at least 54 9th and 10th grade special needs students into 11th grade so that their FCAT scores would not count towards the school's grade (assigned by the state) in 2005 and 2006. As a result of an investigation into the allegations, Principal Backus was immediately removed as principal.

Philadelphia Inquirer, June 25, 2006 (Patrick & Eichel, 2006).

Edison Schools fires Jayne Gibbs, principal at Parry Middle School in Chester, Pennsylvania for allegedly changing student test answers in 2005. Eighth graders at the school said the principal had given them the answers to questions on the Pennsylvania System of School Assessment. Gibbs is also accused of exempting special-education students from testing, violating state and federal rules. Edison Schools also asks the state and district to investigate exemplary test results at Showalter Middle School, where Gibbs served as principal from 2003-04.

A third way in which educators manipulate test scores is by manipulating the test administration. These methods involve the most blatant forms of cheating educators can use to increase test scores, such as by changing student answers, giving students hints to test questions, or pointing out incorrect answers on the test. These manipulations also include any changes educators make to the test administration instructions, such as providing students with extra time to complete the test, rephrasing test items for students, giving students inappropriate accommodations, allowing students to work together on the test, or allowing students to use forbidden reference materials such as calculators or dictionaries on the test. Each of these methods serves to increase test scores without a corresponding increase in test scores, so each method is a manipulation.

The following three news report summaries provide examples of recent incidents in which educators have been caught allegedly manipulating the test administration:

Herald Tribune, August 10, 2007 (Morris, 2007).

Mary Cropsey, a third-grade teacher at Mills Elementary School in Manatee, Florida, is accused of tampering with student answer sheets on the Florida Comprehensive Assessment Test (FCAT). One student reports that Cropsey helped students on the test; another student reported hearing that the teacher gave students extra time to complete the exam. An investigation began after yet another student reported that she had not finished the exam, but the next day all the bubbles had been filled-in. If the allegations are proven true, Cropsey could lose her teaching certificate and even be charged with a crime.

Newsday, June 24, 2007 (Hildebrand, 2007ab; Marcus, 2007ab)

The entire Uniondale school district is placed on academic probation due to evidence of tampering with Regents Math A and

B high school exams and the State Mathematics Assessments for grades 3-8 in 2005 and 2006. Faculty members reportedly allowed students to use calculators, which were not allowed on the exam (Marcus, 2007a). The New York Department of Education reports that complaints of test fraud have more than doubled over the past five years, with the department receiving 37 complaints in 2006. One dozen teachers and administrators accused of test fraud have faced hearings in front of the New York Professional Standards and Practices Board. Of those twelve cases, six cases resulted in revocation of professional certifications, two cases were cleared, and the remaining four cases remain under investigation. The number of complaints verified by the state has remained relatively steady, with between 9-16 in each of the past five years (Hildebrand, 2007a). An analysis of Uniondale's test scores found that 333 answers on the Regents Math A exam were altered, and 97% of the time they were changed to the correct answer. On the Regents Math B exam, 198 answers were changed, with 97% again being changed to the correct answer. On the 2005 8th grade math assessment, Uniondale students scored below average on 11 of the 14 easiest questions, but higher than average on 12 of the 13 most difficult items (Hildebrand, 2007b; Marcus, 2007b).

San Francisco Chronicle May 13, 2007 (Asimov & Wallack, 2007) Teachers in at least 123 public schools have reportedly cheated for students on California's high-stakes tests between 2004-2006. In two-thirds of these cases, the schools admit that they had cheated. The cheating behaviors included (a) allowing students to use reference materials such as maps and flow charts during the test, (b) allowing students to use calculators, (c) helping students answer questions, and (d) erasing and changing student answers. California currently identifies potential misconduct by scanning answer sheets for suspicious erasures. Cheating is virtually ignored in schools in which cheating impacts less than 5% of tests are given. Schools in which cheating impacts more than 5% of the tests are not ranked and receive a note stating "adult irregularity in testing procedure" occurred.

The final way in which educators manipulate test scores is by manipulating score reports or performance standards. The most blatant manipulations in this category involve educators changing or removing student scores from official score records. Less obvious manipulations include educators changing demographic data so that scores from higher-ability students are added to specific lower-scoring subgroups or changing student identification numbers so that score trends cannot be calculated for lower-scoring students. A more nebulous manipulation in this category occurs when educators misrepresent test scores in score reports to make the scores appear better than they

actually are. This includes changes made to lower the difficulty of the test or lower the cut-score for proficiency in order to make it appear as though student achievement has increased. Because these methods cause test scores to increase without an increase in underlying student achievement, each of these methods is a manipulation.

The following three news report summaries provide examples of recent incidents of educators manipulating score reports or performance standards:

San Francisco Gate, June 30, 2006 (Sturrock, 2006).

According to researchers with Policy Analysis for California Education (PACE), California and eleven other states have inflated test outcomes by lowering the achievement standard students need to meet to be proficient in reading and math under the federal No Child Left Behind Act. The study describes large differences in results from state and national tests and outlines several reasons for these large differences. One of the reasons is that states sometimes lower their standards for what they deem proficient.

MSNBC, April 17, 2006 (MSNBC, 2006)

With permission from the federal government, nearly two million students' test scores are not counted when schools report progress by subgroups under the No Child Left Behind requirements. This is due to states being able to define the minimum number of students needed in a subgroup before scores are reported. In the past two years, almost half of all states have successfully petitioned the U.S. Department of Education to increase these minimums. An investigation concludes that about 1 out of every 14 test scores are not being counted under appropriate racial categories. The scores from more than 24,000 students in Missouri, 257,000 in Texas, and 400,000 in California are not being counted.

San Antonio Express News, September 17, 1998 (Stinson, 1998)

The Austin School District manipulated test results last spring to make it appear as if several schools performed better than they did, the Texas Education Agency (TEA) says. Commissioner of Education Mike Moses explained the trickery in this August 14 letter to the school district, stating, "... student identification number changes were submitted for students tested at (the schools), which resulted in the exclusion of those students from the accountability subset of TAAS results used to determine the 1998 accountability ratings." Administrators knew that by changing student identification numbers, the TEA would eliminate those students' scores from ratings calculations (Nichols & Berliner, 2004, p. 27).

Prevalence of Test Score Manipulations

In addition to providing examples of manipulations used by educators, the published news reports also provide a crude estimate of the prevalence of each manipulation method. Figure 2.1 shows the number of published news reports about test manipulations each year from 1995 until 2007. Figure 2.2 shows the number of reports of manipulations by each state over that time period. These figures show that reports of manipulations are widespread and generally increasing over time, especially since the introduction of NCLB in 2002. Any conclusions made from these figures should be made cautiously, however. On one hand, published news reports can only be expected to represent the most interesting and, therefore, blatant incidents of test score manipulations. Also, news reports can only describe incidents in which educators were caught manipulating test scores. Therefore, these news reports may grossly underestimate the prevalence of manipulations. On the other hand, many of these news reports describe *allegations* of manipulations. Since follow-up reports are rarely written about these incidents, it is unknown how many of the reported allegations were eventually proven to be untrue. Therefore, these news reports may actually overestimate the prevalence of manipulations. Since it is unknown whether news reports over- or underestimate the actual prevalence of test score manipulations made by American public school educators, these reports can only provide evidence that educators do manipulate test scores and it appears that manipulations are becoming more prevalent.

Recognizing the limitations in using published news reports, researchers have developed other methods to estimate the prevalence of test score manipulations. These methods include administering surveys to teachers and state officials; directly observing test administration procedures in classrooms; statistical detection; and other targeted research methods. Because each method has advantages and disadvantages, no one method provides the single best estimate of the prevalence of test score manipulations in American schools. Therefore, in order to best estimate the prevalence of manipulations,

Table 2.2 synthesizes the estimates from research using each method. The estimated prevalence of each manipulation method is displayed using the taxonomy previously developed.

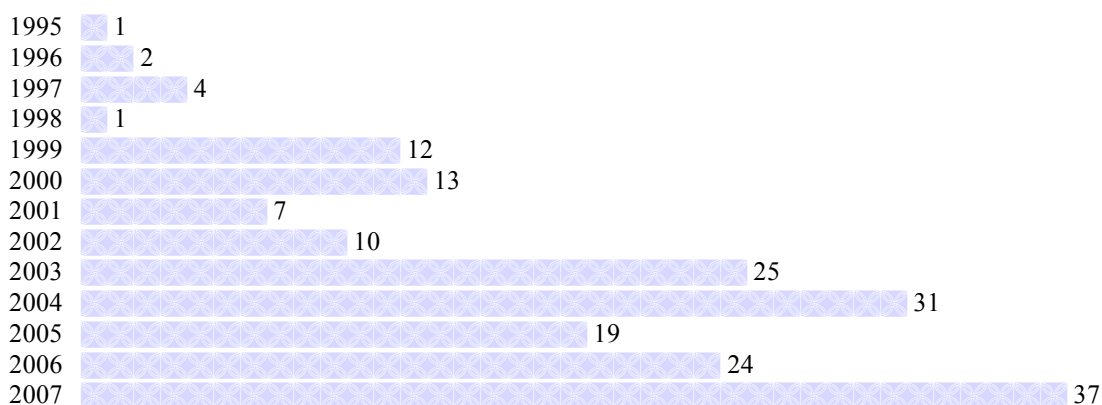


Figure 2.1 The number of published news reports on manipulations from 1995 – 2007.

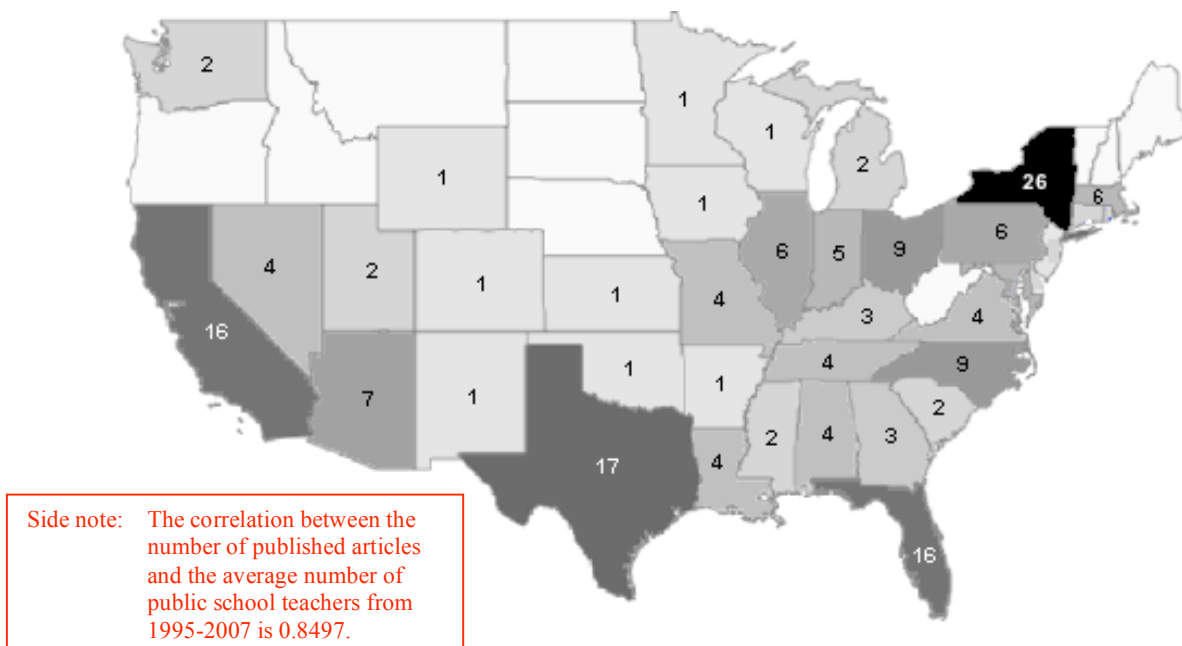


Figure 2.2 The number of published news reports from each state (1995 – 2007).

Table 2.2 Estimated Prevalence of Manipulations

	Teacher Surveys				State Survey	Stat Detect
	Shepard & Dougherty (1991)	KD-LG ¹ (1992)	Nolen et al. ² (1992)	Pedulla et al. (2003)		
Sample size	N=360	N=74	N=1881	N=4195	N=46	JL ⁶ (2003)
Sample location	2 districts	LA	AZ	47 states	States	Chicago
Manipulate Teaching Philosophy or Process						
Making copies of test prior to test administration	11%	8%	10%; 9%	32%	10%	---
Practicing with items identical or similar to the test	---	66%	12%; 9%	---	33%	---
Practice with last year's (alternate form) test items	---	---	---	---	---	---
Practice with items of the same format as the test	---	---	---	---	---	---
Use commercial test preparation packages	---	53%	41%; 12%	---	---	---
Teaching test-taking skills; test-wiseness	---	76%	60%; 36%	76%	---	---
Teaching content from specific test items	23%	41%	23%	6%	---	---
Focusing resources on students closest to proficiency	---	---	---	13%	---	---
Primarily teaching content found on the test	---	39%	66%; 22%	69%	18-22%	---
Manipulate Examinee Pool						
Excluding students from testing	8-13%	0%	---	---	---	---
Bribing or paying students to increase test scores	---	---	---	14% ³	---	---
Having high-scoring students take test multiple times	---	---	---	---	---	---
Providing inappropriate special education placement	---	---	---	---	---	---
Increasing caloric content of school meals	---	---	---	---	---	---
Manipulate Test Administration						
Altering a student's answer sheet (changing answers)	2%	6%	1.4%	2%	---	4-5%
Sanitizing answer sheets (cleaning before scoring)	---	88% ⁵	---	---	---	---
Giving students answers	---	8%	---	---	---	---
Checking or pointing out incorrect answers	10%	---	---	11%	---	---

Sample size	Shepard & Dougherty (1991)	KD-LG ¹ (1992)	Nolen et al. ² (1992)	Pedulla et al. (2003)	Lai & Waltman (2007)	Mehrens et al. ⁴ (1993)	JL ⁶ (2003)
Sample location	NC	LA	AZ	47 states	IA	States	Chicago
N	168	74	1881	4195	1338	46	
Sample location	NC	LA	AZ	47 states	IA	States	Chicago
Manipulate Test Administration (continued)							
Giving students (non)verbal hints on test items	14%	0%	8%; 6%	11%			
Not following test administration procedures exactly	4%						
Rephrasing test items for students	18%	10%					
Allowing students to work together during testing	2%						
Ignoring students who are cheating							
Giving students additional examples			28%; 16%				
Providing students extra time	14%	1%	8%; 3%	15%			
Reading items that are to be read by students	14%						
Answering questions about test content	12%						
Providing students with reference materials or tools							
Having students fill-in unanswered items							
Providing inappropriate accommodations to students	1%						
Review skills that will be on tomorrow's test			44%; 30%				
Manipulate Score Reports or Standards							
Changing student test scores on official records							
Providing false IDs so scores won't count							
Misrepresenting data							
Changing criteria for proficiency; making test easier							

- Notes:
- 1 Kher-Durlabhji & Lacina-Gifford (1992) sampled *preservice* teachers about what manipulations they planned on using
 - 2 Numbers represent separate responses from elementary; secondary school teachers
 - 3 Actual survey item was worded "give prizes to reward students"
 - 4 Numbers represent the percent of states reporting specific incidents during the 1989-1990 academic year
 - 5 Actual survey item was worded "check student's completed answer sheets"
 - 6 Jacob and Levitt (2003) attempted to detect the percentage of teachers who manipulate answer sheets for their students

Prevalence Estimates Based on Teacher Surveys

The most frequently used method to estimate the prevalence of test score manipulations involves administering a survey to teachers or school administrators. These surveys, which are usually byproducts of larger research into test preparation activities or the impact of high-stakes testing on instruction, typically ask teachers to indicate which manipulation methods they use to increase test scores.

Due to concerns that teachers will be reluctant in admitting to using several of the more blatant forms of manipulation, many surveys also ask teachers to report if they are aware of other teachers in their schools manipulating test scores. Table 2.2 displays the results of six of these teacher surveys. The numbers in the table represent the percentage of teachers responding to each survey who report that either they or other teachers in their schools use each method of manipulation.

The first teacher surveys that asked about the use of a wide range of manipulations were administered in the early 1990s. Gay (1990) administered a survey to 168 North Carolina teachers in grades 3 through 8. Gay found 35% of respondents reported participating in or being aware of testing irregularities in their schools (p. 4). These testing irregularities, defined by eight specific examples, all fit within the definition and taxonomy of test score manipulations. According to the results, 23% of teachers reported manipulating the teaching process by copying the test and teaching its contents to students prior to test administration. Fewer respondents reported manipulating the test administration, with 15% adding extra time to the test publisher's time limits, 14% coaching students on the test by giving verbal or nonverbal hints, 10% calling attention to incorrect student answers, 4% changing the publisher's test administration directions, and 2% leaving students unsupervised during testing. The most blatant form of manipulation, changing student answers, was reported by 2% of the respondents. Gay also noted that one respondent reported an incident in which a teacher

encouraged students to use reference materials during a writing test and another respondent admitted to checking student responses to “be certain that her students answered as they had been taught” (p. 4). In reporting the survey results, Gay suggested that the estimated prevalence of each test score manipulation “may only be the tip of the iceberg” (p. 3) and reported that 43% of respondents reported a belief that test manipulations were increasing among teachers.

The next year, Shepard and Dougherty (1991) administered a similar survey to 360 teachers from two American school districts as part of a larger study to determine the effect of the high-stakes tests on instruction and student learning. In this survey, teachers were asked to report the frequency with which “controversial testing practices” happened in their schools. These testing practices all fit within the definition of test score manipulations. Supporting the results from Gay’s survey, the researchers found that manipulating the teaching process was the most prevalent form of manipulation, with 41% of respondents occasionally or frequently giving highly similar items to students for practice and 11% practicing with items from the actual test. Fewer teachers reported manipulating the test administration, with 23% providing hints to students during testing, 20% giving students more time than the test directions call for, 18% rephrasing test questions for students, 14% reading test questions that were supposed to be read by students, 12% answering questions about test content during test administration, and 8% giving answers to students. Shepard and Dougherty also found that 6% of respondents reported the most blatant manipulation of changing incorrect answers to correct ones on student answer sheets. Shepard and Dougherty also found evidence of educators manipulating the examinee pool. 13% of respondents did not administer the test to students who would have trouble and 8% encouraged lower-ability students to be absent on the days of the test.

The following year, Hall and Kleine (1992) surveyed 220 Oklahoma public school teachers and found that 55% reported awareness of fellow teachers cheating on

tests for their students. While the term *cheating* was not defined on the survey, the results would certainly provide another estimate of the prevalence of activities designed to increase test scores without a corresponding increase in student achievement.

Rather than asking teachers to report manipulations they currently use, Kher-Durlabhji and Lacina-Gifford (1992) surveyed 74 pre-service teachers from Louisiana to determine the types of test preparation and administration activities they plan to use in teaching. Once again, the most prevalent forms of manipulation were methods used to manipulate the teaching process or philosophy. 76% of these pre-service teachers intended to spend instructional time on teaching test-taking skills, 66% intended to practice with previous years' test questions, 53% intended to use commercial test preparation packages, 39% intended to teach only content found on the test, and 8% intended to practice with items from the actual test. Also supporting previous survey results, fewer pre-service teachers intended to manipulate the test administration. 10% indicated that they plan to rephrase test items for their students and 1% planned to provide students extra time to complete the test. Not a single pre-service teacher intended to give students hints on test items. Also, no respondents intended to manipulate the examinee pool in order to increase test scores. It should be noted that while 88% of pre-service teachers intended to "check students' completed answer sheets," this does not necessarily imply they will take any action after checking the answer sheets. Thus, this does not describe a manipulation of test scores.

While the previous surveys were administered to small samples of educators, Nolen, Haladyna, and Haas (1992) administered a series of large-scale surveys to Arizona elementary and secondary teachers to determine their uses of test scores. As part of this research, 1,881 Arizona teachers were asked to indicate the test preparation practices they always or usually employ and which test administration activities were common or very common in their classrooms. In addition to finding that elementary school teachers were more likely to report manipulations than secondary teachers, the researchers once again

found that manipulations of teaching philosophy or process were most frequently reported. The prevalence of specific manipulations in this category ranged from 66% of elementary school teachers teaching only content found on the test to 9% of secondary teachers giving their students items from the actual test as practice. As further evidence of the prevalence of manipulations of the teaching process, 23% of elementary teachers responding to the survey indicated a belief that administrators required them to spend class time on test preparation activities and 33% reported spending more time on test preparation than was required. Also, 7% of elementary teachers, 10% of secondary teachers, and 3% of school administrators surveyed reported that teachers are encouraged to raise test scores by teaching items from the actual test. Manipulations of the test administration were less prevalent, with 44% of elementary school teachers admitting to reviewing the tested content and skills immediately before testing and 3% of secondary teachers providing extra time for students to complete the test.

Lai and Waltman (2007) administered a similar survey to a sample of 1,338 teachers from 125 public schools in Iowa. The estimated prevalence of manipulations of the teaching process or philosophy ranged from a median of 82% of teachers within a school teaching test-taking skills to their students to a median of 11% of teachers within a school giving actual items from the test to their students as practice. The researchers found that as many as 50% of teachers within a school admitted to using actual items from the test as practice and as many as 67% of teachers within a school allowed students to practice with the alternate form of the test. More troubling was the finding that “unexpectedly high percentages of teachers rated practicing with exactly the same test that will be administered this year as being ‘very ethical’” (p. 12).

While the surveys from Nolen, Haladyna, and Haas (1992) and Lai and Waltman (2007) were administered to large samples of teachers, the fact that their surveys were administered to teachers in only one state may limit the generalizations that can be made. To address this problem, the National Board on Educational Testing and Public Policy

funded a survey of 4,195 teachers from every state except Iowa, Oregon, and Idaho (Pedulla et al., 2003). The survey, designed to measure teacher attitudes towards state testing programs, asked teachers how their state's testing program impacted their classroom instruction and test preparation activities. Results from this survey once again showed that manipulations of the teaching process or philosophy were the most prevalent form of manipulation used by teachers. More than two-thirds of the respondents reported teaching test-taking strategies to their students and teaching only content found on the actual test, and one third reported practicing with items identical or similar to those found on the actual test. Less prevalent ways of manipulating the teaching process or philosophy included 13% of teacher focusing instructional resources on students who were closest to achieving a proficient score and 6% teaching content from specific test items. To manipulate the test administration, 15% give students extra time to complete the test, 11% give hints to students, and 11% point out incorrect answers to students. Only 2% admitted to the most blatant manipulation of changing answers on student answer sheets. In an effort to manipulate the examinee pool, 14% of teachers reported giving prizes to reward students for higher test scores. As a general estimate of the prevalence of test score manipulations, 38% of the respondents indicated that teachers in their schools have found ways to raise state-mandated test scores without really improving learning.

While the results from this method of administering surveys to teachers seem to converge to provide estimates of the prevalence of each manipulation method, the results should be interpreted carefully. As stated earlier, results from the surveys of Gay (1990), Shepard and Dougherty (1991), Hall and Kleine (1992), and Kher-Durlabhji and Lacina-Gifford (1992) are based on extremely limited sample sizes that may not generalize beyond the sample. While Nolen, Haladyna, and Haas (1992) and Lai and Waltman (2007) administered their surveys to much larger samples of teachers, the fact that each of these surveys was administered to teachers within a single state may also limit any

generalizations. This concern is further supported by the low 42% response rate from the survey administered by Nolen, Haladyna, and Haas (1992, p. 10). Results from the survey from the National Board on Educational Testing and Public Policy (Pedulla et al., 2003), while collected from a much more representative sample of teachers, should also be cautiously interpreted due to the reported 35% response rate. Further encouraging a cautious interpretation, the demographics of the sample did not match national demographics, with the sample containing more males, more experienced teachers, and more English teachers than the national average (Pedulla et al., 2003, p. 138).

Results should also be interpreted carefully due to the impact of social desirability bias on surveys. Social desirability bias is the inclination one has to respond to survey items in a manner that will be viewed favorably by others (Paulhus, 1991). In a study of student cheating on achievement tests, Scheers and Dayton (1987) found this effect, concluding that survey responses underestimated the actual extent of cheating behaviors. In a similar study of student cheating, Nelson and Schafer (1986) found the exact opposite effect. These researchers found that surveyed responses overestimated actual cheating incidents. While Eve and Bromley (1981) developed *culture-conflict theory* to explain these opposite results, the fact that response bias cannot be predicted limits the accuracy of manipulation prevalence estimates from these surveys.

Finally, with the exception of Lai and Waltman (2007), each survey treated individual teachers as independent units of analysis. Since administrative decisions might be made to influence teachers' decisions to manipulate test scores, teachers are not independent units (Lai & Waltman, 2007, p. 14) perhaps a more accurate estimate could be made using schools, school districts, or states as the units of analysis.

Prevalence Estimates Based on State Surveys

Rather than surveying individual educators to estimate the prevalence of manipulations, Mehrens, Phillips, and Schram (1993) surveyed 46 state departments of

education on behalf of the National Council on Measurement in Education. The researchers asked state officials to indicate the number of test security incident reports they received involving their state tests. Of the 41 states with testing programs at the time, 36 (88%) indicated receiving reports of test security breaches, with 28 states (68%) receiving such reports during the 1989-1990 academic year. Although some states completed the survey incorrectly (p. 5), the researchers found that the 41 states received an average of 6.34 reports of security breaches in that single academic year.

Supporting the results from the teacher surveys, Mehrens, Phillips, and Schram (1993) found that manipulations of the teaching process or philosophy were the most prevalent form of test score manipulation. Four states (10%) received at least one report regarding missing test materials in 1989-1990. The states reported that the test materials went missing due to unauthorized personnel having access to the materials or due to shipping irregularities. Ten states (24%) received reports of teachers using items very similar to the actual test items as practice and 14 states (34%) received reports of teachers providing their students with test questions in advance.

States received fewer reports of manipulations of the test administration during the 1989-1990 academic year. While 18 states (44%) received reports of teachers failing to follow the written test administration procedures, only 7 states (17%) received reports of educators erasing student answers. Only two states (5%) received reports of the most blatant form of manipulation – teachers marking answers on student answer sheets. As further evidence of the prevalence of test score manipulations, 10 states (24%) received reports of “dramatic increases” in average test scores for schools or school districts (p. 7).

In another survey of state departments of education, the test security firm Caveon (Sorensen, 2006) asked states to indicate the number of times in the past two years that they have taken formal action because cheating was either confirmed or suspected. While 6 of the 34 states surveyed (18%) reported no formal actions taking place, 13 states (38%) took formal action up to 3 times, 5 states (15%) took formal action 4 to 6 times, 2

states (6%) of states took formal action 7 to 9 times, and 8 states (24%) took formal action more than 10 times in two years.

As a more indirect estimate of manipulation prevalence, state officials were also asked to rate the importance of various test security threats to their state testing programs. Yet again manipulations of the test administration were perceived to represent a greater threat to state testing programs than manipulations to the teaching process or philosophy. States rated lost test materials, inappropriate administrative pressure put on teachers to increase test scores, and teacher coaching of students based on prior knowledge of test questions as the most important perceived threats. Manipulations of the test administration, including students working together on the test, teachers providing answers to students during testing, and educators changing student answers after testing were perceived to be less important.

While the results from these state surveys do support the results from teacher surveys in showing that manipulations do happen and manipulations of the teaching process or philosophy are more prevalent than manipulations of the test administration, the reported prevalence estimates must once again be interpreted cautiously. These surveys only indicate the number of manipulation incidents reported to state officials. In his survey of 60 teachers, Gay (1990) found that only 20% were willing to report testing irregularities to their school administrators (p. 4). If we assume that school administrators are similarly reluctant to report testing irregularities to district and state officials, then the results from these state surveys could greatly underestimate the actual prevalence of test score manipulations.

Prevalence Estimates Based on Direct Observation

To overcome the potential impact of response bias in teacher surveys and the limitations of state surveys, some researchers have directly observed schools and classrooms to determine the prevalence of test score manipulations. Due to its resource

intensiveness, observational studies of classroom test administration have been rare. In 1981, Horne and Gary observed the administration of a standardized test in 16 elementary school classrooms. The researchers found more than half of the teachers varied from the written test directions with 6 (38%) of the teachers “purposely and consciously [manipulating] the test administration procedures” (Horne & Gary, 1981, p. 12).

White, Taylor, Carcelli, and Eldred (1981) similarly observed the administration of a standardized test in 38 Utah classrooms. Supporting the previous study, the researcher found nearly half of the teachers did not follow the test administration directions exactly. The study found 46% of teachers failed to follow the exact wording in presenting test questions as stated in the test manual, 41% changed the wording of the test directions to a “vocabulary more familiar to students,” and only 50% “refrained from repeating a test question unless the directions specified to do so” (White et al., 1981).

Wodtke, Harper, Schommer, and Brunelli (1989) continued this method of research by observing the test administration practices of ten kindergarten classrooms. The researcher found, once again, that teachers manipulated the test administration. The researchers observed teachers failing to follow the time limit directions found in the test manual in 27% of the testing sessions. The researchers also recorded 21 unauthorized item repetitions, 40 incidents of teachers “cueing correct answers,” and 149 “significant procedural variations,” including rephrasing test questions and failing to disseminate practice test booklets (p. 228). The researchers concluded that administrations of standardized tests are manipulated so much as to render them incomparable (Wodtke et al., 1989).

While the method of direct observation has advantages over teacher and state surveys, the method is not without faults. First, due to its resource-intensiveness, only a small (and usually nonrepresentative) sample of classrooms can be observed. Second, direct observation methods may be subject to both the *Hawthorne Effect* and the *observer-expectancy effect*. The Hawthorne Effect suggests that subjects in direct

observation studies temporarily change their behavior due to their knowledge that they are part of a study (BJA, 2007). Thus, in studies of test administration manipulations, the Hawthorne Effect suggests that estimates based on direct observation may underestimate the actual prevalence. The observer-expectancy effect, on the other hand, is “a cognitive bias that occurs when a researcher expects a given result and therefore unconsciously manipulates an experiment or misinterprets data in order to find it” (Ferguson, 2007). If the researchers expected to observe test administration manipulations and if the observer-expectancy effect is real, then results from these direct observation studies might overestimate the actual prevalence.

Prevalence Estimates Based on Statistical Detection

In an effort to eliminate the response bias in surveys and potential biases in direct observation, statistical detection methods have also been used to estimate the prevalence of manipulations. These methods, detailed in Appendix B and Cizek (1999), use statistical analyses of student answer sheets in order to detect potential manipulations. First developed with the intention of detecting students who cheat on tests, these methods attempt to find examinee answer sheets with unusually large score gains or wild score fluctuations (Perlman, 1985); statistically improbable numbers of erasures (Qualls, 2001); or unusual patterns of answers (Advanced Psychometrics, 1993; Angoff, 1974; Anikeeff, 1954; Bay, 1995; Bellezza & Bellezza, 1989; Bird, 1927, 1929; Cizek, 1999; Drasgow, Levine, & Williams, 1985; Frary, 1977; Hanson & Brennan, 1987; Jacob & Levitt, 2004; Karabatsos, 2003; Kvam, 1996; Levine & Drasgow, 1979, 1988; Meijer & Sijtsma, 2001; Roberts, 1987; Saupe, 1960; Sotaridona & Meijer, 2001, 2002; van der Linden & Sotaridona, 2002; Wesolowsky, 2000; Wollack, 1997).

In 1985, Perlman analyzed student answer sheets from a standardized test administration in Chicago Public Schools. Discovering some schools had unusually large score gains and unusually high numbers of answer sheet erasures, Perlman hypothesized

that educators at these “suspect” schools may have manipulated the test administration process. In an effort to test his hypothesis, Perlman retested 23 of these suspect schools and 17 control schools that were not suspected of manipulating scores. The results of the retesting led him to state “clearly the suspect schools did much worse on the retest than the [control] schools,” and conclude that, “it’s possible that we may have underestimated the extent of cheating [manipulations of the test administration] at some schools” (Perlman, 1985, pp. 4-5).

Almost two decades later, Jacob and Levitt (2003, 2004) conducted another analysis of answer sheets from Chicago Public Schools. The researchers developed a statistical index to detect educators who manipulate student responses on tests by supplying them with answers or erasing and changing student responses. Applying their index to analyze 8 years of answer sheets from the Chicago Public Schools’ administration of the *Iowa Tests of Basic Skills* (ITBS), the researchers concluded, “Empirically, we detect cheating [manipulation of answer sheets] in approximately 4 to 5 percent of the classes in our sample” (2004, p. 846). They also found that between 1.1% and 2.1% of educators in their sample manipulated answer sheets on any particular ITBS subtest and 3.4% to 5.6% of educators manipulated answer sheets on at least one ITBS subtest. Further describing the prevalence, Jacob and Levitt (2003) conclude that educators found to manipulate answer sheets on one test were 10 times more likely to manipulate answer sheets on other tests and educators found to manipulate answer sheets one time were 9 times more likely to do so again in the future (p. 73).

In order to validate their estimate of 4% to 5% of educators manipulating the test administration by changing student answer sheets, the researchers retested 117 Chicago classrooms in 2002. The educators in these classrooms included “cheaters” whose classrooms experienced large score gains and showed evidence of unusual response patterns; “bad teachers who cheat” whose classrooms had unusual response patterns but did not experience large score gains; “anonymous tips” whose classrooms were not

identified by the statistical index but who were accused of cheating; “effective teachers” whose classrooms experienced large score gains with no evidence of unusual response patterns or manipulations; and “randomly selected” educators whose classrooms were not suspected of manipulations. Scores from the students in the “effective teachers” classrooms increased on the retest, while the “randomly selected” classrooms experienced a small decline of 2.3 standard score units. The “cheaters,” “bad teachers who cheat,” and “anonymous tips” classrooms experienced a large decline in score of 16.2, 8.8, and 6.8 standard score units, respectively. One of the classrooms taught by a “cheater” experienced a loss of 54 standard score units on the retest – a loss roughly equivalent to three full grade equivalent units on the ITBS. Based on the results of this retesting, Jacob and Levitt (2004) expressed confidence in their estimate of the prevalence of manipulations of student answer sheets.

Using a different statistical index, Wesolowsky (2000) analyzed answer sheets from the 2005-2006 administration of the *Texas Assessment of Knowledge and Skills* (TAKS). While not providing an overall estimate of the prevalence of manipulations of the test administration, Wesolowsky found that scores from more than 50,000 students showed evidence of irregularities that could include students copying answers from other students or educators doctoring student answer sheets. Additionally, the analysis found 112 schools in which at least 10% of student answer sheets were identified as potentially being manipulated. Expressing confidence in his index’s conservative estimate of manipulations, Wesolowsky stated, “The evidence of substantial cheating is beyond any reasonable doubt” (Benton & Hacker, 2007a, 2007b).

While these statistical detection methods have advantages over the survey and direct observation methods, they do have limitations. First, statistical indices can only detect manipulations of answer sheets. They cannot detect other manipulations of the test administration or manipulations of the teaching process or philosophy, examinee pool, or score reports or standards. Second, these methods can only detect *possible* manipulations

of answer sheets. Some students, classrooms, and schools can experience legitimate large gains in test scores. Likewise, an unusual response string from an examinee or group of examinees does not necessarily mean that the answer sheets have been manipulated. For these reasons, statistical detection indices might overestimate the actual prevalence of test administration manipulations.

On the other hand, many of the statistical indices have been shown to be ineffective in detecting simulated manipulations of answer sheets (Chason & Maller, 1996; Iwamoto, Ningester, & Luecht, 1996). Demonstrating this ineffectiveness, test security firm Caveon analyzed a simulated data set using six of their indices to detect unusual response strings. The data set was simulated so that 3,283 answer sheets had been manipulated by teachers (teachers changing answers from incorrect to correct). The data were also simulated to include 5 “schools” of answer sheets had been manipulated by administrators (principals or school personnel changing answers for all students). With this data set, the firm’s indices were only able to detect 41 (1.2%) of the simulated manipulated answer sheets and none of the five simulated schools (Impara, Kingsbury, Maynes, & Fitzgerald, 2005). Due to their poor ability in detecting manipulations, statistical indices might actually underestimate the actual prevalence of test administration manipulations.

Prevalence Estimates Based on Targeted Research

While the survey, direct observation, and statistical detection methods provide estimates of the prevalence of manipulations of the teaching process or philosophy and administration, they rarely provide information regarding the prevalence of manipulations of the examinee pool. To fill this gap, researchers have designed targeted research studies. These studies find that schools, school districts, and states do manipulate the examinee pool in order to increase test scores.

Educators who manipulate the examinee pool usually do so by excluding lower-ability students from testing. One way in which educators do this is by suspending or otherwise punishing lower-ability students during the test administration period so they are not able to participate. Figlio (2005) hypothesized that during the test administration period, schools with accountability systems would give low-ability students harsher penalties (longer suspensions) than they would give to higher-ability students. Figlio supported his hypothesis by citing evidence that students who receive suspensions of at least one week in length were twice as likely to miss the test administration and the make-up testing dates as students who receive shorter suspensions (p. 3). To test his hypothesis, Figlio analyzed test and discipline data from 41,803 students in Florida school districts during the four years following the introduction of the state's high-stakes *Florida Comprehensive Assessment Test* (FCAT). In his analysis, Figlio compared the lengths of suspensions given to at least two students for the same incident. By classifying students as low- or high-ability based on previous years' test scores, Figlio was able to compare the suspension lengths to see if ability had an influence. Figlio found that, "While schools always tend to assign harsher punishments to lower-ability students than to higher-performing students throughout the year, this gap grows substantially during the testing window. Moreover, this testing window-related gap is only observed for students in testing grades" (pp. 4-5). Figlio also found that given two students suspended for the same incident during the test administration period, low-ability students were 12.3% less likely to take the FCAT than higher-ability students (p. 19).

In stating his conclusions, Figlio did address potential concerns. The first concern is that maybe low-ability students are more likely to be suspended during the test administration period because they want to avoid testing. The second concern is that perhaps low-ability students are more likely to cause the incident or be worse offenders, so therefore they are more likely to receive longer suspensions. The researcher addresses these concerns by reporting that low-ability students tend to get suspended at similar

rates, relative to high-ability students, during the test administration period as in other times of the year.

Another way in which educators can manipulate the examinee pool is by inappropriately classifying low-ability students as disabled. Before NCLB, states such as Florida had rules in which disabled students were exempt from taking the state test (Figlio & Getzler, 2002). States could, then, improve their test scores by simply classifying their lowest-scoring students as disabled. To test the hypothesis that states do, in fact, manipulate the examinee pool in this way, Figlio and Getzler (2002) analyzed 9 years of test data from six school districts in Florida. The researchers found that in the five years before the introduction of the state high-stakes accountability system, between 7.3% and 8.8% of students were classified as disabled. In the three years following the introduction of the accountability system, the classification rate increased each year from 9.4% to 9.6% to 10.8%. Controlling for this nearly linear increase in classification rates, the researchers found that the introduction of the high-stakes accountability system led to a 5.6% increase in the likelihood that a student was classified as disabled. According to the researcher, “the introduction of FCAT testing is associated with a more than 50% higher rate of disability classification” (p. 9).

To address the concern that perhaps Florida is unique in its manipulations of the examinee pool, Cullen and Reback (2002) and Jacob (2007) conducted similar studies in Texas to determine if lower-ability students were inappropriately classified as disabled in order to increase test scores. Both studies concluded that lower-ability students were more likely to be exempt from testing and that educators do, in fact, manipulate the examinee pool through inappropriate disability classification.

Another way in which educators manipulate the examinee pool is by disproportionately focusing instructional resources on students who have the best chance to improve the school’s overall test scores. If school performance is determined by the percent of students earning a proficient score on a test, then educators may be tempted to

focus their attention on students who earned scores just below proficient the previous year at the expense of students who scored extremely high or low on the previous year's test. Neal and Schanzenbach (2007) tested this hypothesis in a study of test data from Chicago Public Schools. The data came from 1998, following the introduction of a high-stakes accountability system. This accountability system evaluated schools by examining the percentage of students in each school who earned a proficient score. The researchers found that students near the middle of the achievement distribution achieved at a higher level after the accountability system was introduced. They also found that students at the low-end of the achievement distribution achieved at the same or even a lower level after the accountability system was introduced. The researchers found that "for at least the bottom 20% of students, there is little evidence of significant gains and a possibility of lower than expected scores" (p. 27) following the introduction of the accountability system. This seems to support the hypothesis that educators manipulate the examinee pool in Chicago.

Reback (2007) conducted a similar analysis on test data from Texas. The researcher hypothesized that a state accountability system "increases incentives for schools to improve the performance of students who are on the margin of passing but does not increase short-run incentives for schools to improve other students' performance" (p. 1). While Reback found that accountability systems do improve overall student achievement, most of these gains were realized for students whose achievement levels were closest to the cut-scores. The researchers found that "other students only make greater than expected gains in this situation if their own performance is particularly important for their schools' rating" (p. 33). These conclusions support the belief that educators in Texas manipulate the examinee pool in order to increase test scores.

An unusual way in which educators have been shown to manipulate the examinee pool is through, surprisingly enough, the school lunch program. Figlio and Winicki (2003) designed a study to target this specific manipulation method. After claiming "the

link between nutrition and cognitive ability has been well established” (p. 382), the researchers examined the school lunch menus from Virginia public schools during the 1999-2000 academic year. From this data, the researchers were able to conclude, “schools threatened with accountability sanctions increase the caloric content of their lunches on testing days in an apparent attempt to boost short-term student cognitive performance” (p. 381). Moreover, the researchers found evidence that this manipulation of the school lunch program was effective in raising test scores. The researchers found that schools that increased the caloric content of their lunches on testing days by 100 experienced a 7% increase in the pass rate of students on the mathematics test (p. 392).

Another targeted research method used to determine the prevalence of manipulations of score reports or score standards is through analyses of state testing programs and their cut-scores for proficiency. By lowering their proficiency standards or making their state tests easier, states can inflate test scores (or test score comparisons) without actually increasing student achievement. Some studies that attempt to document these manipulations simply compare the tests and proficiency standards from each state. In a report for *CBS News*, Wallace (2007) concluded that the large variability in state proficiency rates was due, primarily, to differences in the difficulty of state tests and the cut scores used for proficiency. Sturrock (2006), reporting for the *San Francisco Gate*, suggested that differences in results from California’s state test and the National Assessment of Educational Progress (NAEP) provided evidence that California had manipulated scoring standards to make it appear as though student achievement had increased. Although the methods and conclusions from these analyses are oftentimes questionable (as will be discussed later), they can provide evidence of educators manipulating test scores

Why Do Educators Manipulate Test Scores?

The previously discussed surveys, observational studies, statistical detection studies, and other targeted research results provide evidence that educators are manipulating test scores and that these manipulations can have a significant impact on test scores. In order to prevent this behavior, it must first be understood why educators manipulate test scores. Research into cheating and inappropriate test preparation activities suggests at least four reasons why educators would manipulate test scores: (1) educators are former students, (2) pressure from high-stakes test accountability systems, (3) a lack of understanding of what behaviors are inappropriate, and (4) a lack of oversight and policies to deter manipulations. While these reasons are neither mutually exclusive nor exhaustive, they can provide insight into how manipulations could possibly be deterred.

Educators Are Former Students

Over the last century, researchers have published more than one hundred studies on the prevalence of student cheating on exams (Cizek, 1999, 2003). Estimates have ranged from 5% of examinees cheating on any particular occasion (Bellezza & Bellezza, 1989) to 75% of students admitting to some form of cheating before graduating high school (Impara, Kingsbury, Maynes, & Fitzgerald, 2005) to more than 80% of American undergraduate students admitting to cheating during college (Passow, Mayhew, Finelli, Harding, & Carpenter, 2006). A 2006 survey of 36,000 students by the Josephson Institute (2006) found that 60% of examinees cheated on a test during the past year; 35% cheated two or more times; 33% admitted using the internet to plagiarize an assignment; and 27% admitted to lying on at least one question on this survey.

Educators are former students. If students cheat on tests and students become educators, then educators may continue that cheating behavior to manipulate test scores in their classrooms. In a study of 5,280 students across nine academic majors, Bowers

(1964) found that 52% of undergraduate education majors reported cheating on a test during college. Cizek (1999) described an observational study conducted in the 1920s in which 110 women about to begin student teaching were allowed to score their own tests in a college-level education course. The researchers found that 30 (27%) of the women cheated by changing their answers during the self-scoring, with 4 (3.6%) of the women changing more than 10 answers. Based on these studies, it is safe to assume that many educators have cheated on a test at least once in their academic careers. Because research has found that the decision to cheat in college is correlated with the decision to later engage in other unethical behaviors in the workplace (Crown and Spiller, 1998), educators who cheated as students may choose to manipulate test scores in their classrooms.

In replicating a 30-year old large-scale study of undergraduate student cheating, McCabe, Trevino, and Butterfield (2001) found that while the overall prevalence of cheating increased only modestly, the prevalence of the most blatant forms of test cheating “increased significantly” (p. 221). In explaining possible causes of cheating, the researchers found that contextual factors (peer cheating, peer disapproval of cheating, or perceived severity of penalties for cheating) were significantly more influential than individual factors (age, gender, academic ability, or participation in extracurricular activities). The researchers explained:

Students who might otherwise complete their work honestly observe cheating by others and convince themselves they cannot afford to be disadvantaged by students who cheat and go unreported or unpunished. Although many find it distasteful, they too begin cheating to level the playing field. The strong influence of peers’ behavior may suggest that academic dishonesty not only is learned from observing the behavior of peers, but that peers’ behavior provides a kind of normative support for cheating. The fact that others are cheating may also suggest that, in such a climate, the non-cheater feels left at a disadvantage. Thus cheating may come to be viewed as an acceptable way of getting and staying ahead (pp. 220-222).

This reasoning is supported by a 1997 report from *Who's Who Among American High School Students* on academic cheating (Newberger, 1997). According to the report, 92% of confessed cheaters declared that they had never been caught. Newberger suggested that students cheat not only to keep up with other cheating students, but also because they think they can get away with it. Crown and Spiller (1998) also found that while lower-ability students were slightly more likely to cheat, contextual factors such as a school's lack of an honor code or weak penalties for cheating increase the likelihood of cheating. The researchers also concluded, "The amount of unflattering attention the popular press gives to the students reporting high percentage levels of collegiate cheating could lead many students to the conclusion that they must cheat just to keep up with their peers" (pp. 684-695).

Cizek (2003) suggests that these contextual factors also influence an *educator's* decision to manipulate test scores. As he describes, "Because so much of that cheating went undetected and unpunished, and because they can easily put themselves in the position of examinees desperate to pass a test, those who give tests may often be tempted to turn a blind eye to cheating (pp. 6-7). Jacob and Levitt (2003) found evidence of these contextual influences in their study of educators who change student answers on tests. In addition to finding that younger educators were more likely to cheat than older educators, the researchers found that educators in classrooms that performed poorly on the previous year's exam, and educators in classrooms with higher poverty rates and more minority students were more likely to cheat. These contextual factors along with many educators' histories of cheating as students may explain why many educators manipulate test scores.

Pressure From State Accountability Systems

Other researchers suggest that high-stakes accountability systems are the reason why educators manipulate test scores. Some studies have found that the pressure felt by educators, whether real or perceived, to improve test scores causes them to manipulate

test scores. Hatch and Freeman (1988) found this when they interviewed kindergarten teachers in Ohio. 67% of the interviewed teachers reported “implementing instructional practices in their classrooms that they considered to be antithetical to the learning needs of young children; they did this because of the demands of parents and the district and state accountability systems” (Hatch & Freeman, 1988, p. 146). Hamilton and Stecher (2006) also found this in their survey of 2,628 math teachers and 262 principals from elementary and middle schools in California, Georgia, and Pennsylvania. The researchers found that 79-92% of teachers felt a great deal of pressure to improve scores on the state mathematics test. Because of this pressure, 19% - 78% of teachers manipulated the teaching philosophy or process by: (a) focusing more on topics emphasized on the state test, (b) emphasizing the formats and styles of test items in instruction, (c) spending more time teaching general test-taking strategies, (d) focusing more effort on students who are close to proficient, or (e) offering more assistance outside of school to help students who are not proficient (p. 22). More than half of the principals also responded to this pressure by encouraging teachers to manipulate the teaching philosophy or process by: (a) distributing commercial test preparation materials, (b) encouraging or requiring teachers to spend more time on tested subjects and less time on other subjects, or (c) encouraging teachers to focus their efforts on students close to meeting the standards (p. 24).

Nolen, Haladyna, and Haas (1992) found that the perception of pressure due to the accountability system was enough to cause educators to manipulate test scores. In a survey of Arizona educators, the researchers found that more than 43% of teachers believed that administrators and school boards used test scores to evaluate teacher effectiveness. This would be perceived as a great deal of pressure, since only 7% of teachers believed test scores should be used to evaluate teacher effectiveness. When the researchers interviewed administrators, they found that only 15% actually used test scores in the evaluation of teacher performance. The teachers further perceived pressure from

administrators to manipulate test scores. The researchers found that 7% of teachers believed they were encouraged to teach actual test items to their students. Furthermore, more than one-third of teachers believed they were encouraged to use more class time than required for test preparation activities and more than two-thirds believed they were encouraged to teach test-taking skills, focus on skills from the test, and use the item format from the test on classroom tests (pp. 11-12).

Survey results from the National Board on Educational Testing and Public Policy (Pedulla et al., 2003) further supports the notion that a state's accountability systems is the reason why educators manipulate test scores. Although obtaining a response rate of only 35%, the researchers found that 72% of the 4,195 teachers responding to the survey agreed to the statement, "The state-mandated testing programs lead some teachers in my school to teach in ways that contradict their own ideas of good educational practice" (p. 31). Jacob and Levitt (2003) also reached this conclusion, finding that "a high-stakes testing environment increases probability that a teacher would cheat" (p. 17).

In reviewing published news reports of educators manipulating test scores, Nichols and Berliner (2005) concluded that the reports provided evidence of Campbell's Law. Campbell's Law states, "The more any quantitative social indicator is used for social decision-making, the more subject it will be to corruption pressures and the more apt it will be to distort and corrupt the social processes it is intended to monitor" (Campbell, 1975, p. 35). The researchers argued that pressures from state accountability systems lead educators to feel justified in manipulating test scores, stating, "It is plausible that teachers and administrators are trying to resist a system they see as corrupt or unfair, as do tax, religious, and civil rights protestors across this nation" (Nichols & Berliner, 2004, p. 24).

Gay (1990) found further evidence of this feeling of justification. Of the 161 respondents to his survey, Gay found 60% rationalized the use of manipulations in order to improve the image of the teacher and 11% justified manipulations to help students.

Shepard (1990) found that educators believe “a district is at a disadvantage if it plays fair by teaching to a broad curricular domain and by avoiding more than one-time practice on test formats” (p. 21). Cizek (2003) also reached the conclusion that “educators appear to be growing increasingly indifferent toward cheating and even increasingly to feel that cheating is a justifiable response to externally mandated tests” (pp. 375-376).

Educators Unaware of Manipulations & Their Impact

Another reason why educators manipulate test scores is because they are unaware of the definition and impact of manipulations. Educators simply do not know which behaviors are manipulations and why they should not manipulate test scores. Cizek (2003) claims:

There is an abundance of information to guide test takers and test administrators in how to avoid inappropriate testing practices. For their part, test developers usually produce carefully scripted directions for administering their tests and provide clear guidelines as to which kinds of behaviors on the part of examinees and administrators are permissible and which are not. Acceptable and unacceptable behaviors are sometimes formalized in state administrative codes or statutes.... Numerous professional organizations have published statements on cheating. (pp. 364-365).

To support Cizek’s claim, Appendix C displays some of the codes and standards endorsed by the American Counseling Association, the American Educational Research Association, the American Psychological Association, the American Speech-Language-Hearing Association, the Joint Committee on Testing Practices, the National Association of School Psychologists, the National Association of Test Directors, the National Council on Measurement in Education, and the National Education Association. These codes and standards all provide guidance to educators in determining which testing behaviors are appropriate or inappropriate. The existence of these codes and standards led Cizek (2003) to conclude:

... there has not generally been a dissemination problem regarding what constitutes integrity in testing or cheating on tests. Virtually

everyone involved in testing knows how to administer tests that yield credible, accurate results. (p. 365).

Evidence from teacher surveys, however, contradict Cizek's conclusion. Surveys have shown that educators do not agree with administrators or testing specialists as to which behaviors are appropriate or inappropriate in preparing for or administering tests. Kher-Durlabhji and Lacina-Gifford (1992) asked pre-service teachers to determine which testing activities are appropriate or inappropriate to use. Table 2.3 displays the percent of respondents who rated each testing practice as appropriate. While some of the listed activities are appropriate, the fact that some teachers believed changing completed answer sheets or presenting actual test items for practice are appropriate demonstrates that teachers do not understand what behaviors are manipulations.

Lai and Waltman (2007) similarly surveyed a large sample of Iowa public school teachers to determine their perceptions of the ethicality of testing behaviors. Table 2.4 displays the percent of respondents who believed each testing activity was ethical. The researchers found "unexpectedly high percentages of teachers rated practicing with exactly the same test that will be administered this year as being 'very ethical'" (p. 11) and only a median 80% of teachers within schools viewed this practice as being unethical. Another surprising result was that only 75% of teachers believed that "providing instruction without checking the content of the test" was an ethical practice. These results provide further evidence that educators are unaware of which testing behaviors are appropriate and which behaviors are manipulations.

Lai and Waltman (2007) similarly surveyed a large sample of Iowa public school teachers to determine their perceptions of the ethicality of testing behaviors. Table 2.4 displays the percent of respondents who believed each testing activity was ethical. The table provides further evidence that educators are unaware of which testing behaviors are appropriate and which behaviors are manipulations. The researchers found "unexpectedly high percentages of teachers rated practicing with exactly the same test that will be administered this year as being 'very ethical'" (p. 11) and only a median 80%

of teachers within schools viewed this practice as being unethical. Another surprising result was that only 75% of teachers believed that “providing instruction without checking the content of the test” was an ethical practice. The researchers found that the results of the survey “raise questions about the extent to which teachers understand the testing procedures used in Iowa and are aware of existing professional standards dictating appropriate versus inappropriate testing practices” (p. 49).

Moore (1994) presented 42 elementary school teachers and 10 testing specialists from the Midwest a list of 40 test administration and preparation practices. The subjects were asked to rate the inappropriateness of each activity on a five-point scale (1 = appropriate; 5 = inappropriate). Table 2.5 shows the mean inappropriateness ratings assigned by teachers and testing specialists to eight categories of test administration and preparation practices. The table shows significant differences between teacher and testing specialist inappropriateness ratings in six of the eight categories of practices. Testing specialists rated practices as being more inappropriate and teachers rated the practices as being more appropriate. Popham (1991) conducted a similar analysis and also found discrepancies between teachers’ and administrators’ perceptions of the relative appropriateness of testing behaviors. These discrepancies in perceived appropriateness once again demonstrate that educators are unaware of which testing behaviors are inappropriate manipulations.

Table 2.3 Percentage of pre-service teachers rating each activity as appropriate

98.6%	Encourage students to do their best
96.0%	Teach test-taking skills
86.5%	Check student's completed answer sheets
79.7%	Send note home to parents to elicit cooperation
66.2%	Use commercial test preparation materials
37.9%	Teach according to test objectives
37.8%	Develop curriculum based on test
37.4%	Practice alternative forms of test
36.5%	Teaching objectives based on standardized test
23.4%	Rephrasing wording of questions
8.1%	Present actual test items for practice
2.7%	Allow more time than allocated for testing
1.4%	Change completed answer sheets
1.4%	Give hints or clues
1.4%	No special test preparation
0%	Dismiss low-achieving students from test taking
0%	Change answers of low-achieving students

Source: Kher-Durlabhji and Lacina-Gifford, 1992, p. 13

Table 2.4 Median percent of Iowa teachers within schools rating each activity as ethical or unethical

Ethical	Behavior	Unethical
100%	Teach test-taking skills	0%
100%	Use previous year's test data to inform instruction	0%
80%	Use practice tests	0%
75%	Provide instruction without checking the content of the test	0%
75%	Structure all/most classroom tests like the ITBS (standardized test)	0%
60%	Review tested content/skills prior to testing	20%
20%	Practice with last year's questions	67%
17%	Routinely provide instruction only on tested content areas	67%
11%	Practice with the same test questions	80%

Source: Lai & Waltman, 2007

Table 2.5 Teacher and testing specialist mean appropriateness ratings

Testing Behavior	Teacher Mean Rating	Testing Specialist Mean Rating
Generalized Test-Taking	1.2	1.1
Motivational Activities	1.6	2.3*
Same Format Preparation	1.7	2.4*
Pretest Intervention	2.4	2.8*
Previous Form Preparation	2.5	3.6*
Posttest Intervention	3.7	3.8
Current Form Preparation	2.9	3.8*
During Test Intervention	2.9	3.7*

* $p < .05$

Source: Moore, 1994

Educators do have access to an abundance of information to guide them in selecting testing practices. If dissemination of this information is not the problem, perhaps the reason why educators manipulate test scores is because of an *overabundance* of professional codes and standards. In addition to the codes and standards listed in Appendix C, many states and school districts have rules and laws regarding test security. Additionally, several researchers have provided guidelines for teachers to evaluate the appropriateness of test preparation activities. These guidelines are displayed in Tables 2.6 and 2.7. The overwhelming number of guidelines, codes, standards, rules, and evaluative criteria may overwhelm educators who try to learn which testing behaviors are appropriate or inappropriate. Furthermore, since these guidelines are not requirements, educators may simply ignore them. Finally, educators such as Kilian (1992) have found these guidelines and rules to be, at best, vague, and at worst, contradictory. The lack of clear guidance has led to the lack of understanding of which testing behaviors are inappropriate which may be a major reason why educators manipulate test scores.

Table 2.6 Examples of appropriate and inappropriate test preparation practices

Source	Appropriate	Uncertain	Inappropriate
Frederickson (1984)	<ol style="list-style-type: none"> Physical/emotional/intellectual preparation Time-use skills Error-avoidance skills Guessing strategies 	<ol style="list-style-type: none"> Teach objectives from organizations that scan objectives from many tests Teach objectives matching test objectives Use the format of the test questions 	<ol style="list-style-type: none"> Teach from published parallel form of test Instruction from the actual test
Mehrens & Kaminski (1989)	<ol style="list-style-type: none"> General instruction on objectives (no reference to test) Teach test-taking skills 	<ol style="list-style-type: none"> Teach objectives from organizations that scan objectives from many tests Teach objectives matching test objectives Use the format of the test questions 	<ol style="list-style-type: none"> Developing curriculum based on test content Preparing objectives based on test items to teach Presenting items similar to test items Using commercial test prep. packages Dismissing low-achieving students from testing^b Presenting items verbatim from the test^b
Haladyna et al. (1991)^a	<ol style="list-style-type: none"> Training in test-wiseness skills Sanitizing answer sheets (for all students) Increasing motivation via appeals to parents & students 	<ol style="list-style-type: none"> Unauthorized access to test materials Failure to keep materials in locked storage Tampering with sealed test booklets Administration of parallel forms prior to testing Using specific content of test items in instruction Failure to sign and return security forms 	<ol style="list-style-type: none"> Practicing with actual test items prior to testing Providing answers or coaching students Altering answer documents prior to scoring Photocopying test materials
Mehrens et al. (1993)^c	<ol style="list-style-type: none"> Teach test-taking skills Use previous year's data to inform instruction Provide instruction without checking test content Use practice tests Structure all/most classroom tests like the actual test Review tested content/skills prior to testing 	<ol style="list-style-type: none"> Practice with last year's questions 	<ol style="list-style-type: none"> Practice with the same test questions Provide instruction only on tested content
Lai & Waltman (2007)^d	<ol style="list-style-type: none"> Ensure the curriculum is being taught effectively Ensure students are ready physically & psychologically Ensure students are using appropriate time management and deductive reasoning strategies Ensure testing environments are conducive to optimal test performance Ensure test administrators are knowledgeable/prepared Ensure teaching and learning climates in classrooms and schools are positive and productive 	<ol style="list-style-type: none"> Provide students a review of tested content Developing classroom tests with items of the same format as the test Special efforts before test administration to teach test-taking skills or increase student motivation for the test 	<ol style="list-style-type: none"> Practice with current form of the test Practice with next year's form of the test Practice with previous forms of the test Develop practice tests locally that are similar in content or format to the actual test form Drill students for short-term retention Any activity designed primarily to fit the test before the test is scheduled to be administered Any activity that needs to be implemented just before the test is scheduled to be administered Instruct students to blindly guess or use testing "tricks"
CEA (2007)^e			

a – the researchers rated the ethicality of practices; skills rated as “unethical” or “highly unethical” are listed as inappropriate in this table.

b – these activities were described as being “highly unethical”

c – actions deemed unacceptable by the vast majority of states are classified as “inappropriate” and actions deemed unacceptable by many but not all states are classified as “uncertain”

d – based on median school-level ethicality rating from Iowa public school teachers

e – these test preparation activities were classified in regards to their use in preparing for the *Iowa Tests of Basic Skills* in Iowa

Table 2.7 Guidelines and criteria to select appropriate test preparation activities

Study	Guidelines / Criteria
Ligon & Jones (1982)	An appropriate test preparation activity “contributes to students’ performing near their true achievement levels, and contributes more to their scores than would an equal amount of regular classroom instruction”
Matter (1986)	Inappropriate activities are, “Any additional activities not incorporated into regular ongoing instruction”
Popham (1991)	<p>Test preparation activities should be evaluated through reference to two evaluative standards:</p> <p>Professional ethics: No test preparation activity should violate the ethical standards of the education profession. Practices can be considered unethical if they violate general ethics of theft, cheating, and lying. Practices such as changing student answers on the answer sheet, providing practice items from the actual test, and excluding lower ability students from testing would be examples of violations.</p> <p>Educational defensibility: No test preparation practices should increase students’ test scores without simultaneously increasing student mastery of the content domain tested. Violations of this standard would include all practices that attempt to improve test performance through test-taking skills, testwiseness skills, motivational strategies, and inappropriate practices such as providing additional examples to students during testing.</p>
Crocker (2006)	Criteria for assessing the appropriateness of proposed preparation activities: Validity, Academic ethics, Fairness, Educational value, Transferability
Iowa Testing Programs (2005)	<p>The appropriateness of any proposed practice should meet either of the two following standards:</p> <ul style="list-style-type: none"> • It will promote the learning and retention of important knowledge and content skills that students are expected to learn. • It will decrease the chance that students will score lower on the test than they should due to inadequate test-taking skills or limited familiarity with the item formats used on the test. <p>Activities that do not meet one of these criteria are more likely to be unethical, to promote only temporary learning, or to waste instructional time.</p> <p>Test preparation activities must meet three criteria:</p>
CEA (2007)	<p>Academic Ethics: The action should not contribute to the misrepresentation or falsification of information The action should not be perceived by students, parents, or the community as being dishonest The action should not result in a violation of district policy or copyright (e.g., an illegal act)</p> <p>Score Meaning & Use: Test scores should accurately represent student learning related to the specific set of content and skills covered by the test Test scores should not be influenced by a inadequate test-taking skills or limited familiarity with test item formats Test scores should allow users to make accurate inferences related to the larger domain of content and skill areas</p> <p>Educational Value: The action should promote learning and long-term retention of content/skills defined by district standards/curriculum The action should provide students with knowledge/skills that apply to a broad range of situations/context The amount of instructional time dedicated to test preparation should be warranted in light of educational opportunities lost The actions should be matched with the needs of individual students</p>

Lack of Oversight and Policies

A fourth reason why educators manipulate test scores is because many states have not developed or implemented high-quality policies to prevent these behaviors. Surveys of state departments of education from the National Council on Measurement in Education (Mehrens, Phillips, & Schram, 1993) and Caveon (Sorensen, 2006) along with a survey of Iowa public school districts (Thiessen, 2007) have found that many states and school districts have not implemented simple policies to deter educators from manipulating test scores. Table 2.8 displays the results of these surveys.

Mehrens, Phillips, and Schram (1993) surveyed 46 state departments of education and found the majority claimed to have written test security policies. 65% of states claimed to have written policies addressing test preparation activities, 88% claimed to have policies addressing the security of test materials, and 95% claimed to have policies addressing test administration activities. More than a decade later, after the accountability systems under NCLB had been implemented, Sorensen (2006) conducted a similar survey and found fewer states claimed to have test security policies. While 77% claimed to have policies addressing test preparation activities, only 63% had policies addressing the security of test materials, and less than half (47%) had policies addressing test administration behaviors. It is not known why fewer states claimed to have test security policies in 2006 than in 1993.

Because Iowa does not have a state test security policy, Thiessen (2007) surveyed 154 Iowa public school districts to determine the existence and quality of test security policies at the district level. One year after the Iowa Department of Education disseminated a sample policy and guidance to develop their own policies to school districts, only 27% reported adopting a test security policy. 73% of Iowa public school districts had not yet developed a test security policy and 65% had no plans to adopt a policy in the near future.

Table 2.8 State and school district test security policies

	Mehrens, Phillips, & Schram, 1993	Sorensen, 2006	Thiessen, 2007
Sample	46 state departments of education	34 state departments of education	154 Iowa public school districts
Have written policies addressing:			
Test preparation activities	65%	77%	27%
Test administration activities	95%	47%	27%
Security of test materials	88%	63%	27%
Have no plans to adopt a written policy in the near future	---	---	65%
Policies:			
Require test proctors to be trained	---	78%	78%
Require independent monitoring of test administration	---	---	31%
Identify individuals responsible for responding to incidents	---	65%	90%
Provide a separate budget for test security	---	5%	---
Have policies to specify that:			
Test materials must be sealed before administration	48%	---	---
New test forms must be used each year	54%	---	---
Teachers cannot examine tests before administration	62%	---	---
Routinely run statistical analyses of answer sheets to check for:			
Unusual number or pattern of erasures	49%	---	---
Unusual score fluctuations	64%	---	---
If cheating or manipulations are suspected:			
The incident is investigated	64%	81%	---
A written policy guides the investigation	24%	47%	---
Percentage of suspected incidents that are later confirmed	50%	---	---
If cheating or manipulations are confirmed:			
No sanctions are imposed (or only a letter is sent)	52%	11% - 18%	---
The guilty party is given a stern warning or reprimand	19%	61%	---
The guilty party is suspended	12%	34%	---
The guilty party is dismissed	8%	45%	---

The policies that are adopted by states and school districts are not all of high quality. In fact, some policies ignore fundamental methods to prevent test score manipulations. For example, Sorensen (2006) and Thiessen (2007) both found that only 78% of states and Iowa school districts required their test proctors to be trained prior to test administration. Survey results also show that less than two-thirds of states have identified an individual or group of individuals to be responsible for answering questions

about test security or responding to test score manipulations, and only 5% of states provide a separate budget for test security. In Iowa, less than one-third of the school districts surveyed had policies requiring independent monitoring of the test administration. Of those districts that did have policies, Thiessen found that only 8% had policies that could be effective in deterring test score manipulations.

Mehrens, Phillips, and Schram (1993) found that many states did not require basic test materials security precautions. They found that less than two-thirds of state policies specify that teachers cannot examine test questions before the test administration and only 54% of state policies require new test questions to be used annually. They also found that less than half of the state policies require test materials to be sealed prior to administration.

Even fewer states had policies that called for routine statistical analyses of answer sheets to detect potential manipulations. 49% of states reported routinely analyzing answer sheets for unusual patterns of erasures and 64% reported routinely checking for unusual score gains in the 1993 survey (Mehrens, Phillips, & Schram, 1993). A 2006 poll conducted by the Philadelphia Inquirer found that fewer than half of all states attempt to detect manipulations on their state tests (Patrick & Eichel, 2006). The 2006 survey from Caveon found that 52% of states claimed to routinely run statistical analyses to check for evidence of manipulations (Sorensen, 2006). This survey also found that 25% of states had no plans to implement statistical analysis methods to detect manipulations. Thiessen (2007) found that only 6% of Iowa public school districts routinely conducted these analyses.

The surveys of state departments of education also found that many state policies are weak when it comes to investigating reports of manipulation incidents. Mehrens, Phillips, and Schram (1993) found that when educators are suspected of cheating on state tests, only 64% of those incidents are investigated. The researchers also found that when states do decide to investigate reported incidents, only 24% of states have a written policy

to guide the investigation. Sorensen (2006) found that the situation may have improved slightly since 1993, finding that 81% of states claim to investigate suspected incidents of cheating and 47% have a written policy that prescribes actions to be taken when cheating is suspected.

According to the 1993 survey, half of all suspected incidents of cheating are later confirmed (Mehrens, Phillips, & Schram, 1993). The states admit, however, that the sanctions imposed on confirmed manipulators are not standardized. The 1993 survey found that in more than half of all confirmed incidents; the guilty party was either not penalized at all or only sent a letter from the state department of education. By 2006, only 11% of states did not penalize confirmed cheaters (Sorensen, 2006). In 1993, a stern warning or reprimand was given in 19% of confirmed cases; the guilty party was suspended in 12% of confirmed cases; and the guilty party was dismissed in 8% of confirmed cases of educator cheating. By 2006, the sanctions increased in magnitude, with 61% of confirmed cheaters receiving a stern reprimand, 34% receiving a suspension, and 45% being dismissed.

The lack of policies providing basic considerations of test materials security, test preparation activities, and test administration behaviors might not encourage educators to manipulate test scores, but it certainly does nothing to deter or prevent educators from engaging in these behaviors. As will be discussed in the next section, high quality, enforceable policies can deter educators from manipulating test scores.

How to Prevent Manipulations: Evaluation of State Policies

Researchers have provided several suggestions to deter or prevent educators from manipulating test scores. In their research on student cheating, Aiken (1991), Burns (1988), Cizek (1999), and Singhal and Johnson (1983) suggested that test developers can prevent manipulations by modifying the tests used to make high-stakes decisions. The researchers recommend test publishers develop constructed-response test items,

suggesting that it would be more difficult for educators to manipulate scores from these constructed-response tests than it would be for multiple-choice tests. The researchers also suggested that test developers develop new test forms with new items each time the test is administered to make it more difficult for educators to manipulate scores.

Impara and Foster (2006) suggest that while these methods may be effective in preventing students from cheating on tests, “item and test development strategies do little to reduce” educator score manipulations (p. 93). Also, developing new forms of constructed-response tests for each administration would be labor-intensive and inefficient. Test developers should focus on developing the best items to measure the construct of interest; not developing items that are most resistant to manipulations.

Another method to prevent test score manipulations might be to catch and punish educators who manipulate test scores. This would require states to implement methods to detect manipulations, such as statistical analyses of answer sheets or surveys after test administration, and to strengthen the sanctions imposed upon educators found to manipulate scores.

Ignoring the facts that statistical analyses have been shown to be ineffective in detecting manipulations (Chason & Maller, 1996; Impara, Kingsbury, Maynes, & Fitzgerald, 2005; Iwamoto, Ningester, & Luecht, 1996) and that states have been reluctant to punish educators who manipulate test scores, suppose states could accurately detect manipulators. Even if states punished these manipulators, this *after-the-fact* approach to deter manipulations would be labor-intensive and, if used as the only deterrent, would most likely be ineffective. In their study on student cheating, Bunn, Caudill, and Gropper (1992) found that both the expectation and severity of punishment had no effect on reducing cheating behaviors in students.

A third method to reduce the number of educators who manipulate test scores would be by developing, implementing, and disseminating high quality policies that both discourage manipulations and encourage honesty and integrity. In their study on student

cheating, McCabe and Trevino (1993) found that students were less likely to cheat if their schools had severe penalties for cheating coupled with high quality policies or honor codes on student cheating. Based on a decade of research from more than 14,000 students, the researchers found that neither sanctions nor honor codes alone reduced cheating, but that the combination of the two was effective at reducing student cheating by as much as 20% (McCabe & Trevino, 2002). In order to be effective, the honor code must explain why academic integrity is important, describe which behaviors are appropriate or inappropriate, and clearly show the school's commitment to academic integrity. The sanctions described in the policy must be significant and consistently applied to those caught cheating. The researchers also found that in order to be effective, the policies must be developed with input from students and supported by top administrators.

Cizek (2003) suggested that the combination of policies and sanctions might also work to deter educators from manipulating test scores. While educators have had an overabundance of professional codes and standards and test administration manuals to guide their behavior, these guidelines have not been enforceable and educators have not been held accountable for following them. Policies developed by state boards of education, on the other hand, may be effective because educators would be required to follow them. Cizek (1999, 2001) recommended that states bear responsibility for developing policies to prevent educators from manipulating test scores.

Some states have developed specific policies and regulations to address test score manipulations, but many others have left this task up to individual school districts (Cizek, 1999; Mehrens, Phillips, & Schram, 1993; Patrick & Eichel, 2006). Unfortunately, very little research has been conducted to determine the existence or evaluate the quality of these state and district developed policies. In 1999, Cizek wrote, "Only one study has been conducted to investigate the existence of policies at the elementary and secondary school level" (1999, p. 171); "Unfortunately, no research has actually examined the

content of cheating policies” (p. 174); and that it was not even known if schools, school districts, or states had any policies to address educator manipulations (p. 171). Since that time, Thiessen (2007) conducted an evaluation of test security policies in Iowa public school districts and found that 73% of districts had no policy, 22% had policies that were inadequate to deter test score manipulations, and 5% had policies that could effectively prevent test score manipulations.

Due to a lack of research, it is not currently known how effective test security policies are in preventing educators from manipulating test scores. To address this, a framework must be developed to evaluate the quality of state test security policies. Then, analyses should be conducted to determine the effectiveness of those policies in reducing the prevalence and impact of test score manipulations.

Test Security Policy Content: Evaluative Framework

As Cizek (1999) noted, little research exists to evaluate the content of state policies to deter educators from manipulating test scores. In their surveys of state departments of education, Mehrens, Phillips, and Schram (1993) and Sorensen (2006) provides general guidelines such as recommending states conduct statistical analyses of answer sheets and outline sanctions for those caught manipulating scores. In his discussion of policies and honor codes, Cizek (1999, 2001) recommended that states describe specific activities in defining what testing behaviors are appropriate or inappropriate. Professional codes and standards, example policies from test publishers (Harcourt Assessment, 2006, 2007; Iowa Testing Programs, 2005; Riverside, 2006), and state departments of education also provide guidance as to the content of effective test security policies. Finally, Thiessen (2007) provided content recommendations in an evaluative framework for the development, adoption, and implementation of district test security policies in Iowa. These sources all provide guidance as to what content a test security policies needs to effectively prevent test score inflation by deterring educators

from manipulating the teaching philosophy or process, examinee pool, test administration, or score reports or standards. The recommendations, along with the specific manipulations they are intended to deter, are summarized in Table 2.9.

If state test security policies are to foil educators from manipulating test scores, these policies must:

- (F) Formalize beliefs of state educators regarding the role of testing and practices
- (O) Oversee test preparation, administration, and scoring activities
- (I) Inform educators about why some behaviors and activities are unacceptable
- (L) Limit opportunities for educators to manipulate test scores

The policy content recommendations can be reclassified to fit into this framework to FOIL test score manipulations. Table 2.10 displays the recommendations under this FOIL framework.

The first way in which a state test security policy can foil test score manipulations is by formalizing both state educators' beliefs about the role of testing in education and current state testing practices. This formalization begins with state educators providing input into the content of the written test security policy. In their review of research on student cheating, McCabe and Trevino (2002) found that school cheating policies were more effective if students were encouraged to assist in the development of the policy content. Similar results could be found if educators assist in the development of state test security policies.

The formalization of a state test security policy also requires states to disseminate their policies and ensure all educators understand the policies. Cizek (2003) made a similar recommendation after examining several published news reports on test score manipulations:

Reports of cheating are often accompanied by protestations from the guilty parties that they did nothing wrong. Every implementation of high-stakes tests should be accompanied by dissemination of clear guidelines regarding permissible and impermissible behaviors. Such reminders should be clearly

worded, pilot tested, distributed, and signed by all who handle testing materials, including test site supervisors, proctors, and examinees (pp. 377-378).

If educators do not understand the content, the policy should identify individuals responsible for answering questions about policy content. Cizek (1999) recommended identifying an individual in charge to prevent student cheating and test publisher Harcourt Assessment (2007) requires the identification of such an individual before shipping tests to a customer.

An effective test security policy should also formalize the state's current testing practices. The policy should outline procedures for handling testing materials and testing irregularities, for re-testing students, and for correcting possible scoring errors. The policy should also formalize the state's belief that test score manipulations are unacceptable. To do this, the policy should require mandatory reporting of all incidents of manipulations, while, at the same time, providing protection for those reporting suspected incidents. These recommendations are based on the finding of Gay (1990) that only 20% of educators were willing to report testing irregularities to their school administrators (p. 4). Based on research into effective policies to prevent student cheating (McCabe & Trevino, 2002), state test security policies should outline the procedures that will be used to investigate suspected cases of manipulations and specify the sanctions imposed on those found guilty of manipulating test scores. Cizek (2003) made similar recommendations, stating:

Enforce penalties for cheating and change the system of investigation. ...there are strong disincentives for educational personnel to report cheating; and in most jurisdictions, the responsibility for investigating cheating involves personnel at the school or district level and agencies such as boards of education with an inherent conflict of interest when it comes to ferreting out inappropriately high apparent student achievement. Revised procedures should include... increased protection for whistleblowers; more streamlined procedures and stiffer penalties for cheating, including permanent disqualification from teaching within a state and more coordinated sharing of information regarding educators who have had their licenses revoked; and delegation of responsibility for investigating incidents of cheating to an independent agency (p. 381).

The next way in which a state test security policy can foil test score manipulations is by overseeing all aspects of test security. As recommended by the test security firm Caveon (Sorensen, 2006) and Cizek (2003), states should regularly audit the security of their current testing programs. Regular audits serve to evaluate the effectiveness of a state's test security policies and encourage states to refocus efforts on test security.

Effective state policies should also provide for oversight of the administration of tests. As Cizek notes, many tests are administered behind closed doors with little independent oversight (2003, p. 381). Jacob and Levitt (2003) found that teachers who administered exams to their own students without independent oversight were 50% more likely to cheat. By requiring independent monitoring of the test administration in a random sample of classrooms, state test security policies can ensure that test administration directions are followed and that educators will not give answers or hints to students, provide students with more time to complete the test, or provide students with inappropriate reference materials or tools.

Effective policies should also provide for oversight in the form of statistical analyses of test scores and answer sheets. State policies should require all student answer sheets to be analyzed for unusual patterns of erasures (Qualls, 2001), unusual patterns of responses (Cizek, 2003; Jacob & Levitt, 2003, 2004; Sorensen, 2006; Wesolowsky, 1990), and unusual score fluctuations (Cizek, 1999; Jacob & Levitt, 2003). Bellezza and Bellezza (1989) found that when examinees were made aware that statistical analyses would be used to identify cheaters, the incidence of cheating declined from 5% to 1%. The provision for statistical analyses of student answer sheets could have a similar effect of reducing educator manipulations of test scores.

Table 2.9 Content Recommendations for Test Security Policies

	Policy Recommendations
<p>General</p>	<p>Educators should have input into the development of policy content (Cizek, 1999, 2003; McCabe & Trevino, 2002)</p> <p>Policy content should be clearly worded and signatures should be obtained to ensure information is disseminated and understood (Cizek, 1999, 2003; McCabe & Trevino, 2002)</p> <p>The policy should identify individuals who are in charge of answering questions about policy content (Cizek, 1999; Harcourt Assessment, 2007)</p> <p>The policy should require mandatory reporting of all incidents of manipulations (Cizek, 2003; Gay, 1990)</p> <p>The policy should provide protections for those reporting incidents of manipulations (Cizek, 2003)</p> <p>The policy should outline due process and procedures to investigate and handle incidents of suspected manipulations. Ideally, the policy should call for an independent agency to handle investigations. (Cizek, 1999; McCabe & Trevino, 2002)</p> <p>The policy should outline the sanctions imposed on those found to have manipulated test scores (McCabe & Trevino, 2002)</p> <p>The policy should provide for regular audits of test security (Cizek, 2003; Sorensen, 2006)</p>
<p>Manipulate Teaching Philosophy or Process</p> <p>Making copies of test prior to test administration</p>	<p>The policy should explain copyright laws and penalties for violating copyright laws (Harcourt Assessment, 2007; Iowa Testing Programs, 2005; Riverside Publishing, 2006)</p> <p>The policy should specify who has access to test materials and how to document who handles test materials</p>

	<p>(Cizek, 2003; Sorensen, 2006)</p> <p>The policy should require test materials to be sealed prior to test administration (Cizek, 1999, 2003; Harcourt Assessment, 2007; Iowa Testing Programs, 2005; Sorensen, 2006)</p> <p>The policy should limit the amount of time educators have access to tests before and after administration (Cizek, 2003; Sorensen, 2006)</p> <p>The policy should require the use of multiple test forms. Ideally, new test forms would be administered each year (Cizek, 2003; Crocker, 2003, 2006)</p> <p>The policy should provide examples of specific appropriate and inappropriate test preparation activities (Cizek, 2003; Moore, 1994; Popham, 1991; McCabe & Trevino, 2002)</p> <p>The policy should provide guidelines as to how much instructional time should be spent on test preparation activities (Lai & Waltman, 2007; Moore, 1994; Popham, 1991)</p> <p>The policy should explain the importance of validity and test scores generalizing to a broader domain (Lai & Waltman, 2007; McCabe & Trevino, 2002)</p> <p>Explain the uses of test scores beyond accountability (for example, to make instructional improvements) (McCabe & Trevino, 2002)</p>
<p>Practicing with items identical or similar to the test</p> <p>Practice with last year's (alternate form) test items</p> <p>Practice with items of the same format as the test</p> <p>Use commercial test preparation packages</p> <p>Teaching test-taking skills; test-wiseness</p> <p>Teaching content from specific test items</p> <p>Focusing resources on students closest to proficiency</p> <p>Teaching only content found on the test</p> <p>Changing curricula to better match the test</p>	
<p>Manipulate Examinee Pool</p> <p>Excluding students from testing</p> <p>Having high-scoring students take test multiple times</p> <p>Providing inappropriate special education placement</p> <p>Bribing or paying students to increase test scores</p> <p>Increasing caloric content of school meals</p> <p>Manipulate Test Administration</p>	<p>No Child Left Behind test participation requirements should be explained, including the testing of disabled students and English language learners. (Cullen & Reback, 2006; Figlio & Getzler, 2002, Jacob, 2007)</p> <p>The policy should explain why accommodations are used in test administration and provide examples of appropriate and inappropriate testing accommodations (or make reference to materials that provide these examples). (McCabe & Trevino, 2002)</p> <p>The policy should provide examples of appropriate and inappropriate school- or classroom-level activities on the day of testing. (McCabe & Trevino, 2002)</p>

<p>Altering a student’s answer sheet (changing answers) Giving students answers Checking or pointing out incorrect answers Giving students (non)verbal hints on test items Not following test administration procedures exactly Allowing students to work together during testing Ignoring students who are cheating Giving students additional examples Providing students extra time Rephrasing test items for students Reading items that are to be read by students Answering questions about test content Providing students with reference materials or tools Having students fill-in unanswered items Providing inappropriate accommodations to students</p>	<p>The policy should provide for monitoring of the test administration. Ideally, independent monitors would be used to oversee test administration in a randomly selected sample of classrooms or schools. (Cizek, 2003)</p> <p>The policy should require all test proctors to be trained. (Lai & Waltman, 2003; Cizek, 2003)</p> <p>The policy should provide specific examples of appropriate and inappropriate test administration behaviors (or make reference to materials that provide examples), including how to respond to student questions and what materials are allowed during testing. (Cizek, 2003; McCabe & Trevino, 2002)</p> <p>The policy should explain the importance of standardization and following test administration procedures. (McCabe & Trevino, 2002)</p> <p>The policy should provide for statistical analyses of answer sheets to check for unusual erasure patterns, response patterns, or score fluctuations. (Bellezza & Bellezza, 1989; Cizek, 2003; Jacob & Levitt, 2003, 2004; Qualls, 2001; Sorensen, 2006; Wesolowsky, 1990)</p> <p>The policy should provide specific examples of how to handle testing irregularities, including how to clean student answer sheets following testing (Cizek, 2003; McCabe & Trevino, 2002)</p>
<p>Sanitizing answer sheets (cleaning before scoring) Review skills that will be on tomorrow’s test</p> <p>Manipulate Score Reports or Standards</p> <p>Changing student test scores on official records Providing false IDs so scores won’t count</p> <p>Misrepresenting data Changing criteria for proficiency, making test easier</p>	<p>The policy should outline procedures to be followed if test scores are suspected of being incorrect, including procedures for re-testing. (Cizek, 2003)</p> <p>The policy should provide a system (barcodes, for example) to ensure accurate student information is matched to test scores (Cizek, 2003)</p> <p>The policy should provide examples of appropriate and inappropriate interpretations and uses of test scores (McCabe & Trevino, 2002)</p>

Table 2.10 FOIL Framework for Evaluating Test Security Policy Content

<p>Formalize beliefs of state educators regarding the role of testing and practices</p> <ul style="list-style-type: none"> • Educators should have input into the development of the content of the written state policy • Policy content should be clearly worded and signatures should be obtained to ensure information is disseminated and understood • The policy should identify individuals who are in charge of answering questions about policy content • The policy should require mandatory reporting of all incidents of manipulations • The policy should provide protections for those reporting incidents of manipulations • The policy should outline due process and procedures to investigate and handle incidents of suspected manipulations • The policy should outline the sanctions imposed on those found to have manipulated test scores • The policy should outline procedures to be followed if test scores are suspected of being incorrect, including procedures for re-testing • The policy should provide a system (barcodes, for example) to ensure accurate student information is matched to test scores • The policy should provide specific examples of how to handle testing irregularities, including how to clean student answer sheets following testing
<hr/> <p>Oversee test preparation, administration, and scoring activities</p> <ul style="list-style-type: none"> • The policy should provide for regular audits of test security • The policy should provide for monitoring of the test administration. Ideally, independent monitors would be used to oversee test administration in a randomly selected sample of classrooms or schools • The policy should provide for statistical analyses of answer sheets to check for unusual erasure patterns, response patterns, or score fluctuations.
<hr/> <p>Inform educators about why some behaviors and activities are unacceptable</p> <ul style="list-style-type: none"> • The policy should explain copyright laws and penalties for violating copyright laws • The policy should provide examples of specific appropriate and inappropriate test preparation activities • The policy should provide guidance as to how much time should be spent on test preparation activities • The policy should explain the importance of validity and test scores generalizing to a broader domain • Explain the uses of test scores beyond accountability (for example, to make instructional improvements) • No Child Left Behind test participation requirements should be explained, including the testing of disabled students and English language learners • The policy should explain why accommodations are used in test administration and provide examples of appropriate and inappropriate accommodations (or make reference to materials that provide these examples) • The policy should provide examples of appropriate and inappropriate school- or classroom-level activities on the day of testing • The policy should require all test proctors to be trained • The policy should provide specific examples of appropriate and inappropriate test administration behaviors (or make reference to materials that provide examples), including how to respond to student questions and what materials are allowed during testing • The policy should explain the importance of standardization and following test administration procedures • The policy should provide examples of appropriate and inappropriate interpretations and uses of test scores
<hr/> <p>Limit opportunities for educators to manipulate test scores</p> <ul style="list-style-type: none"> • The policy should specify who has access to test materials and how to document who handles test materials • The policy should require test materials to be sealed prior to test administration • The policy should limit the amount of time educators have access to tests before and after administration • The policy should require the use of multiple test forms. Ideally, new test forms would be administered each year <hr/>

The third way in which state test security policies can deter manipulations is by informing educators as to what behaviors are acceptable and why other behaviors or activities are unacceptable. Recall that a major reason why educators manipulate test scores is that they are simply unaware as to what behaviors or activities constitute manipulations. Policies that explain copyright laws, NCLB requirements, the role of test preparation, the purposes of testing, the uses of test scores, the importance of validity, and the importance of standardized test administration procedures may help reduce test score manipulations.

McCabe and Trevino (1993, 2002) found that when schools informed students about the importance of testing and the seriousness of cheating through honor codes, student cheating reduced. They found that, in order to be effective, these honor codes must also provide specific examples of appropriate and inappropriate behaviors. Likewise, state test security policies should provide specific examples of appropriate and inappropriate test preparation, administration, and scoring behaviors and activities. If state policies were developed with input from educators, then these policies would represent state educators' collective beliefs as to the appropriateness of each testing activity. By codifying these collective beliefs, the state policies would be more effective in reducing manipulations than the overabundance of professional guidelines, standards, and codes currently available to educators.

State policies should also inform educators by requiring all test proctors to be trained regularly. Cizek (2003) noted:

Too often, the qualifications for proctoring exams are only faintly spelled out, the training provided is minimal if any, and no incentives exist to heighten proctors' vigilance or pursuit of instances of cheating. Proper training must include instruction on methods examinees use to cheat and effective procedures for documenting on-site testing irregularities (pp. 380-381).

Training would ensure that educators are aware of which behaviors are unacceptable and the sanctions they will face if they manipulate test scores. This combination of

information about the importance of testing, specific examples of appropriate and inappropriate activities, and sanctions was found to be effective in reducing student cheating at the collegiate level (McCabe, Trevino, Butterfield, 2001).

The fourth way in which state test security policies can deter manipulations is by limiting educators' opportunities to manipulate scores. This can be done by limiting access to test materials and by administering new test forms annually. As many test publishers recommend or require (Harcourt Assessment, 2007; Iowa Testing Programs, 2005; Riverside Publishing, 2006) test materials should be sealed prior to administration. Policies should also outline the handling of test materials to limit the amount of time educators have access to materials before and after the test administration. Sorensen (2006) and Cizek (2003) both recommend that states specify who has access to materials and document anyone who has been given access. This would deter educators from manipulating the teaching process or philosophy through inappropriate test administration and from manipulating the test administration by changing student answers. Cizek (2003) further recommends that states administer new test forms annually (or that test disclosure laws be revised) in order to prevent educators from manipulating test scores through inappropriate test preparation activities.

Relationship Between Test Security Policies and Score

Trend Discrepancies

In order to determine if test security policies are effective in deterring educators from manipulating test scores, the impact of test score manipulations must first be estimated. Unfortunately, it is extremely difficult to do this. While the retesting experiments conducted in Chicago (Jacob and Levitt, 2003, 2004; Perlman, 1985) provide evidence of the impact of manipulations on test scores for individual students or classrooms, they do not provide evidence that manipulations significantly impact test scores on a district or state level.

In order to even begin to address the impact of test score manipulations at the state-level, the definition of *manipulation* must again be considered. Recall that the term *manipulation* is defined as any practice used by educators to increase test scores without an equal, corresponding increase in student performance on the underlying construct. This suggests that in order to provide evidence of the impact manipulations, state test scores can be compared to a different measure of the same construct. If scores from the state test align with scores from the other measure, then one could conclude that manipulations have very little impact on state test results. If, on the other hand, a large discrepancy exists between state test scores and scores on the other measure, then one possible explanation for this discrepancy could be that educators manipulated state test scores.

As will be discussed later, this conclusion would be only one of many possible explanations for the discrepancy between two measures of the same construct. When a large discrepancy exists between two measures of the same construct, all alternative plausible rival hypotheses should be ruled-out before drawing a conclusion of causality (Koretz, 1991). This study will only attempt to find if a relationship exists between the quality of a state's test security policy and discrepancies in score trends as measured by two tests.

Due to the requirements of NCLB, every state has already implemented an accountability system using a state test. These state tests all measure student performance in, at least, reading and mathematics and provide a percentage of students who score at or above a proficient level in each subject. In order to attempt to estimate the impact of manipulations on state level test scores, researchers must choose an appropriate second measure of the same constructs of reading and mathematics. In addition to measuring the same construct, this second measure must be designed, administered, and scored so that its scores cannot be manipulated in the same way as state

test scores. An obvious choice for this second measure would seem to be the National Assessment of Educational Progress (NAEP).

NAEP is a congressionally mandated assessment administered by the National Center for Education Statistics (NCES) ([Chadwick, 2006](#)). While states and school districts are required to participate in the testing, not all students are tested on the NAEP each year. As part of the State NAEP program, representative samples of students from grades 4 and 8 are selected from each state to take the test. Instead of testing each student in reading and mathematics, each student is administered a portion of the entire test. These results are then combined to provide average scale scores and percentages of students meeting basic, proficient, and advanced performance standards at the state level.

NAEP is designed and administered in a way that has made it potentially more robust against educator manipulations. This means that it can provide a good comparison to the results from state tests that can be subject to manipulations. First, results from the NAEP are not used to make high-stakes decisions regarding the performance of an individual educator, school, or school district. Because of this, educators should feel no pressure to manipulate test scores on the NAEP. Second, although some items are publicly released after testing, educators do not have access to the items on the NAEP before it is administered. This virtually eliminates the possibility that educators will manipulate the teaching process or philosophy through inappropriate practice or coaching to inflate NAEP scores. Third, by forcing make-up testing for classrooms with less than 90% attendance and by comparing sample demographics to state demographics, the NAEP provides some level of protection against manipulations of the examinee pool ([NCES, 2007](#)). Finally, because the U.S. Department of Education hires staff to administer the NAEP and classroom teachers can monitor the administration, educators would have difficulty manipulating the test administration in order to inflate NAEP scores ([Massachusetts Department of Education, 2007](#)). Thus, the NAEP provides results that can be more robust against educator manipulation.

Single-Year Comparisons of State and NAEP Results

Because manipulations may have a smaller impact on NAEP scores than state test scores, researchers have developed methods to compare the results from state and NAEP tests. One simple method involves researchers making single-year comparisons of state and NAEP test proficiency results. Studies that have used this method include research from the Civil Rights Project at Harvard University (Lee, 2006), The Education Trust (Hall & Kennedy, 2006), and The Hoover Institution (Peterson & Hess, 2005, 2006).

The logic behind these studies is this: (a) if the state tests and NAEP measure the same constructs of reading and math, and (b) if the state tests and NAEP both define *proficiency* in these constructs, and (c) if the definitions of proficiency are similar, then (d) the percentages of proficient students provided by both tests should be similar. If the reported proficiency rates from the two tests do not provide similar results in a given year, then that discrepancy provides possible evidence that the state tests, which are more susceptible to manipulation, may have been inflated through manipulation. Again, these conclusions are based on strong assumptions that will be discussed later.

Table 2.11 displays an example of the results from these single-year analyses. The second column of Table 2.11 displays the percentage of 8th grade students who scored at or above proficient in mathematics during the 2005 administration of the state test. The third column shows the percentage scoring at or above proficient on the 8th grade NAEP test. As the table shows, the proficiency rates obtained from the state tests are higher than proficiency rates from the NAEP for 46 states. The median state percentage of students scoring at or above a proficient level on state tests was 62% in 2005. The state median percentage of students scoring proficient or above on the NAEP was 30%. Therefore, the median proficiency rate reported from state tests is 2.07 times larger than the median proficiency rate reported from the NAEP. The fifth column of the table displays this information for each state.

Table 2.11: Results from 2005 state and NAEP tests of 8th grade mathematics

State	% at or above proficient on...		% at or above basic on...	Ratio of % proficient or above on state test to...	
	State Test	NAEP	NAEP	% at or above proficient on NAEP	% at or above basic on NAEP
Alabama	63	15	53	4.20	*1.19*
Tennessee	87	21	61	4.14	*1.43*
West Virginia	71	18	60	3.94	*1.18*
Mississippi	53	14	52	3.79	*1.02*
Oklahoma	69	21	63	3.29	*1.10*
Louisiana	51	16	59	3.19	0.86
Georgia	69	23	62	3.00	*1.11*
North Carolina	84	32	72	2.63	*1.17*
Virginia	81	33	75	2.45	*1.08*
Utah	73	30	71	2.43	*1.03*
Arizona	63	26	64	2.42	0.98
Indiana	71	30	74	2.37	0.96
Colorado	75	32	70	2.34	*1.07*
Idaho	70	30	73	2.33	0.96
Nevada	49	21	60	2.33	0.82
Florida	59	26	65	2.27	0.91
Iowa	74	34	75	2.18	0.99
Connecticut	76	35	70	2.17	*1.09*
Alaska	62	29	69	2.14	0.90
Michigan	62	29	68	2.14	0.91
Nebraska	72	35	75	2.06	0.96
Wisconsin	73	36	72	2.03	0.88
Pennsylvania	63	31	76	2.03	0.96
Kansas	68	34	77	2.00	0.88
Texas	61	31	72	1.97	0.85
South Dakota	69	36	80	1.92	0.86
Ohio	63	33	74	1.91	0.85
Oregon	64	34	72	1.88	0.89
North Dakota	65	35	68	1.86	0.79
Illinois	54	29	81	1.86	0.80
New York	56	31	70	1.81	0.80
Minnesota	76	43	79	1.77	0.96
Delaware	53	30	72	1.77	0.74
Montana	63	36	80	1.75	0.79
Maryland	52	30	66	1.73	0.79
New Jersey	62	36	74	1.72	0.84
New Mexico	24	14	53	1.71	0.45
California	37	22	57	1.68	0.65
Rhode Island	39	24	63	1.63	0.62
New Hampshire	56	35	77	1.60	0.73
Vermont	60	38	78	1.58	0.77
Kentucky	36	23	64	1.57	0.56
Arkansas	33	22	64	1.50	0.52
Washington	51	36	75	1.42	0.68
Wyoming	38	29	76	1.31	0.50
Hawaii	20	18	56	1.11	0.36
Maine	29	30	74	*0.97*	0.39
Massachusetts	39	43	80	*0.91*	0.49
South Carolina	23	30	71	*0.77*	0.32
Missouri	16	26	68	*0.62*	0.24
Median	62	30	71	2.07	0.42

Source: Hall & Kennedy, 2006

The table shows 46 states reported larger proficiency rates than what were reported by the NAEP. Alabama reported the greatest discrepancy, with a state test proficiency rate 4.20 times larger than the NAEP proficiency rate. Only four states reported a discrepancy in the opposite direction, with Missouri's state test proficiency rate being only 0.62 times the proficiency rate reported by the NAEP.

Based on similar results from state and NAEP testing in 4th and 8th grade reading and mathematics in 2003 and 2005, researchers have concluded that state test scores are inflated. Reports from The Brookings Institution (Ravitch, 2005), The Civil Rights Project (Lee, 2006), The Education Trust (Hall and Kennedy, 2006), and The Hoover Institution (Peterson and Hess, 2006) all conclude that the single-year discrepancies in proficiency rates between state tests and NAEP are due to states manipulating their score standards. Ravitch (2005) concluded that states lower their proficiency standards "for fear of alienating the public and embarrassing public officials responsible for education" (p. 2). Peterson and Hess (2005, 2006), in using similar data to rate each state's accountability system, concluded that state were "tempted to race to the bottom, lowering expectations to ever lower levels so that fewer schools are identified as failing, even when no gains are being made" (p. 1). Lee (2006) found a positive correlation between the strength of a states high-stakes accountability system and the size of the discrepancy in proficiency rates, concluding that states make tests easier and "water down [their] own performance standards" (p. 51) in order to inflate test scores. Hall and Kennedy (2006) reached a similar conclusion, stating, "most state standards for proficiency are closer to the *basic* level on the NAEP" (p. 19) than they are to the NAEP proficiency level.

Other researchers argue that state proficiency standards are closer to the NAEP basic standards because the tests use different definitions of *proficiency*. In a 2007 report, Idaho's NAEP State Coordinator reviewed the literature and developed six guidelines on the proper use of NAEP scores in confirming results from state tests (Stoneberg, 2007ab). Among the guidelines, Stoneberg noted that state and NAEP

definitions of *proficient* were not the same and that “NAEP’s percentage at or above *basic* is the most directly comparable statistic for confirming state results” (2007a, p. 7). Stoneberg noted that under NCLB, the U.S. Department of Education required states to define proficiency in terms of grade-level expectations (p. 3). A student scoring proficient according to state standards should represent a student who is achieving at or above grade-level expectations. Stoneberg notes that NAEP, on the other hand, does not consider grade-level expectations in defining proficiency. The National Assessment Governing Board printed the following in a booklet designed to inform the public about interpretations of NAEP scores:

Achievement levels define performance, not students. Notice that there is no mention of “at grade level” performance in these achievement goals. In particular, it is important to understand clearly that the Proficient achievement level does not refer to “at grade” performance. Nor is performance at the Proficient level synonymous with “proficiency” in the subject. That is, students who may be considered proficient in a subject, given the common usage of the term, might not satisfy the requirements for performance at the NAEP proficiency level. Further, Basic achievement is more than minimal competency. Basic achievement is less than mastery but more than the lowest level of performance on NAEP. (Loomis & Bourque, 2001).

Further supporting the argument that proficiency rates from state tests should not be compared to proficiency rates from the NAEP, the 2004 NAEP Validity Studies Panel (Mosquin & Chromy, 2004) recommended that “of the various statistics that might be used for measuring a gap on the NAEP scale – proportion at or above the basic, proficient, or advanced achievement level, or mean standardized score – the proportion at or above the *basic* achievement level will both have the greatest correlation with the adequate yearly progress statistic and also be the most directly comparable” (p. 12).

Table 2.11 illustrates the impact of comparing state proficiency to the NAEP basic level instead of the NAEP proficient level. The fifth column shows the ratio of the percentage of students scoring proficient on the state test to the percent of students scoring at or above the basic level of achievement on the NAEP. Whereas 46 states had

proficiency rates higher than NAEP proficiency rates, only 11 states had proficiency rates higher than NAEP basic rates. The number of students scoring proficient on the Tennessee state test, for example, was 1.43 times larger than the number of students scoring at a basic level on the NAEP. Missouri, on the other hand, had a NAEP basic rate 4.25 times larger than the state reported proficiency rate of 16%, implying that state standards have been set higher than NAEP standards.

Due to the fact that conclusions from these studies on the impact of manipulations on state test scores differ depending on which NAEP standard (basic or proficient) is used, these single-year state and NAEP comparisons are limited. In an attempt to address this limitation, researchers from the U.S. Department of Education (USDE) (Braun & Qian, 2007), the American Institutes of Research (AIR) (McLaughlin et al., 2000, 2002) have developed another group of single-year comparison methods to compare state test and NAEP results. These methods, which were developed through a series of 12 studies beginning in 1993, involve linking state scores or proficiency standards onto the NAEP score scale (Buckley, 2007). Kolen and Brennan (2004) provide detailed descriptions of the linking methods used in those 12 studies.

The goal of these linking methods is to put scores from each state test on the same NAEP scale. While each method is unique, the methods used by the AIR the USDE both involve a three-step process. In this process, as described by Ho and Haertel (2006b), researchers first examine state test scores from the sample of students and schools that were administered the NAEP. In the second step, the researchers calculate the percentage of these students who are proficient or above on the state test. The final step is to find the NAEP cut score that sets the same percentage of students as proficient. This NAEP cut score then represents the state proficiency standard mapped onto the NAEP scale.

Ho and Haertel (2006) note that, “all else being equal, states that report greater percents proficient will have lower mapped standards” (pp. 2-3) and that, “higher scoring NAEP states will have higher mapped standards” (p. 3). The researchers concluded that,

“The mapping essentially penalizes state performance standards for reporting high percents of proficient students without commensurately high NAEP performance” (p. 3). Thus, mapped proficiency standards that are relatively low on the NAEP score scale may represent states that have manipulated their scoring standards (lowering the standard for proficiency or making the test easier) in order to inflate their test scores.

The AIR (McLaughlin et al, 2000) and USDE (Braun & Qian, 2007) methods both found large differences among the proficiency standards used by states. In employing their method to analyze 2005 data, Braun and Qian (2007) found that state proficiency standards varied widely. For grades 4 and 8 in reading and mathematics, the state proficiency standards spanned 60 to 80 score points on the 500-point NAEP scale. The researchers also found:

a strong negative correlation between the proportions of students meeting the states’ proficiency standards and the NAEP score equivalents to those standards, suggesting that the observed heterogeneity in states’ reported percents proficient can be largely attributed to differences in the stringency of their standards (p. iii)

Thus, if these linking methods provide valid results, it appears as though states reporting higher percentages of proficient students are manipulating score standards in order to inflate their scores.

Stoneberg (2007) provides a guideline that suggests that *none* of the single-year comparisons described in this section should be used to estimate the impact of manipulations. Stoneberg noted that in 2002, an Ad Hoc committee from the National Assessment Governing Board (NAGB) recommended that comparisons between state test and NAEP results, “should not be conducted on a ‘point-by-point’ [single-year] basis” (p. 6) because of the potential impact of the differences between state and NAEP testing programs (to be discussed later). Because of these flaws, single-year methods are not ideal methods to use to compare state and NAEP results.s

Trend Comparisons of State and NAEP Results

As an alternative to single-year comparisons, an Ad Hoc Committee convened by the NAGB recommended that NAEP achievement levels be used as evidence to confirm *trends* in state test scores in 4th and 8th grade reading and mathematics (Ad Hoc Committee, 2002). The National Academy of the Sciences also recommended comparing trends in scores on state and NAEP tests, suggesting that comparisons should focus on changes in the percentages of students scoring proficient rather than focusing on results from a single year (Pellegrino, Jones, & Mitchell, 1998). As will be discussed, discrepancies between state and NAEP results from a single year can be influenced by differences in test content and the motivational levels of examinees. Some researchers believe these differences are not as problematic when comparing score trends. Linn, Baker, & Betebenner (2002) noted that, “Despite differences in the stakes attached to the results of state tests and measures such as NAEP in content coverage, it is relevant to ask the degree to which gains on a state test generalize to gains on other measures of achievement” (p. 6). Klein, Hamilton, McCaffrey, and Stecher (2000) also believed that trend comparisons addressed some of the problems with single-year comparisons, noting that “any reduction in student effort or performance that may stem from NAEP being a relatively low-stakes test should be fairly consistent over time and therefore not bias our measurement of score improvements across years” (p. 4). Linn (2000) further justified the trend comparison method, noting that, “Divergence of trends does not prove that NAEP is right and the state assessment is misleading, but it does raise important questions about the generalizability of gains reported on a state’s own assessment, and hence about the validity of claims regarding student achievement” (p. 14).

Reports from The Thomas B. Fordham Foundation (2005), RAND (Klein, Hamilton, McCaffrey, & Stecher, 2000), and other researchers (Koretz & Barron, 1998; Linn, 2000; Linn, Baker, & Betebenner, 2002) have employed this method of trend comparison to estimate the impact of test score manipulations on state test results. In

comparing scores on the Kentucky Instructional Results Information System (KIRIS) to NAEP results, Koretz and Barron (1998) found that the large KIRIS gains reported for fourth grade reading and mathematics from 1992 to 1994 were more than four times larger than the gains in NAEP results over the same time period. The researchers concluded that the large KIRIS gains were due to teachers manipulating the teaching process to teach only the content of previous tests.

Klein, Hamilton, McCaffrey, and Stecher (2000), on behalf of RAND, and Linn, Baker, and Betebenner (2002) similarly compared trends from the Texas Assessment of Academic Skills (TAAS) to the NAEP. Both sets of researchers found that score gains reported from the TAAS were significantly larger than the gains reported from the NAEP. Figure 2.3 illustrates the findings. Figure 2.3a displays the score trends on the TAAS compared to trends on the NAEP in 8th grade mathematics from 1990 until 2001. The slopes of the test scores over time represents test score trends. Whether using the NAEP proficient or basic standards, Figure 2.3a shows that trends in TAAS pass rates outpace the NAEP trends. As a counterexample, Figure 2.3b shows a similar comparison between NAEP results and scores from the Maryland School Performance Assessment Program (MSPAP) in 8th grade mathematics. Maryland experienced trends similar to Texas on the NAEP over this time period. While TAAS results showed much greater growth than the NAEP, the trends from the MSPAP appear to support NAEP trends. Klein, Hamilton, McCaffrey, and Stecher (2000) concluded that the discrepancy in score trends in Texas could be attributed to:

- (1) students being coached to develop skills that are unique to the specific types of questions that are asked on the statewide exam (i.e., as distinct from what is generally meant by reading, math, or the other subjects tested);
- (2) narrowing the curriculum to improve scores on the state exam at the expense of other important skills and subjects that are not tested;
- (3) an increase in the prevalence of activities that substantially reduce the validity of the scores

In other words, discrepancies in score trends could be due to manipulations of the teaching philosophy or process or manipulations of the test administration.

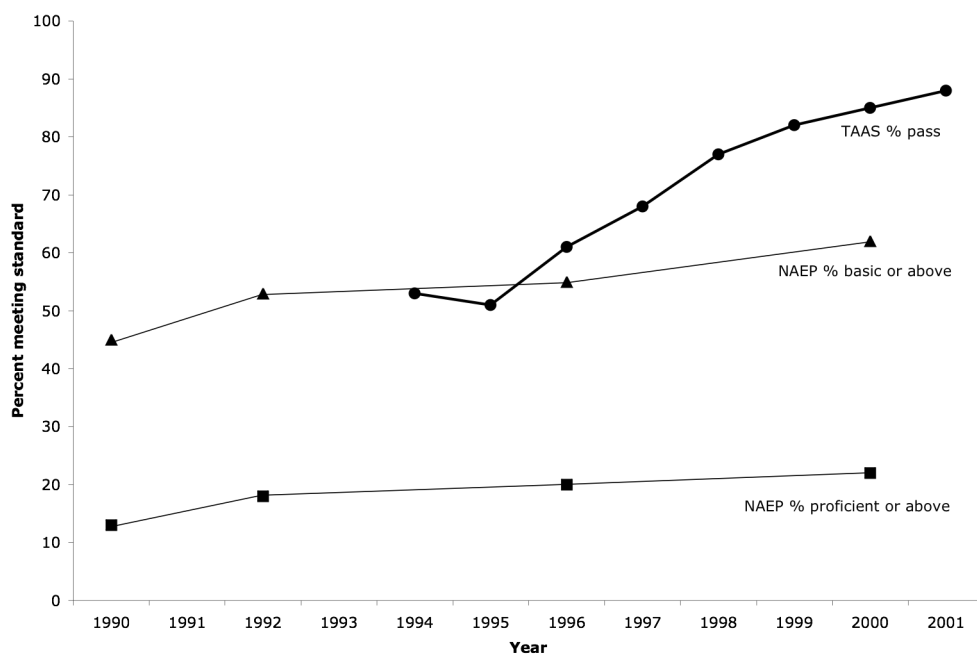


Figure 2.3a Score trends in 8th grade mathematics measured by the TAAS and NAEP.

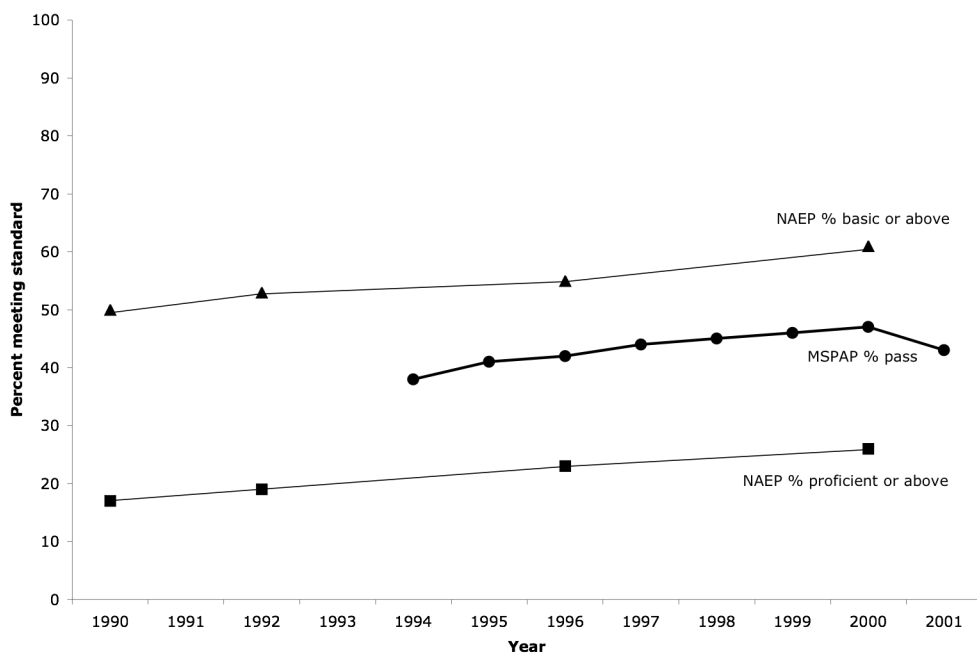


Figure 2.3b Score trends in 8th grade mathematics measured by the MSPAP and NAEP.

Source: Linn, Baker, & Betebenner (2002)

Jacob (2007) extended the research into discrepancies between score trends on the TAAS and NAEP. Jacob found that between 1996 and 2000, TAAS scores increased at a much higher rate than NAEP scores. He found, for example, that math performance increased by more than 0.5 standard deviations on the TAAS compared to a 0.1 standard deviation increase on the NAEP. Using item-level data, Jacob concluded that the discrepancy in score trends could not be explained by changes in the demographic composition of examinees or differences in test item formats. The fact that the NAEP is administered with a time limit whereas the TAAS is not timed also did not explain the differences in score trends. After conducting the analysis, Jacob expressed concern about the generalizability of student achievement gains under state accountability systems.

The Thomas B. Fordham Foundation (2005) conducted an analysis to compare trends on other state tests to the NAEP. Tables 2.12 and 2.13 illustrate a sample of this data. Table 2.12 shows that of the 29 states with reported score trends in 8th grade reading achievement from 2003-2005, 19 reported increases in the percent scoring proficient on the state tests. Of those 19 states, none experienced an increase in the percent proficient on the NAEP. In fact, 14 of those states experienced declines in proficiency on the NAEP. Table 2.13 shows that only three states experienced trends of the same direction in state test and NAEP results. Similar results in reading for grades 4 and 8 led the researchers to conclude that manipulations of the scoring standards inflated state test results. Fordham Foundation president Chester E. Finn, Jr. stated, "If states ease their standards, construct simple-minded tests, or set low passing scores, they can mislead their own citizens and educators into thinking that just about everyone is proficient" (p. 1).

Table 2.12: Discrepant trends in 8th grade reading achievement (2003-2005)

State	Change in % proficient from 2003-2005 on...		Change in % basic from 2003-2005 on...	Change in % proficient on state test is greater than change in % ___ on NAEP	
	State Test	NAEP	NAEP	Proficient	Basic
Alabama	11	0	-2	*	*
California	9	-1	-1	*	*
Idaho	9	0	0	*	*
Arizona	8	-2	-1	*	*
Delaware	8	-1	3	*	*
Tennessee	8	0	2	*	*
Maryland	6	-1	-2	*	*
Virginia	6	0	-1	*	*
Kentucky	5	-3	-3	*	*
Indiana	3	-5	-4	*	*
Iowa	3	-2	0	*	*
New York	3	-2	0	*	*
North Dakota	3	-1	2	*	*
Oregon	3	0	-1	*	*
Georgia	2	-1	-2	*	*
North Carolina	2	-2	-3	*	*
Oklahoma	2	-5	-2	*	*
Missouri	1	-3	-3	*	*
South Dakota	1	-4	0	*	*
Colorado	0	-4	-3	*	*
Mississippi	0	-4	-5	*	*
Wisconsin	0	-2	0		*
Wyoming	0	2	2		
Hawaii	-1	-4	-3	*	*
Maine	-1	0	2		
Connecticut	-2	-4	-3	*	*
Florida	-5	-2	-2		
Texas	-5	0	-2		
Montana	-6	0	0		
Median	2	-2	-1		

Source: Thomas B. Fordham Foundation, 2005

Table 2.13: Discrepant trends in 8th grade reading achievement (2003-2005)

		2003-2005 State Test Trend		
		Decline in % proficient	No change in % proficient	Growth in % proficient
2003-2005 NAEP Trend	Decline in % proficient	Hawaii Florida Connecticut	Colorado Wisconsin Mississippi	California Arizona Delaware Indiana Maryland Iowa Kentucky S Dakota Georgia Oklahoma Missouri N Carolina New York N Dakota
	No change in % proficient	Montana Maine Texas		Alabama Idaho Oregon Tennessee Virginia
	Growth in % proficient		Wyoming	

Unfortunately, these trends comparison methods are limited due to both substantive and technical issues. In explaining the technical issues, Ho (2007) describes “the act of comparing state and NAEP results as the act of comparing the height of two children on pogo sticks” (p. 2). When researchers measure trends in the percentage of students scoring above a cut-score, the magnitude and sign of those Percent Above Cut (PAC-based) trends depend on the selection of cut-score. As Ho explains:

The interpretive problems with PAC-based statistics may be simply explained by their interaction with unimodal distributions. If a unimodal distribution of test scores shifts in the positive direction, the rate at which examinees cross a cut-score will not be constant. As the mode of the distribution approaches the cut-score, more and more examinees will cross in equal units of time. After the mode of the distribution passes the cut-score, fewer and fewer examinees will cross in equal units of time. If the cut-score were different under this model, the trend would be different. In this sense, PAC-based trends may be described as *pliable* under the choice of cut-score. (p. 4)

Ho goes on to demonstrate the pliability in PAC-based trends for state 4th grade reading results on the NAEP from 2003 to 2005 by calculating these PAC-based trends from the Basic, Proficient, and Advanced cut-scores on the NAEP. He finds, for example, that Arizona experienced a 1% gain in students scoring above the NAEP Advanced cut-score from 2003 to 2005. Using a different cut-score, Arizona experienced a 2% decline in students scoring above the NAEP *Basic* cut-score. Other states showed similar results in that the choice of cut-score changes the magnitude and sign of the trend in students scoring above that cut-score. Because conclusions from trend comparisons change depending on the selection of cut-score, these trend comparisons should be interpreted cautiously.

Scale-Invariant Trend Comparison Methods

To address the pliability of PAC-based trends, Ho (2007) introduced a scale-invariant trend statistic based on the Probability-Probability (PP) plot of score distributions from a test given at two times. The V statistic (Ho & Haertel, 2005) is

described as a scale-neutral effect size – a measure of the change in test scores from one time to the next that does not change depending on the selection of a cut-score. A detailed explanation of the V statistic will be provided in the next chapter.

Ho (2007) estimated the V statistic for 82 combinations of test results from 4th and 8th grade state and NAEP reading and mathematics tests from 2003 and 2005. The researcher found that the average trend in state test scores was significantly more positive than the average trend in NAEP scores, with 76% of the state trends being more positive than NAEP trends. After cautioning readers that these findings could be influenced by content differences, examinee motivation, examinee sampling, or other reasons, Ho concluded, “These results are consistent with the hypothesis that increased attention to state test content leads to improved performance on state tests but not on NAEP” (p. 13).

State and NAEP Trend Discrepancies: Plausible Rival

Hypotheses

While researchers have concluded that manipulations may have caused discrepancies between state and NAEP results (Hall & Kennedy, 2006; Jacob, 2007; Kleine, Hamilton, McCaffrey, & Stecher, 2000; Koretz & Barron, 1998; Lee, 2006; Peterson & Hess, 2005, 2006; Ravitch, 2005), it must be noted that the existence of discrepancies does not prove the existence or impact of test score manipulations. Jacob (2007) notes, “... there has been little research on reasons *why* student performance differs between NAEP and local assessments” (p. 11). This research is important because, as Koretz (1991) states, in order to conclude that a discrepancy between state and NAEP results, “reflects specific policies or practices, one needs to be able to reject with reasonable confidence other plausible explanations...” (p. 20).

Hill (1998), Ho and Haertel (2007), the Iowa Department of Education (2007), Jacob (2007), and Koretz (1999) all address plausible rival hypotheses that may explain

any discrepancies between state and NAEP results. Synthesizing this research, some of these plausible rival hypotheses include:

- Differences in content coverage or sequence or opportunity to learn
- Differences in item formats or administration mode (paper- or computer-based)
- Differences in test difficulty
- Differences in score standards or standard-setting procedures
- Differences in test administration procedures/environment or administration date
- Differences in accommodations allowed during testing
- Differences in examinee populations or subgroup definitions
- Differences in examinee motivation or effort

The first four plausible rival hypotheses address differences between state and NAEP tests and scoring procedures. If a state test differs from NAEP in content coverage or sequence, then it would be expected that students would score higher on the state test (due to educators focusing on state content standards). Likewise, differences in item formats, test administration mode, or test difficulty may have a significant impact on score discrepancies between state tests and NAEP. Also, as was discussed previously, score standards may also impact state-NAEP discrepancies.

The next two rival hypotheses for score discrepancies address differences in test administration procedures. If state test administration procedures significantly differ from NAEP procedures (in terms of testing time, use of accommodations, or use of materials such as calculators during testing), then discrepancies in results between the two tests would not be completely unexpected.

The final three possible explanations for score discrepancies deal with potential differences in the examinees being tested under state tests and NAEP. While NCLB requires at least 95% of students to be tested annually and NAEP sets its standard at 85% (Hill, 1998, p. 3), this means that up to 20% of examinees could have been excluded from at least one of the tests. Clearly, these potential differences in the examinee pool could

impact discrepancies between results from the two tests. Also, since state tests under NCLB are high-stakes and NAEP remains a relatively low-stakes test, differences in examinee motivation or effort could have an impact on discrepancies.

These plausible rival hypotheses are not exhaustive, but they do provide a reminder that discrepancies between state and NAEP score trends do not automatically mean that educators have manipulated test scores. In order to have confidence in a causal relationship, strong assumptions must be made that the discrepancies are not due to the above plausible rival hypotheses. While several studies have concluded that differences in test content (Wei, Shen, Lukoff, Ho, & Haertel, 2006), examinee motivation (Klein, Hamilton, McCaffrey, & Stecher, 2000; Linn, Baker, & Betebenner, 2002), examinee demographics, test item formats, and test administration time limits (Jacob, 2007) cannot explain the discrepancies between state test and NAEP results, the existence of these differences should at least temper expectations about the comparability of state test and NAEP score trends.

Another assumption implicitly made in comparing state and NAEP score trends is that results from NAEP are somehow the “gold standard.” **While NAEP may not be the gold standard, it may be the only available standard with which to compare the performance of states in reading and mathematics achievement.** While NAEP scores have been more robust against educator manipulations, Hill (1998) notes, “As more and more states see the need for increased NAEP scores, practices will evolve that will virtually ensure gains on NAEP” (p. 10). Thus, Hill suggests that if NAEP results are used to validate state test results, NAEP results will become high-stakes and NAEP will become subject to the same manipulations as state tests.

While discrepancies between state and NAEP score trends cannot be attributed to educator manipulations, the relationship between the quality of a state’s test security policy and the magnitude of state-NAEP discrepancies is still interesting. If an inverse relationship is found between policy quality and discrepancy magnitude, future research

could be targeted to determine if those discrepancies could possibly have been caused by manipulations.

Summary

This literature review has shown that educators do manipulate test scores by manipulating the teaching philosophy or process, manipulating the examinee pool, manipulating the test administration, or manipulating score reports or standards. While the exact prevalence of each form of manipulation is unknown, the evidence suggests that reported incidents of manipulations are widespread and increasing.

This literature review has also shown that educators might manipulate test scores because of a lack of effective test security policies at the state level. Using research into student cheating and honor codes along with test security survey results, this literature review suggests that states can foil test score manipulations by formalizing testing beliefs and practices, overseeing test activities, informing educators about what behaviors are appropriate and inappropriate, and limiting opportunities for educators to manipulate test scores.

Finally, while cautioning against causal interpretations and providing a list of some plausible rival hypotheses, this literature review makes the case that the relationship between the quality of a state's test security policy and the magnitude of discrepancies in trends between the state test and NAEP is of interest.

METHODOLOGY

The purpose of this study is to determine if a relationship exists between the existence/quality of state test security policies and discrepancies between state test and NAEP score trends. Specifically, this study attempts to address the following research questions:

1. What kinds of manipulations do educators use to increase test scores? Why do educators manipulate test scores? What is the estimated prevalence of each type of manipulation?
2. What test security policies and practices do states implement in an attempt to deter educators from manipulating test scores? What is the quality of each state's test security policy?
3. What is the relationship between the quality of a state's test security policy and any discrepancies between score trends on state and NAEP tests? Which aspects of a state's test security policy seem to have the strongest relationship with score trend discrepancies? What are some potential explanations for the discrepancies between state test and NAEP score trends?

To address the first set of research questions, studies and news reports were synthesized to develop a taxonomy of test score manipulations. While the prevalence of each manipulation method could only be roughly estimated, the evidence suggests that incidents of manipulations are widespread and growing.

To address the second set of research questions, a FOIL framework was developed to evaluate the existence and quality of four aspects of test security policies. The quality of test security policies serves as the independent variable in this study.

To address the third set of research questions, a scale-invariant framework is used to compare trends between state tests and NAEP scores. Test score discrepancies (the dependent variable) and test security policies are then compared among states to

determine if states with higher quality test security policies experience smaller score trend discrepancies than states with lower-quality test security policies. When possible, comparisons are also made within states that recently adopted new test security policies or significantly modified existing policies to determine the relationship of these policies with score trend discrepancies.

Independent Variable: Test Security Policy Quality

The independent variable in this study is the quality of test security policies implemented by states to deter educators from manipulating test scores. First, information regarding state test security policies was collected and organized. Then, the quality of the test security policy content and implementation was evaluated using the FOIL framework developed in the previous chapter (displayed in Table 2.10). Each state's adopted policy was evaluated holistically and with regards to each of the four aspects of the FOIL framework.

In order to collect information about each state's test security policies, the term *policy* must be defined. According to the Stanford Policy Repository (2007), a policy is:

a statement of *principles* and/or *values* that mandate or *constrain the performance of activities used in achieving institutional goals*. A policy is general in nature, has broad application and helps to *ensure compliance with: applicable laws and regulations*; contract requirements; and *delegation of authority*.... Policies promote operational efficiencies and *reduce institutional risk*. ... Directives, processes, procedures, work instructions, and the like flow from policies... (emphasis added).

Since the institutional goal of a state's public education system, according to the mandates of NCLB, is to increase student achievement as measured by test scores, a state test security policy is a written plan of action to guide educators' decisions and actions in testing. The test security policy guides these decisions and actions to ensure they comply with the state's values, principles, laws, and regulations.

A state test security policy must be contrasted with test administration manuals and the national testing codes and standards outlined in Appendix C. As explained in the

previous chapter, evidence indicates that administration manuals and national standards are ineffective in deterring educators from manipulating test scores. This ineffectiveness may be due, in part, to the overabundance of these materials (which sometimes offer conflicting guidance) and/or to the fact that these materials are often perceived as suggestions rather than mandates.

State test security policies differ from test administration manuals or national standards in that policies are mandates. For this study, a test security policy must be connected with state rules, regulations, laws, or sanctions to clearly demonstrate that the policy is a mandate. Also, for this study, a state test security policy must be issued or endorsed by a state Department of Education (DOE), Board of Education (BOE), or State Legislature (SL). This will differentiate state security policies from guidelines developed by schools, school districts, or local education agencies.

Data Collection and Verification

State test security policy information was collected from publicly available information published on state DOE, BOE, and SL websites. These sites were navigated to find information regarding state assessment programs. Also, these sites were searched for key phrases such as *test security*, *test policy*, *test guidelines*, and *ethics codes*. Any information regarding the state's testing principles, requirements, laws, and regulations was collected. The data were then entered into a standardized form (see Figure 3.1). If data from any state were missing, incomplete, or contradictory, state DOE officials were called for verification.

Sampling

While it was not necessary to collect test security policy information from states excluded from the analysis (see Table 3.6), policy information was collected for all 50 states. When available, information about changes to policies in each of the years 2003, 2005, and 2007 were collected for each state.

State test security policy evaluation form

State name: _____ Test Name: _____

Year policy adopted/modified: _____ Item format: _____

Rate each of the following according to the scale: 0 = Missing 1 = Meets 2 = Exceeds

Formalize beliefs of state educators regarding the role of testing and testing practices

• Prominence / Availability of information

_____ The state has a separate test security office or budget

_____ The state has a separate web page for information about test security

_____ The policy is mentioned in test administration manuals

_____ Percentage of test administration manual pages dedicated to test security = _____%

_____ Number of clicks to navigate from the front page to test security information = _____

• Content

_____ Teachers provided input into content (evidence of a committee or documentation of development)

_____ Clarity of test security policy information

_____ Availability of FAQs regarding test security

_____ The policy requires educator signatures to indicate understanding

_____ Amount of information available (number of documents = _____)

• Implementation

_____ Identifies an individual or group in charge of security: ___ individual or ___ group

_____ Identifies individuals responsible for security at both state and district levels

_____ Evidence of a dissemination plan being followed

_____ Policy content is updated regularly

_____ Policy provides for barcodes (or another system) to automate test score identification

_____ Availability of forms and checklists for districts/schools to use to aid in test security practices

• Requirements and sanctions

_____ Mandatory reporting requirements (for suspected incidents of manipulation)

_____ Provides standard forms (or online reporting) for suspected incidents

_____ Explains the protections for individuals who report suspected incidents

_____ Due process is explained (procedures for investigating suspected incidents)

_____ Sanctions for confirmed cases of manipulation are outlined

_____ Sanctions include suspension and dismissal of confirmed manipulators

• Other

_____ Explains the importance of test security: ___ positive message or ___ negative message

_____ Laws, regulations, rules are mentioned (to indicate that security is a mandate)

_____ N/A Test security policies are to be developed at: ___ state-level or ___ district-level

Oversee test preparation, administration, and scoring activities

• Test Security Audits

_____ Implementation of test security policy is audited regularly (evidence of improvements)

• Test administration oversight

_____ The policy provides for independent monitoring of test administration

_____ Teachers are not to administer the test to their own students

• Statistical Analyses

_____ The policy provides for statistical analyses of answer sheets to detect possible manipulations

_____ The policy provides for erasure analysis

_____ The policy provides for aberrant response analysis

_____ The policy provides for analysis of score fluctuations

_____ Evidence of security analysis reports

• Score Reports

_____ The policy outlines procedures to follow before making any changes to test scores

Inform educators about why some behaviors and activities are unacceptable

• Principles & Rules

_____ The policy explains copyright laws and penalties for violating copyright

_____ The policy explains the importance of validity & generalizing from test scores

_____ The policy explains that all students must be tested (as required by NCLB)

_____ The policy explains the importance of standardized test administration

_____ The policy refers to an honor code or code of ethics

_____ The policy describes the uses of test scores

• Examples of appropriate and inappropriate behaviors

_____ Examples of test preparation activities: _____ appropriate and _____ inappropriate

_____ Examples of test administration activities: _____ appropriate and _____ inappropriate

_____ Examples of accommodations: _____ appropriate and _____ inappropriate

_____ Examples of school/class activities on test day: _____ appropriate and _____ inappropriate

_____ Examples of uses/interpretation of scores: _____ appropriate and _____ inappropriate

• General Guidance

_____ The policy limits the amount of time spent on test preparation activities

_____ The policy specifically states that educators cannot change student answers

_____ The policy specifically states that educators cannot give students hints or answers

_____ The policy specifically states that educators cannot read certain sections aloud to students

<input type="checkbox"/>	The policy lists what materials teachers can or cannot provide to students during testing			
<input type="checkbox"/>	The policy outlines procedures for retesting students or sanitizing answer sheets			
• Training				
<input type="checkbox"/>	The policy provides for regular training of district-level testing coordinators			
<input type="checkbox"/>	The policy provides for regular training of school-level testing coordinators			
<input type="checkbox"/>	The policy provides for regular training of all test proctors			
<input type="checkbox"/>	The quality of training materials available online			
<input type="checkbox"/>	Regularity / amount of training			
Limit opportunities for educators to manipulate test scores				
• Materials security				
<input type="checkbox"/>	The policy specifies who has access to test materials (and at what times)			
<input type="checkbox"/>	The policy specifies the amount of time materials are available			
<input type="checkbox"/>	The policy provides for a tracking system of test materials			
<input type="checkbox"/>	The policy requires test materials to remain sealed until testing			
• Test Forms				
<input type="checkbox"/>	The policy requires new test forms (different test items) to be used annually			
<input type="checkbox"/>	The policy requires multiple test forms (which may be reused)			
Overall categorization of test security policy:				
A) <input type="checkbox"/>	Clear and accessible	vs.	<input type="checkbox"/>	Ambiguous or difficult-to-find
B) <input type="checkbox"/>	State-level mandate	vs.	<input type="checkbox"/>	District- or school-level responsibility
C) <input type="checkbox"/>	Punitive or law-focused	vs.	<input type="checkbox"/>	Instructive/informative
D) <input type="checkbox"/>	Independent monitoring	vs.	<input type="checkbox"/>	No independent monitoring
E) <input type="checkbox"/>	Investigative	vs.	<input type="checkbox"/>	Preventative
F) <input type="checkbox"/>	Example-based	vs.	<input type="checkbox"/>	Not many examples
G) <input type="checkbox"/>	Positive message	vs.	<input type="checkbox"/>	Negative message
Additional Information:				

Figure 3.1: Evaluation form for state test security policies.

Analysis

Once policy information had been collected, the quality of each policy was evaluated according to the FOIL framework. Figure 3.1 shows the form that was used to evaluate the quality of more than 60 features of each state's policy according to a 3-point scale. A score of 0 was given to any state policy that was missing the feature; a score of 1 indicates the state policy contained the feature; a score of 2 indicates the state policy clearly exceeded the majority of states in that feature.

The evaluation form contains additional information about each state's test security policy, including the year in which the policy was developed or modified, and the format of the items on the state test. The year of policy development/modification was used, when possible, for longitudinal analyses of policy effectiveness. Item format information was used to test the hypothesis that states administering tests similar in format to the NAEP experience smaller trend discrepancies between the two tests.

The evaluation form was also used to collect the following information:

- Percentage of pages in state test administration manuals that are focused on test security information
- The number of test security policy documents available online
- Whether the policy identifies an individual or a group of individuals (committee) responsible for test security
- Whether the policy explains the importance of test security via a positive or negative (punitive) message
- Whether the policy is developed at the state- or district-level
- The number of examples of appropriate and inappropriate testing activities

The results for each feature along with this additional information were used to dichotomize state policies in 7 ways:

- Clear and accessible (specific information is readily available) vs. ambiguous or difficult to find

- State-level mandates vs. district- or school-level responsibility (states that require districts to develop test security policies)
- Punitive (focused on sanctions) vs. instructive (focused on informing educators on the importance of test security)
- Policies that require independent monitoring of test administration vs. policies that allow teachers to administer tests to their own students
- Investigative (focused on reporting and investigating potential manipulation incidents) vs. preventative (focused on preventing manipulations)
- Example-based (provides many examples of appropriate and inappropriate behaviors) vs. general (provides general information without specific examples)
- Positive message (provides examples of positive testing behaviors) vs. negative message (provides examples of negative testing behaviors).

A short narrative provides additional information collected from test security policy documents.

Once all state policies have been examined, the data were summarized. One way to summarize this information was to report the percentage of states receiving scores of 0, 1, and 2 for each feature. These scores were also summed to produce composite scores for each section (formalize, oversee, inform, and limit) and subsection (prominence, content, implementation, etc.) displayed in Figure 3.1. Then the mean (or median) of these composite scores were reported for each section and subsection. These composite scores were also displayed for each state. Figure 3.2 displays the composite scores that were calculated.

Based on the actual data collected, the number of independent features identified from state policies might be reduced. For example, if every state with feature X also has feature Y, then those features will be combined. Based on the information collected thus far (712 documents total from all 50 states), it does not appear that this will be a big issue.

State test security policy evaluation form

State name: _____	Test Name: _____			
Year policy adopted/modified: _____	Item format: _____			
Composite Scores				
_____ Formalize beliefs of state educators regarding the role of testing and testing practices				
_____	Prominence / Availability of information			
_____	Content			
_____	Implementation			
_____	Requirements and sanctions			
_____	Other			
_____ Oversee test preparation, administration, and scoring activities				
_____	Test Security Audits			
_____	Test administration oversight			
_____	Statistical Analyses			
_____	Score Reports			
_____ Inform educators about why some behaviors and activities are unacceptable				
_____	Principles & Rules			
_____	Examples of appropriate and inappropriate behaviors			
_____	General Guidance			
_____	Training			
_____ Limit opportunities for educators to manipulate test scores				
_____	Materials security			
_____	Test Forms			
Overall categorization of test security policy: Check one for each row				
A) _____	Clear and accessible	vs.	_____	Ambiguous or difficult-to-find
B) _____	State-level mandate	vs.	_____	District- or school-level responsibility
C) _____	Punitive or law-focused	vs.	_____	Instructive/informative
D) _____	Independent monitoring	vs.	_____	No independent monitoring
E) _____	Investigative	vs.	_____	Preventative
F) _____	Example-based	vs.	_____	Not many examples
G) _____	Positive message	vs.	_____	Negative message

Figure 3.2: Example of composite scores calculated for each state.

Technical Quality of Policy Evaluation Data

As stated in the previous chapter, the components of the test security policy evaluation are based on suggestions from test developers and publishers (Harcourt Assessment, 2007; Iowa Testing Programs, 2005; Riverside Publishing, 2006); research into test preparation activities (Crocker, 2003, 2006; Gay, 1990; Lai & Waltman, 2007;

Moore, 1994; Popham, 1991), research into statistical detection of aberrant responders (Bellezza & Bellezza, 1989; Jacob & Levitt, 2003, 2004; Qualls, 2001; Sorensen, 2006; Wesolowsky, 1990), and research into student and teacher cheating (Cizek, 1999, 2003; Cullen & Reback, 2006; Figlio & Getzler, 2002; Jacob, 2007; McCabe & Trevino, 2002).

Through discussions with a district assessment director from an Iowa public school district, one university assessment coordinator, and two (faculty) experts in educational measurement, the evaluation form was further refined.

Assumptions and Limits

The primary assumption with regards to the independent variables in this study is that the quality of a state's test security policy can be inferred from the quality of the materials and information made publicly available. This might be problematic, especially when trying to evaluate the quality of policy implementation based on available information. Just because a state has an exceptional policy (and materials describing the policy) does not mean that the policy was implemented well. An in-depth case study of a single state (or small group of states) would need to be conducted to evaluate the quality of implementation.

The main limitation with regards to the independent variables is that some information may not be available. It may be that some states lack information about certain aspects of their test security policies. With over 700 policy documents collected from all 50 states, this does not appear to be a big problem. Another limitation is the lack of an interval scale on the data collection form. To create a common scale, section and subsection composite scores were standardized when aggregating the data. A third limitation is the subjective nature of some of the ratings assigned to state policies. While many of the 60-plus features can be objectively scored, features such as the clarity of a policy or the quality of materials to train test proctors were subjectively rated. **To address**

this, the rubric used to rate states in these subjective features was more clearly specified once the data have been collected.

Dependent Variable: Scale-Invariant State Test and NAEP

Score Trend Discrepancies

The discrepancy in trends between state test and NAEP scores is used as the dependent variable in this study. Ideally, traditional effect sizes (mean differences divided by a pooled standard deviation) for both state and NAEP score trends would be compared to estimate the discrepancies. Unfortunately, as Jacob (2007) also discovered, “many states do not even publish state level averages of the underlying raw or scaled score, but rather report student performance in terms of the percent meeting various proficiency levels” (p. 14). Because of this limitation in the data (and because of the limitations in simply comparing proficiency rates between the two tests), a scale-invariant method (to be described later in this chapter) was used to estimate the discrepancies in effect-sizes between state and NAEP score trends.

First, results from state and NAEP testing in 2003, 2005, and 2007 were collected. Then, scale-invariant effect sizes, V_{state} and V_{NAEP} (to be described later in this chapter), were used to measure trends in state and NAEP results from 2003-2005, 2005-2007, and 2003-2007. The simple difference between V_{state} and V_{NAEP} define the state-NAEP score trend discrepancies. Table 3.1 displays the data that was collected and estimated for each state. With 50 states, 2 subjects (reading and mathematics), 2 grade levels, and 3 trends, a maximum of 600 scale-invariant discrepancy estimates could have been collected. The final sample in this study consists of 215 discrepancy estimates (36% of the maximum possible) from 32 states.

Table 3.1 Scale-Invariant Trend Discrepancy Data

			2003-05 score trends	2005-07 score trends	2003-07 score trends
State	Reading	4 th Grade	$V_{state} - V_{NAEP}$	$V_{state} - V_{NAEP}$	$V_{state} - V_{NAEP}$
		8 th Grade	$V_{state} - V_{NAEP}$	$V_{state} - V_{NAEP}$	$V_{state} - V_{NAEP}$
	Mathematics	4 th Grade	$V_{state} - V_{NAEP}$	$V_{state} - V_{NAEP}$	$V_{state} - V_{NAEP}$
		8 th Grade	$V_{state} - V_{NAEP}$	$V_{state} - V_{NAEP}$	$V_{state} - V_{NAEP}$

Data Collection and Verification

Online resources were searched to collect state and NAEP test scores. To collect results from state tests, the website of the Council of Chief State School Officers (CCSSO) was first explored to collect general information about each state's testing program. For each state, the following information for grades 4 and 8 in reading and mathematics for the years 2003, 2005, and 2007 was collected: the name of the test used for accountability requirements of NCLB, the test format (norm-referenced, criterion-referenced, or other), test administration dates (fall or spring), and the number of cut-scores and performance-levels reported. The years 2003, 2005, and 2007 were chosen because they correspond with NAEP administration years. After entering this information into a text file, any changes made to a state's testing program in 2005 or 2007 were flagged.

Figure 3.3 displays an example of the data collected from the CCSSO website. From this figure, one can see that Alabama changed from a norm-referenced to a criterion-referenced test between 2003 and 2005 and that the state testing program classifies students into four performance levels (two of which represent proficient levels of performance). Because Alabama changed tests, this state was flagged to indicate that no valid trend comparisons exist from 2003-2005 or from 2005-2007.

After collecting and entering information from the CCSSO website for all state testing programs, this information was verified by examining the websites of each state's Department of Education (State DOE). From these State DOE websites (and web

searches) additional data about each state's testing program was entered. The general test information collected is displayed in Table 3.2. Once again, any changes made to a state's testing program were documented so that inappropriate trends (trends in scores on different tests or from tests administered at different times in the year) were not computed.

State: Alabama			
Overview: CCSSO Database: http://accountability.ccsso.org/state_profiles.asp			
Test Information			
	2003	2005	2007
4 th Grade Reading	Stanford 10	ART	ART
4 th Grade Math	Stanford 10	AMT	AMT
8 th Grade Reading	Stanford 10	ART	ART
8 th Grade Math	Stanford 10	AMT	AMT

Notes: Stanford 10 was a norm-referenced test administered in April of 2003
 ART = Alabama Reading Test
 AMT = Alabama Mathematics Test
 ART and AMT were administered in April of 2005 and 2007

The state uses 3 cut-scores to classify students into the following categories:

Level I	Level II	Level III	Level IV
(Below proficient)		(Proficient)	

Figure 3.3 Example of data collected from CCSSO.org website.

Table 3.2 General test information collected from each state DOE website

Information	Notes
Test name and format	When both are available, results from <i>reading</i> tests are preferred to <i>English</i> or <i>Language Arts</i> tests. When available, item types (selected-response, constructed-response, or mixed) are recorded. Also, when possible, results from alternate assessments were excluded from analysis.
Test Administration Date	Fall or Spring
Number of Cut-Scores	Cut-scores are used to classify students into performance levels
Proficiency Levels	Which performance levels that represent proficiency
Changes	Any changes in test name, format, administration date, cut-scores, or performance levels are flagged.

Once general information about state testing programs had been collected and entered, state test results were collected. At a minimum, the percentage of students scoring in each performance level in reading and mathematics in grades 4 and 8 in 2003, 2005, and 2007 were collected. With this data, the percentage of students scoring at or above proficient in each subject and grade level in each year were calculated. If this minimum amount of data was not readily available online, calls or emails were made to state DOE officials to request the information. When available online, additional information was collected for each grade and subject, such as the number of students tested each year and the scale score means and standard deviations (to calculate traditional effect sizes). Table 3.3 shows an example of the test score data that was collected for each state.

Table 3.3 Example of test score data to be collected

		Alabama		
		2003	2005	2007
4 th Grade Reading	Test Format	Stanford 10	ART	ART
	Items	NRT	CRT	CRT
	Administration	Constructed	Mixed	Mixed
	Cut-Scores	Spring	Spring	Spring
	Decimals	3	3	3
	# Tested	--	2	2
	Avg. Score	--	N/A	56,083
	Std. Deviation	--	N/A	N/A
	% in PL 1	--	N/A	N/A
	% in PL 2	--	0.31%	0.50%
	% in PL 3	--	16.35%	14.49%
	% in PL 4	--	33.18%	31.86%
	% Proficient	--	50.16%	53.15%
	Adjusted?	--	83.34%	85.0%
	Scaled-up?	--	No	No
	Cat Shift?	1	2	2

Notes: N/A = information not available online

PL = Performance level

Adjusted = Were scores adjusted in any way from what's reported online?

Scaled-up = Were changes made to ensure the percentages sum to 100%?

Cat Shift = Equal values represent years in which trends can be compared

The table shows, once again, that Alabama changed tests between 2003 and 2005. This change means that the 2003 test scores were used to measure trends from 2003-2005 or 2005-2007. Because of this, the 2003 test score data was collected. If a state changed tests after 2003 and again after 2005, then only the 2007 data was collected as no appropriate trends could be measured.

The table also shows, once again, that the Alabama testing program used 3 cut-scores to place students into one of 4 performance levels. The *Decimals* variable indicates with what precision the state reported its test results. In this example, Alabama reported the percentage of students scoring in each performance level with two decimals of precision. The table shows that while the number of students tested in 2005 was not found online, a total of 56,083 4th grade students were administered the ART in 2007. The table also shows that the mean scale scores and standard deviations were not readily available online, but that the percentages of students scoring in each performance level were available. The percentage of students scoring at or above proficient was calculated by adding the percentage of students scoring in the third and fourth performance levels. Finally, the table also shows that scores were not adjusted in any way (from the reports available online) and that trends can only be made from the 2005 and 2007 data.

To aid in data entry, a spreadsheet was developed and state-specific data collection issues were discussed with Educational Measurement and Statistics faculty. When available, the data entered into the spreadsheet were verified by other test score reports available online. The percentages of students scoring in each performance level were also verified by checking to see if the percentages sum to 100%. Some states had percentages that summed to either 99% or 101% due to rounding in the score reports. For these states, the percentages of students scoring within each performance level were divided by the sum of the percentages (scaled-up) to ensure all states had a sum of 100%.

Any changes in the testing program or problems in data entry were flagged for further investigation. As stated earlier, any changes in the test, test administration date,

or cut-scores were flagged so that inappropriate trends would not be computed. When available, the numbers of students tested each year (for subgroups defined by race, disability status, socio-economic status, and English proficiency) were compared to check for any major changes in the testing population from year-to-year. The percentages of students scoring within each performance level and average scale scores were also visually examined to check for unusually large fluctuations.

For the flagged states, other online resources were searched and DOE officials were contacted to determine the explanation for the unusual data. These flagged states were also discussed with faculty during regular data collection meetings. From these meetings, decisions regarding which state test results to include or exclude from the analysis were made.

Sampling

Ideally, the data set would have included the percentage of students scoring within each performance level in grades 4 and 8 in reading and mathematics during the 2003, 2005, and 2007 administrations of the state tests. Due to incomplete or missing data, or due to idiosyncrasies in some states' testing programs, cut scores, or score reports, some state data had to be excluded from analysis. The following rules were used to determine if data should be included in the analysis:

- All pairs of results (the 2003-2005, 2005-2007, and 2003-2007 pairs) were included if:
 - The state tested grades 4 and 8 in reading and mathematics both years.
 - The same test, or parallel forms, was administered both years.
 - Cut-scores and performance levels were not changed in either year.
 - The state administered the test during the same season (fall or spring) each year.

- The population of students tested each year remained relatively stable (no wild fluctuations in the number of students within subgroups tested each year).
- The state uses at least 3 cut-scores to place students into at least 4 non-overlapping performance levels.

Pairs of results were excluded from the analysis if the state changed the test or cut-scores (without equating the new scale to the old). Pairs were also excluded when the test administration date changed (from fall to spring or spring to fall) or when the data indicated that the tested student population had changed in some significant way. These pairs of results were excluded, because the test results may have different meanings in each year and, therefore, no appropriate trend could be computed. The last criteria requiring states to have at least 3 cut-scores and 4 non-overlapping performance levels was used only due to the requirements of the nonparametric estimation procedure that was used to analyze the data (explained later).

NAEP Data Collection

With the state test results collected, NAEP results were then collected from the official NAEP results website, *The Nation's Report Card* (2007ab). The *Trends in Achievement Levels by States* reports display the percentage of students scoring in each of the four NAEP performance levels: below basic, basic, proficient, and advanced. Through these reports, the percent of students scoring in each performance level were collected for each state in reading and mathematics in 2003, 2005, and 2007. The sums of the percentages were calculated as a check of the accuracy of data entry.

Analysis

Once state test and NAEP data had been collected, discrepancies in score trends for the two types of tests were estimated. As previously discussed, inadequacies in many state test score reports did not allow traditional effect sizes to be estimated. Also, as

discussed in the previous chapter, many common methods to calculate discrepancies between state test and NAEP results have technical limitations. Specifically, single-year and trend comparisons of the percentage of students scoring at or above proficient are troublesome because of their pliability under different choices of cut-scores (to be illustrated later). In this study, a scale-invariant framework developed by Ho and Haertel (2006) was used with state test and NAEP results to estimate the discrepancy in score trends on the two tests.

Technical Limitations of Comparing Changes in Percentages of Proficient Students (PPS)

As has been mentioned previously, trend comparisons based on changes in the percentage of students scoring above a cut-score are known to be dependent on the choice of cut-score (Holland, 2002; Ho & Haertel, 2005; Koretz & Hamilton, 2006), which makes them of limited usefulness in comparing state-NAEP score trends. To illustrate this, Figure 3.4 illustrates score distributions from two simulated administrations of the same test. The data were simulated so that from Time 1 to Time 2, the mean score increased from 550 to 600 and the standard deviation decreased from 150 to 100 from the first to the second administration. The data were simulated this way not only to provide a clear example, but also because the goal of NCLB is to both increase achievement for all students (increasing the overall mean) and decrease gaps in student achievement (perhaps decreasing the standard deviation).

Suppose the simulated data in Figure 3.4 come from a test with a cut-score of 500. The figure shows that 63% of students at Time 1 and 84% of students at Time 2 scored above this cut-score of 500. If this cut-score were defined as the proficiency standard, the state producing these results would be lauded for increasing proficiency by 21%. If, instead, a cut-score of 700 defined proficiency, the figure shows that the state would be viewed as ineffective in increasing achievement (26% of students scored above this cut-

score at both Time 1 and Time 2). Using a cut-score of 800 (possibly reflecting higher expectations), a comparison of the percentage of proficient students (PPS) would lead to a conclusion that score trends were actually negative (proficiency dropped from 5% to 2%). Thus, this figure illustrates that the choice of cut-score can impact the conclusions drawn from PPS-based trend statistics.

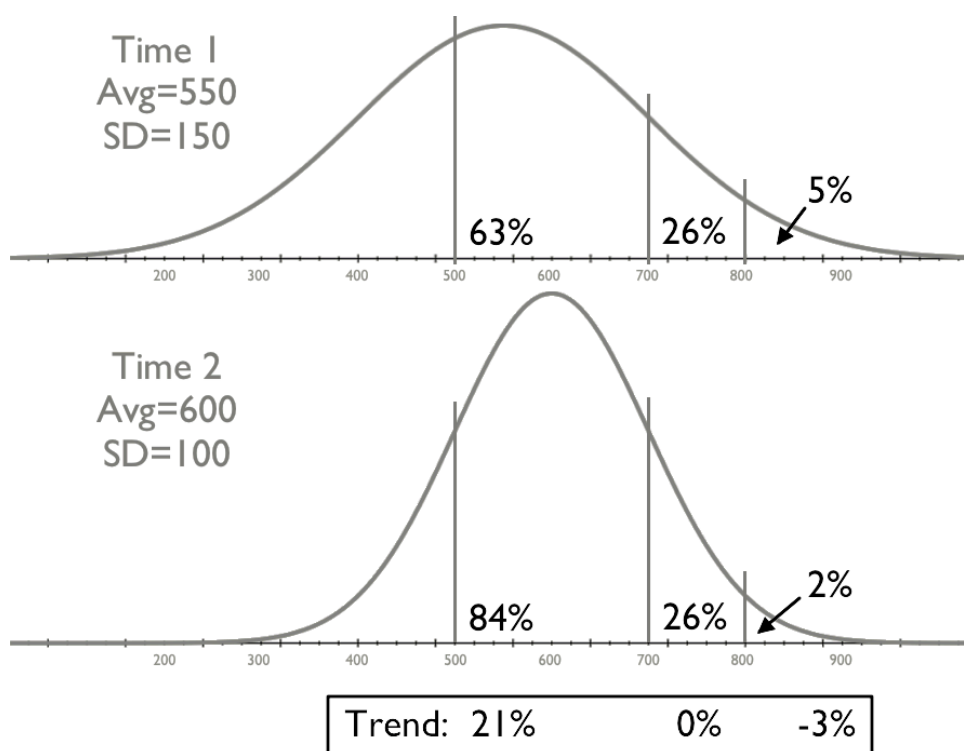


Figure 3.3 Simulated test score distributions to illustrate pliability of PPS-based trends.

Holland (2002) recommends using cumulative distribution functions to display the gap between two test score distributions. As explained by Wilk and Gnanadesikan (1968) CDFs provide a graphical display of a distribution's location, spread, and shape; and CDFs lend themselves to smoothing and interpolation. For a test administered at Time 1 and Time 2, a CDF provides a visual display of both:

$$F_1(x) = \% \text{ of students scoring at or below cut-score } x \text{ at Time 1} \quad (1)$$

and

$$F_2(x) = \% \text{ of students scoring at or below cut-score } x \text{ at Time 2.} \quad (2)$$

Figure 3.4 displays CDFs for the same simulated distributions described earlier.

The gap between the CDFs displays the trend in scores from Time 1 to Time 2. As Holland (2002) explains, this gap could be measured in several ways, most obviously by measuring the gap either horizontally or vertically. In his article, Holland recommended measuring the gaps between CDFs horizontally to represent the difference in percentiles from each distribution. The figure shows, for example, that the 50th percentile at Time 1 is equal to a scale score of 550. The 50th percentile at Time 2 is equal to a scale score of 600. Therefore, this horizontal gap displays a general positive trend in scores from Time 1 to Time 2 (at the 50th percentile).

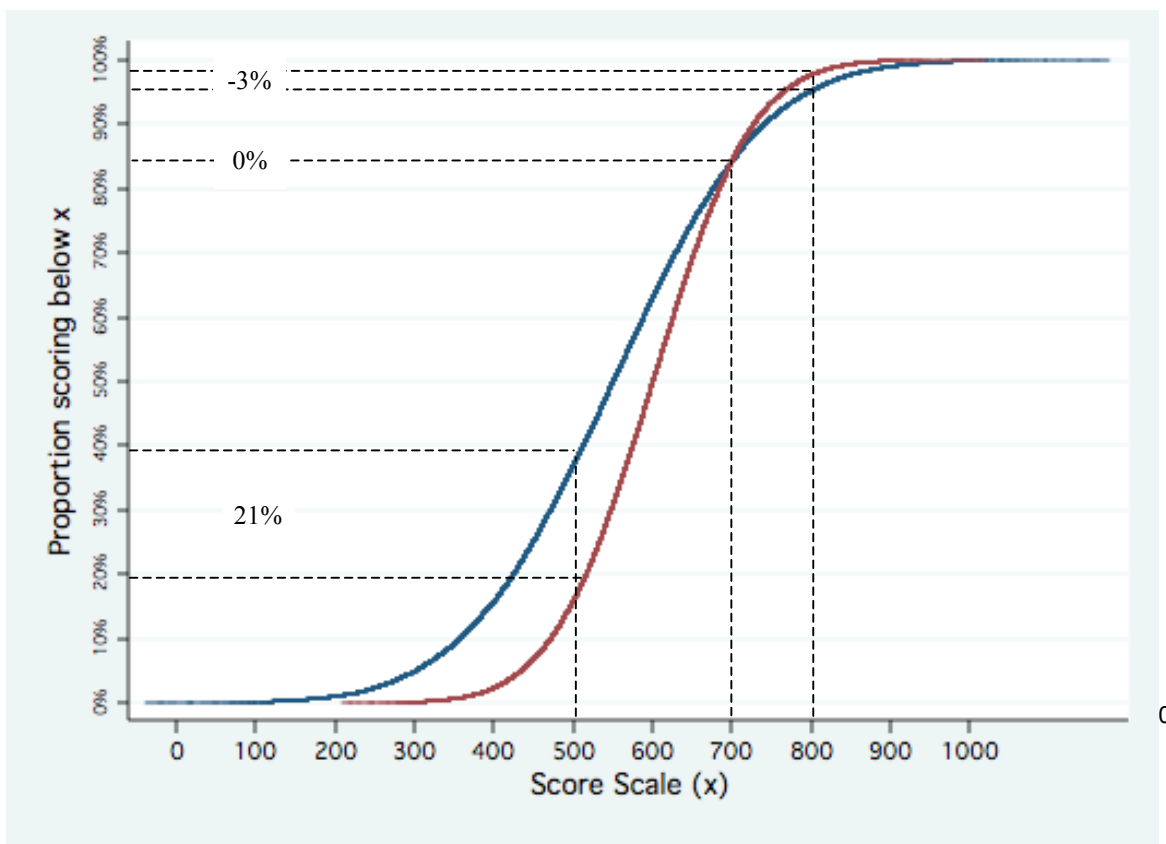


Figure 3.4: CDFs from a test administered at Time 1 and Time 2 with cut-scores of 500, 700, and 800

In the same article, Holland (2002) defended the use of vertical gap measurements in standards-based testing situations by explaining:

The use of cut-scores is common among those interested in “standards-based” assessments. What could be more natural than to measure just how many in a group of examinees meet or exceed the cut-off for a standard with a name such as “Master” or “Proficient”? When educational reform is coupled with standards-based assessment, reformers are naturally led to ask if the percents of certain groups of students meeting or achieving a standard are increasing, possibly as a result of the reforms. From this point of view, the fact that score distributions often change over time by slow but steady shifts upwards is of little interest. The ultimate goal is to implement reforms in such a way as to get as many students as possible to meet or exceed the standards (p. 16).

Since NCLB is focused on standards-based assessment, the vertical gaps between CDFs are of interest in this study. The vertical gaps between the CDFs in Figure 3.4 once again show how the choice of cut-score impacts the conclusions drawn from measuring the change in percentage of students scoring above a cut-score (PPS-based trends).

Scale-Invariant Framework: P-P Plots

The scale-invariant framework developed to address this pliability in PAC-based trends is based on the Probability-Probability (P-P) plot (Haertel, Thrash, & Wiley, 1978; Ho & Haertel, 2006a; Livingston, 2006; Spencer, 1983). A P-P plot is, “a comparative plot of sample cumulative probabilities” (Fisher, 1983, p. 31) “constructed solely from vertical slices across CDFs” (Ho, 2007, p. 13). Thus, P-P plots display the vertical gaps between CDFs of test scores administered at Time 1 and Time 2. As Ho (2007) notes, “a monotone transformation of scale may contort the CDFs horizontally, but will not change the vertical relationships between the cumulative proportions” (p. 13). Thus, P-P plots, and any statistics derived from them, are invariant to transformations of the score scale.

P-P curves, which increase monotonically from the origin to the point (1,1), display the percentiles of one distribution versus the percentiles of another distribution (Holmgren, 1995). When the distributions represent scores from the same test

administered twice, the P-P curve shows the proportion of students scoring at or below a given cut-score at each time. In other words, for a given percent p ,

$$F_2^{-1}(p_2) = \text{the } p^{\text{th}} \text{ percentile from Time 2} \quad (3)$$

(which represents the test score at which $p\%$ of students scored below at Time 2), the P-P plot displays

$$p_1 = F_1[F_2^{-1}(p_2)], \quad (4)$$

the percentage of students scoring at Time 1 scoring below given percentiles of Time 2.

Figure 3.5 displays the P-P plot for the simulated data set from Figures 3.3 and 3.4.

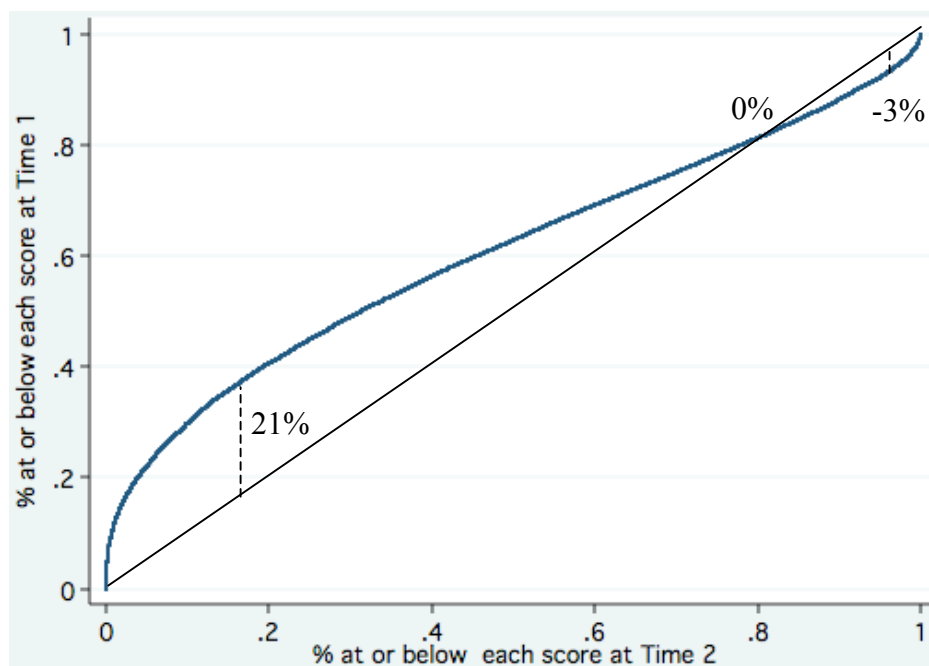


Figure 3.5: P-P plot from the simulated data displayed in Figures 3.3 and 3.4

The diagonal line in Figure 3.5 is shown for reference. A P-P function that lies on the diagonal would represent identical score distributions at Time 1 and Time 2, whereas a P-P function that lies mainly above the diagonal would indicate a positive score trend. The vertical lines drawn in Figure 3.3 show this. The point (.26, .37) on the P-P curve

shows that only 26% of students at Time 2 scored below the 37th percentile from the Time 1 distribution (the same 21% “gain” displayed in Figures 3.3 and 3.4).

Summary Statistics: V Coefficients

The P-P plot is a scale-independent representation of the difference between two distributions, as horizontal scale distortions of the CDFs do not change the P-P plot. Statistics generated from P-P plots are likewise scale-invariant. Since vertical deviations from the P-P plots to the diagonal represent score trends, one useful and interpretable statistic of interest would be the area under the P-P curve.

The area under the P-P curve (AUC):

$$\text{AUC} = \int_0^1 F_1(F_2^{-1}(p_2)) dp_2 = P(X_2 > X_1), \quad (4)$$

represents the probability that a randomly chosen test score from the Time 2 distribution is greater than a randomly chosen test score from the Time 1 distribution (Ho, 2007, p. 14). For identical score distributions at Time 1 and Time 2, the area under the P-P curve (which would fall on the diagonal) would be 0.50 (representing the chance probability). When scores improve from Time 1 to Time 2, the P-P curve would fall above the diagonal and the area would be greater than 0.50. As Ho notes, “the usefulness of this statistic is that it is invariant to discretionary choices such as cut-scores, percentile, and score scale” (p. 8). Thus, the area under the P-P curve addresses the problem of pliability of PAC-based trend comparisons under choice of cut-scores.

For the P-P plot in Figure 3.5, the area under the curve is approximately 0.611. This positive value represents the positive trend in scores from Time 1 to Time 2. It also indicates that a randomly chosen test score from Time 2 has a 61% probability of being greater than a randomly chosen test score from Time 1.

Wilk and Gnanadesikan (1968) note that the nonparametric Kolmogorov-Smirnov statistic is represented as the maximum deviation from the 45-degree diagonal to a point

on the P-P plot (p. 11). Likewise, Ho (2007) notes that P-P plots have conceptual ties to Receiver Operator Characteristic (ROC) curves and Lorenz Curves; and that the nonparametric Mann-Whitney U statistic is a simple linear transformation of the $P(X_2 > X_1)$ statistic. Another useful transformation of $P(X_2 > X_1)$ is found by assuming that the distribution from Time 1 has a standard normal distribution and the distribution from Time 2 has a normal distribution with unit variance. Under these assumptions, the area under the P-P curve defines the mean for the distribution from Time 2 that can be interpreted in terms of standard deviation units. Thus, these assumptions lead to a transformed summary statistic:

$$V = \sqrt{2}\Phi^{-1}(P(X_2 > X_1)) = \sqrt{2}\Phi^{-1}\left(\int_0^1 F_1(F_2^{-1}(p_2))dp_2\right), \quad (4)$$

where Φ^{-1} represents an inverse normal transformation. Ho and Haertel (2006a) describe V as a scale-free effect size of the trends in scores from Time 1 to Time 2. Unlike traditional effect sizes, the V statistic cannot be distorted by scale transformations, yet it may still be loosely interpreted as a distance in terms of standard deviation units.

For the P-P plot in Figure 3.5, $V \approx \sqrt{2}\Phi^{-1}(.611) \approx .40$. This indicates that the Time 2 scores increased by 0.40 standard deviation units over the Time 1 scores. This is supported by the fact that the distributions were simulated to have an effect size of approximately 0.40.

Calculating V From Reported Performance Level Data

P-P plots and V statistics are calculated from test score distributions. In this study, score distributions are not known. As Table 3.4 shows, the collected data simply show the percentage of students scoring at or below specific cut-scores for each state test and NAEP. If a state administers the same test twice and cut-scores do not change, then the corresponding percentages of students scoring below each cut-score at each time define points on a P-P plot.

Table 3.4 Example of collected test data to interpolate P-P curves

South Carolina	State Test			NAEP		
	2003	2005	2007	2003	2005	2007
4 th Grade Reading	18.9	21.4	21.9	20.8	18.5	20.3
	65.6	59.5	58.6	68.2	64.1	64.2
	86.0	85.8	80.3	96.1	95.3	95.3
4 th Grade Math	23.6	20.4	17.3	40.6	42.6	41.1
	67.1	63.6	57.8	74.3	74.4	74.2
	97.7	97.1	95.9	94.7	94.2	94.6
8 th Grade Reading	32.9	33.7	32.1	32.2	28.6	29.1
	80.2	76.8	80.2	73.7	70.1	68.1
	93.7	92.0	93.2	95.2	93.3	92.6
8 th Grade Math	*33.2*	*25.3*	28.7	*30.6*	*33.0*	31.3
	79.3	*70.3*	75.4	*75.8*	*75.3*	75.4
	97.7	*94.1*	96.6	*98.3*	*98.1*	98.3

Notes: Numbers represent percentages of students scoring below 3 cut-scores.

As an example, consider the 2003-2005 8th grade math (highlighted) data reported in Table 3.4. The data show that from 2003-2005, the percentage of students scoring below the first cut-score on the state test decreased from 33.2% to 25.3%. From this information, the point (.253, .332) was placed on the P-P plot. Likewise, the points (.703, .793) and (.941, .977) were plotted on the P-P curve for the state test. For the NAEP data, the points (.330, .306), (.753, .758), and (.981, .983) were plotted.

Thus, states reporting data from at least 3 cut-scores provided 3 points for the P-P plot. The theoretical points (0, 0) and (1, 1) were then added to the P-P plot to yield five data points. Figure 3.6 displays these points plotted for the example data in Table 3.4.

Using these five points, the smoothed curve algorithm implemented in Microsoft Excel was then used to plot a curve from a cubic Bezier-based interpolation function with control points (Ho, 2007). Figure 3.6 displays the smoothed (interpolated) P-P curves for the example data in Table 3.4. According to Ho and Haertel (2006), “Simulation studies suggest that three P-P points is the minimum number of points necessary for the interpolation function to obtain a reasonable approximation of the P-P curve” (p. 34).

Thus, data from states reporting fewer than three cut-scores were eliminated from this study

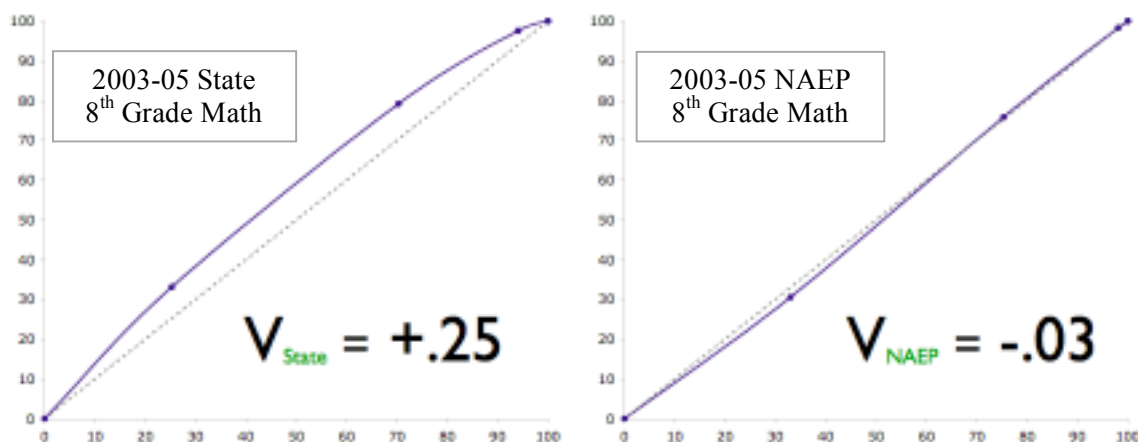


Figure 3.6: Smoothed (interpolated) P-P plots for the example data in Table 3.4.

The cubic spline macro for Microsoft Excel (SRS1 Software, 2007) was then used to obtain interpolated points from this smoothed P-P curve. With these points, numerical integration procedures were used to estimate the area under the P-P curve. From these estimated areas, values of the V statistic were calculated. In this study, Simpson's Rule was used to estimate the area under the smoothed P-P curves using 10,000 interpolated subdivisions. These area estimates were identical (to at least 6 decimal places) to the area estimates using 50,000 interpolated subdivisions.

Applying these procedures to the example data in Figure 3.6, V statistic values were estimated to be 0.25 for the state test and -0.03 for NAEP. Thus, for 8th grade math in South Carolina from 2003-2005, state test trends indicate an increase of 0.25 standard deviation units, while NAEP trends indicate a decline of 0.03 standard deviation units. Thus, the discrepancy in score trends was estimated to be $V = 0.28$.

Since all test scores contain measurement error, the effect of this error on the scale invariant effect sizes is of concern. As Ho (2007) explains:

Effect size-based trend statistics are generally attenuated by measurement error, but NAEP reports statistics that are corrected for this effect. In contrast, [state test] effect sizes are biased towards zero due to measurement error. If [state test] V statistics are treated like traditional effect sizes, they can be corrected by disattenuating them in inverse proportion to the square root of the reliability of the test (Hedges & Olkin, 1985). As reliabilities for State assessments are not always reported, the uncorrected [state test] statistics are used. As the results will show, if the reliabilities of state tests are taken into account, disattenuation will increase the degree of average State-NAEP trend discrepancy (p. 16).

Final Sample of Test Score Data

Tables 3.5 and 3.6 display the final sample of state test data that was included in, or excluded from, analysis. For each state, the scale-invariant effect size V for trends in state test and NAEP results were computed. Recall that, ideally, the data set would include results from each state in 2003, 2005, and 2007 in reading and mathematics in grades 4 and 8 for a total of 600 state test trend effect sizes and 600 NAEP trend effect sizes. Due to changes in tests or cut-scores, a lack of available data, or the use of fewer than three cut-scores, the final data set included 215 state test score distributions.

Once the data had been entered, the P-P plots generated through the interpolation function were visually inspected to check the accuracy of data entry. P-P plots with extreme deviations from the diagonal or that cross the diagonal were flagged for further investigation. Estimated values of AUC and V were also inspected and outliers were checked again for accuracy.

Visual displays of the V estimates were also examined to determine the accuracy of data entry and analysis. Since the V estimates should tend to agree with the proficiency trends reported by states, a scatterplot was inspected. Any observations in which the V estimates and proficiency trends differ were flagged for further investigation. A scatterplot of the V estimates and traditional effect sizes (calculated from the states that report means and standard deviations) were also inspected and outliers were flagged. A similar scatterplot using NAEP V estimates and traditional effect sizes was likewise inspected to determine the accuracy of data entry and analysis.

Table 3.5 Data included in the analysis

	Number of cut-scores reported on state tests											
	Grade 4						Grade 8					
	Reading			Mathematics			Reading			Mathematics		
	03-05	05-07	03-07	03-05	05-07	03-07	03-05	05-07	03-07	03-05	05-07	03-07
Alabama		3	3		3	3		3	3		3	3
Alaska												
Arizona	3	3	3	3	3	3	3	3	3	3	3	3
Arkansas		3			3			3			3	
California	4	4	4	4	4	4	4	4	4	4	4	4
Colorado	3	3	3		3		3	3	3	3	3	3
Connecticut	4			4			4			4		
Delaware							4	4	4	4	4	4
Florida	4	4	4	4	4	4	4	4	4	4	4	4
Georgia												
Hawaii		3			3		3	3	3	3	3	3
Idaho	3	3	3	3	3	3	3	3	3	3	3	3
Illinois							3			3		
Indiana												
Iowa												
Kansas				4	4	4	4	4	4			
Kentucky		3									3	
Louisiana	4	4	4	4	4	4	4	4	4	4	4	4
Maine	3	3	3	3	3	3	3	3	3	3	3	3
Maryland												
Massachusetts		3			3						3	
Michigan	3	3	3	3	3	3				3	3	3
Minnesota												
Mississippi	3	3	3	3	3	3	3	3	3	3	3	3
Missouri				3						3		
Montana	(3)	(3), 3	3	(3)	(3), 3	3	(3)	(3), 3	3	(3)	(3), 3	3
Nebraska												
Nevada												
New Hampshire												
New Jersey												
New Mexico												
New York	3			3			3			3		
North Carolina	3			3			3			3		
North Dakota		3			3			3			3	
Ohio		4						4			4	
Oklahoma		3			3		3	3	3	3	3	3
Oregon												
Pennsylvania							3	3	3	3	3	3
Rhode Island												
South Carolina	3	3	3	3	3	3	3	3	3	3	3	3
South Dakota												
Tennessee												
Texas	4			4			4			4		
Utah												
Vermont												
Virginia												
Washington	3	3	3	3	3	3						
West Virginia		4			4			4			4	
Wisconsin	3	3	3	3	3	3	3	3	3	3	3	3
Wyoming	3			3			3			3		

Notes: Blank cells represent data excluded from the analysis

Montana administered 2 tests in 2003 & 2005. The (ITBS) changed from high- to low-stakes in 2005.

Table 3.6 Data excluded from the analysis

	Number of cut-scores reported on state tests											
	Grade 4						Grade 8					
	Reading			Mathematics			Reading			Mathematics		
	03-05	05-07	03-07	03-05	05-07	03-07	03-05	05-07	03-07	03-05	05-07	03-07
Alabama	T			T			T			T		
Alaska	TS	C	TCS	TS	C	TCS	TS	C	TCS	TS	C	TCS
Arizona												
Arkansas	T		T	T		T	T		T	T		T
California												
Colorado				D		D						
Connecticut		S	S		S	S		S	S		S	S
Delaware	D	D	D	D	D	D						
Florida												
Georgia	C	C	C	C	C	C	C	C	C	C	C	C
Hawaii	D		D	D		D						
Idaho												
Illinois	D	D	D	D	D	D		D	D		D	D
Indiana	CD	C	CD	CD	C	CD	C	C	C	C	C	C
Iowa	C	C	C	C	C	C	C	C	C	C	C	C
Kansas	D	D	D								D	D
Kentucky	S		S	D	D	D	D	D	D	S		S
Louisiana												
Maine												
Maryland	D	C	CD	D	C	CD	C	C	C	C	C	C
Massachusetts	D		D	D		D	D	D	D	D		D
Michigan							D	D	D			
Minnesota	D	D	D	D	D	D	D	D	D	D	D	D
Mississippi												
Missouri	D	TS	TSD		TS	TS	D	TSD	TSD		TS	TS
Montana												
Nebraska	D	D	D	D	D	D	D	D	D	D	D	D
Nevada	D	D	D	D	D	D	D	D	D	D	D	D
New Hampshire	D	D	D	D	D	D	D	D	D	D	D	D
New Jersey	C	CD	CD	C	CD	CD	C	CD	CD	C	CD	CD
New Mexico	T	D	T	T	D	T	T	D	T	T	D	T
New York		D	D		D	D		D	D		D	D
North Carolina		D	D		D	D		D	D		D	D
North Dakota	S		S	S		S	S		S	S		S
Ohio	T		T	T	D	T	D		D	D		D
Oklahoma	T		T	T		T						
Oregon	C	CS	CS	C	CS	CS	C	CS	CS	C	CS	CS
Pennsylvania	D	D	D	D	D	D						
Rhode Island	T	T	T	T	T	T	T	T	T	T	T	T
South Carolina												
South Dakota	C	C	C	C	C	C	C	C	C	C	C	C
Tennessee	T	C	TC	T	C	TC	T	C	TC	T	C	TC
Texas		C	C		C	C		C	C		C	C
Utah	C	C	C	C	C	C	C	C	C	C	C	C
Vermont	D	D	D	D	D	D	D	D	D	D	D	D
Virginia	D	D	D	D	D	D	C	C	C	C	C	C
Washington							D	D	D	D	D	D
West Virginia	T		T	T		T	T		T	T		T
Wisconsin												
Wyoming		T	T		T	T		T	T		T	T

Notes: Blank cells represent data that were included in the analysis.

T = data were excluded due to a change in state tests

C = data were excluded due to the state reporting fewer than 3 cut-scores

S = data were excluded due to a change in the scoring standards (or number of cut-scores reported)

D = data were not reported

South Dakota reports 3 cut-scores, but extremely few students scored in one performance level

Reporting

Once values of V are estimated for state test and NAEP trends, the values for each state are displayed on a scatterplot to show the discrepancies. The centroid of the scatterplot represents the average discrepancy between state test and NAEP trends. The axes for the scatterplot are displayed in Figure 3.7. Because the values of V arrive from a normality assumption, a matched-pairs t-test is used to test the null hypothesis of equal trends in both state and NAEP tests.

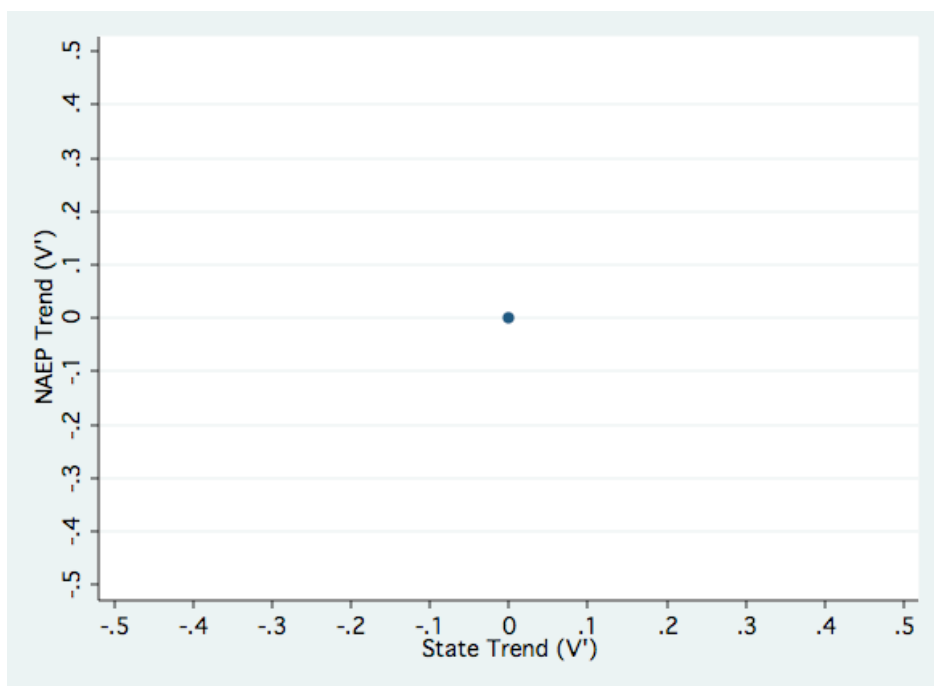


Figure 3.7: Scatterplot to display scale-invariant trend effect sizes

With a sample size of 215, the power of this t-test can be estimated via G*Power 3 (Faul, Erdfelder, Lang, & Buchner, 2007) by making some assumptions. Based on the 2003-05 discrepancy results reported by Ho (2007), assume the true discrepancy between state and NAEP trend effect sizes is 0.80 with a standard deviation of 0.15 for state trends, a standard deviation of 0.10 for NAEP trends, and a correlation of 0.60 between

state and NAEP trends. From these assumptions, a one-tailed t-test with an alpha of 0.01 would have estimated power above 0.99. These assumptions lead to an assumed effect size of 0.70 in the difference between state test and NAEP trends. With a sample size of 215, an effect size of at least 0.25 will be required for power to be 0.90. If alpha is set at 0.05, then an effect size of at least 0.20 will lead to a power above 0.90.

To further summarize the discrepancies in score trends, the percentages of paired observations (dots on the scatterplot) falling above and below the 45-degree diagonal are reported. Dots below the diagonal represent state-subject-grade-year combinations in which state score trends were more positive than NAEP trends. The percentages of dots falling in each quadrant of the scatterplot are also reported to show how many state-subject-grade-year trend combinations show a difference in sign (positive state trends with negative NAEP trends or vice versa).

The above analyses are also computed for each subject, each grade level, and each pair of years (2003-05, 2005-07, and 2003-07) to see if any different conclusions are reached. Finally, for each state-subject-grade-year combination, simple differences in V values are computed. These differences in V values will be used as a single estimate of the discrepancy between state test and NAEP results.

Assumptions and Limitations

One methodological limitation of this study involves the interpolation function used to generate the P-P plots from reported state test data. P-P plots and values of V are estimated from the limited number of cut-scores reported by each state. The estimations for states with fewer cut-scores probably contain more error than states that report larger numbers of cut-scores. Ho and Haertel (2006) note that, “the cubic Bezier interpolation method provides an estimate whose error has not been assessed” (p. 39). The interpolation procedure also requires the elimination of data from any state test with

fewer than 3 cut-scores. The elimination of this data may bias the results in some unknown way.

An analysis was conducted to investigate the impact of the choice of 3 cut-scores as a limitation of this method. For this analysis, the V estimates for states reporting 4 cut-scores were set aside. Then, for each of these states, the first cut-score was eliminated and V was re-estimated using the remaining 3 cut-scores. Likewise, V was estimated from 3 cut-scores after eliminating the second, third, and fourth reported cut-score. Comparisons were then made among the V estimates from 4 and 3 cut-scores. A significant discrepancy in the estimates would make the choice of 3 cut-scores suspect. Even if no significant discrepancy is found, the results of this study depend, in small part, to the interpolation method chosen to construct the smoothed P-P curves.

Another possible limitation is that V estimates are calculated by applying a normal transformation to the area under the P-P curve estimates. The value of these V estimates would change under different distributional transformations. The results in this study, therefore, are dependent upon this choice of transformation to normal distributions.

The usefulness of the V estimates was also analyzed through comparisons with traditional effect sizes (the ratio of the mean difference to the pooled standard deviation) from states reporting score means and standard deviations. Although this resulted in a small sample of data, it provides evidence as to the usefulness of the V estimates. To do this, a scatterplot (and the standard error of estimate) of the relationship between V and traditional effect size estimates were examined. Outliers were investigated in an attempt to provide an explanation as to why the V estimate would differ from the traditional effect size estimate.

The quality of the data also could have negatively impacted this study. The estimated NAEP percentiles and published state test results (which are oftentimes subject to rounding) contain error that impacted the V estimates. Also, some states only reported English or Language Arts test scores instead of reading test scores (as is reported

by NAEP). Trend discrepancies in these cases could simply be due to significant content differences in the tests. Also, some states administered their tests in the fall, whereas the NAEP is administered in the spring. This could have impacted results, since there is little reason to expect fall-to-fall score trends to be similar to spring-to-spring score trends.

Analysis

As explained earlier, the data collected from this analysis consists of estimates of the discrepancy in trends between state tests and the NAEP ($V_{\text{state}} - V_{\text{NAEP}}$) and various measures of the quality of state test security policies. From this data, the following questions can be addressed:

- What is the relationship between the quality of a state's test security policy and any discrepancies between score trends on state and NAEP tests? This question is addressed by regressing ($V_{\text{state}} - V_{\text{NAEP}}$) on the overall composite score (possibly standardized) from the test security policy evaluation form. This analysis is conducted separately for each grade level and subject.
- Which aspects of a state's test security policy seem to have the strongest relationship with score trend discrepancies? This question is addressed by regressing ($V_{\text{state}} - V_{\text{NAEP}}$) on (standardized) composite scores for each of the four sections (F-O-I-L) of the policy evaluation form. A comparison of the standardized beta weights would indicate which aspect had the greatest impact on trend discrepancies.
- Do states with higher numbers of reported test score manipulations experience greater trend discrepancies? To address this question, the correlation between the numbers of published news reports and the trend discrepancies for each state is calculated. An analysis is conducted to determine if states with higher numbers of published reports prior to 2003 experienced smaller discrepancies (evidence that the published reports caused states to focus more on test security).

The seven dichotomizations of state policies are used to conduct mean comparison tests (t-tests or nonparametric analyses) in order to address the following questions:

- Do states with clear and accessible policies experience smaller trend discrepancies than states with ambiguous or difficult-to-find policy information? **Explain this**
- Do states with state-level mandates experience smaller trend discrepancies than states with district-level policies? **Explain this**
- Do states with punitive policies experience smaller trend discrepancies than states with instructive policies? **Explain this**
- Do states requiring independent test administration monitoring experience smaller trend discrepancies than states that allow teachers to administer tests to their own students? **Explain this**
- Do states with investigative policies experience smaller trend discrepancies than states with preventative policies? **Explain this**
- Do states with example-based policies experience smaller trend discrepancies than states with more general policies? **Explain this**
- Do states with positive-message policies experience smaller trend discrepancies than states with negative-message policies? **Explain this**

After these analyses, all possible longitudinal analyses for states that adopted or modified test security policies between 2005-2007 were conducted. To do this, trend discrepancies from 2003-2005 were compared with discrepancies from 2005-2007. It is unknown how many observations are available for these analyses.

Finally, the narrative (additional information) collected from each state policy are used to briefly describe some case studies. For example, Montana administered two tests in 2003, 2005, and 2007. One test, the ITBS, was high-stakes in 2003 and low-stakes in 2005 and 2007. The other test became high-stakes in 2005. Trends on these tests can be compared to see if the change from high-stakes to low-stakes on the ITBS caused score

trends to decline. Also, simple score trends for states that changed tests sometime during 2003-2007 can be calculated and compared. For example, Alabama changed tests in 2004. From 2004-2005 (the first two years of administration), test scores increased at a high rate. After that, score trends appeared to slow down. Because Alabama did not administer the test in 2003, the increase in scores from 2004-2005 will not be “captured” by this study even if that increase in scores was caused by educator manipulations.

Scope

This study analyzes data from 2003, 2005, and 2007 for grades 4 and 8 in reading and mathematics in states with tests having at least 3 cut-scores. It only analyzes results from state tests and the NAEP in those years, subjects, and grades. Conclusions are cautiously made due to the presence of confounding variables.

Assumptions, Limitations, and Confounding Variables

The biggest limitation of this study is that it cannot attempt to find a causal relationship between test security policy quality and score trend discrepancies. A causal relationship cannot be inferred due to the influence of confounding variables (plausible rival hypotheses, discussed in the previous chapter). Differences in test content, test administration, examinee motivation, examinee populations, and other factors could influence the relationship between the independent and dependent variables in this study. Likewise, this study can only evaluate the quality of test security policy materials. No causality can be concluded because the quality of a state’s policy materials does not necessarily represent the quality of the state’s policy implementation.

For example, one alternate hypothesis for any significant relationship between test security policy quality and trend discrepancies is that states with higher-quality policies simply have higher-quality testing programs. Higher-quality testing programs could, in turn, simply represent states with higher-quality educational systems (curriculum development, teacher training, etc.). States with higher-quality educational systems

might be expected to score higher on state tests (trained teachers focusing on state-developed curriculum) than on NAEP. This is just one of many plausible hypotheses. This study only attempts to discover if a relationship exists between the quality of state test security policies and score trend discrepancies.

A second limitation is that some states were excluded from analysis. Data from states experiencing changes in tests or cut-scores, or from states with fewer than three cut-scores, were excluded. The exclusion of this data may bias the results in some unknown way and limit the generalization of the results.

A third limitation is that the trend discrepancy statistics are calculated at the state-level. The decision to manipulate test scores may be made at the level of individual teachers, schools, or school districts. If, in fact, manipulations have a relationship with score trend discrepancies, a state-level trend estimate may not be able to detect this relationship.

Summary

This study has been designed to determine if a relationship exists between the existence and quality of state test security policies and discrepancies in state-NAEP score trends. The quality of state test security policies is evaluated based on a framework derived from analyses of newspaper reports, educator surveys, state surveys, direct observation studies, statistical analyses, and targeted research into test score manipulations. State-NAEP trend effect size discrepancies are estimated via a scale-invariant framework that is not impacted by choice of cut-scores.

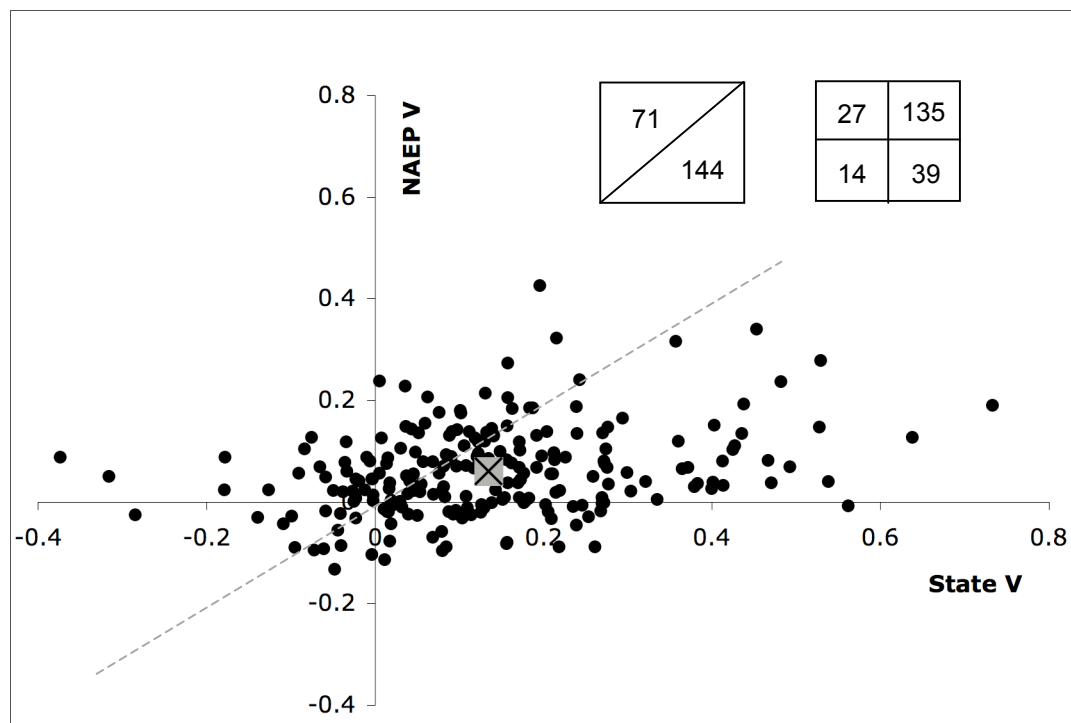
The study is designed in such a way as to determine which aspects of test security policies might be more strongly related to score trend discrepancies. The study also provides information regarding the content and quality of test security policies adopted by states to deter educators from manipulating scores.

The study extends the work of researchers who have focused on inappropriate testing practices, test score pollution, and detecting cheating on achievement tests. It also extends the work of researchers who have focused on comparing state test results to NAEP results. It extends the work of Ho (2005, 2007) and Haertel (2006a) in their development of scale-invariant measures of trend discrepancies and furthers the research of Cizek (1999) and McCabe and Trevino (1993, 2002) in evaluating the impact of test security practices and honor codes in educational organizations.

Results of this research could be used by states to develop, improve, or audit their current test security policies and practices. Results could also be used in professional development to train teachers in appropriate test preparation and administration activities. Finally, this study will contribute to the debate over the effectiveness of accountability systems and sanctions in improving student achievement.

COMPARE TO TRADITIONAL EFFECT SIZES (Center for Ed. Policy)

CEP – more states plan on reporting means & SDs



APPENDIX A: PUBLISHED NEWS SUMMARIES

- 12/21/2007 USA Today:
Doris Alvarez, principal of Preuss high school in San Diego, resigns in connection with a case of alleged cheating and grade-tampering. An audit by the University of California at San Diego found that three-fourths of reviewed Preuss transcripts had one or more grades changed — most of them to benefit students. Just last month, the school, which prepares low-income and minority students for college, ranked 10th out of 18,000 U.S. high schools in new rankings by U.S. News & World Report. Preuss ranked second in the USA among charter high schools. Last May, Preuss ranked 10th among 1,200 high schools deemed the nation's best by Newsweek.
(Toppo, 2007)
- 12/16/2007 Orlando Sentinel:
In June, a Pasco County teacher was fired for rewording questions while proctoring the math portion of the FCAT. The Caveon Test Security company, which investigates cheating allegations and has a contract with Florida, estimates that adult-led cheating is probable in some 1 percent or 2 percent of the schools it has investigated.
(Julian, 2007)
- 12/13/2007 New York Post:
A former Staten Island school administrator ordered teachers to change scores on state Regents exams to enable students to graduate - and even hiked some of the test scores herself at home. Former Wagner HS Assistant Principal Mary Incantalupo, of Staten Island, was recommended for termination for her role in a 2006 test-tampering scandal that's been under investigation by the Department of Education since October 2006. The school's principal, Gary Giordano - who was dating Incantalupo at the time of the grading and who married her this year - was cleared of serious wrongdoing, but will be disciplined for not keeping tabs on the exams.
(Gonen, 2007)
- 11/11/2007 KSTP Eyewitness News (Minnesota):
A Channel 5 Eyewitness News investigation finds teachers help their students cheat on standardized tests. According to reports from the Minnesota Department of Education, teachers may find it tempting to point out a correct answer on a multiple choice test or correct a composition prior to grading. According to the Minnesota Board of Teaching, one teacher had their license suspended for three years for altering student's answers on a test. Another teacher had their license suspended for nine months for administering a test improperly. Even though the test booklets are supposed to be sealed and kept locked up, the reports showed that teachers were found to be taking the tests, sharing them with other teachers, and in one case, a teacher's pre-test lesson plan included a math problem which was "strikingly similar" to one appearing on the actual test.
(Muehlhausen, 2007)
- 11/01/2007 Boston.com:
20 Massachusetts teachers are accused of improperly helping students during the MCAS exams in 2007. This compares to 15 accusations in 2006 and 3 in 2005. A New Bedford elementary school teacher briefed students on the subjects of the reading passages they would encounter on the reading MCAS. Other teachers provided dictionaries or other forbidden tools, or made mistakes in administering the exam, such as forgetting to remove helpful material from a visible place in the classroom.
(Jan, 2007)
- 10/15/2007 The Columbus Dispatch:

An analysis of investigations in Ohio's 10 largest school districts in 2006 finds that 5 of the 8 educators accused of cheating are either under investigation by the state or have been punished already. Some districts prefer to let the problem teachers move on quietly to another school, a practice known as "passing the trash." Some districts have a "basics-only" policy for reference checks that prevents them from telling potential employers anything except when the teacher worked for the district.

(Smith Richards & Riepenhoff, 2007)

- 10/12/2007 The Detroit News:
Thousands of fifth and sixth grade students in Michigan will be forced to retake the writing portion of the Michigan Educational Assessment Program test after a newspaper published sensitive information about the test. The Jackson Citizen Patriot published two of the writing topics before the test administration period. The decision to retest students was made after the Michigan Department of Education learned a reporter was allowed into Jackson Public Schools during the administration of the test, a violation of the state's testing ethics.
(Mrozowski, 2007)
- 09/07/2007 The Press-Enterprise (California):
An English teacher at Citrus Hills Intermediate School in Corona, CA and a fourth-grade teacher at Red Maple Elementary School in Moreno Valley, CA are caught allegedly giving their students copies of old exams as practice for the current year's test. The two schools where these teachers worked were among 15 California schools that are currently under investigation for testing irregularities. This is not the first report of cheating in these school districts. A math teacher at Sierra Middle School in Riverside Unified District resigned in 2004 after being caught changing student answers on dozens of state tests.
(Parsavand, 2007)
- 08/23/2007 Recordnet.com (California):
A school district in Calaveras reports allegations of teacher cheating to state officials. One teacher in the district allegedly read portions of the state test or taught content during the test administration. Another teacher allegedly used actual test questions to prepare students for the test.
(Johnson, 2007)
- 08/23/2007 CBS 13 (California):
One teacher at August Elementary school and another at McKinley Elementary school in Stockton, California are accused of cheating to increase their students' scores on the state tests. The McKinley school teacher was caught writing down questions from the test, while the August teacher is accused of giving a student "punctuation tips." District official Dianne Barth is quoted as saying, "I don't believe it's widespread. I don't believe there is cheating in Stockton Unified [school district]."
(Cannata, 2007)
- 08/23/2007 Washington Post:
Severna Park High School in Washington D.C. is put on probation after allegations of cheating on an Advanced Placement history exam. A College Board investigation found that the test proctor failed to follow test directions and allowed students to talk and use cell phones. 42 students were forced to re-take the test and the test proctor was banned from ever administering the AP tests again.
(Wan, 2007)
- 08/18/2007 The Modesto Bee:
Twelve schools in San Joaquin Valley, California submit 13 reports of testing irregularities to state education officials. One teacher at Oak View School allegedly made a practice worksheet of

items “almost identical” to the items on the CAT-6 exam. The school’s principal suggested the behavior was due to the pressure felt by teachers under No Child Left Behind to increase test scores. Another teacher allegedly allowed her students to use calculators on the math test. A third alleged incident involves a teacher who helped students during the writing exam. Another teacher was caught copying answers from student test booklets.
(Balassone, 2007)

- 08/16/2007 Dayton Daily News:
Test scores from the City Day Community charter school in Dayton, Ohio plummet when independent officials monitor the 2007 test administration. In 2006, the school had produced extraordinary score gains under suspicion that students were given practice tests that were identical to the actual test.
(Elliott, 2007b)
- 08/10/2007 Herald Tribune (Florida):
Mary Cropsey, a third-grade teacher at Mills Elementary School in Manatee, Florida, is accused of tampering with student answer sheets on the Florida Comprehensive Assessment Test (FCAT). One student reports that Cropsey helped students on the test; another student reported hearing that the teacher gave students extra time to complete the exam. An investigation began after yet another student reported that she had not finished the exam, but the next day all the bubbles had been filled-in. If the allegations are proven true, Cropsey could lose her teaching certificate and even be charged with a crime.
(Morris, 2007)
- 07/16/2007 Miami Herald:
Hollywood Hills Elementary, an A-rated school for the past five years, receives an incomplete grade because FCAT scores for 16 students were flagged as irregular. The students in Katie Steinberg-Lessard’s third grade class do not have their test scores months after taking the test. The teacher spent three mornings a week before school with students who wanted extra practice for the FCAT. Five other Dade schools with incomplete grades are currently being investigated.
(Shah, 2007)
- 07/16/2007 San Francisco Chronicle:
The California Department of Education concludes that for the second consecutive year, educators at University Preparatory Charter High School in San Francisco interfered with state-mandated testing. State investigators seized illegal copies of the 2005 form of the test that was used to prepare students for the exams. Eight former teachers at the school assert the existence of a culture of cheating at the school. According to those former teachers, student grades are frequently falsified and low-scoring students are excluded from state-mandated testing. Last year, the state found that hundreds of answers on the ninth-grade English and math tests had been changed from wrong to right. A counselor from Oakland’s Skyline High school reports that a student earning D’s and F’s transferred to University Preparatory Charter High School and received A’s and B’s while taking 16 classes in a single semester. When the student returned to Skyline High, he once again earned D’s and F’s. Last year, investigators concluded that educators at the school changed hundreds of test answers before they were sent for scoring. Former testing coordinator Mike Schwartz is suing school founder and director Isaac Haqq for breach of contract, claiming Haqq was responsible for the altered answer sheets.

Eight days after the original story was published, Isaac Haqq resigned as principal from University Preparatory Charter High School.
(Asimov, 2007ab)
- 07/13/2007 The Dallas Morning News:

A state investigation finds that David Tamez, an elementary school teacher in Amarillo, Texas, leaked the fourth-grade writing test prompt on the spring TAKS writing test to colleagues before the test administration. Tamez reportedly leaked the test information because he believed educators in other districts were doing it as well. The teacher obtained the test information by volunteering to serve on the committee that selects questions for the final form of the TAKS. He alleges that committee members “regularly smuggle out secret TAKS information to share in the home districts.” Another teacher interviewed by investigators signed a statement indicating that Tamez “bragged that the source of his insider test information was... a person he had sex with who works for a company that helps build the TAKS.” The Amarillo Independent School District concluded that the teacher obtained the information from an unidentified employee at Pearson Educational Measurement. Tamez resigned from his position, but will retain his teaching certificate if he cooperates with the investigation. (Benton, 2007b)

- 07/03/2007 News 10 Now:
An investigation concludes that a third grade teacher at West Leyden Elementary in the Adirondack Central School District in Boonville, NY cheated for students on the New York State Mathematics Exam. According to parents, the teacher told students how she was going to help them and then tapped on any incorrect answers during the test administration. According to superintendent Frederick Morgan, the teacher remains employed by the district – she was “... moved to a grade where there is no state exam given.” (Ohler, 2007)
- 06/27/2007 New York Times:
A 23-month investigation into cheating allegations at Cobble Hill High School of American Studies in Brooklyn, NY concludes that the whistle-blower wrongly accused both the principal and assistant principal. Teacher Philip Noble accused principal Lennel George and assistant principal Theresa Capra of ordering teachers to cheat on the scoring of Regents exams. The investigation alleges that Mr. Noble was “a sub-par teacher with poor evaluations who wrongly accused [them] of engineering a cheating scheme because [they] had given him a negative review that could have led to his firing.” The investigation does not explicitly rule out the possibility of cheating at the school; only that the principal and assistant principal did nothing wrong. (Bosman, 2007)
- 06/24/2007 New York Times:
A 2006 investigation concludes that wrong answers were erased and changed to correct answers on state-mandated English tests in four New York City elementary schools. 2007 test scores for the schools in which cheating allegedly took place fell substantially, providing evidence that cheating had inflated the 2006 scores. (Fessenden, 2007)
- 06/24/2007 Newsday.com:
State officials blame “adult interference” for suspiciously high test scores in eight schools in Camden, NJ. Faculty members reportedly allowed students to use calculators, which were not allowed on the exam. (Marcus, 2007a)

The entire Uniondale school district is placed on academic probation due to evidence of tampering with Regents Math A and B high school exams and the State Mathematics Assessments for grades 3-8 in 2005 and 2006. The New York Department of Education reports that complaints of test fraud have more than doubled over the past five years, with the department receiving 37 complaints in 2006. One dozen teachers and administrators accused of test fraud have faced hearings in front of the New York Professional Standards and Practices Board. Of those twelve cases, six cases resulted in revocation of professional certifications, two cases were cleared, and

the remaining four cases remain under investigation. The number of complaints verified by the state has remained relatively steady, with between 9-16 in each of the past five years. (Hildebrand, 2007a)

An analysis of Uniondale's test scores found that 333 answers on the Regents Math A exam were altered, and 97% of the time they were changed to the correct answer. On the Regents Math B exam, 198 answers were changed, with 97% again being changed to the correct answer. On the 2005 8th grade math assessment, Uniondale students scored below average on 11 of the 14 easiest questions, but higher than average on 12 of the 13 most difficult items. (Hildebrand, 2007b; Marcus, 2007b)

- 06/21/2007 KLTU-7:
Lincoln Intermediate schoolteacher Bernice Martin allegedly changed answers on 17 math TAKS tests in 2006. The investigation into the alleged test fraud included erasure and handwriting analyses. The teacher, a counselor who served as test coordinator, and a retired principal could lose their credentials over the incident. (McCollum, 2007)
- 06/15/2007 The Dallas Morning News:
The Texas Education Agency could begin proceedings to close Theresa B. Lee Academy, a Fort Worth charter school, do to alleged tampering with the 2005 administration of the Texas Assessment of Knowledge and Skills (TAKS). The school was identified through a statistical analysis of test scores conducted by test security firm Caveon. When the state followed-up on this analysis, the school repeatedly refused to provide information to investigators. When asked for testing paperwork, principal William Powell reportedly replied that the paperwork had been, "lost in the flood." The details of this flood story changed with each telling and state investigators were unable to confirm the existence of any flood. The academy's vice principal, Shirley Dukes, reportedly changed answers on student answer sheets, informed teachers of the test's essay topics before test administration, and even wrote essays for the students. Dwaine Guyton, a former teacher at the academy, alleges that vice principal Dukes changed student answer sheets after test administration. Jakobus Wolf, a former science teacher at the academy, alleges that vice principal Dukes copied test answers onto the chalkboard for students and later asked if Mr. Wolf would be interested in being paid to manipulate student answer sheets. Mr. Wolf is reported to have said, "Kids know that if they go to Theresa B. Lee, somebody else will pass the TAKS for them." (Benton, 2007a)
- 06/14/2007 Texas Education Agency News:
The Texas Education Agency is recommending sanctions against three educators in three schools because of cheating on the TAKS. An analysis found excessive erasures and evidence of tampering with student answer sheets at Winona High School. In San Augustine Intermediate School, a student complained that someone had changed his answers on the 7th grade math test. An analysis later found evidence of tampering on the answer sheets of 17 out of 25 students in that class. (TEA, 2007)
- 06/12/2007 WREG-TV Memphis:
Connie Smith is fired over a cheating scandal in Tunica County, Mississippi. The teacher at Robinsonville Elementary witnessed the school's principal leaking test answers and reported the misconduct to the district superintendent. The school board terminated Smith's contract even though two other teachers confirmed her account of the misconduct. (Turner, 2007)
- 06/06/2007 Newsday:

New York state auditors discover 11 schools opened exam materials prematurely. Under state rules, the exam materials must be stored in steel safes or concrete vaults and are not to be unsealed until the day of testing. At 14 other New York schools, auditors found that exams had been removed from their locked boxes before being stored in safes. The state Education Department has recently revoked the rights of more than 20 schools to store exams after finding security breaches.

(Hildebrand, 2007c)

- 06/04/2007 The Dallas Morning News:
A conservative statistical analysis of 2005-2006 TAKS answer sheets conducted by Dr. George Wesolowsky, a professor at McMaster University in Canada, finds that the scores from more than 50,000 students show evidence of cheating that could include students copying answers from other students or educators doctoring student answer sheets. The analysis found 112 schools in which at least 10% of answer sheets were flagged for cheating. Many of the suspicious scores were found on the 11th grade test – the test students must pass to graduate. The schools with the strongest evidence of cheating include Forest Brook High School (North Forest ISD); Worthing High School and Sam Houston High School (Houston ISD); and South Oak Cliff High School (Dallas ISD). Based on the analysis, it appears as though cheating was more than 3 times as common in Dallas and Houston as it was in other large Texas districts. Professor Wesolowsky is quoted as saying, “The evidence of substantial cheating is beyond any reasonable doubt.” (Benton & Hacker, 2007a, 2007b)
- 05/26/2007 Visalia Times-Delta:
Three Visalia Unified School District teachers in Visalia, CA are reported for misconduct during the 2006 administration of tests for the Academic Performance Index (API). One La Joya Middle School teacher and one Crestwood Elementary School teacher are reported to have read portions of the test that were designed to be read by students. A teacher at Mineral King Elementary reviewed questions with students after administering the test. District Superintendent Stan Carrizosa reportedly views the incidents as innocent mistakes. (Garcia, 2007)
- 05/22/2007 WMC-TV Memphis:
A teacher at Germanshire Elementary in Memphis, Tennessee allegedly cheated for students on the TCAP test. The district’s evidence of this misconduct includes erasure marks on test booklets, stray marks on answer sheets, and statements from students. (Rhodes, 2007)
- 05/20/2007 Inside Bay Area:
Three high school teachers in Oakland resign after being caught cheating on the California High School Exit Exam. Two of the teachers reportedly clarified a test question on the math portion of the exam while the third teacher proctored the exam. Dale Brodsky, an attorney hired by the Oakland Education Association, is quoted to have said cheating is “a non-issue in this whole debate about testing,” and then questioned the term cheating by saying, “What is ‘cheating’?” (Murphy, 2007)
- 05/13/2007 San Francisco Chronicle:
Teachers in at least 123 public schools have reportedly cheated for students on California’s high-stakes tests between 2004-2006. In two-thirds of these cases, the schools admit that they had cheated. The cheating behaviors included (a) allowing students to use reference materials such as maps and flow charts during the test, (b) allowing students to use calculators, (c) helping students answer questions, and (d) erasing and changing student answers. California currently identifies potential misconduct by scanning answer sheets for suspicious erasures. Cheating is virtually ignored in schools in which cheating impacts less than 5% of tests are given. Schools in which cheating impacts more than 5% of the tests are not ranked and receive a note stating “adult

irregularity in testing procedure” occurred. Since 2005, the following San Francisco area schools have confirmed testing irregularities: Mission Elementary (Antioch Unified); Hidden Valley Elementary, Cambridge Elementary, and Glenbrook Middle (Mount Diablo Unified); Forty-Niners Academy and Ravenswood City Elementary (East Palo Alto); John Muir Elementary (San Francisco Unified), Scott Lane Elementary (Santa Clara Unified), Bay Farm Elementary (Alameda City Unified), Chavez Middle and Treeview Elementary (Hayward Unified), Los Paseos Elementary (Morgan Hill Unified), Petaluma Junior High, Petaluma Joint Union High, Fair Oaks Elementary, and Williams Elementary (San Jose Unified). (Asimov & Wallack, 2007)

- 05/11/2007 Associated Press:
A teacher in Bloomington, Ohio is placed on paid leave after allegedly helping students cheat on the state graduation test. Nine students will be forced to retake the exam instead of graduating with their classmates.
(Associated Press, 2007a)
- 05/02/2007 St. Petersburg Times:
Barbara Heggaton, a special education teacher at Moon Lake Elementary School in Pasco County Florida, is accused of giving answers to three students during administration of the FCAT.
(Solochek, 2007)
- 05/01/2007 Tyler Morning Telegraph:
An investigation by the Texas Education Agency Office of Inspector General reports that test administration misconduct occurred in the 2005 administration of the TAKS in Winona High School. While the district was cleared of any wrongdoing in three previous investigations, the report alleges misconduct on the basis of suspicious erasure patterns and a possible breach of security caused by a missing key to the cabinet in which tests were kept.
(Waters, 2007)
- 02/04/2007 Dayton Daily News:
A newspaper investigation found that students at City Day Elementary School in Dayton, Ohio were given 44 practice questions that were identical or “substantially the same” as questions from the actual state exam. In some questions on the practice test, only names or small details were changed from the real test questions. The investigation was launched due to the suspiciously large amount of improvement shown by the school. In 2005, no sixth grade student in the school passed the math subtest of the Ohio Achievement Test. One year later, 100% of these students (now in 7th grade) passed the math test.
(Elliott, 2007)
- 12/03/2006 Deseret News:
The Utah State Board of Education accepted a test protocol pamphlet that defines cheating on the U-PASS exams. According to state testing director Judy Park, the ethics policy comes as a response to, “an unusually high volume of calls to the state office from testing directors, parents, teachers, and superintendents with ethical questions on the way tests are given.” According to Park, the state receives about five reported testing protocol violations each year. The new policy defines the following behaviors as cheating: (a) changing student answers in any way, for any reason, (b) looking at a test beforehand and altering lessons, (c) using inflections or gestures to help students in any way, (d) leaving helpful materials on classroom walls, (e) reclassifying student ethnicity, (f) letting students or parents supervise other students taking a test, (g) suggesting a student rethink his or her answer, (h) letting students take testing materials away from the testing site.
(Toomer-Cook, 2007)
- 11/20/2006 New York Daily News:

City officials are investigating teachers from Millenium Art Academy in Castle Hill for allegedly coaching 35 students during testing and inflating student scores.
(Einhorn & Melago, 2006)

- 11/06/2006 Staten Island Advance:
Seventeen Staten Island teachers inform the United Federation of Teachers of tampering with the Regents exam. The vice principal at Wagner High School allegedly re-scored student tests at home while teachers added points to student test scores. The teachers claim they were told to change test answers in their classrooms. The informants also claim the principal said he would make them pay for coming forward. Other Staten Island teachers suggest this behavior is a system-wide practice. According to Frank DeSantis, a teacher in St. George High School, “A lot of teachers get that feeling that all [schools] are looking for is statistics, and [they’re] lying and cheating to get them.”
(Gonen, 2006; W-CBS TV, 2006)
- 10/22/2006 The Columbus Dispatch:
Of the 28 Ohio school districts analyzed by The Columbus Dispatch, 15 had instances of educators cheating on standardized tests. Barbara Oaks, a teacher in the Coventry district, looked through the test and wrote out a geometry problem she thought her students would have trouble with. Winifred Shima, a teacher from the Parma district, used a copy of the test to create a study guide for students that included 45 of the 46 actual test questions. Brian Wirick (East Knox) and Heather Buchanan (Wapakoneta) both used the test to create study guides for students. Judy Wray, a veteran teacher in Marietta, made copies of the actual state test to help students prepare. Wray is reported to have said that teachers cheat more than administrators know.
(Richards, 2006a)
- 10/11/2006 The Indianapolis Star:
Two Corpus Christi Catholic School teachers in South Bend, Indiana are found to have cheated on statewide exams. Beth Troyer and Sandra Ernst were suspended for one week without pay for allegedly sending questions and answers (from an older version of the test) home with the students. State officials have received about a dozen reports of testing violations this year, but only half are suspected cheating incidents.
(Hupp, 2006)
- 10/01/2006 The Dallas Morning News:
5 months after being found guilty for cheating on the Texas Assessment of Knowledge and Skills (TAKS), at least 10 of the 22 Wilmer-Hutchins teachers are now working in other North Texas Public Schools. More than two years after the cheating took place, none of the teachers ever faced official sanction. Several of the school districts that now employ these teachers were unaware that these teachers have cheated in the past.
(Benton, 2006b)
- 09/25/2006 The Indy Channel.com:
A fifth-grade teacher from Wayne Township, Indiana receives a one-week suspension without pay for allegedly giving four students extra time to complete the math portion of the Indiana State Test of Educational Progress. Tom Langdoc, the district’s Director of School Community Services, believes the teacher was aware that she was cheating.
(The Indy Channel, 2006)
- 08/31/2006 Fairtest.org:
A teacher’s aide and a guidance counselor at Morton Elementary in Franklin City, VA are suspended after allegedly changing student answers on state exams. School officials report the answer changes would have actually resulted in more student failures on the exam.
(Fairtest, 2006)

- 08/20/2006 The Boston Globe:
The Massachusetts Department of Education documents 15 cases of inappropriate educator behaviors on the 2006 administration of the MCAS (compared to 3 allegations in 2005). A sixth-grade teacher from Andover West Middle School is reprimanded for reviewing a student's test and returning it to the student for revision. A fifth-grade test booklet at Pentucket Lake Elementary School was stolen and mailed to a local newspaper. Teachers in New Bedford and Peabody allowed students to use dictionaries during the test.
(Jan, 2006)
- 07/30/2006 Houston Chronicle:
Two Houston fifth-grade teachers resign after being accused of giving test answers to their students. Sheryle Douglas and Shawn Manning, the teachers once praised by President Bush and Oprah Winfrey, admit to giving students answers to an old version of the Stanford 10 Achievement Test as practice for this year's test. Scores from this test are used to award pay bonuses to teachers. The teachers worked at Wesley Elementary, which was also under investigation in 2003 when a former teacher accused school administrators of pressuring teachers to give test answers to students.
(Tresague & Viren, 2006)
- 07/28/2006 Dallas Star-Telegram:
The Texas Education Agency announces it will investigate testing irregularities at 609 schools from the 2005 administration of the Texas Assessment of Knowledge and Skills. Four types of irregularities were reported in Texas: patterns of similar responses, multiple marks on answer sheets, large score gains compared to previous years, and unusual response patterns. State-appointed monitors will oversee future test administrations.
(Brock, 2006)
- 07/04/2006 Baltimore Examiner:
Officials revoke the certificates of two fourth-grade teachers in Carroll County after they were accused of cheating on the Maryland School Assessments. One of the teachers admitted to copying questions from a previous test in order to create a practice worksheet for students.
(Johnson, 2006)
- 06/30/2006 Brevard School District web site:
Lori Backus, principal of Cocoa High School in Brevard, FL is accused of moving at least 54 9th and 10th grade special needs students into 11th grade so that their FCAT scores would not count towards the school's grade (assigned by the state) in 2005 and 2006. As a result of an investigation into the allegations, Principal Backus was immediately removed as principal.
(Brevard School District, 2006)
- 06/25/2006 Philadelphia Inquirer:
Edison Schools fires Jayne Gibbs, principal at Parry Middle School in Chester, Pennsylvania for allegedly changing student test answers in 2005. Eighth graders at the school said the principal had given them the answers to questions on the Pennsylvania System of School Assessment. Gibbs is also accused of exempting special-education students from testing, violating state and federal rules. Edison Schools also asks the state and district to investigate exemplary test results at Showalter Middle School, where Gibbs served as principal from 2003-04.
(Patrick & Eichel, 2006)
- 06/09/2006 Abilene Reporter-News:
An elementary school in the Big Spring district in Texas is flagged for testing irregularities. Third-graders at Marcy Elementary were found to have too many erasure marks on the reading test in the 2005 Texas Assessment of Knowledge and Skills.

(Levesque, 2006)

- 05/23/2006 The Dallas Morning News:
According to Caveon, a test security firm hired by the Texas Education Agency, almost 9% of schools had unusual scores on the Texas Assessment of Knowledge and Skills. Using statistical analyses, the firm found suspicious scores from 702 classrooms in 609 Texas schools in 2005. In one elementary school, 45 of the 262 students had identical answer sheets. An additional 29 students had perfect scores on the test. The chances of this happening naturally would be less than 1 in 1 trillion trillion trillion trillion trillion trillion (a 1 followed by 27 zeros).
(Benton, 2006a)
- 05/21/2006 St. Louis Post-Dispatch:
(excerpt) “Principal instructed teachers to encourage children to retry specific questions if the teachers thought the children knew the answer but had missed it on their first try.”
(CEA, 2007, p. 13)
- 04/17/2006 MSNBC:
With permission from the federal government, nearly two million students’ test scores are not counted when schools report progress by subgroups under the No Child Left Behind requirements. This is due to states being able to define the minimum number of students needed in a subgroup before scores are reported. In the past two years, almost half of all states have successfully petitioned the U.S. Department of Education to increase these minimums. An investigation concludes that about 1 out of every 14 test scores are not being counted under appropriate racial categories. The scores from more than 24,000 students in Missouri, 257,000 in Texas, and 400,000 in California are not being counted.
(Associated Press, 2006b)
- 04/11/2006 The Columbus Dispatch:
The Ohio Department of Education is investigating possible security breaches on the 2006 state tests. According to the department, 11 districts are investigating security breaches. The allegations include opening sealed boxes of test booklets early and teachers helping students cheat on the exams. Lora DeCarlo, a teacher at Franklin Middle School, was suspended without pay for 10 days. According to the teacher, she reviewed some student answer sheets and returned their tests to them with pages open to the items they needed to review. Other Ohio teachers accused of helping students cheat on tests in 2006 have resigned. Two years ago, a Hilliard teacher and a Reynoldsburg administrator resigned after acknowledging they broke test rules.
(Richards, 2006b)
- 03/28/2006 The Baltimore Sun:
(excerpt) “Teacher took notes based on the test administered last year and created worksheets for her pupils for this year’s test. She also shared the worksheet with other teachers. Some of these other teachers, no knowing the origin of the questions on the worksheet, alerted the principal to similarities between the worksheets and this year’s test.”
(CEA, 2007, p. 13)
- 03/08/2006 – 06/16/2006 Philadelphia Inquirer:
Joseph Carruth, principal of Charles Brimm Medical Arts High in Camden, New Jersey, is fired after accusing Assistant Superintendent Luis Pagan of pressuring him to alter student answers on the 2005 High School Proficiency Exams. Carruth was allegedly told to create his own answer key and change answer sheets after the test was administered. The test scores from the high school significantly dropped the following year. The state also investigated two elementary schools for alleged cheating. Michael Mimms, principal of Sumner Elementary, is put on administrative leave after it is discovered that he possessed opened copies of the 2006 TerraNova exam and distributed it to teachers.

(Kummers & Burney, 2006abc)

- 02/07/2006 Memphis Eyewitness News:
Teachers in Memphis schools are being investigated for test irregularities. According to the Tennessee Department of Education, an unusually high number of erasure marks were found on student exams. In many cases, incorrect answers were changed to correct answers.
(Memphis Eyewitness News, 2006)
- 01/12/2006 New York Daily News:
Fifth-grade students in Brooklyn were allegedly given actual copies of an exam to use as practice. Some students at Public School 58 in Cobble Hill reported that they recognized passages and questions from the test. Joyce Plus-Saly, the school principal, allegedly gave the materials to teachers to share with students, not knowing the questions would be used on the actual test.
(Lucadamo, 2006)
- 12/23/2005 WCBS-TV New York:
Ross Rosenfeld, a teacher at Junior High School 14 in Sheepshead Bay, was fired from his job after secretly recording conversations with the school principal. According to Rosenfeld, the recordings show that administrators ignored cheating on a state social studies exam. Rosenfeld was allegedly told to ignore a student who was found to have a cheat sheet during an exam.
(Lyon, 2005)
- 09/29/2005 Pittsburgh Post-Gazette:
Beth Boysza, a fourth-grade teacher in Pittsburgh, is suspended after being accused of helping her students on a math test in 2003. Boysza allegedly put Post-It® notes in the test booklets, providing students with special test instructions. She also is alleged to have re-read test questions to students. Boysza argues that she was simply providing accommodations to students, following directions provided by the district and test developer.
(Ove, 2005)
- 09/19/2005 The Courier-Journal in Louisville, Kentucky:
Following two cheating scandals, the Indiana Professional Standards Board increased the consequences for teachers who are caught helping their students cheat on tests. A teacher in Muncie, IN allegedly tapped her students on the shoulder to notify them of incorrect answers. A principal at Shakamak Elementary School in Jacksonville was found to have modified test questions and give them to students before the test administration. Both educators were caught after parents or state education officials noticed unusually large increases in school test scores.
(Hupp, 2005)
- 08/29/2005 Union-Tribune in San Diego, CA:
A teacher in Vista, CA was transferred to another school after allegations that she cheated on the California Standards Test. The teacher had allegedly put helpful materials on the classroom walls. Nearly half the students in the classroom reported that they had been told correct answers. The teacher was caught after a student reported the unusual behavior to her parents.
(Jenkins, 2005)
- 06/28/2005 Free Republic:
Isben Jeudy, a Long Island high school assistant principal, is arraigned after allegedly giving his son the answers to a history Regents exam. An official caught Jeudy's son with blue writing on his hand – the writing reportedly had answers to 35 of the exam questions.
(Eltman, 2005)
- 05/24/2005 St. Louis Post-Dispatch:
(excerpt) "Teachers prompted students with hand signals and pointed to answers."

(CEA, 2007, p. 13)

- 05/16/2005 Seattle Post Intelligence:
Lisa Poitras alleges that her daughter's teachers at Lake Dolloff Elementary have cheated on exams for two consecutive years. The teachers allegedly check student answers, give assistance, and urge students to make corrections on the Washington Assessment of Student Learning. Poitras is reported to have said "her daughter was made to erase and rewrite her answer to a question so many times that she wore a hole through the booklet page and had to reinforce it with scotch tape."
(Blanchard, 2005)
- 05/09/2005 Honolulu Advertiser:
The Hawaii Department of Education is investigating reports of cheating on the Hawaii State Assessment. Eighth-grade students were allegedly given test questions and answers to prepare for the test administration. An anonymous school employee notified the newspaper that teachers were given review sheets with actual test items on them.
(Shapiro, 2005)
- 05/04/2005 WHO TV in Des Moines, Iowa:
Gene Zwiefel, a seventh-grade teacher in the Adel district, resigns after allegations were made that he quizzed students on materials found in the actual Iowa Tests of Basic Skills. According to David Frisbie, director of the Iowa Testing Programs, similar incidences have occurred at four other Iowa Schools.
(WHO TV, 2005)
- 05/05/2005 Houston Chronicle:
(excerpt) "Teachers signaled students by tapping them on their shoulders to let them know an answer was wrong."
(CEA, 2007, p.13)
- 05/03/2005 Atlanta Journal-Constitution:
Following an investigation of cheating in Texas, Georgia begins an investigation of its own test results. While no high-profile cheating case emerged in Georgia, 159 educators were sanctioned for test administration problems in the past five years.
(Ghezzi, 2005)
- 05/03/2005 Star-Telegram in Texas:
Two teachers at A.M. Pate Elementary School are no longer working after allegedly giving students answers to the Texas Assessment of Knowledge and Skills. One of the teachers, Georgia Johnson (a 25-year veteran), had 18 of the 19 students in her class pass the test. Six of her students had perfect scores. The other teacher, Mildred Lawrence-Medearis (17 years experience), had all 29 of her students pass the reading and math exams.
(Garza, 2005)
- 04/13/2005 Rockford Register Star:
The Illinois Department of Education is investigating Tiffany Parker, principal of Lewis Lemon Elementary School in Rockford, for allegedly altering student answers in 2003.
(Watters, 2005)
- 04/30/2005 St. Louis Post-Dispatch:
"After the allotted time for testing, a teacher told students to fill-in answers for questions they had left blank."
(CEA, 2007, p. 13)

- 04/13/2005 NBC 6 in Miami, Florida:
The Florida Department of Education has reassigned Nicholas Emmanuel, principal of West View Middle School, after he allegedly helped students cheat on the Florida Comprehensive Assessment Test.
(NBC 6, 2005)
- 03/24/2005 Philadelphia Inquirer:
Shirley Neeley, Pennsylvania State Education Commissioner, moves to dissolve the Wilmer-Hutchins Independent School District board after 22 educators were found to have cheated on the Texas Assessment of Knowledge and Skills. The teachers allegedly ordered students who finished the test early to fix answers on other students' answer sheets.
(Mezzacappa et. al, 2005)
- 02/18/2005 The Ithaca Journal in Ithaca, New York:
Robert Blair, a fourth grade teacher with 19 years experience at Palmer Elementary School, resigns after administrators discover altered answer sheets on his students' state English Language Arts tests. Based on an analysis of erasures, 17 or 18 of the 22 students in his class had their answer sheets altered. The report states that there were 14 proven cases of teacher cheating in 2003-04 in New York.
(Associated Press, 2005)
- 01/31/2005 WRAL Raleigh-Durham, North Carolina:
Following rumors of test misconduct at Sallie B. Howard School for the Arts and Education, North Carolina administrators report there have been at least 10 investigations into testing irregularities. In that time, two teachers had their licenses revoked and a third case is in litigation.
(Carlson, 2005)
- 01/11/2005 Christian Science Monitor:
The Houston Independent School District launches an investigation into suspicious results on the 2004 administration of the TAKS. Recent examples of reported educator cheating include: (a) a third grade teacher in Indiana was suspended for allegedly tapping students on the shoulder when they marked wrong answers, (b) a fifth grade teacher in Mississippi was fired for allegedly helping students on the writing portion of a test, (c) nine Arizona school districts discarded test results because teachers allegedly read the test to students and gave students extra time.
(Axtman, 2005)

APPENDIX B: STATISTICAL DETECTION INDICES

Statistical methods to detect cheating do not, in fact, detect cheating. These methods, first developed in the 1920s to detect student cheating on multiple-choice tests, measure the likelihood of observing score gains, erasures, or answer patterns from student answer sheets. While the methods can identify unlikely large score gains, an improbably number of erasures, or unusual patterns of answers to items, they cannot determine if these events were due to cheating or simply due to chance.

Furthermore, the methods cannot identify all forms of cheating. They can only attempt to detect cheating due to educators giving answers to students or changing student answer sheets.

While most statistical detection methods were developed to detect student cheating, they can also be used to identify possible educator cheating. After all, if multiple students within a class or school are flagged as having unusual patterns of answers or erasures, the educator in charge of that class or school may have cheated.

Early Developments

Saupe (1960) summarizes the development of statistical detection methods to detect students who copy answers from other students. Bird (1927, 1929) developed three empirical approaches to detect possible copying in which the number of matching incorrect answers between two student tests is compared to the distribution of identical incorrect answers observed from a large random sample of answer sheet pairs (Saupe, 1960, p. 476). Because the number of incorrect answers depends upon the ability level of the student, the empirical distribution was based on random samples of test pairs from students with similar total scores to the suspected cheater. If the tests from the suspected cheater and the source student (from whom the cheater allegedly copied) were found to have an unusually large number of identical incorrect responses in comparison to this empirical distribution, then the suspected cheating could be verified.

In an application of his method, Bird describes a test administration in which the test proctor observed suspicious behaviors from four examinees. Bird calculated an average of 4.0 identical incorrect answers from a random sample of pairs of tests from examinees not suspected of cheating. The suspected cheaters had 17, 25, 28, and 31 identical incorrect answers on the 149-item test. As validation of his method, Bird reports that three of the suspected cheaters “confessed guilt when confronted with the evidence” (Bird, 1927, p.261).

From Empirical to Chance Models

Rather than taking the time to develop the empirical distribution, Dickenson (1945) developed a method to determine the likelihood of identical answers occurring by chance. This method simply compares the actual number of identical incorrect answers on a pair of answer sheets to the expected number based on the number of possible responses to each item. Under this method, it is assumed that each incorrect item response is equally likely to be chosen by students. If k is the number of possible responses to each item, then $(k - 1) / k^2$ is the expected proportion of incorrect answers on one test that are identical to another test. Dickenson suggested that if the observed proportion of identical incorrect answers is more than twice the expected proportion, then copying is implied (Saupe, 1960, p. 476).

Anikeeff developed another chance model using the binomial distribution to determine the likelihood of observing a specific number of identical incorrect answers between two tests. The number of observed identical incorrect answers between a pair of tests is compared to a binomial distribution with a mean of N , and a standard deviation of $\sqrt{Np(1 - p)}$, where N is the number of wrong responses by the suspected cheater and p is the reciprocal of the number of possible responses to each item (Saupe, 1960, p.

476). A low likelihood of observing that number of matching incorrect answers may indicate copying.

In an application of his method, Anikeeff concludes that his method is not effective at detecting copying. He concluded that his method would be useful in situations in which an examinee copies more than 16% of the answers from another examinee (Anikeeff, 1954).

Bellezza and Bellezza (1989) developed a method similar to Anikeeff's method called Error Similarity Analysis (ESA). This method, used by the *Scrutiny!* software package (Advanced Psychometrics, 1993), calculates the total number of times all pairs of examinees chose identical incorrect answers for each item. The probability of observing a given number of identical incorrect answers is estimated by the binomial distribution:

$$\frac{w!}{c!(w-c)!} P^c (1-P)^{w-c}$$

where c is the number of common items answered incorrectly by a pair of examinees, w is the number of items for which the pair of examinees had identical incorrect responses, and P is the estimated probability of two examinees selecting an identical incorrect answer. Using this equation, or a method based on a standard normal approximation, Bellezza and Bellezza were able to determine the likelihood of observing a specific number of identical incorrect answers between two examinees.

Holland (1996) describes another popular method to detect possible cheaters called the K-Index. This index, used by the Educational Testing Service (ETS), may be the most popular method currently used (Cizek, 1999, p. 141). Although limited information about this index exists, Holland describes it as a method used to "assess the degree of unusual agreement between the incorrect multiple-choice answers of two examinees" based on an estimate of the probability two examinees would agree on a response by chance (Holland, 1996, p. 5). The index uses the binomial distribution to

model this probability. Sotaridona (2001) developed two indices, S1 and S2, similar to the K-Index except using the poisson distribution to model the probability of observing a specific number of identical examinee responses.

Incorporating More Information

Acknowledging the limitation in methods that only analyze matching incorrect answers, Saupe (1960) developed his method to detect copying on multiple-choice tests.

In this method, the total number of items on the test, K , is partitioned into

$K = R_i + R_j - R_{ij} + W_{ij}$, where R_i and R_j are the number of correct responses for students i and j , respectively; R_{ij} is the number of items both students answered correctly; W_{ij} is the number of items both students answer incorrectly; and w_{ij} is the number of items in which both students gave matching incorrect answers.

Under chance conditions, the expected number of items answered correctly by both students would be the proportion of all items answered correctly by student i multiplied by the number of items answered correctly by student j : $ER_{ij} = \frac{1}{k} R_i R_j$. Thus, the regression of R_{ij} on the product $R_i R_j$ is of interest. This regression line can be written as: $\hat{R}_{ij} = b_{r1} R_i R_j + b_{r0}$.

The distance of an observed point $(R_i R_j, R_{ij})$ from the regression line can be used to evaluate the observed degree of correspondence between the items answered correctly by a pair of students. If this distance exceeds ts_r , where t is the appropriate value from the t-distribution and s_r is an appropriate estimate of the standard error of estimate of R_{ij} , then the assumption of chance correspondence can be rejected at a specified confidence level (assuming a bivariate normal distribution of R_{ij} and $R_i R_j$). A correspondence index can be written as:

$$CI = \frac{R_{ij} - b_{r1} R_i R_j - b_{r0}}{ts_r}.$$

A correspondence index greater than 1.00 is equivalent to rejecting the null hypothesis of chance correspondence between the items answered correctly by a pair of students.

The same logic is used to determine the correspondence between the incorrect answers from two students. If each item has k possible responses, the expected number of matching incorrect answers due to chance is:

$$Ew_{ij} = \frac{1}{k-1} W_{ij}.$$

Using the regression of w_{ij} on W_{ij} , the correspondence index would be:

$$CI = \frac{w_{ij} - b_{w1}W_{ij} - b_{w0}}{ts_w}$$

Saupe suggests an advantage to analyzing the correspondence of correct and incorrect answers separately is that the evidence provided by both indices is non-overlapping and, therefore, complementary. In applying his model to a random sample of 150 pairs of tests, Saupe's correspondence indices identified 6 suspicious pairs. In an attempt to validate the results, Saupe examined seating charts and discovered that 5 of the 6 suspicious pairs came from students in adjacent seats. Saupe admits the main disadvantage of his method is its use of a chance model – it is not reasonable to assume students randomly answer test questions (Saupe, 1960).

Attempting to overcome this disadvantage, Angoff (1974) developed 8 more indices to detect examinees copying on tests. Angoff's methods were all based on developing distributions of identical responses made by pairs of non-cheating examinees. The methods only differ in the combinations of independent and dependent variables used to develop the bivariate distributions. The degree to which an examinee's observed value on the dependent variable, conditioned on the observed value of the independent variable, deviates from the mean of the dependent variable from the distribution provides an index of cheating.

Angoff found that six of his indices were not effective in detecting cheating. Of the remaining indices, Angoff favored the method called the B Index. To use this index, the bivariate distribution of $W_i W_j$ and Q_{ij} is estimated from all examinees, where $W_i W_j$ is the product of the number of incorrect answers from two examinees and Q_{ij} is the number of identical incorrect answers for both examinees. For a pair of examinees, A and B , the observed values $W_a W_b$ and Q_{ab} are calculated. The following test statistic can then be used to determine whether the observed value of Q_{ab} is significantly different from the mean value of Q_{ij} :

$$t = \frac{Q_{ab} - \bar{Q}_{ij}}{S_{Q_{ij} W_i W_j}}$$

While Saupe and Angoff used information from both incorrect and correct responses, Frary (1977) developed two indices based on estimating the probability of an examinee choosing a correct response, choosing each incorrect response, or choosing to omit each item. After dismissing his first index, Frary developed the following formula for his g_2 index:

$$g_2 = \frac{C - \sum_i \hat{P}(k_{ia} = k_{ib})}{\sqrt{\sum_i \hat{P}(k_{ia} = k_{ib}) \left[1 - \sum_i \hat{P}(k_{ia} = k_{ib}) \right]}}$$

where C is the number of identical answers for a pair of examinees and $\hat{P}(k_{ia} = k_{ib})$ is the probability that an examinee would choose the identical response of another examinee. Frary used piecewise linear functions of total test scores to estimate this probability.

After applying his method to actual test data and recommending its use to prevent cheating, Frary acknowledged three limitations. First, in order to use his method, one examinee must be identified as “the copier” and another examinee must be identified as “the source.” This will not always be practical in large-scale testing situations. Second,

the g_2 index assumes that the probabilities of an examinee choosing each response to an item are constant, regardless of examinee ability. Third, Frary found that his method decreased in effectiveness for easier tests, stating, “If no examinees can answer as many as 90% correctly, the potential for detection is greatly enhanced” (Frary, 1977, p.253).

Hanson and Brennan (1987) continued to compare responses between pairs of examinees in their development of two more indices to detect possibly copying. The first method, Pair 1, uses the number of identical incorrect responses between a pair of examinees along with the length of the longest string of identical responses. The second method, Pair 2, uses the same information along with the percentage of maximum possible identical incorrect responses between two examinees.

In comparing their methods to the methods developed by Angoff (1974) and Frary (1977) on a simulated data set, Hanson and Brennan conclude that “it might not make a great deal of difference which of the statistical methods of investigating copying considered here are used” (p. 21). They do, however, recommend their method based on the interpretability of their indices.

Controlling for False Positives

In evaluating the effectiveness of the previously developed indices, Post (1994) concludes that while the indices may be used to scan for potential cheaters, “many existing statistical tests designed to detect copying on multiple-choice exams understate the Type I [false positive] error” (p. 140). Because this Type 1 error may be higher than specified, Post discourages using the indices to make accusations of cheating. Post attributes this inflated Type I error rate to the difficulty in estimating the probability of an examinee choosing each possible response to an item.

In an attempt to improve the estimation of item response probabilities and reduce the Type I error rate, Wesolowsky (2000) made a slight modification to Frary’s method. Whereas Frary used piecewise linear functions of raw scores to estimate probabilities,

Wesolowsky uses smooth distance iso-contours from location theory for estimation (p. 912). Also, while previous methods made assumptions about which examinee copied from another, Wesolowsky's method simply examines the number of matching answers and ignores other suspicious patterns such as strings of identical answers. In developing a computer program to analyze answer sheets and employing a Bonferroni adjustment to control for overall Type I error rate, Wesolowsky recommends his method as an effective way to screen for potential cheaters. This method, used in 2007 to scan for cheaters on the Texas Assessment of Knowledge and Skills, flagged more than 50,000 examinees as potentially having cheated (Benton & Hacker, 2007a, 2007b).

Incorporating Item Response Theory

Other researchers improved the estimation of the probability of an examinee choosing each possible response to an item by developing indices based on item response theory (IRT) models. In IRT models, the probability of an examinee choosing each response to an item is a function of the examinee's latent ability, θ , and characteristics of each possible item response. The item response characteristics of interest depend on the IRT model being used.

For the three-parameter logistic model, the probability of examinee a correctly answering dichotomously scored item i can be expressed as:

$$P_{ia} = P_{ia}(\theta_a) = P_{ia}(u_i = 1 | \theta) = c_i + \frac{1 - c_i}{1 + e^{-1.7a_i(\theta_a - b_i)}},$$

where

a_i = the item discrimination parameter,

b_i = the item difficulty parameter,

c_i = the guessing parameter,

θ_a = the latent ability of examinee a ,

and μ_i = the examinee's scored response to the item.

Using maximum likelihood or Bayesian methods, values for the item response and examinee ability parameters can be estimated under the assumptions of the specified IRT model. Under the assumption of local independence, the probability of observing a string of n item responses from an examinee with ability θ_a is equal to the product of the probabilities of the individual item responses:

$$P(u_1, u_2, \dots, u_n | \theta) = \prod_{i=1}^n P_i^{u_i} (1 - P_i)^{(1-u_i)}.$$

Thus, given an observed string of responses from an examinee, the above formula can be used to estimate the likelihood of observing that response string or the probability of observing a different string of responses.

Rather than simply estimating the probability of an examinee answering dichotomously scored items correctly or incorrectly, IRT models can be used to estimate the probability of an examinee choosing each possible choice to a multiple-choice item. Bock's Nominal Model calculates the probability of choosing response u on multiple-choice item g with m possible responses as:

$$P_u = \frac{\pi_u}{\sum_{h=1}^m \pi_h},$$

where π_h represents a scale value directly related to the probability that response h is chosen on a specific test item. This model can be reparameterized as:

$$P_u(\theta) = \frac{e^{(a_u \theta + c_u)}}{\sum_{h=1}^m e^{(a_h \theta + c_h)}},$$

from which item response discrimination and difficulty parameters, a and c , and examinee ability parameter θ can be estimated. Again, under the local independence assumption, the probability of observing a specific string of item responses can be estimated from these item response and examinee ability estimates.

Person-Fit and Aberrant Response Indices

The application of IRT models to detect possible cheating has been theorized via person-fit and aberrant response indices. These indices measure the extent to which the observed pattern of responses from an examinee with ability level θ deviates from the response pattern expected under the chosen IRT model. For example, an examinee whose ability exceeds the difficulty level (b) of an item would have a high probability of answering that item correctly. Likewise, an examinee whose ability is less than the difficulty of an item would have a high probability of answering that item incorrectly. When an examinee's response string fits this pattern across most or all items on the test, the model "fits" the person. Aberrant response strings (high-ability examinees incorrectly answering easy items, low-ability examinees correctly answering difficult items, or examinees choosing unusual responses on a multiple-choice test) indicate poor model fit. Person-fit and aberrant response indices measure the degree to which the chosen IRT model fits the observed responses from an examinee.

More than fifty person-fit indices have been developed to detect aberrant responders (Karabatsos, 2003; Meijer & Sijtsma, 2001; Thiessen, 2004). These indices, displayed in Figure B1, attempt to detect students who provide unusual responses due to luck, language deficiencies, random guessing, low-motivation, misaligned answer sheets, or cheating (Meijer, 1996).

Person-fit indices can be classified into three categories: deviation-based, covariance-based, and likelihood-based. Deviation-based indices, such as Wright and Stone's (1979) Outfit Mean Square index, sum the squared standardized differences between an examinee's scored response to an item and the expected probability of that correct response. A large difference would indicate a disagreement between the model and the examinee.

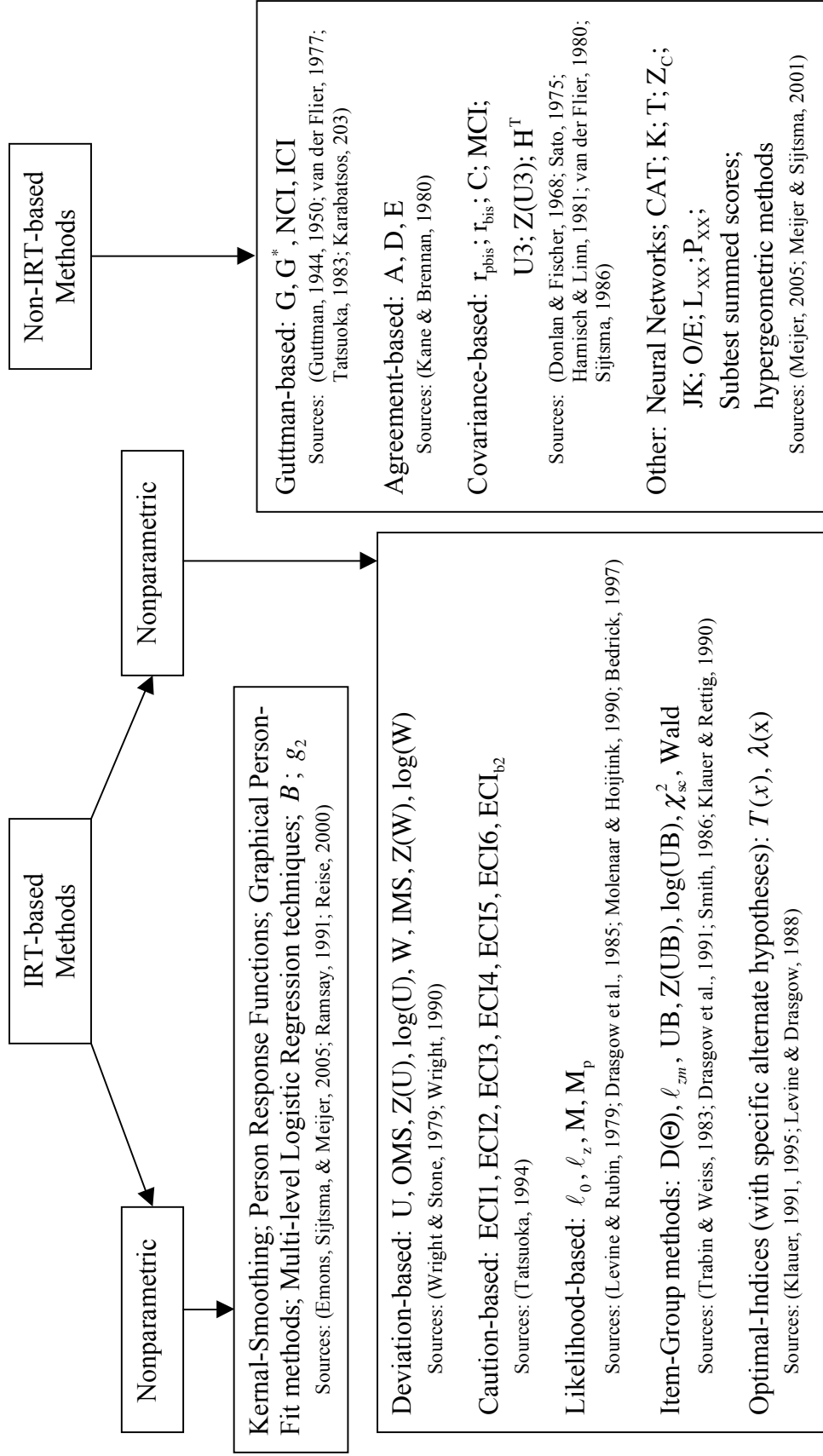


Figure B1 Aberrant Response Detection Methods and Indices

Note: See (Karabatsos, 2003; Meijer, 1996; Meijer & Sijtsma, 2001; Thiessen, 2004) for formulas and discussion of these methods.

Covariance-based indices, such as Tatsuoka's (1984) Extended Caution Indices, measure the degree to which an examinee's item responses deviate from the Guttman Perfect Pattern. Specifically, these indices calculate the ratio of the covariance between an examinee's responses and item difficulty estimates to the covariance between the average probability of correct responses across all examinees (estimated via IRT models) and the item difficulty estimates.

Likelihood-based indices, such as Levine and Rubin's (1979) log-likelihood function ℓ_0 evaluate the shape of a likelihood function. Given an examinee's responses to a set of test items with known item parameters (a, b, c), the ability level of an examinee, θ_a , can be estimated by maximizing the likelihood function:

$$\ell_0 = \ln[L(u|\theta; a, b, c)] = \ln\left[\prod_{i=1}^n P_i^{u_i} Q_i^{(1-u_i)}\right] = \sum_{i=1}^n [u_i \ln(P_i) + (1-u_i) \ln(Q_i)].$$

The value of θ_a that maximizes this function represents the examinee's most likely ability level given the observed item responses and the estimated item parameters. Examinees whose item responses conform to the IRT model produce likelihood functions with relatively high maximum values. Examinees whose responses deviate from what is predicted by the IRT model produce low maximum values of the likelihood function (Davey, et. al, 2003, p.6). Thus, the relative magnitude of ℓ_0 can be used to identify examinees whose responses are aberrant.

Aberrant Response Indices to Detect Examinee Cheating

While person-fit indices detect many forms of aberrance, some indices were developed specifically to detect examinee cheating. The first set of these indices is referred to as optimal person-fit statistics. Levine and Drasgow's (1988) $\lambda(x)$ index was developed to test the null hypothesis of normal examinee responses (based on a chosen IRT model) against the alternative hypothesis that an examinee's responses are consistent with a specific aberrant response model. Thus, if a researcher can specify and estimate

an alternative hypothesis of examinee cheating, optimal person-fit statistics can provide the maximum detection rate for aberrance (Karabatsos, 2003, p. 282).

Unfortunately, studies have shown that aberrant response indices are ineffective at detecting cheating (Chason & Maller, 1996; Iwamoto, Nungester, & Luecht, 1996). Demonstrating this ineffectiveness, test security firm Caveon analyzed a simulated data set using their six aberrant response indices. With the data set simulated so that examinees cheat on 10% of the test items, the firm's six aberrant response indices were only able to detect 41 (1.2%) of the 3,283 simulated cheaters and none of the five simulated cheating schools (Impara, et al., 2005).

Wollack argues that IRT-based aberrant response indices are inadequate in detecting cheating (specifically, student copying) because the statistical significance of these indices "does not depend on the similarity between the suspected copier's responses and those of a neighboring examinee" (p. 307). He stated that previous attempts to detect cheating, based on Classical Test Theory (CTT), are also inadequate. Wollack argues that since CTT item statistics are dependent on the sample of examinees tested, the expected similarity between a pair of examinees depends "largely on the performance of the other examinees in the sample, rather than only on the two examinees of interest" (p. 307). In order to overcome the apparent limitations of CTT-based indices and IRT-based aberrant response indices, Wollack developed his ω -index to specifically detect copying.

The ω -index is similar to Frary's g_2 index in that it attempts to estimate the probabilities associated with each possible item response. It differs from the g_2 index in that it uses Bock's Nominal Model to model these probabilities. The ω -index also attempts to control for Type I errors by examining the seating chart used in test administration and analyzing only pairs of examinees who sit physically close enough to make copying possible.

In applying the ω -index, each examinee is analyzed as a possible copier. If the test is administered in a classroom with a rectangular seating chart, the examinees sitting

to the copier's left, right, front left, front center, or front right are analyzed as potential sources for copying. Assuming item responses are locally independent, as is necessary to use IRT models, the expected number of identical responses and its associated variance can be modeled with the binomial distribution. Thus, an index similar to the g_2 can be compared to the standard normal distribution to evaluate the statistical significance of the degree of similarity between two examinees' item responses:

$$\omega = \frac{C - \sum_i \hat{P}(k_{ia} = k_{ib} | \theta_a, K_b, \xi)}{\sqrt{\sum_i \hat{P}(k_{ia} = k_{ib} | \theta_a, K_b, \xi) \left[1 - \sum_i \hat{P}(k_{ia} = k_{ib} | \theta_a, K_b, \xi) \right]}}$$

where C is the number of identical answers for a pair of examinees and

$\sum_i \hat{P}(k_{ia} = k_{ib} | \theta_a, K_b, \xi)$ is the probability that examinee a would choose the identical response of examinee b given the ability θ_a of examinee a , the item responses K_b of examinee b , and the matrix of item parameters, ξ .

Wollack (1997) compared his index to Frary's g_2 index on two data sets with three types of simulated copiers. The first type of copying was random copying, in which randomly selected items were simulated as being copied from another examinee. The second type was difficulty-weighted copying, in which the more difficult items were simulated as being copied. The third type was random-strings copying, in which strings of 4 consecutive items were simulated as being copied. Between 10- 40% of items were simulated as being copied by 5% of examinees for each type of copying. Wollack ensured that in this simulation, simulated copiers copied answers from examinees with higher ability estimates sitting close to them, as would be expected in reality. Based on these simulations, Wollack concludes that the ω -index is more effective in both detecting copying and controlling Type 1 error rates than the g_2 index when a seating chart is available (Wollack, 1997, p.312).

Adjacent Seating Methods

Similar to Wollack's ω -index, Kvam (1996), Roberts (1987), and the National Board of Medical Examiners (Cizek, 1999) developed methods that require knowledge about which examinees sat adjacent to others. In applying Kvam's method, examinees in a classroom are randomly administered two forms of the exam. Minor changes are made to the questions so that the two multiple-choice forms have different answers. After administering the exams, maximum likelihood methods are used to estimate the probability that an examinee copies an item response from an adjacently-seated source given that the examinee does not formulate an answer to the question (Kvam, 1996, p.239).

The score-difference method developed by Roberts (1987) also relies on two test forms randomly administered to examinees. After administering the test forms, each answer sheet is scored using the answer keys from both test forms. The difference between the scores obtained from the two answer keys provides an indication of possible cheating. While conceptually simple, Roberts concludes that his method is "seriously flawed and has little to recommend" (p. 77).

In administering its exams, the National Board of Medical Examiners (NBME) uses two methods to detect potential cheating. The first method, adjacent-nonadjacent analysis, requires the test to consist of at least two parallel parts. Each part must be administered in separate testing sessions in which test takers are randomly assigned to seat locations. Under these conditions, it is expected that the response patterns from any two examinees would have the same degree of similarity regardless of where the examinees are seated. If two examinees are seated adjacently during at least one test session and it is discovered that their responses are more similar in those sessions than would be expected if they were seated apart, then potential copying has been detected.

The NBME method uses a simple 2x2 chi-square test for independence to detect potential cheaters. The test statistic is calculated from Table B.1 using the formula given below, and then evaluated for significance:

$$\chi_1^2 = \frac{(a+b+c+d)(ad-bc)^2}{(a+b)(a+c)(b+d)(c+d)}.$$

Table B.1: Results from 2005 state and NAEP tests of 8th grade mathematics

Seating Location	Number of items answered incorrectly		Total
	Identical Responses	Different Responses	
Adjacent	<i>a</i>	<i>b</i>	<i>a+b</i>
Nonadjacent	<i>c</i>	<i>d</i>	<i>c+d</i>
Total	<i>a+c</i>	<i>b+d</i>	<i>a+b+c+d</i>

Methods to Detect Educator Cheating

The previously discussed methods all attempt to detect cheating (usually in the form of copying) by students on tests. Only two serious attempts have been made to detect potential educator cheating on tests. The first method, erasure analysis, is currently used by several states in auditing their testing programs. The second method, developed by Jacob and Levitt (2003), was used to detect cheating educators in Chicago Public Schools.

Some erasure analysis methods attempt to detect cheating by identifying answer sheets with an unusual number or pattern of erasures. It is assumed that a large number of erased answers might indicate an educator who is manipulating answer sheets. Other erasure analysis methods only look at wrong-to-right erasures – items in which an incorrect answer was erased and changed to a correct answer.

The Louisiana Educational Assessment Program has written and implemented erasure analysis procedures to audit its state testing program. These procedures require the test scoring contractor to scan every answer sheet for wrong-to-right erasures and

compute the mean and standard deviation for each subject at each grade level. The scoring contractor must then flag examinees whose wrong-to-right erasures exceed the state mean by more than four standard deviations as potential cheaters. The number of erasures for each potential cheater, along with the proportion of wrong-to-right erasures found on their answer sheet, are then reported to the State Superintendent of Education and the scores from the potential cheaters are voided from official records (Louisiana Educational Assessment Program, 2003).

One problem with using erasure analysis to flag potential cheaters is determining how many erasures indicate potential cheating. To address this problem, Qualls (2001) examined 4,553 answer sheets from the 1995-96 administration of the Iowa Tests of Basic Skills to determine the typical erasure behavior of examinees in a low-stakes testing environment. Qualls found that more than 90% of examinees changed 3 or fewer answers per test and only 2% of examinees changed more than 6 responses. While the erasure behavior depended on the test subject, less than 7% of item responses were erased, on average. Furthermore, the results indicate that it would be rare for an examinee to erase and change more than 15% of items on a single test. Qualls found that the first or second items on the test were most frequently erased and that students with one erasure had about a 50% chance of changing an incorrect response to a correct response. In a focus on wrong-to-right erasures, Qualls found that examinees gained an average of between 0.167 and 0.494 points per erasure on the tests. Perhaps these results could be used to develop new and improve upon current erasure analysis methods for high-stakes tests.

The other attempt to detect educator cheating on large-scale standardized tests was developed by Jacob and Levitt (2003). Jacob and Levitt set out to create a statistical index to identify educators, specifically teachers, who manipulate answersheets.

The method used by Jacob & Levitt to detect potential cheating teachers is actually a combination of two indicators: (1) unexpected test score fluctuations and (2)

unexpected patterns in student answers. Rather than focus on individual students, these indices are calculated at the classroom level.

Index #1: Unusual Test Score Fluctuations

On a vertically scaled test like the ITBS, the expected gain in student test scores from year-to-year can be estimated. The test score gains depend on a variety of factors (student ability, curriculum, teacher quality, etc.), but most students will experience growth in achievement each year. An unexpected test score fluctuation would occur when many students in a classroom experience large score gains one year followed by small gains (or even negative growth) the next year. To calculate this unexpected test score fluctuation index, Jacob & Levitt do the following:

- If the interest is in determining if a teacher manipulated answersheets during year t , then test scores must be collected for year t , the previous year ($t-1$), and the next year ($t+1$).
- Find the average test score gains (in grade equivalent units) for all students in each classroom (the gain from year $t-1$ to year t along with the gain from year t to year $t+1$).
- Find the percentile rank of each classroom's average test score gains relative to all other classrooms in that same subject, grade, and year. The percentile rank of growth from year $t-1$ to year t will be called rank gain_t while the rank of growth from year t to year $t+1$ will be called rank gain_{t+1} .
- The index is calculated as: $\text{Index \#1} = (\text{rank gain}_t)^2 + (1 - \text{rank gain}_{t+1})$. As the formula shows, classrooms with relatively large score gains followed by low gains the next year will yield large values of this index. The percentile ranks are squared in computing the index to give relatively more weight to big score gains and big score declines.

- Teachers whose classrooms yield values in the top 95th percentile of this index are identified as having unusual test score fluctuations.

Table B.2 shows an example of this index. The test score gains for three classrooms are displayed. Classroom 1 represents an average classroom, growing from 2.9 to 4.1 to 5.2 grade equivalent units over the course of two years. Classroom 2 represents a classroom with an excellent teacher (during year t). This excellent teacher was able to drastically increase test scores from year $t-1$ to year t . The teacher of Classroom 3 also was able to get drastic increases in test scores. The difference between Classrooms 2 and 3 can be seen by the changes in test scores once the students moved on to the next grade (with a different teacher). Whereas the students from Classroom 2 were able to build upon their score gains (as we would expect from students who were taught by an excellent teacher), the students in Classroom 3 experienced test score declines. This decline in test scores may indicate that the gains experienced the previous year did not represent genuine gains in achievement. Thus Classroom 3 earns a high score on Index #1 and the teacher of that classroom is identified as a potential cheater.

Table B.2: Example data for Index #1

Year	Avg. Classroom Grade Equivalent			Change in Grade Equivalent Units		Index #1 (Percentile Rank)
	$t-1$	t	$t+1$	From $t-1$ to t (percentile rank)	From t to $t+1$ (percentile rank)	
Classroom 1	2.9	4.1	5.2	1.2 (59)	1.1 (56)	0.5400 (20)
Classroom 2	3.4	5.5	6.8	2.1 (92)	1.3 (77)	0.8993 (70)
Classroom 3	3.1	5.2	4.7	2.1 (92)	-0.5 (1)	1.8265 (99)
Avg. Classroom	3.3	4.4	5.3	1.1	0.9	0.5000 (50)

Jacob & Levitt experienced with other measures of unusual test score fluctuations (such as regressing test score gains on previous gains and student demographics). They concluded that these other, more complicated measures yielded information similar to their simple index.

Index #2: Unexpected Patterns in Student Answers

Unusual fluctuations in test scores do not prove that a teacher manipulated student answersheets. In fact, since the 95th percentile is used as a cut-off, Index #1 will always identify 5% of the classrooms as having unusual fluctuations. To determine which of these classrooms cheated (and which just experienced improbable fluctuations), Jacob & Levitt developed a second index to identify unexpected patterns in student answers.

The logic is this: The quickest way for a teacher to cheat is to alter the same block of consecutive items for students in the class (or instruct students in the classroom to change their answers to the same set of items). Thus, if a classroom experiences unusual test score fluctuations and the students in the classroom have unusual answer patterns (identical answers to the same block of items or unexpected correct answers to difficult items), then we have more reason to believe the teacher cheated.

To identify unexpected answer patterns, the researchers combine four measures of suspicious answer strings to calculate Index #2. These four measures will be briefly discussed.

The first measure focuses on identifying the most unlikely block of identical answers given by students on consecutive items. Using a multinomial logit model, the likelihood of each student choosing each possible answer on every item is calculated. This likelihood is based on the student's past test scores, future test scores, and demographic (gender, race, etc). All combinations of students and consecutive items are

searched to find the block of identical answers that were least likely to have arisen by chance (controlling for classroom size).

Given student s in classroom c with answer j on item i , the model is:

$$P(Y_{isc} = j) = \frac{e^{\beta_j X_s}}{\sum_{j=1}^J e^{\beta_j X_s}},$$

Where X represents the vector of past test scores, future test scores, and demographics. The likelihood of a student's answer string for item m to item n is calculated as:

$$P_{sc}^{mn} = \prod_{i=m}^n P_{isc}.$$

This likelihood is multiplied across students in the class with identical responses in the string:

$$\tilde{P}_{sc}^{mn} = \prod_{\text{Students with identical strings}} P_{sc}^{mn}.$$

If each student in the classroom has unique responses from item m to item n , then there will be a distinct value of this index for each student in the class. If all students in the classroom have identical responses across these items, then there will only be one value of this index (and the value will be extremely small). The calculations are repeated for all strings from a length of 3 items to a length of 7 items.

Notice that the values yielded by these calculations will be smaller as: (1) the number of students with identical responses increase, (2) the length of the string of identical responses increase. Thus, smaller values are associated with more improbable answer strings within a classroom.

The minimum value for each classroom is recorded as measure #1:

$$\text{Measure \#1} = \min_s (\tilde{P}_{sc}^{mn}).$$

The second measure calculates the degree of correlation in student responses across the test, especially for unexpected answers. The logic is that teachers who cheat will have students with highly correlated answers. To calculate this measure, the residuals for each item choice are calculated:

$$e_{ijsc} = \begin{cases} 0 - P(Y_{isc}), & \text{for the unchosen options} \\ 1 - P(Y_{isc}), & \text{for the chosen answer} \end{cases}.$$

Then, the residuals for each option are summed across students within the classroom:

$$e_{jc} = \sum e_{ijsc}.$$

The option residuals for the classroom are then summed for each item. At the same time, the residuals are (1) squared to accentuate outliers and (2) divided by the number of students in the class to normalize for class size (n):

$$v_{ic} = \frac{\sum_j e_{jic}^2}{n}.$$

Measure #2 is simply the average of these item residual values:

$$\text{Measure \#2} = \bar{v} = \frac{\sum_i v_{ic}}{n}$$

Higher values indicate classrooms with highly correlated answers.

The third measure calculates the variance in the degree of correlation across test items. With Measure #2, we might expect high correlations among student answers if a teacher emphasizes certain topics during the school year. If a teacher cheats by changing answers for multiple students on selected questions, the within-class correlation on those particular questions will be extremely high, while the degree of within-class correlation on other questions is likely to be typical. Thus, a teacher who changes answers on selected

items will have a classroom with a large degree of variance in the correlation of responses across items.

This measure is calculated as the variance of item residuals from Measure #2:

$$\text{Measure \#3} = \sigma_v = \frac{\sum_i (v_{ic} - v_c)^2}{ni}.$$

The fourth measure compares the answers of students within a classroom to the answers from other equally-able students in the sample. This measure can then detect students who miss easy items while answering difficult items correctly. Students whose answers follow this pattern may have had their answers influenced by a cheating teacher.

To calculate this measure, students are grouped by their total number correct scores on the test. Let A_s represent a specific total correct score. Let $q_{ic} = 1$ if a particular student answers item i correctly and zero otherwise. Then determine the proportion of students with total score A_s who answered each item correctly (call this quantity \bar{q}_A).

The deviations between a student's item score and the expected item score (based on equally-abled students) are squared and summed across items:

$$Z_{sc} = \sum (q_{isc} - \bar{q}_A)^2.$$

This deviation between this Z-value for each student and the average Z-value for all equally-abled students is then summed for all students within a classroom:

$$\text{Measure \#4} = \sum (Z_{sc} - \bar{Z}_A)$$

High values of this index indicate the answers from a large number of students in the classroom deviated from equally-abled students in other classrooms.

After completing the calculations, the classrooms are ranked on each of the four measures. The percentile ranks for each classroom on each measure are then combined to form the second index:

$$\text{Index \#2} = (\text{Measure 1 rank})^2 + (\text{Measure 2 rank})^2 + (\text{Measure 3 rank})^2 + (\text{Measure 4 rank})^2$$

Classrooms falling above the 95th percentile on this index are identified as having unusual answer patterns.

Combining Indices to Detect Cheating Classrooms

Jacob & Levitt argued that taken individually, the above two indices do not detect teachers who manipulate answersheets. After all, there are always going to be (innocent) classrooms with unexpected score fluctuations and there are going to be (innocent) classrooms with improbable answer patterns. The key is to identify classrooms that yield high values on both indices.

In non-cheating classrooms, there is no reason to believe that the two indices would have a strong correlation. If a teacher manipulates student answersheets, we would expect a strong correlation between the two indices. Therefore, educators whose classrooms appear above the 95th percentile on both indices are identified as potential cheaters.

APPENDIX C: NATIONAL TESTING CODES AND STANDARDS

Source: *The Standards for Educational and Psychological Testing* (AERA et al., 1999)

Developer(s): American Educational Research Association, American Psychological Association, National Council on Measurement in Education

Statements related to inappropriate testing behaviors:

Validity, the degree to which evidence and theory support the interpretations of test scores entailed by proposed uses of tests, is the most fundamental consideration in developing and evaluating tests (p. 9).

The usefulness and interpretability of test scores require that a test be administered according to the developer's instructions. Maintaining test security also helps to ensure that no one has an unfair advantage (p. 61).

- 5.1 – Test administrators should follow carefully the standardized procedures for administration and scoring specified by the test developer
- 5.2 – Modifications or disruptions of standardized test administration procedures or scoring should be documented.
- 5.6 – Reasonable efforts should be made to assure the integrity of test scores by eliminating opportunities for test takers to attain scores by fraudulent means.
- 5.7 – Test users have the responsibility of protecting the security of test materials at all times (pp. 63-64).

Fairness requires all examinees to be given a comparable opportunity to demonstrate their standing on the construct(s) the test is intended to measure. Just treatment also includes such factors as appropriate testing conditions and equal opportunity to become familiar with the test format, practice materials, and so forth. Fairness also requires that all examinees be afforded appropriate testing conditions. Careful standardization of tests and administration conditions generally helps to assure that examinees have comparable opportunity to demonstrate the abilities or attributes to be measured (pp. 74-75).

Ideally, examinees would also be afforded equal opportunity to prepare for a test. Examinees should in any case be afforded equal access to materials provided by the testing organization and sponsor which describe the test content and purpose and offer specific familiarization and preparation for test taking. In addition to assuring equity in access to accepted resources for test preparation, this principle covers test security for nondisclosed tests (p. 75).

- 7.12 – The testing or assessment process should be carried out so that test takers receive comparable and equitable treatment during all phases of the testing or assessment process (p. 84)
- 11.7 – Test users have the responsibility to protect the security of tests, to the extent that developers enjoin users to do so
- 11.8 – Test users have the responsibility to respect test copyrights (p. 115)
- 13.11 – In educational settings, test users should ensure that any test preparation activities and materials provided to students will not adversely affect the validity of test score inferences.
- 13.12 – In educational settings, those who supervise others in test [administration] should have received education and training in testing necessary... (p. 148).
- 15.9 – The integrity of test results should be maintained by eliminating practices designed to raise test scores without improving performance on the construct or domain measured by the test (p. 168).

Source: *Code of Professional Responsibilities in Educational Measurement*, (Schmeiser et al., 1995)

Developer(s): National Council on Measurement in Education, National Association of Test Directors

Statements related to inappropriate testing behaviors:

Those who prepare individuals to take assessments and those who are directly or indirectly involved in the administration of assessments as part of the educational process, including teachers, administrators, and assessment personnel, have an important role in making sure that the assessments are administered in a fair and accurate manner. Persons who prepare others for, and those who administer, assessments have a professional responsibility to:

- 4.3: Take appropriate security precautions before, during, and after the administration of the assessment.
- 4.4: Understand the procedures needed to administer the assessment prior to administration
- 4.5: Administer standardized assessments according to prescribed procedures and conditions and notify appropriate persons if any nonstandard or delimiting conditions occur
- 4.6: Not exclude any eligible student from the assessment
- 4.7: Avoid any conditions in the conduct of the assessment that might invalidate the results
- 4.11: Avoid actions or conditions that would permit or encourage individuals or groups to receive scores that misrepresent their actual levels of attainment

Conducting research on or about assessments or educational programs is a key activity in helping to improve the understanding and use of assessments and educational programs. Persons who engage in the evaluation of educational programs or conduct research on assessments have a professional responsibility to:

- 8.3: Preserve the security of all assessments throughout the research process as appropriate

Source: *Code of Fair Testing Practices in Education*, (JCTP, 2004)

Developer(s): Joint Committee on Testing Practices, American Psychological Association, National Council on Measurement in Education, American Counseling Association, American Educational Research Association, American Speech-Language-Hearing Association, National Association of School Psychologists, National Association of Test Directors

Statements related to inappropriate testing behaviors:

Fairness implies that every test taker has the opportunity to prepare for the test (p. 2).

Test users should administer and score tests correctly and fairly.

- B-1: Follow established procedures for administering tests in a standardized manner.
- B-3: Provide test takers with an opportunity to become familiar with test question formats and any materials or equipment that may be used during testing (contradicts Popham)
- B-4: Protect the security of test materials, including respecting copyrights and eliminating opportunities for test takers to obtain scores by fraudulent means
- D-1: Inform test takers in advance of the test administration about the coverage of the test, the types of question formats, the directions, and appropriate test-taking strategies. Make such information available to all test takers

Source: *Standards for Teacher Competence in Educational Assessment of Students*, (NEA, 1990)

Developer(s): National Education Association

Statements related to inappropriate testing behaviors:

Requires teachers to recognize unethical, illegal, and inappropriate methods of assessment. Fairness, the rights of all concerned, and professional ethical behavior must undergird all student assessment activities, from the initial planning for and gathering of information to the interpretation, use, and communication of the results. Teachers must be well-versed in their own ethical and legal responsibilities in assessment. In addition, they should also attempt to have the inappropriate assessment practices of others discontinued whenever they are encountered. Teachers should also participate with the wider educational community in defining the limits of appropriate professional behavior in assessment.

Teachers who meet this standard will have the conceptual and application skills that follow. They will know those laws and case decisions which affect their classroom, school district, and state assessment practices. Teachers will be aware that various assessment procedures can be misused or overused resulting in harmful consequences such as embarrassing students, violating a student's right to confidentiality, and inappropriately using students' standardized achievement test scores to measure teaching effectiveness.

REFERENCES

(RED = citation from published news report)

(BLUE = source used in the paper)

(GREEN = source used only in statistical detection)

(BLACK = source not yet used in this paper)

Ad Hoc Committee on Confirming Test Results (2002). Using the National Assessment of Educational Progress to confirm state test results. Washington D.C.: National Assessment Governing Board. Retrieved September 14, 2007 from: http://www.nagb.org/pubs/color_document.pdf

Advanced Psychometrics (1993). *Scrutiny!* [Computer software], St. Paul, MN.

Aiken, L. R. (1991). Detecting, understanding, and controlling for cheating on tests. *Research in Higher Education*, 32(6), 725-736

American Educational Research Association, American Psychological Association, National Council on Measurement in Education (1999). *Standards for educational and psychological testing*. Washington DC: American Psychological Association

Angoff, W.H. (1974). The development of statistical indices for detecting cheaters. *Journal of the American Statistical Association*, 69(345), 44-49.

Anikeeff, A.M. (1954). Index of collaboration for test administrators. *The Journal of Applied Psychology*, 38, 174-177.

Asimov, N. (2007a). School mired in claims of cheating: former university prep teachers issue scathing report – state clamps down. Retrieved July 8, 2007 from *San Francisco Chronicle* web site: <http://www.sfgate.com/cgi-bin/article.cgi?f=/c/a/2007/07/08/UPREP.TMP&tsp=1>

Asimov, N. (2007b). Oakland principal in cheating stink quits. Retrieved July 16, 2007 from *San Francisco Chronicle* web site: <http://www.sfgate.com/cgi-bin/article.cgi?file=/c/a/2007/07/13/MNGIBR00831.DTL&type=printable>

Asimov, N. & Wallack, T. (2007). Cheating the test system. Retrieved May 13, 2007 from *San Francisco Chronicle* web site: <http://www.sfgate.com/cgi-bin/article.cgi?f=/c/a/2007/05/13/CHEATERS.TMP>

- Associated Press (2007a). Teacher on leave amid probe of alleged cheating. Retrieved May 11, 2007 from *ABC News* web site:
http://abclocal.go.com/wtv/story?section=nation_world&id=5296839&ft=print
- Associated Press (2005). Changed test answers lead to teacher resigning. Retrieved February 18, 2007 from *The Ithaca Journal* web site:
<http://www.theithacajournal.com/news/stories/20050218/localnews/2003002.html>
- Axtman, K. (2005). When tests' cheaters are the teachers: probe of Texas scores on high-stakes tests is the latest case in series of cheating incidents. Retrieved on January 11, 2005 from *Christian Science Monitor* web site:
<http://www.csmonitor.com/2005/0111/p01s03-ussc.html>
- Balassone, M. (2007). Teachers stumble, cheat on state tests. Retrieved August 20, 2007 from *The Modesto Bee* web site: <http://www.modbee.com/local/story/45979.html>
- Bay, L. (1995). *Detection of cheating on multiple-choice examinations*. Paper presented at the annual meeting of the American Educational Research Association.
- Belleza, F.S., & Belleza, S.F. (1989). Detection of cheating on multiple-choice tests: an update. *Teaching of Psychology*, 22(3), 180-182.
- Benton, J. (2006a). Analysis suggests cheating on TAKS. Retrieved May 23, 2006 from *The Dallas Morning News* web site:
<http://www.dallasnews.com/sharedcontent/dws/dn/education/stories/052306dnmetcheating.125e559b.html>
- Benton, J. (2006b). Cheating hasn't hurt Wilmer-Hutchins teachers. Retrieved October 1, 2006 from *The Dallas Morning News* web site:
<http://www.dallasnews.com/sharedcontent/dws/dn/education/stories/100106dnmetwilmercheaters.344f113.html>
- Benton, J. (2007a). FW charter school in trouble over TAKS cheating. Retrieved June 22, 2007 from *The Dallas Morning News* web site:
<http://www.dallasnews.com/sharedcontent/dws/dn/latestnews/stories/061507dnmetcheatinglee.3c44589.html>
- Benton, J. (2007b). TEA: teacher leaked parts of TAKS test. Retrieved July 16, 2007 from *The Dallas Morning News* web site:
<http://www.dallasnews.com/sharedcontent/dws/news/texasouthwest/stories/071307dntextaks.3938280.html>
- Benton, J. & Hacker H. (2007a). Analysis shows TAKS cheating rampant. Retrieved June 3, 2007 from *The Dallas Morning News* web site:
<http://www.dallasnews.com/sharedcontent/dws/news/dmn/stories/060307dnmetcheating.433e87c.html>

- Benton, J. & Hacker H. (2007b). Estimated number of cheaters might be low. Retrieved June 3, 2007 from *The Dallas Morning News* web site:
<http://www.dallasnews.com/sharedcontent/dws/dn/education/stories/060307dnmetconservative.3d3c27f.html>
- Bird, C. (1927). The detection of cheating in objective examinations. *School and Society*, 35, 261-262.
- Bird, C. (1929). An improved method of detecting cheating in objective examinations. *Journal of Educational Research*, 29, 341-348.
- Blanchard, J. (2005). Do teachers coach during WASL testing? Retrieved May 16, 2005 from *Seattle Post Intelligence* web site:
http://seattlepi.nwsourc.com/local/224429_was16.html
- Bosman, J. (2007). New report clears school of cheating. Retrieved June 27, 2007 from *New York Times* web site:
http://www.nytimes.com/2007/06/27/nyregion/27schools.html?_r=1&ref=nyregion&oref=slogin
- Bowers, W.J. (1964). *Student dishonesty and its control in college*. New York: Columbia University Bureau of Applied Social Research.
- Braun, H. & Qian, J. (2007). Mapping 2005 state proficiency standards onto the NAEP scales. U.S. Department of Education, NCES 2007-482. Retrieved September 6, 2007 from: <http://nces.ed.gov/nationsreportcard/pdf/studies/2007482.pdf>
- Brevard School District (2006). News release retrieved June 25, 2007 from *Brevard School District* web site: http://www.brevard.k12.fl.us/This_Week/releases/pdf/6-30-06%20cocoa%20high%20investigation.pdf
- Brock, K.C. (2006). Inquiry targets 20 area schools. Retrieved July 28, 2006 from *Dallas Star-Telegram* web site:
<http://www.dfw.com/mld/dfw/news/state/15148285.htm>
- Brumfield, R. (2005). High-stakes cheating spawns new market. Retrieved March 9, 2005 from *ESchool News* web site:
<http://www.eschoolnews.com/news/showstory.cfm?ArticleID=5564>
- Buckley, J. (2007). Mapping state standards to the NAEP scale. Presentation from the National Center for Education Statistics. Retrieved September 14, 2007 from:
<http://www.ccsso.org/content/PDFs/NAEPMappingBuckley.ppt>

- Bunn, D., Caudill, S., & Gropper, D. (1992). Crime in the classroom: an economic analysis of undergraduate student cheating behavior. *Research in Economic Education*, Spring, 197-207.
- Bureau of Justice Assistance: Center for Program Evaluation. (2007). Glossary. Retrieved September 11, 2007 from *BJA* web site: http://www.ojp.usdoj.gov/BJA/evaluation/glossary/glossary_h.htm
- Burns, J. G. (1988). Computers in class. *Teaching Professor*, 2(7), 2.
- Bush, G. (2000). Speech delivered at the 91st annual convention of the National Association for the Advancement of Colored People. Baltimore, Maryland. July 10, 2000. Retrieved September 3, 2007 from Washington Post web site: <http://www.washingtonpost.com/wp-srv/onpolitics/elections/bushtext071000.htm>
- Campbell, D. (1975). Assessing the impact of planned social change. *Social Research and Public Policies: The Dartmouth/OECD Conference*, Hanover, NH; Dartmouth College, The Public Affairs Center.
- Cannata, B. (2007). Two Stockton teachers accused of cheating. Retrieved August 28, 2007 from *CBS 13* web site: http://cbs13.com/topstories/local_story_236005351.html
- Carlson, K. (2005). WRAL investigates teachers accused of changing EOG test scores. Retrieved January 31, 2005 from WRAL web site: <http://www.wral.com/news/4148782/detail.html>
- Carnoy & Loeb (2002). Does external accountability affect student outcomes? A cross-state analysis. *Educational Evaluation and Policy Analysis*, 24(4), 305-331.
- Center for Evaluation and Assessment (2007). *Test Preparation: Considering the Appropriateness of these Activities: A Professional Development Module for Iowa Educators*. Retrieved August 23, 2007 from the Iowa Department of Education web site: <http://www.iowa.gov/educate/blogcategory/497/920/>
- Chadwick, D. (2006). The National Assessment of Educational Progress: The nation's report card. NAEP Short Review ICN. Retrieved September 12, 2007 from the Iowa Department of Education web site: http://www.iowa.gov/educate/index.php?option=com_docman&task=doc_download&gid=2204
- Chason, W. M., & Maller, S. (1996). *Utility of the Rasch person-fit statistic in detecting answer copying: A comparison with traditional cheating indices*. Paper presented at the annual meeting of the American Educational Research Association, New York.

- Cizek, G. (1999). *Cheating on tests: how to do it, detect it, and prevent it*. Mahwah, NJ: Lawrence Erlbaum.
- Cizek, G. J. (2001). *An overview of issues concerning cheating on large-scale tests*. Paper presented at the annual meeting of the National Council on Measurement in Education, Seattle, WA.
- Cizek, G. J. (2003). Educational testing integrity: why educators and students cheat and how to prevent it. *Measuring Up: Assessment Issues for Teachers, Counselors, and Administrators*. Retrieved August 8, 2007 from ERIC web site, No. ED 480 061:
<http://www.eric.ed.gov:80/ERICWebPortal/contentdelivery/servlet/ERICServlet?accno=ED480061>
- Cozin, M. (1999). *Cheating and its vicissitudes*. Research report from Raritan Valley Community College, NJ.
- Crocker, L. (2003). Teaching for the test: Validity, fairness, and moral action. *Educational Measurement: Issues and Practice*, 22(3), 5-11.
- Crocker, L. (2006). Preparing examinees for test taking: guidelines for test developers and test users. In S. Downing & T. Haladyna (Ed.), *Handbook of Test Development* (pp. 115-128). Mahwah, NJ: Lawrence Erlbaum Associates, Inc.
- Crown, D.F. & Spiller, M.S. (1998). Learning for the literature on collegiate cheating: a review of empirical research. *Journal of Business Ethics*, 17: 683-700.
- Cullen, J.B. & Reback, R. (2006). Tinkering toward accolades: School gaming under a performance accountability system. NBER Working Paper No. 12286. Retrieved September 9, 2007 from NBER web site: <http://www.nber.org/papers/w12286>
- Davey T., E. Stone, & R. McKinley (2003). Robust estimation of IRT parameters. Paper presented at the annual meeting of the American Educational Research Association, Chicago, April 21-25.
- Drasgow, F., Levine, M.V., Williams, E. A. (1985). Appropriateness measurement with polychotomous item response models and standardized indices. *British Journal of Mathematical and Statistical Psychology*. 38, 67-86.
- Dwyer, D. & Hecht, J. (1994). Cheating detection: statistical, legal, and policy implications. Illinois State University
- Education Week (2006). *Quality Counts At 10: A Decade of Standards-Based Education*. Vol. 25, Issue 17.

- Einhorn, E. & Melago, C. (2006). Bronx HS in cheat probe. Retrieved November 20, 2006 from *New York Daily News* web site:
<http://www.nydailynews.com/news/story/473012p-398009c.html>
- Elliott, S. (2007). Questions raised about materials that boosted test scores. Retrieved February 4, 2007 from *Dayton Daily News* web site:
<http://daytondaily.printthis.clickability.com/pt/cpt?action=cpt&title=...%2Flocal%2F2007%2F02%2F03%2Fddn020407citydayinside.html&partnerID=528>
- Elliott, S. (2007b). With proctors in class, City Day's test scores fall. Retrieved August 24, 2007 from *Dayton Daily News* web site:
<http://www.daytondailynews.com/n/content/oh/story/news/local/2007/08/16/ddn081607cityday.html>
- Eltman, F. (2005). Dad accused of giving son regents answers (principal helps his boy cheat). Retrieved June 24, 2007 from *Free Republic* web site:
<http://www.freerepublic.com/focus/f-news/1432174/posts>
- Emons, W.H.M., Sijtsma, K., & Meijer, R.R. (2004). Testing hypotheses about the person-response function in person-fit analysis. *Multivariate Behavioral Research*, 39(1), 1-35.
- Eve, R.A. & Bromley, D.G. (1981). Scholastic dishonesty among college undergraduates: Parallel tests of two sociological explanations. *Youth and Society*, 13, 3-22.
- FairTest.org (2006). FairTest Examiner. Retrieved June 25, 2007 from *FairTest.org* web site: <http://www.fairtest.org/examarts/August%202006/cheating806.html>
- Faul, F., Erdfelder, E., Lang, A. G., & Buchner, A. (2007). G*Power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*, 39, 175-191.
- Ferguson, A. (2007). Utah State University methods of psychology. Retrieved September 11, 2007 from USU web site:
<http://www.usu.edu/psycho101/lectures/chp2methods/methods.htm>
- Fessenden, F. (2007). Drop in school test scores raises further questions in Yonkers. Retrieved June 24, 2007 from *New York Times* web site:
<http://www.nytimes.com/2007/06/24/nyregion/nyregionspecial2/24peoplewe.html>
- Figlio, D. & Getzler, L. (2002). *Accountability, ability, and disability: gaming the system?* Working paper, University of Florida, 2002.
- Fisher, N. I. (1983). Graphical methods in nonparametric statistics: A review and annotated bibliography. *International Statistical Review*, 51(1), 25-58.

- Frary, R.B., Tideman, T.N., & Watts, T.M. (1977). Indices of cheating on multiple-choice tests. *Journal of Educational Statistics*, 2, 235-256.
- Frederickson, L. (1984). Teaching test-taking skills. *Social Studies Review*, 23(2), 23-28.
- Garcia, N. (2007). Testing honesty: has standardized testing turned some teachers into cheaters? Retrieved May 26, 2007 from *Visalia Delta-Times* web site: <http://www.visaliatimesdelta.com/apps/pbcs.dll/article?AID=/20070526/NEWS01/705260328>
- Garza, C. (2005). Students may have been helped. Retrieved April 19, 2005 from *Star-Telegram* web site: <http://www.dfw.com/mld/dfw/news/local/11431968.htm>
- Gay, G.H. (1990). Standardized tests: irregularities in administering of tests affect test results. *Journal of Instructional Psychology*, 17(2), 93-103.
- Ghezzi, P. (2005). Test stress and cheating. Retrieved May 3, 2005 from *Atlanta-Journal Constitution* web site: http://www.ajc.com/metro/content/custom/blogs/education/entries/2005/05/03/test_stress_and_cheating.html
- Gonen, Y. (2006). Test-mania fuels cheating at many schools, teachers say. Retrieved November 6, 2006 from *Staten Island Advance* web site: <http://www.silive.com/printer/printer.ssf?/base/news/1162819827228590.xml&coll=1&thispage=2>
- Gonen, Y. (2007). School big in 'test tamper.' Retrieved December 18, 2007 from *New York Post* web site: http://www.nypost.com/seven/12132007/news/regionalnews/school_big_in_test_tamper_892195.htm
- Green, J & Winters, M. (2003). Testing high stakes tests: can we believe the results of accountability tests? Civic Report 33, Center for Civic Innovation, Manhattan Institute, Manhattan, NY
- Hacker, H. (2005). Cause of diving TAKS scores unclear. Retrieved April 26, 2005 from *The Dallas Morning News* web site: <http://www.dallasnews.com/sharedcontent/dws/dn/education/stories/042605dnmetcheating.491c8d28.html>
- Haertel, E., Thrash, W., & Wiley, D. (1978). Metric-free distributional comparisons. Report. ML-Group for Policy Studies in Education, Chicago, IL.

- Haladyna, T.M., Haas, N.S., & Nolen, S.B. (1990). *Test score pollution*. Paper presented at the meeting of the American Educational Research Association, Boston, MA.
- Haladyna, T. M., Nolen, S. B., & Haas, N. S. (1991). Raising standardized achievement test scores and the origins of test score pollution. *Educational Researcher*, 20(5), 2-7.
- Haladyna, T.M. & Downing, S.M. (2004). Construct-irrelevant variance in high-stakes testing. *Educational Measurement: Issues and Practice*, 23(1), 17-27.
- Hall, D., & Kennedy, S. (2006). Primary Progress, Secondary Challenge. Report. Retrieved on September 12, 2007 from <http://www2.edtrust.org/NR/rdonlyres/15B22876-20C8-47B8-9AF4-FAB148A225AC/0/PPSCreport.pdf>
- Hall, J.L. & Kleine, P.F. (1992). Educators' perceptions of NRT misuse. *Educational Measurement: Issues and Practice*, 11(2), 18-22.
- Hamilton, L.S. & Stecher, B.M. (2006). Measuring instructional responses to standards-based accountability. Retrieved August 14, 2007 from RAND web site: https://rand.org/pubs/working_papers/2006/RAND_WR373.pdf
- Hanson, B.A. & Brennan, R.L. (1987). A comparison of several statistical methods for examining allegations of copying. *ACT Research Report Series No. 87-15*, Iowa City, IA: American College Testing
- Hanushek, E.A. & Raymond, M.E. (2004a). Does school accountability lead to improved student performance? NBER Working Paper No. W10591. Retrieved September 1, 2007 from Stanford University web site: <http://edpro.stanford.edu/Hanushek/admin/pages/files/uploads/accountability.jpamjournal.pdf>
- Hanushek, E.A. & Raymond, M.E. (2004b). The effect of school accountability systems on the level and distribution of student achievement. *Journal of the European Economic Association*, 2(2-3), 406-415.
- Hanushek, E.A. & Raymond, M.E. (2006). School accountability and school performance. Retrieved September 1, 2007 from Stanford University web site: <http://edpro.stanford.edu/hanushek/admin/pages/files/uploads/HanushekRaymond.pdf>
- Harcourt Assessment (2006). Legal policies. Retrieved September 30, 2007 from *Harcourt Assessment* web site: <http://harcourtassessment.com/hai/Templates/Generic/NoBoxTemplat...fTermsandConditionsofSale%2ehmt&NRCACHEHINT=NoModifyGuest#tcpu16>

- Harcourt Assessment (2007). Terms and conditions of sale. Retrieved September 30, 2007 from *Harcourt Assessment* web site:
<http://harcourtassessment.com/haiweb/Cultures/en-US/Harcourt/General/LegalPolicies.htm>
- Harrington-Lueker, D. (2000). When educators cheat. *The School Administrator*, December 2000.
- Hatch, J.A., & Freeman, E.B. (1988). Who's pushing whom? Stress and kindergarten. *Phi Delta Kappan*, 69, 145–147.
- Hedges, L.V., & Olkin, I. (1985). *Statistical Methods for Meta-Analysis*. Academic Press Orlando, San Diego.
- Hildebrand, J. (2007a). Complaints of test fraud rise. Retrieved June 25, 2007 from *Newsday* website: <http://www.newsday.com/search/ny-litest245268180jun24,0,4086000.story>
- Hildebrand, J. (2007b). Details emerge in Uniondale test scandal. Retrieved June 25, 2007 from *Newsday* website: <http://www.newsday.com/search/ny-litest0621,0,2578527.story>
- Hildebrand, J. (2007c). Security worries ahead of regents. Retrieved June 25, 2007 from *Newsday* web site: <http://www.newsday.com/search/ny-liaudi065244073jun06,0,1220786.story>
- Hill, R. (1998). Using NAEP to compare states' data – while it's still possible. Paper presented at the 1998 Annual Meeting of the National Council on Measurement in Education, San Diego.
- Ho, A.D. (2005). *Comparing Score Trends on High-Stakes and Low-Stakes Tests Using Metric-Free Statistics and Multidimensional Item Response Models*. Unpublished Doctoral Dissertation. Stanford University.
- Ho, A.D. (2007). Discrepancies between score trends from NAEP and state tests: A scale-invariant perspective. *Educational Measurement: Issues and Practice*, 26(4), 11-20.
- Ho, A.D. & Haertel, E.H. (2006a). Metric-free measures of test score trends and gaps with policy-relevant examples. CSE Technical Report #665. University of California, National Center for Research on Evaluation, Standards and Student Testing (CRESST), Los Angeles, CA.
- Ho, A.D. & Haertel, E.H. (2006b). *(Over)-Interpreting Mappings of State Performance Standards onto the NAEP Scale*. Retrieved September 6, 2007 from:

<http://www.ccsso.org/content/PDFs/Ho%20Haertel%20CCSSO%20Brief1%20Final.pdf>

Ho, A.D. & Haertel, E.H. (2007). *Apples to apples? The underlying assumptions of state-NAEP comparisons*. CCSSO Policy Brief. Retrieved February 23, 2008 from:
<http://www.ccsso.org/content/PDFs/Ho%20Haertel%20CCSSO%20Brief2%20Final.pdf>

Holland, P.W. (1996). Assessing unusual agreement between the incorrect answers of two examinees using the K-index: statistical theory and empirical support. *ETS Technical Report No. 96-4*. Princeton, NJ: Educational Testing Service

Holland, P.W. (2002). Two measures of change in the gaps between the CDFs of test-score distributions. *Journal of Educational and Behavioral Statistics*, 27(1): 3-17.

Holmgren, E. B. (1995). The p-p plot as a method for comparing treatment effects. *Journal of the American Statistical Association*, 90(429), 360-365.

Horne, L. V., & Gary, M. K. (1981, April). *What the test score really reflects: Observations of teacher behavior during standardized achievement test administration*. Paper presented at the Annual Meeting of the American Educational Research Association, Los Angeles, CA. (ERIC Document Reproduction Service No. ED 202 9 18)

Houser, B. (1982). Student cheating and attitude: a function of classroom control technique. *Contemporary Educational Psychology*, 7(2).

Hupp, S. (2005). Teachers who enable cheating to be fired. Retrieved September 19, 2005 from *The Courier-Journal* web site: <http://www.courier-journal.com/apps/pbcs.dll/article?AID=/20050919/NEWS02/509190353>

Hupp, S. (2006). Teachers gave out ISTEP answers. Retrieved October 11, 2006 from *The Indianapolis Star* web site:
<http://www.indystar.com/apps/pbcs.dll/article?AID=2006610110503>

Impara, J.C & Foster, D. (2006). Item and test development strategies to minimize test fraud. In S. Downing & T. Haladyna (Ed.), *Handbook of Test Development* (pp. 91-114). Mahwah, NJ: Lawrence Erlbaum Associates, Inc.

Impara, J.C., Kingsbury, G., Maynes, D., & Fitzgerald, C. (2005). Detecting cheating in computer adaptive tests using data forensics. Paper presented at the 2005 Annual Meeting of the National Council on Measurement in Education and the National Association of Test Directors, Montreal, Canada.

- Improving America's Schools Act of 1994, Pub. L. No. 103-382. (1994). Retrieved September 3, 2007 from The U.S. Department of Education web site: <http://www.ed.gov/legislation/ESEA/toc.html>
- Indy Channel (2006). Teacher suspended over ISTEP cheating allegations. Retrieved September 25, 2006 from *The Indy Channel.com* web site: <http://www.theindychannel.com/news/9932044/detail.html>
- Iowa Board of Educational Examiners (2004). *Code of Professional Conduct and Ethics*. Chapter 25 of the Licensure Rules (Iowa Administrative Code). Retrieved August 8, 2005 from the Iowa Board of Educational Examiners web site: <http://www.legis.state.ia.us/Rules/Current/iac/282iac/28225/28225.pdf>
- Iowa Department of Education (2007). Comparing NAEP and ITBS results. Retrieved September 18, 2007 from the Iowa Department of Education web site: http://www.iowaccess.org/educate/ecese/nclb/doc/comparing_naep_itbs.pdf
- Iowa Department of Education (2005). Sample board policy. Retrieved August 8, 2005 from the Iowa Department of Education web site: <http://www.education.uiowa.edu/itp/documents/DESAMPLEPOLICYSTATEMENT.PDF>
- Iowa Testing Programs (2005). *Guidance for developing district policy and rules on test use, test preparation, and test security for the Iowa Tests*. Retrieved August 8, 2005 from the Iowa Testing Programs web site: <http://www.education.uiowa.edu/itp/documents/ITPGuidanceDocument.pdf>
- Iwamoto, C.K., Nungester, R.J., & Luecht, R.M. (1996). *Power of similarity methods and person-fit analysis to detect copying behavior*. Paper presented at the annual meeting of the American Educational Research Association, New York.
- Jacob, B.A. (2007). Test-Based Accountability and Student Achievement: An Investigation of Differential Performance on NAEP and State Assessments. Working paper 12817, National Bureau of Economic Research. Retrieved September 1, 2007 from NBER web site: <http://www.nber.org/papers/w12817>
- Jacob, B. & Levitt, S. (2003). To catch a cheat. Retrieved August 8, 2007 from *Education Next* web site: <http://pricetheory.uchicago.edu/levitt/Papers/JacobLevittToCatchACheat2004.pdf>
- Jacob, B. & Levitt, S. (2004). Rotten apples: an investigation of the prevalence and predictors of teacher cheating. Retrieved August 8, 2007 from *Education Next* web site: <http://pricetheory.uchicago.edu/levitt/Papers/JacobLevitt2003.pdf>
- Jan, T. (2006). More teachers accused of cheating. Retrieved August 20, 2006 from *The Boston Globe*, web site:

http://www.boston.com/news/local/articles/2006/08/20/more_teachers_accused_of_cheating?mode=PF

Jan, T. (2007). Cheating on MCAS doubles. Retrieved November 5, 2007 from *Boston.com* web site:
http://www.boston.com/news/local/articles/2007/11/01/increased_cheating_reported_on_mcas/

Jeffrey, J. & Frisbie, D. (2005). Letter addressed to school administrators in Iowa, 8/15/2005. Retrieved August 8, 2005 from the Iowa Department of Education web site: <http://www.education.uiowa.edu/itp/documents/DE-ITPLetter.pdf>

Jenkins, L. (2005). With such high stakes, cheating is no surprise. Retrieved August 29, 2005 from the *Union-Tribune* web site:
<http://www.signonsandiego.com/news/northcounty/jenkins/20050829-9999-1m29jenkins.html>

Johnson, T.W. (2006). Elementary school teachers lose certificates. Retrieved July 4, 2006 from the *Baltimore Examiner* web site:
http://www.examiner.com/a-167297~Elementary_school_teachers_lose_certificates.html

Johnson, Z. K. (2007). Teacher 'help' crosses the line. Retrieved August 27, 2007 from the *Recordnet.com* web site:
http://www.recordnet.com/apps/pbcs.dll/article?AID=/20070823/A_NEWS/708230324

Joint Committee on Testing Practices (2004). *Code of Fair Testing Practices in Education (revised)*. Retrieved August 29, 2007 from APA web site:
<http://www.apa.org/science/fairtestcode.html>

Josephson Institute (2006). 2006 Josephson Institute report card on the ethics of American youth: Part one – integrity summary data. Retrieved September 18, 2007 from: <http://www.josephsoninstitute.org/reportcard/>

Julian, L. (2007). Performance-pay anxiety: teachers unions may actually do a kind deed for Florida schools. Retrieved December 18, 2007 from the *Orlando Sentinel* web site: <http://www.orlandosentinel.com/news/opinion/views/orliam1607dec16,0,3096449.story>

Kantrowitz, B. & McGinn, D. (2000). When teachers are cheaters. *Newsweek*, June 19, 2000.

Karabatsos, G. (2003) Comparing the aberrant response detection performance of thirty-six person-fit statistics. *Applied Measurement in Education* 16:4, 277-298.

- Kher-Durlabhji, N. & Lacina-Gifford, L.J. (1992). Quest for success: preservice teachers' views of "high stakes" tests. Retrieved August 11, 2007 from ERIC web site, ERIC # ED353338:
<http://www.eric.ed.gov:80/ERICWebPortal/contentdelivery/servlet/ERICServlet?accno=ED353338>
- Kimmel, E. (1997). *Unintended consequences or testing the integrity of teachers and students*. Paper presented at the annual Assessment Conference of the Council of Chief State School Officers, Colorado Springs, CO, June 1997.
- King, L. (2007). Data suggest states satisfy No Child Left Behind law by expecting less of students. Retrieved June 22, 2007 from *USA Today* web site:
http://www.usatoday.com/news/education/2007-06-06-schools-main_N.htm
- Klein, S. Hamilton, L., McCaffrey, D., & Stecher, B. (2000). What do test scores in Texas tell us? Report, The RAND Corporation, Santa Monica, CA. Retrieved September 5, 2007 from:
http://www.rand.org/pubs/issue_papers/IP202/index2.html
- Kohn, A. (1999). Tests that cheat students. *New York Times*, December 9, 1999.
- Kolen, M.J. & Brennan, R.L. (2004). *Test Equating, Scaling, and Linking: Methods and Practices* (2nd ed.). New York: Springer-Verlag.
- Koretz, D. (1999). Limitations in the use of achievement tests as measures of educators' productivity. Presentation at the *Devising Incentives to Promote Human Capital National Academy of the Sciences Conference*. Irvine, CA, December 18, 1999.
- Koretz, D. (2001). State comparisons using NAEP: large costs, disappointing benefits. *Educational Researcher*, 20(3), 19-21.
- Koretz, D. (2005). Alignment, high stakes, and the inflation of test scores. In J. L. Herman & E. H. Haertel (Eds.), *Uses and misuses of data for educational accountability and improvement* (104th Yearbook of the National Society of Education, Part 2, pp. 99-118). Malden, MA: Blackwell.
- Koretz, D., & Barron, S. (1998). The validity of gains on the Kentucky Instructional Results Information System (KIRIS). Report, The RAND Corporation, Santa Monica, CA.
- Koretz, D., & Hamilton, L. (2006). Testing for accountability in K-12. In R.L. Brennan (Ed.), *Educational Measurement* (4th ed., pp. 531-578). American Council on Education and Praeger Publishers. Westport, Connecticut.
- Koretz, D., McCaffrey, D., & Hamilton, L. (2001). Toward a framework for validating gains under high-stakes conditions. CSE Technical Report #551, University of

California, National Center for Research on Evaluation, Standards and Student Testing (CRESST), Los Angeles, CA.

- Kummer, F. & Burney, M. (2006a). How a test-rigging case came to light. Retrieved March 22, 2006 from *The Philadelphia Inquirer* web site: <http://www.philly.com/mld/inquirer/14155006.htm>
- Kummer, F. & Burney, M. (2006b). Test scores drop at Camden High School amid probe. Retrieved June 7, 2006 from *The Philadelphia Inquirer* web site: <http://www.philly.com/mld/inquirer/living/education/14758236.htm>
- Kummer, F. & Burney, M. (2006c). 2 schools' scores plummet. Retrieved June 16, 2006 from *The Philadelphia Inquirer* web site: http://www.philly.com/mld/philly/entertainment/family_guide/14831499.htm
- Kvam, P. H. (1996). Using exam scores to estimate the prevalence of classroom cheating. *The American Statistician*, 50(3), 238-242.
- Laffont, J.J. & Martimort, D. (2001). The theory of incentives: The principal-agent model. Princeton: Princeton University Press.
- Lai, E. & Waltman, K. (2007). High stakes testing and test preparation: examining teacher beliefs and practices. *Center for Evaluation and Assessment*, University of Iowa. Paper presented at the annual meeting of the Iowa Educational Research and Evaluation Association, December 2006. Retrieved June 21, 2007 from the *Center for Evaluation and Assessment* web site: http://www.education.uiowa.edu/cea/documents/Test_Prep_AERA_2007_and_IEREA_2006.pdf
- Lee, J. (2006) Tracking Achievement Gaps and Assessing the Impact of NCLB on the Gaps: An In-Depth Look into National and State Reading and Math Outcome Trends. Civil Rights Project. Harvard University, Cambridge, MA.
- Levesque, S. (2006). Big Spring's TAKS tests flagged. Retrieved June 9, 2006 from *Abilene Reporter-News* web site: http://www.reporter-news.com/abil/nw_ed_elem_secondary/article/0,1874,ABIL_7951_4762000,00.html
- Levine, M.V. & Drasgow, F. (1988). Optimal appropriateness measurement. *Psychometrika*, 53, 161-176.
- Levine, M.V. & Rubin, D.B. (1979). Measuring the appropriateness of multiple choice test scores. *Journal of Educational Statistics*, 4, 269-290.
- Ligon, G.D. (2000). Trouble with a capital t. *School Administrator*, 57(11), 40-44.

- Ligon, G.D. & Jones, P. (1982). *Preparing students for standardized testing: one district's perspective*. Paper presented at the meeting of the American Educational Research Association, New York.
- Linn, R.L. (2005). Adjusting for differences in tests. Paper presented at a Symposium on the Use of School-Level Data for Evaluating Federal Education Programs, Washington DC: The Board on Testing and Assessment, The National Academies, December 9, 2005. Retrieved September 15, 2007 from http://www7.nationalacademies.org/bota/School-Level%20Data_Robert%20Linn-Paper.pdf
- Linn, R.L. (2000). Assessments and accountability. *Educational Researcher*, 23(9), 4-16.
- Linn, R.L. (2003). Performance standards: Utility for different uses of assessments. Education Policy Analysis Archives, 11(31). Retrieved March 15, 2007 from <http://epaa.asu.edu/epaa/v11n31/>
- Linn, R.L., Baker, E.L., & Betebenner, D.W. (2002). Accountability systems: Implications of requirements of the No Child Left Behind Act of 2001. *Educational Researcher*, 31(4): 431-435.
- Linn, R., Graue, M., & Sanders, N. (1990). Comparing state and district results to national norms: The validity of claims that "everyone is above average". *Educational Measurement: Issues and Practice*, 9(3):5-14.
- Livingston, S. A. (2006). Double p-p plots for comparing differences between two groups. *Journal of Educational and Behavioral Statistics*, 31(4), 431-435.
- Loeb, S. & Strunk, K. (2005). Accountability and local control: Response to incentives with and without authority over resource generation and allocation. Retrieved September 1, 2007 from Stanford University web site: <http://www.stanford.edu/~sloeb/Papers/LoebandStrunkAccountability.pdf>
- Loomis, S.C. & Bourque, M.L. (Eds.) (2001). *National Assessment of Educational Progress achievement levels 1992-1998 for reading*. Washington, D.C.: U.S. Department of Education, National Assessment Governing Board. Retrieved September 13, 2007 from <http://www.nagb.org/pubs/readingbook.pdf>
- Lord, F.M. & Novick, M.R. (1968). *Statistical theories of mental test scores*. Reading, MA: Addison Wesley.
- Louisiana Educational Assessment Program (2003). *Erasure Analysis Procedures*. Retrieved August 28, 2007 from Louisiana Department of Education web site: <http://www.doe.state.la.us/lde/uploads/1607.pdf>.

- Lucadamo, K. (2006). An early read on test Ed Dept. probe. Retrieved January 12, 2006 from *New York Daily News* web site:
<http://www.nydailynews.com/news/local/story/382127p-324461c.html>
- Lyon, K. (2005). Ex-teacher blows whistle in cheating scandal. Retrieved December 23, 2005 from *WCBS-TV New York* web site:
http://webstv.com/topstories/local_story_357172153.html
- Magnuson, P. (2000). High-stakes cheating: will the focus on accountability lead to more cheating. *Communicator*, February 2000.
- Marcus, D. (2007a). Why some teachers cheat. Retrieved June 25, 2007 from *Newsday* website: <http://www.newsday.com/news/local/longisland/ny-liince0624,0,5126924.story?coll=ny-top-headlines>
- Marcus, D. (2007b). Experts called to sniff out fraud. Retrieved June 25, 2007 from *Newsday* website: <http://www.newsday.com/search/ny-lifrau0624,0,2648649.story>
- Matter, M.K. (1986). Legitimate ways to prepare students for testing: being up front to protect your behind. *National Association of Test Directors 1986 symposia*. Oklahoma City, Oklahoma City Public Schools.
- McCabe, D., Treviño, L. (1993). Honor codes and other contextual influences. *Journal of Higher Education*, 64, 522-538.
- McCabe, D., Treviño, L. (2002). Honesty and honor codes. Retrieved September 30, 2007 from *AAUP* web site:
<http://www.aaup.org/publications/Academe/2002/02JF/02jfmcc.htm>.
- McCabe, D., Treviño, L. & Butterfield, K. (2001). Cheating in academic institutions: a decade of research. *Ethics and Behavior*, 11(3), 219-232.
- McCollum, D. (2007). Allegations of TAKS cheating. Retrieved June 24, 2007 from *KLTV-7* website: <http://www.kltv.com/Global/story.asp?S=6690284>
- McLaughlin, D., Bandeira de Mello, V., Cole, S., Blankenship, C., Hikawa, H., Farr, K., & González, R. (2002). Comparison of National Assessment of Educational Progress and statewide assessment results: Report to Maryland on 1996 and 1998 assessments. Washington, DC: American Institutes for Research. Retrieved September 14, 2007 from:
http://eric.ed.gov/ERICDocs/data/ericdocs2sql/content_storage_01/0000019b/80/16/8d/94.pdf
- McLaughlin, D., Bandeira de Mello, V., Cole, S., Blankenship, C., Hikawa, H., Farr, K., & González, R. (2002). National Longitudinal School-Level State Assessment

Database: Analyses of 2000/2001 school year scores. Washington, DC: American Institutes for Research.

- Mehrens, W.A. & Kaminski, J. (1989). Methods for improving standardized test scores: fruitful, fruitless, or fraudulent? *Educational Measurement: Issues and Practice*, 8(1), 14-22.
- Mehrens, W.A., Phillips, S., & Schram, C. (1993). Survey of test security practices. *Educational Measurement: Issues and Practices*, 12(4), 5-19.
- Meijer, R.R. (1996). Person-fit research: an introduction. *Applied Measurement in Education*, 9(1), 3-8
- Meijer R.R. (2005) Using patterns of summed scores in paper-and-pencil and CAT to detect misfitting item score patterns. *Law School Admission Council Computerized Testing Report 02-04*, December, 2005.
- Meijer R.R., & Sijtsma K. (2001) Methodology review: evaluating person fit. *Applied Psychological Measurement* 25(2), 107-135
- Memphis Eyewitness News (2006). Superintendent: Teachers changed student test answers. Retrieved February 7, 2006 from *Memphis Eyewitness News* web site: http://www.myeeyewitnessnews.com/news/local/story.aspx?content_id=372B8005-42B9-47E0-8E66-5EECF4585D6
- Messick, S. (1984). The psychology of educational measurement. *Journal of Educational Measurement*, 21, 215-237..
- Mezzacappa, D., Langland, C., & Hardy, D. (2005). Principal accused of cheating on tests. Retrieved April 19, 2005 from *The Philadelphia Inquirer* web site: <http://www.philly.com/mld/inquirer/news/local/states/pennsylvania/11412829.ht>
- Million, J. (2000). When a principal cheats. Retrieved September 14, 2003 from *National Association of Elementary School Principals* web site: <http://www.naesp.org/ContentLoad.do?contentId=231>
- Moore, J.L. & Waltman, K. (2007). Pressure to increase test scores in reaction to NCLB: An investigation of related factors. Paper presented at the annual meeting for the American Educational Research and Evaluation Association, April 2007.
- Moore, W.P. (1994). Appropriate test preparation: can we reach a consensus? *Educational Assessment*, 2(1), 51-68.
- Morris, E. (2007). Manatee teacher's future remains uncertain. Retrieved August 17, 2007 from *Herald-Tribune* web site: <http://www.heraldtribune.com/article/20070810/NEWS/708100515>

- Mosquin, P. & Chromy, J. (2004). Federal sample sizes for confirmation of state tests in the No Child Left Behind Act. Washington, D.C.: American Institutes for Research, NAEP Validity Studies Panel. Retrieved September 13, 2007 from http://www.air.org/publications/documents/MosquinChromy_AIR1.pdf
- Mrozowski, J. (2007). MEAP leak forces retest for thousands of students. Retrieved October 16, 2007 from *The Detroit News* web site: <http://www.detroitnews.com/apps/pbcs.dll/article?AID=/20071012/SCHOOLS/710120406/1026>
- MSNBC (2006). 2 million scores ignored under 'No Child' loophole. Retrieved August 21, 2007 from *MSNBC* web site: <http://www.msnbc.msn.com/id/12357165/from/RSS/>
- Muehlhausen, N. (2007). Investigation details teachers cheating on standardized tests. Retrieved November 12, 2007 from *KSTP Channel 5 Eyewitness News* web site: <http://kstp.com/article/stories/S253627.shtml?cat=5>
- Murphy, K. (2007). Teachers to lose jobs over test help. Retrieved May 20, 2007 from *Inside Bay Area*, web site: http://www.insidebayarea.com/portlet/article/html/fragments/print_article.jsp?articleId=5942288&siteId=181
- Muthen, B.O., Khoo, S.T., & Goff, G.N. (1997). Multidimensional description of subgroup differences in mathematics achievement data from the 1992 National Assessment of Educational Progress. CSE Technical Report #432, University of California, National Center for Research on Evaluation, Standards, and Student Testing (CRESST), Los Angeles, CA.
- Nathanson, C., Paulhus, D. & Williams, K. (2--5). *Predictors of behavioral measure of scholastic cheating: personality and competence but not demographics*. Unpublished.
- National Association of Test Directors (2004). Cleaning up answer sheets. Retrieved August 11, 2007 from http://www.natd.org/Case_3_Cherry_Creek_part_E.pdf
- National Center for Education Statistics (NCES) (2007). Item development process. Retrieved September 12, 2007 from the NCES web site: http://nces.ed.gov/nationsreportcard/contracts/item_dev.asp
- National Education Association (1990). *Standards for Teacher Competence in Educational Assessment of Students*. Retrieved June 22, 2007 from *Buros Institute* web site: <http://www.unl.edu/buros/bimm/html/article3.html>

- NBC6 (2005). Miami-Dade principal faces cheating allegations. Retrieved April 12, 2005 from *NBC 6* web site: <http://www.nbc6.net/education/4370439/detail.html>
- Neal, D. & Schanzenbach, D.W. (2007). Left behind by design: Proficiency counts and test-based accountability. NBER Working Paper No. 13293. Retrieved September 9, 2007 from University of Chicago web site: http://home.uchicago.edu/~n9na/web_ver_final.pdf
- Nelson, T.D. & Schaefer, N. (1986). Cheating among college students estimated with the randomized-response technique. *College Student Journal*, 20, 321-325.
- Newberger, E.H. (2003). Why do students cheat? Retrieved August 12, 2007 from web site: http://www.school-for-champions.com/character/newberger_cheating2.htm
- Nichols, S.L. & Berliner, D.C. (2005). The inevitable corruption of indicators and educators through high-stakes testing. Education Policy Studies Laboratory, Arizona State University, March, 2005. Retrieved March 12, 2005 from <http://www.greatlakescenter.org/pdf/EPSTL-0503-101-EPRU.pdf>
- Nichols, S.L. & Berliner, D.C. (2007). *Collateral damage: how high-stakes testing corrupts America's schools*. Cambridge, MA: Harvard Education Press.
- Niels, G. (1996). Is the honor code a solution to the cheating epidemic. Research report from The Klingenstein Center, Teachers College, Columbia University, NY.
- No Child Left Behind Act of 2001, Pub. L. No. 107-110. (2001). Retrieved September 3, 2007 from The U.S. Department of Education web site: <http://www.ed.gov/policy/elsec/leg/esea02/index.html>
- Nolen, S.B., Haladyna, T.M., & Haas, N.S. (1992). Uses and abuses of achievement test scores. *Educational Measurement: Issues and Practice*, 11(2) 9-15.
- Ohler, A. (2007). Cheating investigation at West Leyden Elementary. Retrieved July 10, 2007 from *News 10 Now* web site: <http://news10now.com/shared/print/default.asp?ArID=111329>
- Ove, T. (2005). Teacher testifies book indicated she could help pupils with test. Retrieved September 29, 2005 from *Pittsburgh Post-Gazette* web site: <http://www.post-gazette.com/pg/05272/579553.stm>
- Parsavand, S. (2007). Corona, Moreno Valley teachers violate state testing rules. Retrieved September 10, 2007 from *The Press-Enterprise* web site: http://www.pe.com/localnews/inland/stories/PE_News_Local_H_cheat08.3df89de.html

- Passow, H.J., Mayhew, M.J., Finelli, C.J., Harding, T.S., & Carpenter, D.D. (2006). Factors influencing engineering students' decisions to cheat by type of assessment. *Research in Higher Education*, 47(6), 643-684.
- Patrick, K. & Eichel, L. (2006). Education tests: Who's minding the scores? Retrieved June 25, 2006 from *The Philadelphia Inquirer* web site: <http://www.philly.com/mld/inquirer/living/education/14898076.htm>
- Paulhus, D.L. (1991). BIDR reference manual, Version 6. Vancouver, Canada: University of British Columbia, Department of Psychology. Retrieved September 10, 2007 from: http://www.ori.dhhs.gov/education/products/n_illinois_u/datamanagement/dmglossary.html
- Pedulla, J.J., Abrams, L.M., Madaus, G.F., Russell, M.K., Ramos, M.A., & Miao, J. (2003). *Perceived effects of state-mandated testing programs on teaching and learning: Findings from a national survey of teachers*. Boston, MA: Boston College, National Board on Educational Testing and Public Policy
- Pellegrino, J.W., Jones, L.R., & Mitchell, K.J. (Eds.) (1998). *Grading the Nation's Report Card: Evaluating NAEP and transforming the assessment of educational progress*. Washington D.C.: National Academy Press. Retrieved September 14, 2007 from <http://www.eric.ed.gov/contentdelivery/servlet/ERICServlet?accno=ED446096>
- Perlman, C. L. (1985, March). Results of a citywide testing program audit in Chicago. Paper presented at the annual meeting of the American Educational Research Association, Chicago, IL. (ERIC Document Reproduction Service No. ED 263 212), pp. 4-5.
- Perlman, C. L. (2003). Practice tests and study guides: Do they help? Are they ethical? What is ethical test preparation practice? *Measuring Up: Assessment Issues for Teachers, Counselors, and Administrators*, ERIC Document Reproduction Service No. ED 480 062
- Peterson, P. & Hess, F. (2005). Johnny can read... in some states. Retrieved June 22, 2007 from *Hoover Institution* web site: <http://www.hoover.org/publications/ednext/3219636.html>
- Peterson, P. & Hess, F. (2006). Keeping an eye on state standards. Retrieved June 22, 2007 from *Hoover Institution* web site: <http://www.hoover.org/publications/ednext/3211601.html>
- Popham, W.J. (1991). Appropriateness of teachers' test-preparation practices. *Educational Measurement: Issues and Practice*, 10(4), 12-15.

- Post, G.V. (1994). A quantal choice model for the detection of copying on multiple choice examinations. *Decision Sciences*, 25(1), 123-142.
- Qualls, A.L. (2001). Can knowledge of erasure behavior be used as an indicator of possible cheating? *Educational Measurement: Issues and Practice*, 20(1), 9-16.
- Ramsey, J.O. (1991). Kernel smoothing approaches to nonparametric item characteristic curve estimation. *Psychometrika*, 56(4), 611-630.
- Ravitch, D. (2005). Every state left behind. *The New York Times*, November 7, 2005. Retrieved September 13, 2007 from The Brookings Institution web site: <http://www.brookings.edu/views/op-ed/ravitch/20051107.htm>
- Reback, R. (2007). Teaching to the rating: School accountability and the distribution of student achievement. NBER Working Paper No. WP0602. Retrieved September 1, 2007 from NBER web site: <http://www.nber.org/papers/wp0602>
- Reise, S.P. (2000). Using multilevel logistic regression to evaluate person-fit in IRT models. *Multivariate Behavioral Research*, 35(4), 543-568.
- Rhodes, S. (2007). WMC-TV Memphis. Germanshire parents sound off over cheating allegations. Retrieved May 22, 2007 from *WMC-TV* web site: <http://www.wmcstations.com/global/story.asp?s=6527799&ClientType=Printable>
- Richards, J.S. (2006a). Cheating is up – among teachers. Retrieved October 22, 2006 from *The Columbus Dispatch* web site: <http://www.columbusdispatch.com/news-story.php?story=dispatch/2006/10/22/20061022-A1-01.html>
- Richards, J.S. (2006b). More districts looking for test violations. Retrieved April 11, 2006 from *The Columbus Dispatch* web site: <http://www.columbusdispatch.com/news-story.php?story=dispatch/2006/04/11/20060411-A1-04.html>
- Riverside Publishing (2006). Sales center: Test security policy. Retrieved September 30, 2007 from Riverside Publishing web site: <http://www.riversidepublishing.com/sales/security.html>
- Roberts, P., Anderson, J., & Yanish, P. (1997). *Academic misconduct: where do we start?* Paper presented at the annual conference of the Northern Rocky Mountain Educational Research Association, Jackson, WY, October 1997.
- Roberts, D.M. (1987). Limitations of the score-difference method in detecting cheating in recognition test situations. *Journal of Educational Measurement*, 24(1), 77-81.
- Rooks, C. (1998). www.2cheat.com. Paper presented to the Teaching in the Community Colleges Online Conference, Honolulu, HI, April 7-9, 1998.

- Rotherham, A. (1999). Toward performance-based federal education funding: Reauthorization of the Elementary and Secondary Education Act. Progressive Policy Institute Policy Report, April 1, 1999. Retrieved September 3, 2007 from the Democratic Leadership Council web site: <http://www.ndol.org/documents/ESEA.pdf>
- Saupe, J.L. (1960). An empirical model for the corroboration of suspected cheating on multiple-choice tests. *Educational and Psychological Measurement*, 20, 475-489.
- Scheers, N.J. & Dayton, C.M. (1987). Improved estimation of academics cheating behavior using the randomized response technique. *Research in Higher Education*, 26, 61-69.
- Schmelkin, L, Kaufman, A., & Liebling, D. (2001). *Faculty assessments of the clarity and prevalence of academic dishonesty*. Paper presented at the annual meeting of the American Psychological Association. San Francisco, CA, August 24-28, 2001.
- Schmeiser, C.B., Geisinger, K.F., Johnson-Lewis, S., Roeber, E.D, & Schafer, W. (1995). *Code of Professional Responsibilities in Educational Measurement*. Retrieved August 11, 2007 from the NATD web site: http://www.natd.org/Code_of_Professional_Responsibilities.html
- Schmidt, E. (2005). Foothills, AZ teacher resigns over cheating incident. Retrieved May 18, 2005 from *Explorer News* web site: <http://www.explorernews.com/articles/2005/05/18/education/education01.txt>
- Shapiro, T. (2005). 'Cheating' probe halts Wai'anae school tests. Retrieved May 9, 2005 from *Honolulu Advertiser* web site: <http://the.honoluluadvertiser.com/article/2005/Apr/09/ln/ln03p.html>
- Shepard, L.A. (1990). Inflated test score gains: is the problem old norms or teaching to the test? *Educational Measurement: Issues and Practice*, 9(3), 15-22.
- Shepard, L. & Dougherty, K. (1991). Effects of high-stakes testing on instruction. Paper presented at the annual meeting of the American Education Research Association, Chicago, IL.
- Shah, N. (2007). 'I' grades leave schools in limbo. Retrieved July 16, 2007 from *Miami Herald* web site: <http://www.miamiherald.com/467/story/171727.html>
- Singhal, A., & Johnson, P. (1983). How to halt student dishonesty. *College Student Journal*, 17(1), 13-19.

- Smith, M.L. (1991). Meanings of test preparation. *American Educational Research Journal*, 28(3), 521-542.
- Smith Richards, J. & Riepenhoff, J. (2007). Schools often keep state in the dark. Retrieved October 16, 2007 from *The Columbus Dispatch*, web site: http://www.columbusdispatch.com/live/content/local_news/stories/2007/10/15/FALTER_new.ART_ART_10-15-07_A9_9A85KGB.html?sid=101
- Solochek, J. (2007). Teacher accused of FCAT cheating. Retrieved May 2, 2007 from *St. Petersburg Times*, web site: http://www.sptimes.com/2007/05/02/Pasco/Teacher_accused_of_FC.shtml
- Sorensen, D. (2006). 2006 state education test security survey results. Retrieved December 17, 2006 from *Caveon* web site: http://www.caveon.com/pr/2006_ed_test_security.pdf
- Sotaridona, L. & Meijer, R. (2001) *Two new statistics to detect answer copying*. Research Report RR-01-07, University of Twente Research, Enschede, Netherlands
- Sotaridona, L. & Meijer, R. (2002) Statistical properties of the K-index for detecting answer copying. *Journal of Educational Measurement*, 39, 115-132.
- Snell, L. (2005). How schools cheat: from underreporting violence to inflating graduation rates to fudging test scores, educators are lying to the American public. Retrieved June 30, 2005 from *Reason* web site: <http://www.reason.com/0506/fe.ls.how.shtml>
- Spencer, B. (1983). On interpreting test scores as social indicators: Statistical considerations. *Journal of Educational Measurement*, 20: 317-333.
- Spiller, S. & Crown, D. (1995). Changes over time in academic dishonesty at the collegiate level. *Psychological Reports*, 76(3)
- SRS1 Software (2007). *SRS1 Cubic Spline for Excel v1.01*. Retrieved December 15, 2007 from the *SRS1 Software* web site: <http://www.srs1software.com/>
- Stanford Policy Repository (2007). Document definitions. Retrieved November 16, 2007 from the *Stanford Policy Repository* web site: <http://www2.slac.stanford.edu/policy/definitions.asp>
- Stinson, R. (1998). TAAS cheaters meet national standards. *San Antonio Express-News*, September 17, 1998.

- Stoneberg, B.D. (2007a). An explanation for the large differences between state and NAEP “proficiency” scores reported for reading in 2005. Paper presented at the 37th Annual National Conference on Large-Scale Assessment, Nashville, TN.
- Stoneberg, B.D. (2007b). The valid use of NAEP achievement level scores to confirm state testing results in the No Child Left Behind Act. Paper presented at the NAEP State Service Center Spring Assessment Workshop, Bethesda, MD.
- Sturrock, C. (2006). States distort school test scores, researchers say. Retrieved June 22, 2007 from *San Francisco Gate* web site: <http://sfgate.com/cgi-bin/article.cgi?file=/c/a/2006/06/30/MNG28JN9RC1.DTL&type=printable>
- Tatsuoka, K.K. (1984). Caution indices based on item response theory. *Psychometrika*, 49, 95-110.
- Texas Education Agency (2007). Sanctions recommended against three schools and three educators because of testing improprieties. Retrieved October 16, 2007 from *Texas Education Agency* web site: <http://www.tea.state.tx.us/press/07iginvestigations.pdf>
- The Nation’s Report Card (2007a). Mathematics report card. Retrieved November 10, 2007 from The Nation’s Report Card web site: http://nationsreportcard.gov/math_2007/data.asp
- The Nation’s Report Card (2007b). Reading report card. Retrieved November 10, 2007 from The Nation’s Report Card web site: http://nationsreportcard.gov/reading_2007/data.asp
- Thiessen, B. (2004). Aberrant response patterns and person-fit statistics. Retrieved August 12, 2007 from <http://homepage.mac.com/bradthiessen/pubs/aberrant.pdf>
- Thiessen, B. (2006). Educator cheating: classification, explanation, and detection. Thesis Equivalency Project, University of Iowa, June 2006. Retrieved June 22, 2007 from web site: <http://homepage.mac.com/bradthiessen/pubs/cheating.pdf>
- Thiessen, B. (2007). Defining and disseminating policies to address educator cheating: case study. Comprehensive exam paper, University of Iowa, January 2007. Retrieved August 8, 2007 from web site: <http://homepage.mac.com/bradthiessen/pubs/format.pdf>
- Thissen, D. (2005). Linking assessments base on aggregate reporting: Background and issues. Paper presented at the ETS Conference, Linking and Aligning Scores and Scales: A Conference in Honor of Ledyard R Tucker’s Approach to Theory and Practice, Princeton, NJ: Educational Testing Services, June 24.

- Thomas B. Fordham Foundation (2005). Gains on State Reading Tests Evaporate on 2005 NAEP. Report. Downloaded in November, 2005 from http://www.edexcellence.net/foundation/about/press_release.cfm?id=19
- Toomer-Cook, J. (2006). School tests under scope. Retrieved July 3, 2007 from Deseretnews.com web site: <http://deseretnews.com/dn/view/0,1249,650212042,00.html>
- Toppo, G. (2007). School test scandal claims decorated principal. Retrieved January 11, 2008 from *USA Today* web site: http://www.usatoday.com/news/education/2007-12-21-high-stakes_N.htm
- Tresague, M. & Viren S. (2006). 2 HSID teachers resign in test-cheating probe. Retrieved July 30, 2006 from *Houston Chronicle* web site: <http://www.chron.com/disp/story.mpl/front/4080406.html>
- Turner, D. (2007). Whistleblower out of a job in apparent cheating scandal. Retrieved June 12, 2007 from *WREG-TV Memphis* web site: <http://www.wreg.com/global/story.asp?s=6642648&ClientType=Printable>
- Turner, P.E. (2005). Cheating viruses and game theory. *American Scientist*, 93, 428-435.
- Turtle, J. (2004). How public schools lie to parents and betray our children. *Public Schools, Public Menace*. Retrieved June 16, 2005 from *Press Method* web site: <http://www.pressmethod.com/releasestorage/547.htm>
- United Press International (2005). Teachers ordered cheating. Retrieved March 24, 2005 from *Washington Times* web site: <http://washingtontimes.com/culture/20050324-114418-7654r.htm>
- U.S. Department of Education Institute of Education Sciences (2007). Mapping 2005 State Proficiency Standards Onto NAEP Scales. Retrieved June 22, 2007 from the *NCES* web site: <http://nces.ed.gov/nationsreportcard/pdf/studies/2007482.pdf>
- van der Linden, W. & Sotaridona, L. (2002). *A statistical test for detecting answer copying on multiple-choice tests*. University of Twente Research Report 02-04. Enschede, Netherlands
- Wallace, K. (2007). "No Child Left Behind" state tests vary. Retrieved May 30, 2007 from *CBS News* web site: <http://www.cbsnews.com/stories/2007/05/30/notebook/main2867441.shtml>
- Wan, W. (2007). School put on probation after students accused of cheating on AP tests. Retrieved August 24, 2007 from *Washington Post* web site:

<http://www.washingtonpost.com/wp-dyn/content/article/2007/08/22/AR2007082202645.html>

- Waters, B. (2007). Winona ISD had possible TAKS security breach. Retrieved May 1, 2007 from *Tyler Morning Telegraph*, web site:
<http://www.tylerpaper.com/apps/pbcs.dll/article?AID=/20070501/NEWS01/705010320/-1/FRONTPAGE>
- Watters, C. (2005). Principal's test scores probed. Retrieved April 13, 2005 from *Rockford Register Star*, web site:
<http://www.rrstar.com/apps/pbcs.dll/article?AID=/20050222/NEWS0107/502220321/1004/NEWS>
- W-CBS TV (2006). New York high school accused of fixing test scores. Retrieved November 2, 2006 from *W-CBS TV* web site:
http://cbs11tv.com/education/local_story_306101643.html
- Wei, X., Shen, X., Lukoff, B., Ho, A.D., & Haertel, E.H. (2006). Using test content to address trend discrepancies between NAEP and California State Tests. ERIC # ED491544. Retrieved September 6, 2007 from:
http://eric.ed.gov/ERICDocs/data/ericdocs2sql/content_storage_01/0000019b/80/29/00/1f.pdf
- Wesolowsky, G.O. (2000). Detecting excessive similarity in answers on multiple choice exams. *Journal of Applied Statistics*, 27(7), 909-921.
- White, K. R., Taylor, C., Carcelli, L., & Eldred, N. (1981). *State refinements to the ESEA Title I evaluation and reporting system: Utah 1979- 80 project*. Logan, UT: Utah State University, Exceptional Child Center. (ERIC Document Reproduction Service No. ED 212 087)
- WHO TV (2005). Teacher resigns over standardized testing procedures. Retrieved May 4, 2005 from *WHO-TV* web site:
<http://www.whotv.com/Global/story.asp?S=3302244>
- Wilk, M.B. & Gnanadesikan, R. (1968). Probability plotting methods for the analysis of data. *Biometrika*, 55(1), 1-17.
- Wodtke, K.H., Harper, F., Schommer, M. & Brunelli, P. (1989). How standardized is school testing? An exploratory study of standardized group testing in kindergarten. *Educational Evaluation and Policy Analysis*, 11(3), 223-235.
- Wolf, R.M. (1992). What can we learn from state NAEP? *Educational Measurement: Issues and Practice*, 11(4), 12.

Wollack, J.A. (1997). A nominal response model approach for detecting answer copying. *Applied Psychological Measurement*, 21(4), 307-320.

Wright, B.D. & Stone, M.H. (1979). *Basic test design: Rasch measurement*. Chicago: MESA Press.