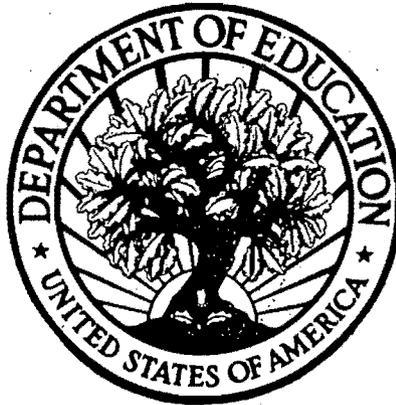


**The Use of Tests When
Making High-Stakes
Decisions for Students:**
*A Resource Guide for
Educators and Policymakers*



**U.S. Department of Education
Office for Civil Rights**

Draft DraftDraft

July 6, 2000

There are few simple or definitive answers to questions about the use of tests for high-stakes purposes. Tests are a means to an end and, as such, can be understood only in the context in which they are used. The education context — in which the relationship (and attendant obligations) of the educator to the student is frequently more complex than that between employer and employee — shows time and again that any decision regarding the legality of a use of a test for high-stakes purposes under federal nondiscrimination law cannot be made without regard to the educational interests and judgments upon which the test use is premised.

Background

Throughout the 1990s, national, state and local education leaders have focused on raising education standards and establishing strategies to promote accountability within the education community. In fact, the promotion of challenging learning standards for all students — coupled with assessment systems that monitor progress and hold schools accountable — has been the centerpiece of the education policy agenda of the federal government as well as many states.

Predictably, the number of states using tests as a condition for high school graduation is on the rise, with (by a recent estimate) 26 states projected to use tests as conditions for graduation by 2003 and six states now using tests as conditions for grade promotion, a significant increase from past years. At the same time, more and more educators and policymakers have requested advice and technical assistance from the U.S. Department of Education regarding test use in the context of standards reforms.

The Department's Office for Civil Rights (OCR) is also addressing testing issues in a more extensive array of complaints of discrimination being filed with our office, most of them in a K-12 setting with implications for high-standards learning. OCR has responsibility for enforcing Title VI of the Civil Rights Act of 1964, Title IX of the Education Amendments of 1972, Section 504 of the Rehabilitation Act of 1973, and Title II of the Americans with Disabilities Act of 1990. These statutes prohibit discrimination on the basis of race, color, national origin, sex, and disability by educational institutions that receive federal funds.

In a similar vein, institutions in the post-secondary community in recent years have engaged in a thoughtful dialogue and analysis regarding merit in admissions and the appropriate use of tests to establish foundations for high-stakes admissions decisions. In some states, the use of tests in connection with admissions decisions has been an important element in public post-secondary education reform.

These trends highlight the salience of two recent conclusions of the National Research Council (NRC) Board on Testing and Assessment. In January of this year, the NRC observed that too many policymakers and educators are not aware of the test measurement standards that should inform testing policies and practices. These standards include the Standards for Educational and Psychological Tests, prepared by a joint committee of the American Psychological Association (APA), the American Educational Research Association (AERA), and the National Council on Measurement in Education (NCME). The NRC also concluded that it "is essential that educators and policymakers

alike be aware of both the letter of the laws and their implications for test takers and test users” [National Research Council, *High Stakes: Testing for Tracking, Promotion and Graduation*, (Heubert and Hauser, eds., 1999)].

The Resource Guide

Toward this end, OCR has prepared this guide in an effort to assemble the best information regarding psychometric standards, legal principles, and resources to help educators and policymakers frame strategies and programs that promote learning to high standards in ways consistent with federal non-discrimination law. Our goal is to inform decisions related to the use of tests that have high-stakes consequences for students when, for instance, they move from grade to grade or graduate from high school. Just as we know that good test use practices can advance high standards for learning and equal opportunity, we know that educationally inappropriate uses of tests do not. If we want this generation of test-taking students and their teachers and schools to meet high standards, then we should insist that the tests they take meet high standards. As foundations for judgments that profoundly shape the lives of students, these tests must be used in ways that accurately reflect educational standards and that do not inappropriately deny opportunities to students based on their race, national origin, sex or disability.

The guide is organized to provide practical guidance related to the use of tests for high-stakes purposes. The Introduction to the guide provides a broad, conceptual overview of relevant principles so that those who are not familiar with test measurement principles or applicable federal law can better understand the kinds of issues that relate to the use of tests in many contexts — from grade-to-grade promotion to college admissions. Chapter one of the guide provides a detailed discussion of the test measurement principles that can provide a foundation for making well-informed decisions related to high-stakes testing. The relevant principles that have been approved by the APA, AERA, and NCME are discussed in detail in this chapter. Adherence to relevant professional standards can help reduce the risk of legal liability when schools are using assessments for high-stakes purposes. Chapter two provides an overview of the existing legal principles that have guided federal courts and OCR when analyzing claims of race, national origin, sex, and disability discrimination related to the use of tests as foundations in high-stakes decisions affecting students. These principles, as applied by the courts and OCR, underscore the importance of adhering to educationally sound testing practices. The Appendix includes a Glossary of Test Measurement Terms, a Glossary of Legal Terms, a Compendium of Federal Nondiscrimination Laws, and a Resources and References section.

Central Principles

There are several central principles reflected in the text of this guide.

First, federal nondiscrimination laws are consistent with the establishment of high standards of learning for all students and educationally sound practices designed to meet that goal. The goals of promoting high educational standards and ensuring nondiscrimination are complementary objectives. Indeed, if the federal courts that have applied civil rights statutes to education cases teach us anything, it is that compliance with federal nondiscrimination standards rests in the first instance upon the school’s

educational judgments, and that those judgments deserve deference. Not surprisingly, the ultimate questions posed by our resource guide on the use of tests for high-stakes purposes also center on educational sufficiency: Is the test valid for the purposes used? Are the inferences derived from test scores, and the high-stakes decisions based on those inferences, accurate and fair? These inquiries are not an effort to dumb down academic standards or alter core education objectives integral to academic admissions or other educational decisions. Rather, they focus the educator and policymaker on ensuring that uses of tests with consequences for students are educationally sound and legally appropriate.

Second, federal nondiscrimination laws support the use of tests, including large-scale standardized tests, when they are used in valid, reliable, and educationally appropriate ways. Importantly, tests can help indicate inequalities in the kinds of educational opportunities students are receiving, and in turn, they may stimulate efforts to ensure that all students have equal opportunity to achieve high standards. When tests accurately indicate performance gaps, our concern should be with the quality of educational opportunities afforded to under-performing students (rather than the integrity of the test itself.) The key question in the context of standards-based reforms and the use of tests as measures of student accountability is: Have all students in certain school districts been provided quality instruction, sufficient resources, and the kind of learning environment that would foster success? — OTL

But rather than just in addition to
Third, a test score disparity among groups of students does not alone constitute discrimination under federal law. The guarantee under federal law is for equal opportunity, not equal results. Test results indicating that groups of students perform differently should be a cause for further inquiry and examination, with a focus upon the relevant educational programs and testing practices at issue. Differences in test scores may result from a range of factors, some of which a school may be able to influence, and others over which it has little control. Federal law recognizes this point, as it must. The legal non-discrimination standard regarding neutral practices (referred to by the courts as the “disparate impact” standard) provides that if the education decisions based upon test scores reflect statistically significant disparities based on race, national origin, or sex in the kinds of educational benefits afforded to students, then questions about the education practices at issue (including testing practices) should be thoroughly examined to ensure that they are in fact non-discriminatory and educationally sound. In short, the goal of the federal legal standards is to help promote accurate and fair decisions that have real consequences for students, not to water down academic standards or deter educators from establishing and applying sensible and rigorous standards.

Conclusion

Recognizing the responsibility that educators and policymakers must shoulder in making the promise of high standards learning a reality, U.S. Secretary of Education Richard Riley in his commemoration of the 45th anniversary of the *Brown v. Board of Education* decision said: “A quality education must be considered a key civil right for the twenty first century.” This is the driving force behind OCR's continuing effort to provide assistance to policymakers and educators as we continue to enforce federal laws that prohibit discrimination against students. Rather than creating false and polarizing “win-

lose” choices on this all-important set of issues, we need to, as Secretary Riley admonishes, “search for common ground” — ground, that is, in this case, expansive.

We have worked with literally dozens of groups and individuals, including educators, parents, teachers, business leaders, policymakers, test publishers, and others, to solicit input and advice regarding the scope, framing, and kinds of resources to include in this guide, and we are grateful for their assistance. In addition, we have contracted with the NRC’s Board on Testing and Assessment, which has reviewed earlier drafts of the guide, to ensure that the guide comports with professional standards. We are grateful for the NRC’s tireless efforts.

Working together with our education partners, we believe that we are providing a useful resource that will serve the education community as it addresses the very complex and important questions that stem from the institution of high standards and accountability systems designed to promote the best schools in the world.

Very truly yours,

DRAFT

Norma V. Cantú

Table of Contents

INTRODUCTION: An Overview of the Resource Guide.....	1
CHAPTER 1. Test Measurement Principles.....	19
CHAPTER 2. Legal Principles.....	46
APPENDIX A: Glossary of Legal Terms.....	63
APPENDIX B: Glossary of Test Measurement Terms.....	67
APPENDIX C: Accommodations Used by States.....	74
APPENDIX D: Compendium of Federal Statutes and Regulations	77
APPENDIX E: Resources and References	80

INTRODUCTION: An Overview of the Resource Guide

I. Introduction

Decisions affecting students' educational opportunities should be made accurately and fairly. When tests are used in making educational decisions for individual students, they should accurately measure students' abilities, knowledge, skills or needs, and they should do so in ways that do not discriminate in violation of federal law on the basis of the students' race, national origin, sex or disability. The U.S. Department of Education's Office for Civil Rights (OCR)¹ has developed this resource guide in order to provide educators and policymakers with a useful, practical tool that will assist in their development and implementation of policies that involve the use of tests in making high-stakes decisions for students. It is intended to facilitate the proper use of tests for those purposes.

Chapter one of this guide provides information about professionally recognized test measurement principles. Chapter two provides the legal frameworks that have guided federal courts and OCR when addressing the use of tests that have high-stakes consequences for students. The test measurement principles described in chapter one are not legal principles. However, the use of tests in educationally appropriate ways — consistent with the principles described in chapter one — can help to minimize the risk of noncompliance with the federal nondiscrimination laws discussed in chapter two.

When tests are used in ways that meet relevant psychometric, legal, and educational standards, students' scores provide important information that, combined with information from other sources, can lead to decisions that promote student learning and equality of opportunity When test use is inappropriate, especially in making high-stakes decisions about individuals, it can undermine the quality of education and equality of opportunity. This lends special urgency to the requirement that test use with high-stakes consequences for individual students be appropriate and fair.

National Research Council, *High Stakes: Testing for Tracking, Promotion and Graduation*, 1999:4.

¹ OCR enforces laws that prohibit discrimination on the basis of race, national origin, sex, disability, and age by educational institutions that receive federal funds. The laws enforced by OCR are: 1) Title VI of the Civil Rights Act of 1964, 42 U.S.C. §§ 2000d, *et seq.* (2000)(Title VI), which prohibits discrimination on the basis of race, color, or national origin; 2) Title IX of the Education Amendments of 1972, 20 U.S.C. §§ 1681, *et seq.* (1999)(Title IX), which prohibits discrimination on the basis of sex; 3) Section 504 of the Rehabilitation Act of 1973, 29 U.S.C. §§ 794, *et seq.* (1999)(Section 504), which prohibits discrimination on the basis of disability; 4) the Age Discrimination Act of 1975, 42 U.S.C. §§ 6101, *et seq.* (1995 and Supp. 1999)(as amended), which prohibits age discrimination; and 5) Title II of the Americans with Disabilities Act of 1990, 42 U.S.C. §§ 12134, *et seq.* (1995 and Supp. 1999)(Title II), which prohibits discrimination on the basis of disability by public entities, whether or not they receive federal financial assistance.

The guide also includes a collection of resources related to test measurement and nondiscrimination principles that are discussed in the guide — all in an effort to help policymakers and educators ensure that decisions that have high-stakes consequences for students are made accurately and fairly.

Educational stakeholders at all levels have approached OCR requesting advice and technical assistance in a variety of test-use contexts, particularly as states and districts use tests as part of their standards-based reforms. Also, increasingly, OCR is addressing testing issues in a broader and more extensive array of complaints of discrimination that have been filed with OCR. These corresponding developments confirm the need to provide a useful resource that captures legal and test measurement principles and resources to assist educators and policymakers. This document does not establish any new legal or test measurement principles.

As used in this resource guide, “high-stakes decisions” refer to decisions with important consequences for individual students. Education entities, including state agencies, local education agencies, and individual education institutions, make a variety of decisions affecting individual students during the course of their academic careers, beginning in elementary school and extending through the post-secondary school years. Examples of high-stakes decisions affecting students include: student placement in gifted and talented programs or in programs serving students with limited-English proficiency; determinations of disability and eligibility to receive special education services; student promotion from one grade level to another; graduation from high school and diploma awards; and admissions decisions and scholarship awards.²

High-stakes decisions in this guide refer to decisions with important consequences for students, such as placement in special programs, promotion, graduation, and admissions decisions.

This guide is intended to apply to standardized tests that are used in making high-stakes decisions affecting individual students and that are addressed in the *Standards for Educational and Psychological Testing (Joint Standards)*. The *Joint Standards* are viewed as the primary technical authority on educational test measurement issues. They have been prepared by a joint committee of the American Educational Research Association, the American Psychological Association and the National Council on Measurement in Education, the three leading organizations in the area of educational test measurement. The *Joint Standards* were developed and revised by these three organizations through a process that involved the participation of hundreds of testing professionals and thousands of pages of written comment from both professionals and the public. The current edition of the *Joint Standards* reflects the experience gained from

² The purpose of this guide is to address tests that are used in making high-stakes decisions for individual students. In addition to using tests for high-stakes purposes for individual students, states and school districts are also using tests to hold schools and districts accountable for student performance. Although using tests for this purpose is not the focus of the guide, we have provided some useful background information about relevant principles and federal statutory requirements.

many years of wide use of previous versions of the *Joint Standards* in the testing community.

The *Joint Standards*, which are discussed in more detail below, apply to standardized measures generally recognized as tests, and also may be usefully applied to a broad range of system-wide standardized assessment procedures.³ For the sake of simplicity, this guide will refer to tests, regardless of the type of label that might otherwise be applied to them. The guide does not address teacher-created tests that are used for individual classroom purposes.

States and school districts are also using another important kind of assessment system for the purpose of promoting school and district accountability. For example, under Title I of the Elementary and Secondary Education Act, states are required to develop content standards, performance standards, and assessment systems that measure the progress that schools and districts are making in educating students to the standards established by the state. Title I explicitly requires that such assessments be valid and reliable for their intended purpose and be consistent with relevant, nationally recognized technical and professional standards.⁴ When educators and policy makers consider using the same test for school or district accountability purposes and for individual student high-stakes purposes, they need to ensure that the test score inferences are valid and reliable for each particular use for which the test is being considered.

When high-stakes decisions are made, test scores are often used in conjunction with other criteria, such as grades and teacher recommendations. A test should not be used as the sole criterion for making a high-stakes decision unless it is validated for this use. The *Joint Standards* state that a high-stakes decision “should not be made on the basis of a single test score. Other relevant information should be taken into account if it will enhance the overall validity of the decision.”⁵ As explained in the *Joint Standards*, “[w]hen interpreting and using scores about individuals or groups of students, considerations of relevant collateral information can enhance the validity of the interpretation, by providing corroborating evidence or evidence that helps explain student performance. ... As the stakes of testing increase for individual students, the importance of considering additional evidence to document the validity of score interpretations and the fairness in testing increases accordingly.”⁶

³ The *Joint Standards* note that the applicability of the Joint Standards to an evaluation device or method is not altered by the label used (e.g., test, assessment scale, inventory). A more complete discussion about the instruments covered by the *Joint Standards* can be found in the introduction section of that document. See *Joint Standards*, Introduction, pp. 3-4.

⁴ 20 U.S.C. 6311(b)(3)(C).

⁵ Standard 13.7 states, “In educational settings, a decision or characterization that will have major impact on a student should not be made on the basis of a single test score. Other relevant information should be taken into account if it will enhance the overall validity of the decision.”

⁶ *Joint Standards*, p. 141.

Although this guide focuses on the use of tests in making high-stakes decisions, policymakers and the education community need to ensure that the operation of the entire high-stakes decision-making process does not result in the discriminatory denial of educational benefits or opportunities to students.⁷ Applicable standards for technical quality set forth in the *Joint Standards* are important principles to consider when other criteria affect high-stakes decisions. Educators should carefully monitor inputs into the high-stakes decision-making process and outcomes over time so that any potential discrimination arising from the use of any of the criteria can be identified and eliminated.

The guide focuses primarily on tests used in making high-stakes decisions at the elementary and secondary education level. However, it is important to recognize that the general principles of sound educational measurement apply equally to tests used at the elementary and secondary education level and at the post-secondary education level, including admissions and other types of test use.⁸ For example, post-secondary admissions policies and practices should be derived from and clearly linked to an institution's overarching educational goals, and the use of tests in the admissions process should serve those institutional goals.⁹

Standardized tests ... offer important benefits that should not be overlooked. ... Both the SAT [I] and ACT cover relatively broad domains that most observers would likely agree are relevant to the ability to do college work. Neither, however, measures the full range of abilities that are needed to succeed in college; important attributes not measured include, for example, persistence, intellectual curiosity, and writing ability. Moreover, these tests are neither complete nor precise measures of 'merit'—even academic merit.

National Research Council, *Myths and Tradeoffs: The Role of Tests in Undergraduate Admissions*, 1999: 21-22.

II. Foundations of the Resource Guide

A. Professional Standards of Sound Testing Practices

Chapter one summarizes the leading professionally recognized standards of sound testing practices

⁷ See Nondiscrimination Under Programs Receiving Federal Financial Education Effectuation of Title VI of the Civil Rights Act of 1964, 34 C.F.R. § 100.3(b)(2) (1999); Nondiscrimination on the Basis of Handicap in Programs Receiving Federal Financial Assistance, 34 C.F.R. §§ 104.4(a), 104.4(b)(1)(i) and (iv), and Basis of Sex in Education Programs and Activities Receiving or Benefiting from Federal Financial Assistance, 34 C.F.R. §§ 106.31(a) and 106.31(b) (1999).

⁸ For additional information regarding testing at the post-secondary level, see *Tradeoffs*, 1999; Messick, S., Validity, in R.L. Linn, ed., *Educational Measurement*, 13-103, 1989; Wigdor, Alexandra K., and Garner, Wendell R., ed., *Assessment Controversies*, chapter 5, National Academy Press, 1982.

⁹ See *High Stakes*, p. 23 and National Research Council, *Placing Children in Special Education: A Strategy for Equity*, 1982.

The proper use of tests can result in wiser decisions about individuals and programs than would be the case without their use and also can provide a route to broader and more equitable access to education ... The improper use of tests, however, can cause considerable harm to test takers and other parties affected by test-based decisions.

Joint Standards, Introduction, at p. 1.

within the educational measurement field. They include those described in the *Joint Standards* (1999), which represent the primary statement of professional consensus regarding educational testing. Other leading professionally recognized standards of sound testing practices within the educational measurement field include the *Code of Fair Testing Practices in Education* (1988), and the *Code of Professional Responsibilities in Educational Measurement* (1995). The guide also cites recent reports from the National Research Council's Board on Testing and Assessment, including *High Stakes: Testing for Tracking, Promotion and Graduation* (High Stakes, 1999), *Myths and Tradeoffs: The Role of Tests in Undergraduate Admissions* (Myths and Tradeoffs, 1999), *Testing, Teaching, and Learning: A Guide for States and School Districts* (Testing, Teaching, and Learning, 1999), *Improving Schooling for Language-Minority Children: A Research Agenda* (Improving Schooling for Language-Minority Children, 1997), and *Educating One & All: Students with Disabilities and Standards-Based Reform* (Educating One & All, 1997).¹⁰ These reports help explain or elaborate principles that are stated in the Joint Standards.

Designed to provide criteria for the evaluation of tests, testing practices, and the effects of test use, the *Joint Standards* recommend that all professional test developers, sponsors, publishers, and users make efforts to observe the *Joint Standards* and encourage others to do so.¹¹ The *Joint Standards* include chapters on the test development process (with a focus primarily on the responsibilities of test developers), the specific uses and applications of tests (with a focus primarily on the responsibilities of test users), and the rights and responsibilities of test takers. Because the *Joint Standards* are the most widely accepted professional standards that are relied upon in developing testing instruments, this guide includes a discussion of specific standards that are contained within the *Joint Standards*, where relevant. Numbered standards that are referenced throughout this guide refer to specific standards that are contained within the *Joint Standards*.

In order to ensure that information presented in the guide is readable and accessible to educators and policymakers, we have paraphrased language from relevant standards. Our goal in paraphrasing is to be concise and accurate. Where we have paraphrased in the text, we have also provided the full text of the relevant standards in the footnotes. Because the *Joint Standards* provide additional relevant discussion, we always encourage readers also to review the full document.

Professional test measurement standards provide important information that is relevant to making determinations about appropriate test use. The *Joint Standards* provide a frame of reference to assist in the evaluation of tests, testing practices, and the effects of test use. The *Joint Standards* caution that the acceptability of a test or test application does

¹⁰ The National Academy of Sciences, which is an independent, private, nonprofit entity, established the Board on Testing and Assessment in 1993 to help policymakers evaluate the use of tests, alternative assessments, and other indicators commonly used as tools of public policy. The Board provides guidance for judging the quality of testing or assessment technologies and the intended and unintended consequences of particular uses of these technologies. The Board concentrates on topics and conducts activities that serve the general public interest.

¹¹ See, e.g., *Joint Standards*, Introduction, p. 2.

not rest on the literal satisfaction of every standard in the *Joint Standards* and cannot be determined by using a checklist.¹² The exercise of professional judgment is a critical element in the interpretation and application of the standards,¹³ and the interpretation of individual standards should be considered in the overall context of the use of the test in question. Failure to meet a particular professional test measurement standard does not necessarily constitute a lack of compliance with federal civil rights laws.

B. Legal Standards

Chapter two of the guide discusses the federal Constitutional, statutory and regulatory nondiscrimination principles that apply to the use of tests for high-stakes purposes. This guide is intended to reflect existing legal principles and does not establish new federal legal requirements. The primary legal focus of the resource guide is an explanation of principles that are clearly embedded in four nondiscrimination laws that have been enacted by Congress: Title VI of the Civil Rights Act of 1964 (Title VI), Title IX of the Education Amendments of 1972 (Title IX), Section 504 of the Rehabilitation Act of 1973 (Section 504), and Title II of the Americans with Disabilities Act of 1990 (Title II).¹⁴ Within the U.S. Department of Education, the Office for Civil Rights has responsibility for enforcing the requirements of these four statutes and their implementing regulations. The due process and equal protection requirements of the Fifth and Fourteenth Amendments to the U.S. Constitution have also been applied by courts to issues regarding the use of tests in making high-stakes educational decisions. Although the Office for Civil Rights does not enforce federal constitutional provisions, a brief overview of these constitutional principles has been included for informational purposes.

Large-scale assessments, used properly, can improve teaching, learning, and equality of educational opportunity. That tests are sometimes used improperly should not discourage policymakers, teachers, and parents. Rather, it should motivate action to ensure that educational tests are used fairly and effectively.

National Research Council, *High Stakes: Testing for Tracking, Promotion, and Graduation*, 1999: 9.

¹² *Joint Standards*, Introduction, p. 4.

¹³ *Joint Standards*, Introduction, p. 4.

¹⁴ Title VI prohibits discrimination on the basis of race, color and national origin in the programs and activities of recipients that receive federal financial assistance. The U.S. Department of Education's regulation implementing Title VI is found at 34 C.F.R. Part 100. Title IX prohibits discrimination on the basis of sex in educational programs and activities of recipients of federal financial assistance. The U.S. Department of Education's regulation implementing Title IX is found at 34 C.F.R. Part 106. Section 504 prohibits discrimination on the basis of disability in the programs and activities of recipients of federal financial assistance. The U.S. Department of Education's regulation implementing Section 504 is found at 34 C.F.R. Part 104. Title II prohibits discrimination on the basis of disability by public entities, regardless of whether they receive federal funding. The U.S. Department of Education's regulation implementing Title II is found at 28 C.F.R. Part 35.

III. Basic Principles

The brief overview of the test measurement and legal principles that follows establishes the framework for more detailed discussions of test quality in chapter one and federal legal standards in chapter two.

A. Test Use Principles

1. Educational Objectives and Context

Tests that are used in educationally appropriate ways and that are valid for the purposes used are important instruments to help educators do their job. Before any state, school district, or educational institution administers a test, the objectives for using the test should be clear: What are the intended goals for and uses of the test in question? As an educational matter, the answer to this question will guide all other relevant inquiries about whether the test use is educationally appropriate. The context in which a test is to be administered, the population of test takers, and the intended purpose for which the test will be used are important considerations in determining which test would be appropriate for a specific use, as illustrated below:

Decisions about tracking, promotion, and graduation differ from one another in important ways. They differ most importantly in the role that mastery of past material and readiness for new material play.

National Research Council, *High Stakes: Testing for Tracking, Promotion, and Graduation*, 1999: 4.

a. Placement Decisions

Placement decisions are by their very nature used to make a decision about the future. Tests used in placement decisions generally determine what kinds of programs, services, or interventions will be most appropriate for particular students. Decisions concerning the appropriate educational program for a student with a disability, placement in gifted and talented programs, and access to language services are examples of placement decisions. The *Joint Standards* state that there should be adequate evidence documenting the relationship among test scores, appropriate instructional programs, and beneficial student outcomes.¹⁵ When evidence about the relationship is limited, the test results should be considered in light of other relevant student information.¹⁶

[At the elementary and secondary education level,] appropriate test use for ... all students requires that their scores not lead to decisions or placements that are educationally detrimental.”

National Research Council, *High Stakes: Testing for Tracking, Promotion, and Graduation*, 1999: 40-41.

¹⁵ Standard 13.9 states, “When test scores are intended to be used as part of the process for making decisions for educational placement, promotion, or implementation of prescribed educational plans, empirical evidence documenting the relationship among particular scores, the instructional programs, and desired student outcomes should be provided. When adequate empirical information is not available, users should be cautioned to weigh the test results accordingly in light of other relevant information about the student.”

¹⁶ See id.

b. Promotion Decisions

Student promotion decisions are generally viewed as decisions incorporating a determination about whether a student has mastered the subject matter or content of instruction provided to date and a determination regarding whether the student will be able to master the content at the next grade level (a placement decision).¹⁷ At present, the focus of most school districts and states with promotion policies has been primarily on assessing mastery of curriculum taught at a given grade level.¹⁸ When a test given for promotion purposes is being used to certify mastery, the use of the test should adhere to professional standards for certifying knowledge and skills for all students.¹⁹ It is important that there be evidence that the test adequately covers only the content and skills that students have actually had an opportunity to learn.²⁰ Educational institutions should have information indicating an alignment among the curriculum, instruction, and material covered on such a high-stakes test. To the extent that a test for promotion purposes is being used as a placement device, it should also adhere, as appropriate, to professional standards regarding tests used for placement purposes.²¹

¹⁷ See *High Stakes*, p. 123.

¹⁸ See American Federation of Teachers, *Passing on Failure: District Promotion Policies and Practices*, 1997.

¹⁹ See Standards 13.5 and 13.6; *High Stakes*, p. 123. Standard 13.5 states, "When test results substantially contribute to making decisions about student promotion or graduation, there should be evidence that the test adequately covers only the specific or generalized content and skills that students have had an opportunity to learn."

Standard 13.6 states, "Students who must demonstrate mastery of certain skills or knowledge before being promoted or granted a diploma should have a reasonable number of opportunities to succeed on equivalent forms of the test or be provided with construct-equivalent testing alternatives of equal difficulty to demonstrate the skills or knowledge. In most circumstances, when students are provided with multiple opportunities to demonstrate mastery, the time interval between the opportunities should allow for students to have the opportunity to obtain the relevant instructional experiences."

²⁰ See Standard 13.5, *supra* note 19.; *High Stakes*, pp. 124-125.

²¹ See Standards 13.2 and 13.9; *High Stakes*, p. 123. Standard 13.2 states, "In educational settings, when a test is designed or used to serve multiple purposes, evidence of the test's technical quality should be provided for each purpose." See Standard 13.9, *supra* note 15.

c. Graduation Decisions

Graduation decisions are generally certification decisions: The diploma certifies that the student has reached an acceptable level of mastery of knowledge and skills.²² When large-scale standardized tests are used in making graduation decisions, there should be evidence that the test adequately covers only the content and skills that students have had an opportunity to learn.²³ Therefore, all students should be provided a meaningful opportunity to acquire the knowledge and skills that are being tested, and information should indicate an alignment among the curriculum, instruction, and material covered on the test used as a condition for graduation.

2. Overarching Principles

The highly contextual and fact-based test measurement analyses applicable to a variety of circumstances ultimately focus upon the following question: Is there sufficient confidence in the test results at issue to allow for informed decisions to be made that will have specified consequences for the students taking the test?

Is it ever appropriate to test [elementary or secondary] students on material they have not been taught? Yes, if the test is used to find out whether the schools are doing their job. But if that same test is used to hold students "accountable" for the failure of the schools, most testing professionals would find such use inappropriate. It is not the test itself that is the culprit in the latter case; results from a test that is valid for one purpose can be used improperly for other purposes.

National Research Council, *High Stakes: Testing for Tracking, Promotion and Graduation*, 1999: 21.

In the elementary and secondary education context, regardless of whether tests are being used to make placement, promotion, or graduation decisions, the National Academy of Sciences' Board on Testing and Assessment has identified three principal criteria, which are based on established professional standards, that can help inform and guide conclusions regarding this issue.²⁴

- (1) *Measurement validity: Is a test valid for a particular purpose, and does it accurately measure the test taker's knowledge in the content area being tested?*

State and local educational agencies and educational institutions should ensure that a test actually measures what it is intended to measure for all students. The inferences derived from the test scores for a given use — for a specific purpose, in a specific type of

²² See *High Stakes*, p. 166.

²³ See Standard 13.5, *supra* note 19.

²⁴ See *High Stakes*, p. 23 and National Research Council, *Placing Children in Special Education: A Strategy for Equity*, 1982.

situation, and with specific types of students — are validated, rather than the test itself. It is important for educators who use the test to request adequate evidence of test quality (including validity and reliability evidence), evaluate the evidence, and ensure that the test is used appropriately in a way that is consistent with information provided by the developers or through supplemental validation studies.

- (2) *Attribution of cause: Does a student's performance on a test reflect knowledge and skills based on appropriate instruction, or is it attributable to poor instruction or to such factors as language barriers unrelated to the skills being tested?*

In some contexts, whether a particular test use is appropriate depends on whether test scores are an accurate reflection of a student's knowledge or skills or whether they are influenced by extraneous factors unrelated to the specific skills being tested. For example, when tests are used in making student promotion or graduation decisions, state and local education agencies should ensure that all students have an equal opportunity to acquire the knowledge and skills that are being tested.²⁵ In some situations, it may be necessary to provide appropriate accommodations for limited English proficient students and students with disabilities to accurately and effectively measure students' knowledge and skills in the particular content area being assessed.²⁶

- (3) *Effectiveness of treatment — Do test scores lead to placements and other consequences that are educationally beneficial?*

"Appropriate test use for ... all students requires that their scores not lead to decisions or placements that are educationally detrimental."

National Research Council, *High Stakes: Testing for Tracking, Promotion, and Graduation*, 1999: 40-41.

The most basic obligation of educators at the elementary and secondary level is to meet the needs of students as they find them, with their different backgrounds, and to teach knowledge and skills to allow them to grow to maturity with meaningful expectations of a productive life in the workforce and elsewhere.²⁷ This elementary and secondary educational obligation is no less present when educators administer tests and evaluate and act on students' test results than it is during classroom instruction. Relying upon the sound premise that tests should be

²⁵ See Standard 7.10, which states, "When the use of a test results in outcomes that affect the life chances or educational opportunities of examinees, evidence of mean test score differences between relevant subgroups of examinees should, where feasible, be examined for subgroups for which credible research reports mean differences for similar tests. Where mean differences are found, an investigation should be undertaken to determine that such differences are not attributable to a source of construct underrepresentation or construct-irrelevant variance. While initially, the responsibility of the test developer, the test user bears responsibility for uses with groups other than those specified by the developer."

²⁶ See *Joint Standards*, p. 143.

²⁷ See *Brown v. Bd. of Educ.*, 347 U.S. 483, 493 (1954) (stating that "[education] is required in the performance of our most basic public responsibilities, ... is the very foundation of good citizenship, ... [and] is [a] principal instrument ... in preparing [the child] for later professional training....").

integral to the learning and achievement of students, one federal court distinguished between testing in the employment and education settings:

If tests predict that a person is going to be a poor employee, the employer can legitimately deny the person the job, but if tests suggest that a young child is probably going to be a poor student, a school cannot on that basis alone deny that child the opportunity to improve and develop the academic skills necessary to success in our society.²⁸

Tests, in short, should be instruments used by elementary and secondary educators to help students achieve their full potential. Test scores should lead to consequences that are educationally beneficial for students. When making high-stakes decisions that involve the use of tests, it is important for policymakers and educators to consider the intended and unintended consequences that may result from the use of the test scores.²⁹

B. Legal Principles

Federal constitutional, statutory, and regulatory principles form the federal legal nondiscrimination framework applicable to the use of tests for high-stakes purposes. Title VI, Title IX, Section 504, and Title II, as well as the equal protection clause of the Fourteenth Amendment to the United States Constitution, prohibit intentional discrimination based on race, national origin, sex, or disability. In addition, the regulations that implement Title VI, Title IX, Section 504 and Title II prohibit intentional discrimination and policies or practices that have a discriminatory disparate impact on students based on their race, national origin, sex, or disability.³⁰ The Section 504 regulation and the Individuals with Disabilities Education Act³¹ contain specific provisions relative to the use of high-stakes tests for individuals with disabilities.³²

²⁸ *Larry P. v. Riles*, 793 F.2d 969, 980 (9th Cir. 1984)(quoting *Larry P. v. Riles*, 495 F. Supp. 926, 969 (N.D. Cal. 1979)).

²⁹ Research indicates that students in low-track classes do not have the opportunity to acquire knowledge and skills strongly associated with future success that is offered to students in other tracks. The National Research Council recommends that neither test scores nor other information should be used to place students in such classes. See *High Stakes*, 1999: 282.

³⁰ 34 C.F.R. § 100.3(b)(2); 34 C.F.R. §§ 106.21(b)(2), 106.36(b), 106.52; 34 C.F.R. § 104.4(b)(4)(i); and 28 C.F.R. § 35.130(b)(3).

The authority of federal agencies to issue regulations with an "effects" standard has been consistently acknowledged by U.S. Supreme Court decisions and applied by lower federal courts addressing claims of discrimination in education. See, e.g., *Lau v. Nichols*, 414 U.S. 563, 568 (1974); *Guardians Ass'n v. City Service Comm'n. of City of N.Y.*, 463 U.S. 582, 584-593 (1983); *Alexander v. Choate*, 469 U.S. 287, 289-300 (1985). See also Memorandum from the Attorney General for Heads of Departments and Agencies that Provide Federal Financial Assistance, "Use of the Disparate Impact Standard in Administrative Regulations under Title VI of the Civil Rights Act of 1964," July 14, 1994.

³¹ The IDEA establishes rights and protections for students with disabilities and their families. It also provides federal funds to local school districts and state agencies to assist in educating students with disabilities. Individuals with Disabilities Education Act, 20 U.S.C. § 1400(1)(c).

³² 34 C.F.R. §§ 104.35, 104.42(b); 20 U.S.C. §§ 1412(a)(17), 1414(b); 34 C.F.R. § 300.138 - .139, 300.530 - .536.

Further discussion of issues regarding testing of limited English proficient students and students with disabilities is provided below.

1. Frameworks for Analysis

a. **Different Treatment**

Under federal law, policies and practices generally must be applied consistently to similarly situated individuals or groups, regardless of their race, national origin, sex, or disability. For example, a court concluded that a school district had intentionally treated students differently on the basis of race where minority students whose test scores qualified them for two or more ability levels were more likely to be assigned to the lower level class than similarly situated white students, and no explanatory reason was evident.³³

In addition, educational systems that were previously segregated by race in violation of the Fourteenth Amendment and have not achieved unitary status have an obligation to dismantle their prior *de jure* segregation. In such instances, when a school district or other educational system uses a test or assessment procedure for a high-stakes purpose that has racially disproportionate effects, the school district or other educational system must show that the disparity is not traceable to prior intentional segregation or that the test or assessment procedure does not perpetuate the adverse effects of such segregation.³⁴ The school district is under “a ‘heavy burden’ of showing that actions that increase[] or continue [] the effects of the dual system serve important and legitimate ends.”³⁵

b. **Disparate Impact**

Discrimination under federal law may also occur where the application of neutral criteria has discriminatory effects and those criteria are not educationally justified. The federal nondiscrimination regulations provide that a recipient of federal funds may not “utilize criteria or methods of administration which have the effect of subjecting individuals to discrimination.”³⁶ (For a further discussion of issues related to testing of students with

³³ See *People Who Care v. Rockford Bd. of Educ.*, 851 F. Supp. 905, 958-1001 (N.D. Ill. 1994), *remedial order rev'd, in part*, 111 F.3d 528 (7th Cir. 1997). On appeal, the Seventh Circuit Court of Appeals stated that the appropriate remedy in this case was to require the district to use objective, non-racial criteria to assign students to classes, rather than abolishing the district's tracking system. 111 F.3d at 536.

³⁴ See also *United States v. Fordice*, 505 U.S. 717, 731-732 (1992); *Debra P. v. Turlington*, 644 F.2d 397, 407 (5th Cir. 1981); *McNeal v. Tate County Sch. Dist.*, 508 F.2d 1017, 1020-1021 (5th Cir. 1975); *GI Forum v. Texas Educ. Agency*, No. SA-97-CA-1278-EP, 2000 U.S. Dist. LEXIS 153, slip op. at 56-57 (W.D. Tex. 2000).

³⁵ *Dayton Bd. of Educ. v. Brinkman*, 443 U.S. at 538 (quoting *Green v. Country School Board*, 391 U.S. 430, 439 (1968)).

³⁶ See 34 C.F.R. § 100.3(b)(2) (Title VI); 34 C.F.R. § 104.4(b)(4)(i) (Section 504); and 28 C.F.R. § 35.130(b)(3)(i) (Title II). See also 34 C.F.R. § 106.31 (Title IX). In *Guardians*, 463 U.S. at 589, the United States Supreme Court upheld the use of the effects test, stating that the Title VI regulation forbids the use of federal funds “not only in

disabilities, see below.)

The disparate impact analysis has been frequently misunderstood to indicate a violation of law based merely on disparities in student performance and to obligate educational institutions to change their policies and procedures to guarantee equal results. Under federal law, a statistically significant difference in outcomes creates the need for further examination of the educational practices in question that have caused the disparities in order to ensure accurate and nondiscriminatory decision making, but disparate impact alone is not sufficient to prove a violation of federal civil rights laws.

"It is ... important to note that group differences in test performance do not necessarily indicate problems in a test, because test scores may reflect real differences in achievement. These, in turn, may be due to a lack of access to a high quality curriculum and instruction. Thus, a finding of group differences calls for a careful effort to determine their cause."

National Research Council, *High Stakes: Testing for Tracking, Promotion, and Graduation*, 1999:5.

Courts applying the disparate impact test have generally examined three questions to determine if the practices at issue are discriminatory: (1) Does the practice or procedure in question result in substantial differences in the award of benefits or services based on race, national origin or sex? (2) Is the practice or procedure educationally justified? (3) Is there an equally effective alternative that can accomplish the institution's educational goal with less disparity?³⁷ Under the regulations implementing Title VI and Title IX, the party challenging the test has the burden of establishing disparate impact. If disparate impact is established, the educational institution must provide sufficient evidence of an educational justification for the practice in question. If sufficient evidence of an educational justification has been provided, the party challenging the test must then demonstrate, in order to prevail, that an alternative with less disparate impact is equally effective in meeting the institution's educational goals or needs.³⁸

2. Principles Relating to Inclusion and Accommodations

a. **Limited English Proficient Students**

programs that intentionally discriminate, but also in those endeavors that have a [racially disproportionate] impact on racial minorities."

³⁷ Courts use a variety of terms when discussing whether an alternative offered by the party challenging the practice is feasible and would also effectively meet the institution's goals. See, e.g., *Georgia State Conf. of Branches of NAACP v. Georgia*, 775 F.2d 1403, 1417 (11th Cir. 1985) (party challenging the practice "may ultimately prevail by proffering an equally effective alternative practice which results in less racial disproportionality"); *Sandoval v. Hagan*, 7 F.Supp.2d 1234, 1278 (M.D. Ala. 1998), *aff'd.*, 197 F.3d 484, 507 (11th Cir. 1999) (plaintiff may prevail by offering a "comparably effective" alternative practice which results in less proportionality). These terms appear to be used synonymously.

³⁸ See *Georgia State Conf.*, 775 F.2d at 1417. See also the Department of Justice's Title VI Legal Manual at p. 2.

The obligations of states and school districts with regard to high-stakes testing of limited English proficient students in elementary and secondary schools must be examined within the overall context of their Title VI obligation to provide equal educational opportunities to limited English proficient students. Under Title VI, school districts have an obligation to identify limited English proficient students and to provide them with a program that enables them to acquire English-language proficiency as well as the knowledge and skills that all students are required to master.³⁹

States or school districts using tests for high-stakes purposes must ensure that, as with all students, the tests effectively measure limited English proficient students' knowledge and skills in the particular content area being assessed. For limited English proficient elementary and secondary students in particular, it may be necessary in some situations to provide accommodations so that the tests provide accurate and valid information about the knowledge and skills intended to be measured.⁴⁰

b. Students with Disabilities

Under Section 504, Title II, and the IDEA,⁴¹ school districts have a responsibility to provide students with disabilities with a free appropriate public education. Providing effective instruction in the general curriculum for students with disabilities is an important aspect of providing a free appropriate public education. Under federal law, students with disabilities must be included in statewide or district-wide assessment programs and provided with appropriate accommodations, if necessary.⁴² There must be an individualized determination of whether a student with a disability will participate in a particular test and the appropriate accommodations, if any, that a student with a disability will need. The individualized determinations of whether a student with a disability will participate in a particular test, and what accommodations, if any, are appropriate must be addressed through the individualized education program (IEP) process or other applicable

³⁹ See Equal Educational Opportunities Act of 1974, P.L. No. 93-380, codified at 20 U.S.C. §§ 1701-1720; *Lau v. Nichols*, 414 U.S. at 568-569; *Castaneda v. Pickard*, 648 F.2d 989, 1011 (5th Cir. 1981); Memorandum to OCR Senior Staff from Michael L. Williams, Former Assistant Secretary for Civil Rights, September 27, 1991 (hereinafter *Williams Memorandum*).

⁴⁰ States and school districts are also required to provide LEP students with "reasonable adaptations and accommodations" in certain situations when using assessments for the purpose of holding schools and districts accountable for student performance under Title I. Title I of the Elementary and Secondary Education Act, 20 U.S.C. § 6311(a)(3)(F)(ii). Moreover, Title I requires States, to the extent practicable, to provide native-language assessments to LEP students for Title I accountability purposes if that is the language and form of assessment most likely to yield accurate and reliable information about what students know and can do. 20 U.S.C. § 6311(a)(3)(F)(iii). For a discussion of comparability issues arising in the testing of LEP students, see pages 38-42 of this guide.

⁴¹ The Section 504 regulation is found at 34 C.F.R. Part 104 (1999). The Title II regulation is found at 28 C.F.R. Part 35 (1999). The IDEA regulation is found at 34 C.F.R. Part 300 (1999).

⁴² States and school districts are also required to provide students with disabilities with "reasonable adaptations and accommodations" in certain situations when using assessments for the purpose of holding schools and districts accountable for student performance under Title I. 20 U.S.C. § 6311(a)(3)(F)(ii).

evaluation and placement processes and included in either the student's IEP or Section 504 plan.⁴³

Under Section 504, post-secondary education institutions may not make use of any test or criterion for admission that has a disproportionate adverse impact on individuals with disabilities unless (1) the test or criterion, as used by the institution, has been validated as a predictor of success in the education program or activity and (2) alternate tests or criteria that have a less disproportionate adverse impact are not shown to be available by the party asserting that the test or criterion is discriminatory.⁴⁴ Admissions tests must be selected and administered so as best to ensure that, when a test is administered to an applicant with a disability, the test results accurately reflect the applicant's aptitude or achievement level, rather than reflecting the effect of the disability (except where the functions impaired by the disability are the factors the test purports to measure).⁴⁵ Admissions tests designed for persons with impaired sensory, manual, or speaking skills must be offered as often and in as timely a manner as are other admissions tests. Admissions tests must be offered in facilities that, on the whole, are accessible to individuals with disabilities.

3. Federal Constitutional Questions Related to Testing of Elementary and Secondary Students For High-Stakes Purposes

interest
The equal protection and due process requirements of the Fifth and Fourteenth Amendments to the U.S. Constitution would apply to ensure that high-stakes decisions by public schools or states based on test use are made appropriately.⁴⁶ The equal protection principles involved in discrimination cases are, generally speaking, the same as the standards applied to intentional discrimination claims under the applicable federal nondiscrimination statutes.⁴⁷ Courts addressing due process claims have examined three questions related to the use of tests as bases for promotion or graduation decisions:

⁴³ Under the IDEA, students with disabilities must be included in state and district-wide assessment programs. See 34 C.F.R. § 300.138(a). However, if the IEP team determines that a student should not participate in a particular statewide or district-wide assessment of student achievement (or part of such an assessment), the student's IEP must include statements of why that test is not appropriate for the student and how the student will be assessed. See 34 C.F.R. § 300.347(a)(5). The IDEA also requires state or local educational agencies to develop guidelines for students with disabilities who cannot take part in state and district-wide assessments to participate in alternate assessments; these alternate assessments must be developed and conducted beginning not later than July 1, 2000. See 34 § C.F.R. 300.138(b).

⁴⁴ See 34 C.F.R. § 104.42(b)(2).

⁴⁵ See 34 C.F.R. § 104.42(b)(3).

⁴⁶ The requirements of Title VI, Title IX and Section 504 apply only to recipients of federal financial assistance. The protections afforded by the Fifth and Fourteenth Amendments to the U.S. Constitution extend to actions by governmental entities that are "state actors" and are not dependent on their receipt of federal financial assistance.

⁴⁷ Federal cases may involve equal protection challenges to a jurisdiction's use of tests in which the claim is not based on intentional race or sex discrimination, but, instead, on the alleged impropriety of the jurisdiction's use of tests to separate out those students who should not be allowed to graduate. As a general matter, courts express reluctance to second guess a state's educational policy choices when faced with such challenges, although they recognize that a state cannot "exercise that [plenary] power without reason and without regard to the United States Constitution." See *Debra P. v. Turlington*, 644 F.2d 397, 403 (5th Cir. 1981). When there is no claim of discrimination based on membership in a

- Is the purpose of the testing program legitimate and reasonable?⁴⁸
- Have students received adequate notice of the test and its consequences?⁴⁹
- Have students actually been taught the knowledge and skills measured by the test?⁵⁰

Federal courts have typically deferred to educators' judgments about the beneficial educational purposes of a testing program, as long as these judgments are not arbitrary or capricious.⁵¹ Improving the quality of education, ensuring that students can compete on a national and international level, and encouraging educational achievement through the establishment of academic standards have been found to be reasonable goals for testing programs.⁵²

Courts have generally required advance notice of test requirements in order to give students a reasonable chance to understand the standards against which they will be evaluated and to learn the material for which they are to be accountable. A reasonable transition period is required between the development of a new academic requirement and the attachment of high-stakes consequences to tests used to measure academic

suspect class, the equal protection claim is reviewed under the rational basis standard. In these cases, the jurisdiction need show only that the use of the tests has a rational relationship to a valid state interest. See *Debra P.*, 644 F.2d at 406; *Erik V. v. Causby*, 977 F. Supp. 384, 389 (E.D. N.C. 1997).

⁴⁸ See *Regents of the Univ. of Mich. v. Ewing*, 474 U.S. 214, 222, 226-27 (1985); *Debra P.*, 644 F.2d at 406; *Anderson v. Banks*, 520 F. Supp. 472, 506 (S.D. Ga. 1981).

⁴⁹ See *Brookhart v. Illinois State Bd. of Educ.*, 697 F.2d 179, 185 (7th Cir. 1983); *Debra P.*, 644 F.2d at 404; *Erik V.*, 977 F. Supp. at 389-90 (E.D. N.C. 1997); *Anderson*, 520 F. Supp. at 1410-12.

⁵⁰ See *Brookhart*, 697 F.2d at 184-87; *Debra P.*, 644 F.2d at 406; *Anderson*, 520 F. Supp. at 509. Insofar as due process cases may involve additional questions regarding the validity, reliability, and fairness of the test used to address the educational institution's stated purposes, these issues are discussed in the portions of the guide addressing discrimination under federal civil rights laws.

⁵¹ See *Ewing*, 474 U.S. at 226-27; *Debra P.*, 644 F.2d at 406; *Anderson*, 520 F. Supp. at 506.

⁵² See *Ewing*, 474 U.S. at 226-27; *Debra P.*, 644 F.2d at 406; *Anderson*, 520 F. Supp. at 506.

achievement. That time period varies, however, depending upon the precise context in which the high-stakes decision is to be made. Relevant inquiries affecting determinations about the constitutionality of notice and timing have included questions about the alignment of curriculum and instruction with material tested, the number of test taking opportunities provided to students, tutorial or remedial opportunities provided to students, and whether factors in addition to test scores can affect high-stakes decisions.

bullets?

Ultimately, in due process cases, federal courts have required, as a matter of “fundamental fairness,” that students have a reasonable opportunity to learn the material covered by the test where passing the test is a condition of receipt of a high school diploma or a condition for grade-to-grade promotion.⁵³ For the test to meaningfully measure student achievement, the test, the curriculum, and classroom instruction should be aligned.

⁵³ See *Brookhart*, 697 F.2d at 184-87; *Debra P.*, 644 F.2d at 406; *GI Forum*, 2000 U.S. Dist. LEXIS 153, slip op. at 50-51; *Anderson*, 520 F. Supp. at 509.

CHAPTER 1. Test Measurement Principles

This chapter explains basic test measurement standards and related educational principles for determining whether tests that are being used to make high-stakes educational decisions for students provide accurate and fair information. As explained in chapter two below, federal court decisions have been informed and guided by professional test measurement standards and principles. Professional test measurement standards, products of the test measurement community, can provide a basis for compliance with federal nondiscrimination laws.⁵⁴ This chapter is intended as a helpful discussion of how to understand test measurement concepts and their use. These are not specific legal requirements, but rather are foundations for understanding appropriate test use.

Educational institutions use tests to accomplish specific purposes based on their educational goals, including making placement, promotion, graduation, admissions, and other decisions. It is only after they have determined the underlying goal they want to accomplish that they can identify the types of information that will best inform their decision making. Information may include test results, as well as other relevant measures, that will be able to effectively, accurately, and fairly address the purposes and goals specified by the institutions.⁵⁵ As stated in the *Joint Standards*, “[w]hen interpreting and using scores about individuals or groups of students, considerations of relevant collateral information can enhance the validity of the interpretation, by providing corroborating evidence or evidence that helps explain student performance....As the stakes of testing increase for individual students, the importance of considering additional evidence to document the validity of score interpretations and the fairness in testing increases accordingly.”⁵⁶

In using tests to make high-stakes decisions, educational institutions should ensure that the test will provide accurate results that are valid, reliable, and fair for all test takers. This includes requesting adequate evidence of test quality, evaluating the evidence, and ensuring that appropriate test use is based on adequate evidence provided by the developers or through supplemental validation studies.⁵⁷ When test results are used to make high-stakes decisions about student promotion or graduation, evidence should be

Perhaps more important than a footnote

⁵⁴ See, e.g., *High Stakes*, p. 59-60.

⁵⁵ Among other considerations, institutions will determine if they want test score interpretations that are norm-referenced or criterion-referenced, or both. Norm-referenced means that the performances of students are compared to the performances of other students in a specified reference population; criterion-referenced indicates the extent to which students have mastered specific knowledge and skills.

⁵⁶ *Joint Standards*, p. 141. See also Standard 13.7, which states, “In educational settings, a decision or characterization that will have a major impact on a student should not be made on the basis of a single test score. Other relevant information should be taken into account if it will enhance the overall validity of the decision.”

⁵⁷ In order to provide educational institutions with tests that are accurate and fair, test developers should develop tests in accordance with professionally recognized standards, and provide educational institutions with adequate evidence of test quality.

available which documents that students have had an adequate opportunity to learn the material being tested.⁵⁸

I. Key Considerations in Test Use

This section addresses the fundamental concepts of test validity and reliability. It will also discuss issues associated with ensuring fairness in the meaning of test scores, and issues related to using appropriate cutscores in high-stakes tests.

A. Validity

Test validity refers to a determination of how well a test actually measures what it says it measures. The *Joint Standards* define validity as “[t]he degree to which accumulated evidence and theory support specific interpretations of test scores entailed by proposed uses of a test.”⁵⁹ The demonstration of validity is multifaceted and must always be determined within the context of the specific use of a test. In order to promote readability, the discussion on validity presented here is meant to reflect this complex topic in an accurate, but concise and user-friendly way. The *Joint Standards* identify and discuss in detail principles related to determining the validity of test scores within the context of their use, and readers are encouraged to review the *Joint Standards*, Chapter 1, Validity, for additional, relevant discussion.⁶⁰

There are three central points to keep in mind:

- The focus of validity is not really on the test itself, but on the validity of the inferences drawn from the test results for a given use.
- All validity is really a form of “construct validity.” — ?
- In validating the inferences of the test results, one must also consider the consequences of the test’s interpretation and use.

⁵⁸ Standards 13.5 and 7.5. Standard 13.5, *supra* note 19.

Standard 7.5 states, “In testing applications involving individualized interpretations of test scores other than selection, a test taker’s score should not be accepted as a reflection of standing on the characteristic being assessed without consideration of alternate explanations for the test taker’s performance on that test at that time.”

⁵⁹ *Joint Standards*, p. 9, 184.

⁶⁰ *Joint Standards*, Chapter 1, Validity, p. 9-24.

1. Validity of the Inferences of the Scores

It is not the test that is validated per se, but the inferences or meaning derived from the test scores for a given use—that is, for a specific purpose, in a specific type of situation, and with specific groups of students. The meaning of test scores will differ based on such factors as how the test is designed, the types of questions that are asked, and the documentation that supports how all groups of students are interpreting what the test is asking and how effectively their performance can be generalized beyond the test.

For instance, in one case, the educational institution may want to evaluate how well students can analyze complex issues and evaluate implications in history. For a given amount of test time, they would want to use a test that measures the ability of students to think deeply about a few selected history topics. The meaning of the scores should reflect this purpose and the limits of the range of topics being measured on the test. In another case, the institution may want to assess how well students know a range of facts about a wide variety of historical events. The institution would want to use a test that measures a broad range of knowledge about many different occurrences in history. The inferences of the scores should accurately reflect how well students know a broad range of historical facts.

6x
Box?
5th

2. Construct Validity

Construct validity refers to the degree to which the scores of test takers accurately reflect the constructs a test is attempting to measure. The *Joint Standards* defines a construct as “the concept or the characteristic that a test is designed to measure.”⁶¹ Test scores and their inferences are validated to measure one or more constructs described in a particular content domain.⁶² In K-12 education, these domains are often explained in state or district content standards in various subject areas.

For instance, in mathematics, constructs of mathematical problem solving and the knowledge of number systems would be among the constructs described in a state’s elementary mathematics content standards. These standards would define the mathematics domain in this situation. Items would be selected for the test that sample from this domain, and are properly representative of the constructs identified within it. The meaning of the test scores should accurately reflect the knowledge and skills defined in the mathematics content standards domain.

Validity should be viewed as the overarching, integrative evaluation of the degree to which all accumulated evidence supports the intended interpretation of the test scores for

⁶¹ Page 173.

⁶² The *Joint Standards* defines a content domain as “the set of behaviors, knowledge, skills, abilities, attitudes or other characteristics to be measured by a test, represented in a detailed specification, and often organized into categories by which items are classified (p.174).” A domain, then, represents a definition of a content area for the purposes of a particular test. Other tests will likely have a different definition of what knowledge and skills a particular content area entails.

Box
Example
Box

a proposed purpose.⁶³ This unitary and comprehensive concept of validity is referred to as “construct validity.” Different sources of evidence may illuminate different aspects of validity, but they do not represent distinct types of validity.⁶⁴

Therefore, “construct validity” is not just one of the many types of validity—it *is* validity. Demonstrating construct validity then means gathering a variety of types of evidence to support the intended interpretations and uses of test scores. All validity evidence and the interpretation of the evidence are focused on the basic question: Is the test measuring the concept, skill, or trait in question? Is it, for example, really measuring mathematical reasoning or reading comprehension for the types of students that are being tested? A variety of types of evidence can be used to answer this question—none of which provides a simple yes or no answer. The exact nature of the types of evidence that needs to be accumulated is directly related to the intended use of the test, which includes information regarding the skills and knowledge being measured, the purpose for which the information will be used, and the population of test takers.⁶⁵

For instance, an educational institution may want to use a test to help make promotion decisions. It may also want to use a test to place students in the appropriate sequence of courses. In each situation, the types of validity evidence an institution would expect to see would depend on how the test is being used.

In making promotion decisions, the test should reflect content the student has learned. Appropriate validation would include adequate evidence that the test is measuring the constructs identified in the curriculum, and that the inferences of the scores accurately reflect the intended constructs for all test takers. Validation of the decision process involving the use of the test would include adequate evidence that low scores reflect lack of knowledge of students after they have been taught the material, rather than lack of exposure to the curriculum in the first place.

In making placement decisions, on the other hand, the test may not need to measure content that the student has already learned. Rather, at least in part, the educational institution may want the test to measure aptitude for the future learning of knowledge or skills that have been identified as necessary to complete a course sequence. Appropriate validation would include documentation of the relationship between what constructs are being measured in the test, and what skills and knowledge are actually needed in the

⁶³ *Joint Standards*, Chapter 1, Validity, pp. 9-11, 184.

⁶⁴ Therefore, construct validity can be seen as an umbrella that encompasses what has previously been described as predictive validity, content validity, criterion validity, discriminant validity, etc. Rather, these terms refer to types or sources of evidence that can be accumulated to support the validity argument. Definitions of these terms can be found in Appendix B, Measurement Glossary.

⁶⁵ Rather than follow the traditional nomenclature (e.g. predictive validity, content validity, criterion validity, discriminant validity, etc.), the *Joint Standards* define sources of validity evidence as evidence based on test content, evidence based on response processes, evidence based on internal structure, evidence based on relations to other variables, and evidence based on consequences of testing. These are discussed in Chapter 1 of the *Joint Standards*, p. 11-17.

future placements. Differential evidence would provide documentation that scores are not significantly confounded by other factors irrelevant to the knowledge and skills the test is intending to measure.

Institutions often think about using the same test for two or more purposes. This is appropriate as long as the validity evidence properly supports the use for the test for each purpose, and properly supports that the inferences of the results accurately reflect what the test is measuring for all students taking the test.

The empirical evidence related to the various aspects of construct validity is collected throughout test development, during test construction, and after the test is completed. It is important for educators and policymakers to understand and expect that the accumulated evidence spans the range of test development and implementation. There is not just one set of documentation collected at one point in time.

Handwritten note: } Weird format

When the empirical database is large and includes results from a number of studies related to a given purpose, situation, and type of test takers, it may be appropriate to generalize validity findings beyond validity data gathered for one particular test use. That is, it may be appropriate to use evidence collected in one setting when determining the validity of the meaning of the test scores for a similar use. If the accumulated validity evidence for a particular purpose, situation, or subgroup is small, or features of the proposed use of the test differ markedly from an adequate amount of validity evidence already collected, evidence from this particular type of test use will generally need to be compiled.⁶⁶ Regardless of where the evidence is collected, educational institutions should expect adequate documentation of construct validity based on needs defined by the particular purposes and populations for which a test is being used.

a. Sources of Validity Error

When considering the types of construct validity evidence to collect, the *Joint Standards* emphasize that it is important to guard against the two major sources of validity error. This error can distort the intended meaning of scores for particular groups of students, situations, or purposes.⁶⁷

One potential source of error omits some important aspects of the intended construct being tested. This is called construct underrepresentation⁶⁸ An example would be a test that is being

Handwritten note: bold, a set off as a definition

⁶⁶ As indicated in the *Joint Standards*, "The extent to which predictive or concurrent evidence of validity generalization can be found in new situations is in large measure a function of accumulated research. Although evidence of generalization can often help to support a claim of validity in a new situation, the extent of available data limits the extent to which the claim can be sustained." *Joint Standards*, Chapter 1, p. 15-16.

⁶⁷ *Joint Standards*, Chapter 1, Validity, p. 10.

⁶⁸ Messick, S. (1989). Validity. In *Educational Measurement, 3rd Edition*, R.L. Linn, ed. New York: Macmillan, p. 13-103.

used to measure English language proficiency. When the institution has defined English language proficiency as including specific skills in listening, speaking, reading, and writing the English language, and wants to use a test which measures these aspects, construct underrepresentation would occur if the test only measured the reading skills.

The other potential source of error occurs when a test measures material that is extraneous to the intended construct, confounding the ability of the test to measure the construct that it intends to measure. This source of error is called construct irrelevance.⁶⁹ For instance, how well a student reads a mathematics test may influence the student's subtest score in mathematics computation. In this case, the student's reading skills are irrelevant when the skill of mathematics computation is what is being measured by the subtest.⁷⁰

ie.
LEP or
systemic

An essential part of the accumulated validity information is collecting evidence not only about what a test measures in particular situations or for particular students, but also evidence that seeks to document that the intended meaning of the test scores is not unduly influenced by either of the two sources of validity error.

3. Considering the Consequences of Test Use

Evidence about the intended and unintended consequences of test use can provide important information about the validity of the inferences of the test results, or it can raise concerns about an inappropriate use of a test where the inferences may be valid for other uses.

For instance, significant differences in placement test scores based on race, gender, or national origin may trigger a further inquiry about the test and how it is being used to make placement decisions.⁷¹ The validity of the test scores would be called into question if the test scores are substantially affected by irrelevant factors that are not related to the academic knowledge and skills that the test is supposed to measure.⁷²

Messick, S. (1995). Validity of psychological assessment: Validation of inferences from persons' responses and performances as scientific inquiry into score meaning. *American Psychologist* 50(9): p.741-749.

⁶⁹ Messick, 1989; 1995.

⁷⁰ On the other hand, if an item is measuring the student's ability to apply mathematical skills in a written format (for instance when an item requires students to fill out an order form), then writing skills may not be extraneous to the construct being measured in this item.

⁷¹ See *Code of Fair Testing Practices in Education*, 1988.

⁷² Standards 7.5, 7.6 and 1.24. Standard 7.5, *supra* note 58.

Standard 7.6 states, "When empirical studies of differential prediction of a criterion for members of different subgroups are conducted, they should include regression equations (or an appropriate equivalent) computed separately for each group or treatment under consideration or an analysis in which the group or treatment variables are entered as moderator variables."

Standard 1.24 states, "When unintended consequences result from test use, an attempt should be made to investigate whether such consequences arise from the test's sensitivity to characteristics other than those it is intended to assess or to the test's failure fully to represent the intended construct."

On the other hand, a test may accurately measure differences in the level of students' academic achievement. That is, low scores may accurately reflect that some students do not know the content. However, test users should ensure that they interpret those scores correctly in the context of their high-stakes decisions.⁷³ For instance, test users could incorrectly conclude that the scores reflect lack of ability to master the content for some students when, in fact, the low test scores reflect the limited educational opportunities that the students have received. In this case, it would be problematic to use the test scores to place low performing students in a special services program for students who have trouble learning and processing academic content. It would be appropriate to use the test to evaluate program effectiveness, however.⁷⁴

Standard 13.1

When educational testing programs are mandated by school, district, state, or other authorities, the ways in which test results are intended to be used should be clearly described. It is the responsibility of those who mandate the use of tests to monitor their impact and to identify and minimize potential negative consequences. Consequences resulting from the uses of the test, both intended and unintended, should also be examined by the test user.

B. Reliability

Reliability refers to the consistency of test results. While no test is ever an "error-free" measure of student performance,⁷⁵ inferences of adequate test reliability refer to estimates which demonstrate that the inconsistency of the scores are minimized over test administrations, forms, items, scorers, and/or other facets of testing.⁷⁶ An example of reliability of test results on different occasions is when the same students, taking the test multiple times, receive similar scores. Consistency over parallel forms of a test occurs

⁷³ Standards 7.5 and 7.10. Standard 7.5, *supra* note 58. Standard 7.10, *supra* note 25.

⁷⁴ *High Stakes*, p. 89-113.

⁷⁵ All sources of assessment information, including test results, include some degree of error. There are two types of error. The first is random error that affects scores in such a way that sometimes students will score lower and sometimes higher than their "true" score (the actual mastery of the students' knowledge and skills). This type of error, also known as measurement error, particularly affects reliability of scores. Therefore, test scores are considered reliable when evidence demonstrates that there is a minimum amount of random measurement error in the test scores for a given group.

The second type of error that affects test results is systematic error. Systematic error consistently affects scores in one direction; that is, this type of error causes some students to consistently score lower or consistently score higher than their "true" (or actual) level of mastery. For instance, visually impaired students will consistently score lower than they should on a test which has not been administered for them in Braille or large print, because their difficulty in reading the items on the page will negatively impact their score. This type of error generally affects the validity of the interpretation of the test results and is discussed in the validity section above. Systematic error should also be minimized in a test for all test takers.

When educators and policy makers are evaluating the adequacy of a test for their local population of students, it is important to consider evidence concerning both types of error.

⁷⁶ Evaluating the reliability of a test includes identifying the major sources of measurement error, the size of the errors resulting from these sources, the indication of the degree of reliability to be expected, and the generalizability of results across items, forms, raters, sampling, administrations, and other measurement facets.

when forms are developed to be equivalent in content and technical characteristics. Reliability can also include estimates of a high degree of relationship across similar items within a single test or subtest that are intended to measure the same knowledge or skill. For judgmentally scored tests, such as essays, another widely used index of reliability ✓ addresses consistency across raters or scorers. In each case, reliability can be estimated in different ways, using one of several statistical procedures.⁷⁷ Different kinds of reliability estimates vary in degree and nature of generalization.

In order to promote readability, the discussion on reliability presented here is meant to reflect this complex topic in an accurate, but concise and user-friendly way. Readers are encouraged to review Chapter 2, Reliability and Errors of Measurement, in the *Joint Standards* for additional, relevant information.⁷⁸

↳ Just say it
is not a complete
technical examination
of what constitutes
"test reliability"
+ rater ratings
of Joint Standards

— maybe
not a footnote?

⁷⁷ These types of reliability estimates are known as test-retest, alternate forms, internal consistency, and inter-rater estimates, respectively. See *Joint Standards*, Chapter 2, Reliability, for some examples of different procedures.

⁷⁸ *Joint Standards*, Chapter 2, Reliability and Errors of Measurement, p. 25-36.

C. Fairness

Tests are fair when they yield score interpretations that are valid and reliable for all students who take the tests. That is, the academic tests must measure the same academic constructs (knowledge and skills) for all students who take them, regardless of race, national origin, gender, or disability. Similarly, the scores must not substantially and systematically underestimate or overestimate the knowledge or skills of members of a particular group. The *Joint Standards* discuss fairness in testing in terms of lack of bias, equitable treatment in the testing process, equal scores for students who have equal standing on the tested construct, and equity in opportunity to learn the material being tested.⁷⁹ In order to promote readability, the discussion on fairness presented here is meant to reflect this complex topic in an accurate, but concise and user-friendly way. Readers are encouraged to review Chapter 7, Fairness in Testing and Test Use, in the *Joint Standards* for additional, relevant information.⁸⁰

"Fairness, like validity, cannot be properly addressed as an afterthought. . . . It must be confronted throughout the interconnected phases of the testing process, from test design and development to administration, scoring, interpretation, and use"

National Research Council, *High Stakes: Testing for Tracking, Promotion and Graduation*, 1999: pages 80-81.

1. Fairness in Validity

Demonstrating fairness in the validation of test score inferences focuses primarily on making sure that the scores reflect the same intended knowledge and skills for all students taking the test. For the most part this means that the test should minimize the measurement of material that is extraneous to the intended constructs and which confounds the ability of the test to accurately measure the constructs that it intends to measure. Rather, a test score should accurately reflect how well each student has mastered the intended constructs. The score should not be significantly impacted by construct-irrelevant influences.

⁷⁹ *Joint Standards*, Chapter 7, Fairness in Testing and Test Use, p. 74-80. In test measurement, the term fairness has a specific set of technical interpretations. Four of these interpretations are discussed in the *Joint Standards*. For instance, bias is discussed in relation to fairness and is defined in the *Joint Standards* in two ways: "In a statistical context, (bias refers to) a systematic error in a test score. In discussing test fairness, bias (also) may refer to construct underrepresentation or construct-irrelevant components of test scores that differentially affect the performance of different groups of test takers (p. 172)." Fairness as equitable treatment in the testing process "requires consideration not only of the test itself, but also the context and purpose of testing, and the manner for which test scores are used (p. 74)." Equal scores for students of equal standing reflects that "examinees of equal standing with respect to the construct the test is intended to measure should on average earn the same test score, irrespective of group membership (p. 74)." For educational achievement tests, "When some test takers have not had the opportunity to learn the subject matter covered by the test content, they are likely to get low scores. . . . low scores may have resulted in part from not having had the opportunity to learn the material tested as well as from having had the opportunity and failed to learn (p. 76)."

⁸⁰ *Joint Standards*, Chapter 7, Fairness in Testing and Test Use, p. 73-84.

The *Joint Standards* identify a number of standards that outline important elements related to validly measuring the intended constructs for all students.⁸¹ The elements span considerations of test development, test implementation, and the proper use of reported test results.

Documenting fairness during test development involves gathering adequate evidence that items and test scores are constructed so that the inferences validly reflect what is intended. For all test takers, evidence should support that valid inferences can be drawn from the scores.⁸² When credible research reports that item and test results differ in meaning across examinee subgroups, then to the extent feasible, separate validity evidence for each relevant subgroup should be collected.⁸³ When items function differently across relevant subgroups, appropriate studies should be conducted, when feasible, so that bias in items due to test design, content, and format is detected and eliminated.⁸⁴ Developers should strive to identify and eliminate language, form, and content in tests that have a different meaning in one subgroup than in others, or that generally have sensitive connotations, except when judged to be necessary for adequate representation of the intended constructs.⁸⁵ Adequate differential analyses should be conducted when evaluating the validity of scores for prediction purposes.⁸⁶

?

⁸¹ *Joint Standards*, Chapter 7, Fairness in Testing and Test Use, p. 80-84.

⁸² Standard 7.2 states, "When credible research reports differences in the effects of construct-irrelevant variance across subgroups of test takers on performance of some part of the test, the test should be used if at all only for those subgroups for which evidence indicates that valid inferences can be drawn from test scores."

⁸³ Standard 7.1 and 7.3. Standard 7.1 states, "When credible research reports that test scores differ in meaning across examinee subgroups for the type of test in question, then to the extent feasible, the same forms of validity evidence collected for the examinee population as a whole should also be collected for each relevant subgroup. Subgroups may be found to differ with respect to appropriateness of test content, internal structure of test responses, the relation of test scores to other variables, or the response processes employed by individual examinees. Any such findings should receive due consideration in the interpretation and use of scores as well as in subsequent test revisions."

Standard 7.3 states, "When credible research reports that differential item functioning exists across age, gender, racial/ethnic, cultural, disability and/or linguistic groups in the population of test takers in the content domain measured by the test, test developers should conduct appropriate studies when feasible. Such research should seek to detect and eliminate aspects of test design, content, and format that might bias test scores for particular groups."

⁸⁴ See Standard 7.3, *supra* note 83.

⁸⁵ Standard 7.3 and Standard 7.4. Standard 7.3, *supra* note 83.

Standard 7.4 states, "Test developers should strive to identify and eliminate language, symbols, words, phrases, and content that are generally regarded as offensive by members of racial, ethnic, gender, or other groups, except when judged to be necessary for adequate representation of the domain." Comment: "Two issues are involved. The first deals with the inadvertent use of language that, unknown to the test developer, has a different meaning or connotation in one subgroup than in others. Test publishers often conduct sensitivity reviews of all test material to detect and remove sensitive material from the test. The second deals with settings in which sensitive material is essential for validity. For example, history tests may appropriately include material on slavery or Nazis. Tests on subjects from life sciences may appropriately include material on evolution. A test of understanding of an organization's sexual harassment policy may require employees to evaluate examples of potentially offensive behavior."

⁸⁶ Standard 7.6, *supra* note 72.

not 2
footnote

Adequate evidence should document the fair implementation of tests for all test takers. The testing process should reflect equitable treatment for all examinees.⁸⁷ Linguistic or reading demands in tests should be kept to a minimum except when these constructs are being measured.⁸⁸

Documentation of appropriate reporting and test use should be available. Reported data should be clear and accurate, especially when there are high-stakes consequences for students.⁸⁹ When tests are used in decisions that have high-stakes consequences for students, evidence of mean score differences between relevant subgroups should be examined, where feasible. When mean differences are found between subgroups, investigations should be undertaken to determine that such differences are not attributable to construct underrepresentation or construct irrelevant error.⁹⁰ Evidence about differences in mean scores and the significance of the validity errors should also be considered when deciding which test to use.⁹¹ In using test results for purposes other than selection, a test taker's score should not be accepted as a reflection of standing on the intended constructs without consideration of alternative explanations for the test taker's performance.⁹² Explanations might reflect limitations of the test, for instance construct irrelevant factors may have significantly impacted the student's score. Explanations may also reflect schooling factors external to the test, for instance lack of instructional opportunities.

The issue of feasibility is discussed in a few of the standards summarized above. In the comments associated with these standards, feasibility is generally addressed in terms of adequate sample size, with continued operational use of a test as a way of accumulating adequate numbers of subgroup results over administrations. When credible research reports that results differ in meaning across subgroups, collecting separate and parallel validity data verifies that the same knowledge and skills are being measured for all test

⁸⁷ Standard 7.12 states, "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."

⁸⁸ Standard 7.7 states, "In testing applications where the level of linguistic or reading ability is not part of the construct of interest, the linguistic or reading demands of the test should be kept to the minimum necessary for the valid assessment of the intended construct."

⁸⁹ Standards 7.8, 7.9, 7.10, 1.24. Standard 7.8 states, "When scores are disaggregated and publicly reported for groups identified by characteristics such as gender, ethnicity, age, language proficiency, or disability, cautionary statements should be included whenever credible research reports that test scores may not have comparable meaning across these different groups."

Standard 7.9 states, "When tests or assessments are proposed for use as instruments of social, educational, or public policy, the test developers or users proposing the test should fully and accurately inform policymakers of the characteristics of the tests as well as any relevant and credible information that may be available concerning the likely consequences of test use."

Standard 7.10, *supra* note 25. Standard 1.24, *supra* note 72.

⁹⁰ Standard 7.10, *supra* note 25.

⁹¹ Standard 7.11 states, "When a construct can be measured in different ways that are approximately equal in their degree of construct representation and freedom from construct-irrelevant variance, evidence of mean score differences across relevant subgroups of examinees should be considered in deciding which test to use."

⁹² Standard 7.5, *supra* note 58.

takers. Particularly in high-stakes situations, feasibility decisions need to include the potential costs to students of using information where the validity of the scores has not been verified.⁹³

2. Fairness in Reliability

Fairness in reliability focuses on making sure that scores are stable and consistently accurate for all students. Two standards discuss issues of fairness in reliability. First, when there are reasons for expecting that test reliability analyses might differ substantially for different subpopulations, reliability data should be presented as soon as feasible for each major population for whom the test is recommended.⁹⁴ Second, "[w]hen significant variations are permitted in test administration procedures, separate reliability analyses should be provided for scores produced under each major variation if adequate sample sizes are available."⁹⁵ Often, continued operational use of a test is a way to accumulate an adequate sample size over administrations.

✓
Good to point out
reliability is messy
through use of test

D. Cutscores

The same principles regarding fairness, validity, and reliability apply generally to the establishment and use of cutscores for the purpose of making high-stakes educational decisions. Cutscores, also known as cut points or cutoff scores, are specific points on the test or scale where test results are used to divide levels of knowledge, skill, or ability. A cutscore may divide the demonstration of acceptable and unacceptable skills, as in placement in gifted and talented programs where students are accepted or rejected. There may be multiple cutscores that identify qualitatively distinct levels of performance. Cutscores are used in a variety of contexts, including decisions for placement purposes or for other specific outcomes, such as graduation, promotion, or admissions.⁹⁶

ex-
Proficient,
Basic,
etc

⁹³ See comment associated with Standard 10.7: "In addition to modifying tests and test administration procedures for people who have disabilities, evidence of validity for inferences drawn from these tests is needed. *Validation is the only way to amass knowledge about the usefulness of modified tests for people with disabilities. The costs of obtaining validity evidence should be considered in light of the consequences of not having usable information regarding the meanings of scores for people with disabilities.* This standard is feasible in the limited circumstances where a sufficient number of individuals with the same level or degree of a given disability is available (italics added)."

⁹⁴ Standard 2.11 states, "If there are generally accepted theoretical or empirical reasons for expecting that reliability coefficients, standard errors of measurement, or test information functions will differ substantially for various subpopulations, publishers should provide reliability data as soon as feasible for each major population for which the test is recommended."

⁹⁵ Standard 2.18.

⁹⁶ In order to promote readability, the discussion on cutscores presented here is meant to reflect this complex topic in an accurate, but concise and user-friendly way. Readers are encouraged to review Chapter 4, Scales, Norms, and Score Comparability, p. 53-54, in the *Joint Standards* for additional, relevant information about cutscores. See also Standards 1.19, 13.9.

Standard 1.19 states, "If a test is recommended for use in assigning persons to alternative treatments or is likely to be so used, and if outcomes from those treatments can reasonably be compared on a common criterion, then, whenever feasible, supporting evidence of differential outcomes should be provided."

Standard 13.9, *supra* note 15.

Many of the concepts regarding test validity apply to cutscores—that is, the cut points themselves must be accurate representations of the knowledge and skills of students.⁹⁷ Further, the validity evidence for cutscores should generally be able to demonstrate that students above the cut point represent or demonstrate a qualitatively greater degree or different type of skills and knowledge than those below the cut point, whenever these types of inferences are made.⁹⁸

Where the results of the [cutscore] setting process have highly significant consequences, ... those responsible for establishing cutscores should be concerned that the process... [is] clearly documented and defensible.

Joint Standards: page 54

Reliability of the cutscores is also important. The *Joint Standards* state that where cutscores are specified for selection or placement, the degree of measurement error around each cutscore should be reported.⁹⁹ Evidence should also indicate the misclassification rates, or percentage of error in classifying students, that is likely to occur among students with comparable knowledge and skills.¹⁰⁰ This information should be available by group as soon as feasible if there is a prior probability that the misclassification rates may differ substantially by group.¹⁰¹ For example, what percentage of students who should be allowed to graduate would not be allowed to do so because of error due to the test rather than differences in their actual knowledge and skills?¹⁰²

There is no single right answer to the questions of when, where and how cutscores should be set on a test with high-stakes consequences for students.¹⁰³ Many experts suggest,

⁹⁷ *Joint Standards*, Chapter 1, Validity, p. 9-16, discusses that the interpretation of *all* scores should be an accurate representation of what is being measured.

⁹⁸ See Standard 4.20's comment section for a discussion on these points. In high-stakes situations, it is important to examine the validity of the inferences that underlie the specific decisions being made on the basis of the cutscores. In other words, what must be validated is the specific use of the test based on how the scores of students above and below the cutscore are being interpreted. What is also at issue is how scores clustered around the cut-off point are interpreted in light of the high-stakes decision.

⁹⁹ Standard 2.14 states, "Conditional standard errors of measurement should be reported at several score levels if constancy cannot be assumed. Where cut scores are specified for selection or classification, the standard errors of measurement should be reported in the vicinity of each cut score."

¹⁰⁰ "Where the purpose of measurement is classification, some measurement errors are more serious than others. An individual who is far above or far below the value established for pass/fail or for eligibility for a special program can be mismeasured without serious consequences: Mismeasurement of examinees whose true scores are close to the cut score is a more serious concern.... The term *classification consistency* or *inter-rater agreement*, rather than *reliability*, would be used in discussions of consistency of classification. Adoption of such usage would make it clear that the importance of an error of any given size depends on the proximity of the examinee's score to the cut score." *Joint Standards*, p. 30.

¹⁰¹ Standard 2.11, *supra* note 94.

¹⁰² Misclassification of students above or below the cutpoints can result in both false positive and false negative classifications, respectively. The example in the text is a false negative classification.

¹⁰³ *High Stakes*, Chapter 7, p. 168.

however, that multiple methods of determining cutscores should be used when determining a final cutscore.¹⁰⁴ Further, the reasonableness of the standard setting process and the consequences for students should be clearly and specifically documented for a given use.¹⁰⁵ Both the *Joint Standards* and *High Stakes* repeatedly state that decisions should not be made solely or automatically on the basis of a single test score, and that other relevant information should be taken into account if it will enhance the overall validity of the decision.¹⁰⁶

¹⁰⁴ *High Stakes*, Chapter 7, p.169.

¹⁰⁵ See Standards 4.19 and 4.21 and their comments. See also *High Stakes*, Chapters 5,6,7.

Standard 4.19 states, "When proposed score interpretations involve one or more cut scores, the rationale and procedures used for establishing cut scores should be clearly documented."

Standard 4.21 states, "When cut scores defining pass-fail or proficiency categories are based on direct judgments about the adequacy of item or test performances or performance levels, the judgmental process should be designed so that judges can bring their knowledge and experience to bear in a reasonable way."

¹⁰⁶ See *High Stakes*, Chapters 5, 6, 7; *Joint Standards*, Standard 13.7. Standard 13.7, *supra* note 56.

Test Measurement Principles: Questions about Appropriate Test Use

In order to determine if a test is being used appropriately in making high-stakes decisions about students, considerations about the context of the test use, and the validity, reliability, and fairness of the scores and their interpretations need to be addressed. In all cases, it is important that the evidence related to the technical merits of the test be based on the current test being proposed.

1. What is the purpose for which the test is being used?
2. What information, besides the test, is being collected to inform this purpose?
3. Based on how the test results are to be used, is there adequate evidence of validity to document that the test score inferences are accurate and meaningful for the students taking the test? That is,
 - Does the evidence support that the inferences accurately reflect the specific knowledge and skills the test says it measures?
 - Does the evidence support that the inferences are valid for the stated purpose, and in the particular type of setting where the test is to be administered?
 - Does the evidence support that the inferences are valid for the specific groups of students who are taking the test?
4. Is there adequate evidence of reliability of the test scores for the proposed use?
5. Is there adequate evidence of fairness in validity and reliability to document that the test score inferences are accurate and meaningful for all students taking the test? That is,
 - Does the evidence support that the inferences are measuring the same constructs for all students?
 - Does the evidence support that the scores do not systematically underestimate or overestimate the knowledge or skills of members of a particular group?
 - Does the evidence demonstrate validity and reliability of the score inferences for each relevant subgroup when a prior probability exists that, across examinee subgroups, test scores may differ in meaning or that the reliability of the scores may vary substantially?
6. Is there adequate evidence that cutscores have been properly established and that they will be used in ways that will provide accurate and meaningful information for all test takers?

II. Accuracy in Testing Limited English Proficient Students and Students with Disabilities

All aspects of validity, reliability, fairness, and cutscores discussed above are applicable to the measurement of knowledge and skills of all students, including limited English proficient students¹⁰⁷ and students with disabilities. This section addresses additional issues related to accurately measuring the knowledge and skills of these two student populations.

Ensuring that test score inferences accurately reflect the intended constructs for all students is a complex task. It involves several aspects of test construction, pilot testing, implementation, analysis, and reporting. The appropriate inclusion of students from these populations in validation and norming samples, and the meaningful inclusion of limited English proficient experts and disability experts throughout the test development process, helps ensure suitable test quality and use for all test takers.

The proper inclusion of all students in testing programs helps to ensure that high-stakes decisions are made on the basis of tests results that are as comparable as possible across all test takers, rather than on the basis of results from assessments that are developed to measure different content domains.¹⁰⁸ The appropriate inclusion of all students can also help to ensure that educational benefits attributable to the high-stakes decisions will be available to all. In some cases, it is appropriate to test limited English proficient students and students with disabilities under standardized conditions, as long as the evidence supports the validity of the scores in a given situation for these students. In other cases, the conditions may have to be accommodated to assure that the scores validly reflect the students' mastery of the intended constructs.¹⁰⁹ The use of multiple measures generally enhances the accuracy of the educational decisions, and these measures can be used to confirm the validity of the test results.

A. General Considerations about Accommodations

Making similar inferences about academic test scores for all test takers, and making appropriate decisions when using these scores, requires measuring the same academic constructs (knowledge and skills in specific subject areas) across groups and contexts. In measuring the knowledge and skills of limited English proficient students and students with disabilities, it is particularly important that the tests actually measure the intended knowledge and skills and not other factors which are extraneous to the intended

¹⁰⁷ These are students who are learning English as a second language. Other documents sometimes refer to these students as English language learners.

¹⁰⁸ *High Stakes*, p. 7, 80.

¹⁰⁹ See *Joint Standards*, Chapter 7, Fairness in Testing and Test Use; Chapter 9, Testing Individuals of Differing Linguistic Backgrounds; Chapter 10, Testing Individuals with Disabilities.

construct.¹¹⁰ For instance, impaired visual capacity may influence a student's test score in science when the student must sight read a typical paper and pencil science test. In measuring science skills, the student's sight is likely not relevant to her knowledge of science. Similarly, how well a limited English proficient student reads English may influence the student's test score in mathematics when the student must read the test. In this case, the student's reading skills are not relevant when the skills of mathematics computation are to be measured.

Typically, accommodations to established conditions are found in three main phases of testing: 1) the administration of tests, 2) how students are allowed to respond to the items, and 3) the presentation of the tests (how the items are presented to the students on the test instrument). Administration accommodations involve setting and timing, and can include extended time to counteract the

Standard 10.1

In testing individuals with disabilities, test developers, test administrators, and test users should take steps to ensure that the test score inferences accurately reflect the intended construct rather than any disabilities and their associated characteristics extraneous to the intent of the measurement.

increased literacy demands or fatigue for a student with learning or physical disabilities. Response accommodations allow students to demonstrate what they know in different ways. Presentation accommodations can include format variations such as fewer items per page, and plain language editing procedures, which use short sentences, common words, and active voice. There is a wide variation in which accommodations are used across states and school districts. (Appendix C lists many of the accommodations used in large scale testing for limited English proficient students and students with disabilities.)

Issues regarding the use of accommodations are complex. When the possible use of an accommodation for a student is being considered, two questions should be examined: 1) What is being measured if conditions are accommodated? 2) What is being measured if the conditions remain the same? The decision to use an accommodation or not should be grounded in the ultimate goal of collecting test information that accurately and fairly represents the knowledge and skills of the student on the intended constructs. The overarching concern should be that test score inferences accurately reflect the intended constructs rather than factors extraneous to the intent of the measurement.¹¹¹

¹¹⁰ This is known as construct irrelevance. See p. 25 above; *Joint Standards*, p. 173-174.

¹¹¹ Standards 9.1, 10.1, Messick, 1989. Standard 9.1 states, "Testing practice should be designed to reduce threats to the reliability and validity of test score inferences that may arise from language differences."

Standard 10.1 states, "In testing individuals with disabilities, test developers, test administrators, and test users should take steps to ensure that the test score inferences accurately reflect the intended construct rather than any disabilities and their associated characteristics extraneous to the intent of the measurement."

Messick (1989), *supra* note 68.

B. Limited English Proficient Students

The *Joint Standards* and several recent measurement publications discuss the population of limited English proficient students and how test publishers and users have handled inclusion in tests to date.¹¹² This section briefly outlines principles derived from the *Joint Standards* and these publications. It addresses two types of testing situations especially relevant for limited English proficient students: the assessment of English language proficiency and the assessment of academic educational achievement.

Interpretation of the scores of limited English proficient students should accurately and fairly reflect the academic knowledge, skills, or abilities that the test intends to measure, minimizing the effect of factors irrelevant to the intended constructs.¹¹³ When credible research evidence reports that scores may differ in meaning across subgroups of linguistically diverse test takers, then, to the extent feasible, the same form of validity evidence should be collected for each subgroup as for the examinee population as a whole.¹¹⁴

Standard 9.1

Testing practice should be designed to reduce threats to the reliability and validity of test score inferences that may arise from language differences.

“When a test is recommended for use with linguistically diverse test takers, test developers and publishers should provide the information necessary for appropriate test use and interpretation;”¹¹⁵ recommended accommodations should be used appropriately and described in detail in the test manual;¹¹⁶ translation methods and interpreter expertise should be clearly described;¹¹⁷ and evidence of the reliability and validity of the

¹¹² For instance, *Joint Standards*, Chapter 9; *High Stakes*, Chapter 9; *Improving Schooling for Language Minority Children: A Research Agenda* (National Research Council, August and Hakuta, 1997); *Ensuring Accuracy in Testing for English Language Learners* (Kopriva, 2000, Washington D.C. Council of Chief State School Officers).

¹¹³ See Standard 9.1, *supra* note 111.

¹¹⁴ Standard 9.2 states, “When credible research evidence reports that test scores differ in meaning across subgroups of linguistically diverse test takers, then to the extent feasible, test developers should collect for each linguistic subgroup studied the same form of validity evidence collected for the examinee population as a whole.”

¹¹⁵ Standard 9.6

Standard 9.5 states, “When there is credible evidence of score comparability across regular and modified tests or administrations, no flag should be attached to a score. When such evidence is lacking, specific information about the nature of the modification should be provided, if permitted by law, to assist test users properly to interpret and act on test scores.”

¹¹⁶ Standard 9.4 states, “Linguistic modifications recommended by test publishers, as well as the rationale for the modifications, should be described in detail in the test manual.”

¹¹⁷ Standards 9.7, 9.11. Standard 9.7 states, “When a test is translated from one language to another, the methods used in establishing the adequacy of the translation should be described, and empirical and logical evidence should be provided for score reliability and the validity of the translated test’s score inferences for the uses intended in the linguistic groups to be tested.”

Standard 9.11 states, “When an interpretation is used in testing, the interpreter should be fluent in both the language of the test and the examinee’s native language, should have expertise in translating, and should have a basic understanding of the assessment process.”

translated test score's inferences should be collected and made available in order to support sound test use by educators and policy makers.¹¹⁸

1. Assessing English Language Proficiency

Issues of validity, reliability, and fairness apply to tests and other relevant assessments that measure English language proficiency. English language proficiency is typically defined as proficiency in reading, writing, speaking, and understanding English.¹¹⁹ Assessments that measure English language proficiency are generally used to make decisions about

who should receive English language acquisition services, the type of programs in which these students are placed, and the progress of students in the appropriate programs. They are also used to evaluate the English proficiency of students when exiting from services, to ensure that they can successfully participate in the regular school curriculum. In making decisions about which tests are appropriate, it is particularly important to make sure that the tests accurately and completely reflect the intended English language proficiency constructs so that the students are not misclassified. It is generally accepted that an evaluation of a range of communicative abilities will typically need to be assessed when placement decisions are being made.¹²⁰

Standard 9.10

Inferences about test takers' general language proficiency should be based on tests that measure a range of language features, and not on a single linguistic skill.

¹¹⁸ Standard 9.7, *supra* note 117.

¹¹⁹ *Improving Schooling for Language Minority Children*, p. 116-118.

¹²⁰ Comment under Standard 9.10, p. 99-100. Standard 9.10 states, "Inferences about test takers' general language proficiency should be based on tests that measure a range of language features, and not on a single linguistic skill."

2. Testing the Academic Educational Achievement
Of Limited English Proficient Students

Several factors typically affect how well the educational achievement of limited English proficient students is measured on standardized academic tests. For all test takers, any test that employs written or oral skills in English or in another language is, in part, a measure of those skills in the particular language. Test use with individuals who have not sufficiently acquired the literacy or linguistic skills in the language of the test may introduce construct-irrelevant components to the testing process. In such instances, test results may not reflect accurately the qualities and competencies intended to be measured.¹²¹ While it is very important that the test score inferences are valid, reliable, and fair, the technical issues associated with developing meaningful achievement tests for this population are complex and difficult to accomplish. Tests must be developed so that they effectively measure the students' knowledge and skills in intended academic achievement constructs rather than factors irrelevant to those constructs, i.e. literacy skills when literacy is not what is being measured. This is particularly important when tests are used to make high stakes decisions for individual students. Reducing the influence of construct irrelevant factors includes minimizing the confounding conditions in the test or the testing process so that the students can access the test requirements.¹²² It also includes providing native language tests where possible, when this approach would yield more accurate results for limited English proficient students.¹²³ In collecting evidence to support the technical quality of a test for these students, the accumulation of data may need to occur over several test administrations to ensure robust sample sizes.

a. **Background Factors for Limited English Proficient Students**

The background factors particularly salient in ensuring accuracy in testing for students with limited English proficiency tend to relate to literacy, culture, and schooling.¹²⁴

Limited English proficient students often bring varying levels of English and home language literacy skills to the testing situation.¹²⁵ These students may be adept in conversing orally in their home language, but unless they have had formal schooling in their home language, they may not have a corresponding level of literacy. Also, while students with limited English proficiency may acquire a degree of oral proficiency in English, literacy in English for many students comes later.¹²⁶ To add to the complexity,

¹²¹ See *Joint Standards*, p. 91.

¹²² See Standard 9.1, *supra* note 111.

¹²³ Standards 9.3 states "When testing an examinee proficient in two or more languages for which the test is available, the examinee's relative language proficiency should be determined. The test generally should be administered in the test taker's most proficient language, unless proficiency in the less proficient language is part of the assessment.

¹²⁴ *Improving Schooling for Language Minority Children*, Chapter 5; *Ensuring Accuracy in Testing for English Language Learners*, Chapter 1.

¹²⁵ See *Joint Standards*, Chapter 9, p. 91-100; *Ensuring Accuracy in Testing for English Language Learners*, Chapter 1.

¹²⁶ *Testing, Teaching and Learning*, p. 61.

oral and literacy proficiency in either the home language or English involves both social and academic components. Thus, a student may be able to write a well-organized social letter in his or her home language, and may not be able to orally explain adequately in that language how to solve a mathematics problem that includes the knowledge of concepts and words endemic to the field of mathematics. The same phenomena may occur in English as well.¹²⁷

Factors Related to Accurately Testing LEP Students

Literacy Issues

- The student's level of oral and written proficiency in English
- The student's literacy in his or her home language
- The language of instruction

Cultural Issues

- Background experiences
- Perceptions of prior experiences
- Value systems

Schooling Issues

- The amount of formal schooling in the student's home country and in U.S. schools
- Consistency of schooling
- Instructional practices in the classroom

Therefore, in determining how to effectively measure the academic knowledge and skills of this population, educators and policymakers should consider how to minimize the influence of literacy issues, except when these constructs are explicitly being measured. Considering the level of linguistic and literacy proficiencies of limited English proficient students in their home language and in English will often affect which achievement tests are appropriate for these students, and which accommodations to standardized testing conditions, if any, might be most useful for which students.¹²⁸

Additionally, diverse cultural and other background experiences, including variations in amount, type and location (home country and U.S.) of formal schooling, as well as interrupted and multi-location schooling (of the type frequently experienced by children of migrant workers), affect language literacy, the contextual content of items, and the academic foundational knowledge base that can be assumed in educational achievement tests. The format and procedures involved in testing can also affect accuracy in test scores, particularly if the test practices differ substantially from ongoing instructional practices in classrooms.¹²⁹

¹²⁷ *Improving Schooling for Language Minority Children*, Chapter 5, p. 113-137.

¹²⁸ *Id.* at Chapter 5.

¹²⁹ *Ensuring Accuracy in Testing for English Language Learners*, Chapters 3,4, 7, and 9.

b. Accommodations for Limited English Proficient Students

Providing accommodations to established testing conditions for some students with limited English proficiency may be appropriate when their use would yield the most valid scores on the intended academic achievement constructs. Deciding which accommodations to use for which students usually involves an understanding of which construct irrelevant background factors would substantially influence the measurement of intended knowledge and skills for individual students, and how the accommodations would impact the validity of the test score interpretations for these students.¹³⁰ Appendix C lists various test presentation, administration, and response accommodations that states and districts generally employ when testing limited English proficient students. Examples of accommodations in the presentation of the test include editing text so the items are in plain language, or providing page formats which minimize confusion by limiting use of columns and the number of items per page. Presenting the test in the student's native language is an accommodation to a test written in English when the same constructs are being measured on both the English and native language versions. Administration accommodations include extending the length of the testing period, permitting breaks, administering tests in small groups or in separate rooms, and allowing English or native language glossaries or dictionaries as appropriate. Response accommodations include oral response and permitting students to respond in their native language.

C. Students with Disabilities

The *Joint Standards* and several recent measurement publications discuss the population of students with disabilities and how test publishers and users have handled inclusion in tests to date.¹³¹ This section briefly outlines principles derived from the *Joint Standards* and these publications. It addresses three types of testing situations especially relevant for students with disabilities: tests used for diagnostic and intervention purposes, the assessment of academic educational achievement, and alternate assessments for K-12 students with disabilities who cannot participate in school-wide tests.

The *Joint Standards* provide that interpretation of the scores of students with disabilities should accurately and fairly reflect the academic knowledge, skills, or abilities that the test intends to measure. The interpretation should not be confounded by the challenges of the students that are extraneous to the intent of the measurement.¹³² Rather, validity

¹³⁰ See *Ensuring Accuracy in Testing for English Language Learners*, Chapters 6 and 8, for a discussion of which accommodations might be most beneficial for students with various background factors.

¹³¹ For instance, *Joint Standards*, Chapter 10; *High Stakes*, Chapter 8; *Educating One and All: Students with Disabilities and Standards-Based Reform* (National Research Council, McDonnell, McLaughlin, and Morison, 1997); *Testing Students with Disabilities* (Thurlow, Elliot, and Ysseldyke, 1998, NY: Corwin Press).

¹³² Standards, 10.1, 10.10. See Standard 10.1, *supra* note 111. Standard 10.10 states, "Any test modifications adopted should be appropriate for the individual test taker, while maintaining all feasible standardized features. A test

evidence should document that the inferences of the scores of students with disabilities are accurate. Pilot testing and other technical investigations should be conducted where feasible to ensure the validity of the test inferences when accommodations have been allowed.¹³³ Feasibility is always a consideration, although the *Joint Standards* comment, “[T]he costs of obtaining validity evidence should be considered in light of the consequences of not having usable information regarding the meanings of scores for people with disabilities”.¹³⁴

1. Tests used for Diagnostic and Intervention Purposes

All issues of validity, reliability, and fairness apply to tests and other assessments used to make diagnostic and intervention decisions for students with disabilities. Tests that yield diagnostic information typically focus in great detail on identifying the specific

Standard 10.12

In testing individuals with disabilities for diagnostic and intervention purposes, the test should not be used as the sole indicator of the test taker’s functioning. Instead, multiple sources of information should be used.

professional needs to consider reasonably available information and individual capabilities that might impact test performance, and document the grounds for the modification.”

¹³³ Several standards discuss the appropriate types of validity evidence, including Standards 10.3, 10.5, 10.6, 10.7, 10.8, and 10.11. Because of the low incidence nature of several of the disability groups, especially when different severity levels and combinations of impairments are considered, this type of evidence will probably need to be accumulated over time in order to have a large enough sample size.

Standard 10.3 states, “Where feasible, tests that have been modified for use with individuals with disabilities should be pilot tested on individuals who have similar disabilities to investigate the appropriateness and feasibility of the modifications.”

Standard 10.5 states, “Technical material and manuals that accompany modified tests should include a careful statement of the steps taken to modify the test to alert users to changes that are likely to alter the validity of inferences drawn from the test scores.”

Standard 10.6 states, “If a test developer recommends specific time limits for people with disabilities, empirical procedures should be used, whenever possible, to establish time limits for modified forms of timed tests rather than simply allowing test takers with disabilities a multiple of the standard time. When possible, fatigue should be investigated as a potentially important factor when time limits are extended.”

Standard 10.7 states, “When sample sizes permit, the validity of inferences made from test scores and the reliability of scores on tests administered to individuals with various disabilities should be investigated and reported by the agency or publisher that makes the modification. Such investigations should examine the effects of modifications made for people with various disabilities on resulting scores, as well as the effects of administering standard unmodified tests to them.”

Standard 10.8 states, “Those responsible for decisions about test use with potential test takers who may need or may request specific accommodations should (a) possess the information necessary to make an appropriate selection of measures, (b) have current information regarding the availability of modified forms of the test in question, (c) inform individuals, when appropriate, about the existence of modified forms, and (d) make these forms available to test takers when appropriate and feasible.”

Standard 10.11 states, “When there is credible evidence of score comparability across regular and modified administrations, no flag should be attached to a score. When such evidence is lacking, specific information about the nature of the modification should be provided, if permitted by law, to assist test users properly to interpret and act on test scores.”

¹³⁴ Comment under Standard 10.7, pg. 106.

challenges and strengths of a student.¹³⁵ These diagnostic tests are often administered in one-to-one situations (test taker and examiner) rather than in a group situation. In many cases they have been designed with standardized adaptations to fit the needs of individual examinees. In making decisions about which tests are appropriate to use, it is important to make sure that the tests accurately and completely reflect the intended constructs, so that the interventions are appropriate and beneficial for the individual students.

2. Testing the Academic Educational Achievement
Of Students with Disabilities

Several factors affect how well the educational achievement of students with disabilities is measured on standardized academic tests. While it is very important that the test score inferences are valid, reliable, and fair, the technical issues associated with developing meaningful achievement tests for this population are complex and difficult to accomplish. To ensure accuracy in testing of students with disabilities, tests must be developed so that they effectively measure the students' knowledge and skills in academic achievement rather than factors irrelevant to the intended constructs of the test. This is particularly important when achievement tests are used to make high-stakes decisions for individual students with disabilities. Reducing the influence of construct irrelevant factors includes minimizing the confounding conditions in the test or the testing process so that the test accurately measures what it is supposed to measure.¹³⁶ In collecting evidence to support the technical quality of the test for these students, the accumulation of data may need to occur over several test administrations to ensure robust sample sizes.

a. **Background Factors for Students with Disabilities**

The background factors particularly important to students with disabilities are generally related to the nature of the disabilities or to the schooling experiences of these students.¹³⁷

¹³⁵ *Joint Standards*, Chapters 10, 12, and 13; *High Stakes*, Chapter 1.

¹³⁶ See Standard 10.1, *supra* note 111.

¹³⁷ *Educating One and All*, Chapter 3; *Testing Individuals with Disabilities*.

Factors Related to Accurately Testing Students with Disabilities

Disability Issues

- Types of impairments
- Severity of impairments

Schooling Experiences

- Overlap of individualized educational goals and general education curricula
- Pace of schooling
- Instructional practices

Within any disability category, the type, number, and severity of impairments vary greatly.¹³⁸ For instance, some students with learning disabilities have a processing disability in only one subject, such as mathematics, while others experience accessing, retrieval, and processing impairments that affect a broad number of school subjects and contexts. For many of these students, one or more of the impairments may be relatively mild, while for others one or more can be significant. Further, different types of disabilities yield significantly different constellations of issues. For instance, the considerations surrounding hearing impaired students overlap significantly with limited English proficient students in some ways and with other students with disabilities in other respects. This complexity poses a challenge not only to educators, but also to test administrators and developers. In general, in determining how to use academic tests appropriately for students with disabilities, educators and policymakers should consider how to minimize the influence of the impairments in measuring the intended constructs.

¹³⁸ *Joint Standards*, Chapter 10, *Testing Individuals with Disabilities*, p. 101-105.

Educating One and All explains that the schooling experiences of students with disabilities vary greatly as a function of their disability, the severity of impairments, and expectations of their capabilities.¹³⁹ Two sets of educational experiences, in particular, affect how educators and policy makers accommodate tests and use them appropriately for this population. First, guidance about the schooling and evaluation of students with disabilities is provided by individualized education program (IEP) teams made up of educators and parents. These teams often recommend testing accommodations that they feel would be appropriate for individual students. Second, classroom instructional techniques affect large scale testing. While special educators have a long history of accommodating instruction to fit student strengths, not all the instructional practices are appropriate in large scale testing. Additionally, some students may not have been exposed routinely to the types of accommodations that would be possible in large scale testing.¹⁴⁰

b. Accommodations for Students with Disabilities

Providing accommodations to established testing conditions for some students with disabilities may be appropriate when their use would yield the most valid scores on the intended academic achievement constructs. Deciding which accommodations to use for which students usually involves an understanding of which construct irrelevant background factors would substantially influence the measurement of intended knowledge and skills for individual students, and how the accommodations would impact the validity of the test score interpretations for these students.¹⁴¹ Appendix C lists various presentation, administration, and response accommodations that states and districts generally employ when testing students with disabilities. Examples of presentation accommodations are the use of Braille, large print, oral reading, or providing page formats which minimize confusion by limiting use of columns and the number of items per page. Administration accommodations in setting include allowing students to take the test at home or in a small group, and accommodations in timing include extended time and frequent breaks. Variations in response format include allowing students to respond orally, point or use a computer.

3. Alternate Assessments

Alternate assessments are assessments for those students with disabilities who cannot participate in state or district-wide standardized assessments, even with the use of appropriate accommodations and modifications.¹⁴² For the constructs being measured, the considerations with respect to validity, reliability, and fairness apply to alternate assessments, as well. Appropriate content needs to be identified, and procedures designed to ensure technical rigor

¹³⁹ See *Educating One and All*, Chapter 3.

¹⁴⁰ See *Educating One and All*, Chapter 5.

¹⁴¹ See *Testing Students with Disabilities* for a discussion of which accommodations might be most beneficial for students with various impairments and other background factors.

¹⁴² The IDEA requires use of alternate assessments in certain areas. See 34 C.F.R. 300.138.

need to be followed.¹⁴³ In addition, strong evidence should show that the test measures the knowledge and skills it intends to measure, and that the measurement is a valid reflection of mastery in a range of contextual situations.

¹⁴³ See *Educating One and All*, Chapter 5, and *Testing Students with Disabilities* for a discussion of the issues and processes involved in developing and implementing alternate assessments.

CHAPTER 2. Legal Principles

It is important for educators and policy makers to understand the test measurement principles and the legal principles that will enable them to ask informed questions and make sound decisions regarding the use of tests for high-stakes purposes. The goal of this chapter is to explain the legal principles that apply to educational testing.

The primary focus of this chapter is four federal nondiscrimination laws, enacted by Congress, and their implementing regulations: Title VI of the Civil Rights Act of 1964 (Title VI), Title IX of the Education Amendments of 1972 (Title IX), Section 504 of the Rehabilitation Act of 1973 (Section 504), and Title II of the Americans with Disabilities Act of 1990 (Title II).¹⁴⁴ Within the U.S. Department of Education, the Office for Civil Rights has responsibility for enforcing the requirements of these four statutes and their implementing regulations. Although the Office for Civil Rights does not enforce federal constitutional provisions, an overview of these constitutional principles, including under the Fifth and Fourteenth Amendments of the U.S. Constitution, has also been included for informational purposes. The discussion of legal principles in this chapter is intended to reflect existing legal principles and does not establish new requirements.¹⁴⁵

Some of the issues that have been considered by federal courts in assessing the legality of specific testing practices for making high-stakes decisions include:

- The use of an educational test for a purpose for which the test was not designed or validated¹⁴⁶
- The use of a test score as the sole criterion for the educational decision¹⁴⁷
- The nature and quality of the opportunity provided to students to master required content, including whether classroom instruction includes the material covered by a test administered to determine student achievement¹⁴⁸
- The significance of any fairness problems identified, including differential predictive validity and possible cultural biases in the test or in test items¹⁴⁹
- The educational basis for establishing passing or cut-off scores¹⁵⁰

¹⁴⁴ Title VI prohibits discrimination on the basis of race, color and national origin in the programs and activities of recipients that receive federal financial assistance. The U.S. Department of Education's regulation implementing Title VI is found at 34 C.F.R. Part 100. Title IX prohibits discrimination on the basis of sex in educational programs and activities of recipients of federal financial assistance. The U.S. Department of Education's regulation implementing Title IX is found at 34 C.F.R. Part 106. Section 504 prohibits discrimination on the basis of disability in the programs and activities of recipients of federal financial assistance. The U.S. Department of Education's regulation implementing Section 504 is found at 34 C.F.R. Part 104. Title II prohibits discrimination on the basis of disability by public entities, regardless of whether they receive federal funding. The U.S. Department of Justice's regulation implementing Title II is found at 28 C.F.R. Part 35.

¹⁴⁵ Consistent with this approach, court decisions are not cited if the case is still on appeal or the time to request an appeal has not ended.

¹⁴⁶ See *Sharif v. New York State Educ. Dep't.*, 709 F. Supp. 345, 354-355, 364 (S.D. N.Y. 1989) (in granting a motion for preliminary injunction, where girls received comparatively lower scores than boys, court found that the state's use of SAT scores as the sole basis for decisions awarding college scholarships intended to reward high school achievement was not educationally justified for this purpose in that the SAT had been designed as an aptitude test to predict college success and was not designed or validated to measure past high school achievement).

I. Discrimination Under Federal Statutes and Regulations

Congress has enacted four statutes prohibiting discrimination based on race, color, national origin, sex, and disability in schools, colleges, and universities. Title VI prohibits discrimination based on race, color, or national origin; Title IX prohibits discrimination based on sex; and Section 504 and Title II of the Americans with Disabilities Act (ADA) prohibit discrimination based on disability. Title VI, Title IX, and Section 504 apply to all educational institutions that receive federal funds. Title II of the ADA applies to public entities, including public school districts and state colleges and universities.¹⁵¹ The Title VI, Title IX, Section 504, and Title II statutes and their implementing regulations as well as the equal protection clause of the Fourteenth Amendment to the United States Constitution, prohibit intentional discrimination, based on race, national origin, sex, or disability. In addition, the regulations that implement Title VI, Title IX, Section 504 and Title II prohibit policies or practices that have a

¹⁴⁷ See *United States v. Fordice*, 505 U.S. 717, 733-738 (1992) (invalidating state's exclusive reliance on ACT scores as a basis for college admissions at historically segregated colleges where the state adopted the ACT for discriminatory reasons and the ACT administering organization recommended that college admissions decisions consider high school grades along with test scores); see also *Sharif*, 709 F. Supp. at 364.

¹⁴⁸ See *Lau v. Nichols*, 414 U.S. at 566-569 (finding a violation of the Title VI regulations where limited English proficient students were taught only in English and not provided any special assistance needed to meet English language proficiency standards required by the state for a high school diploma). See also *Debra P.*, 644 F.2d at 406-408 (holding that use of a graduation test that covered material that had not been taught in class would violate the due process and equal protection clauses and that, under the circumstances of the case, immediate use of the diploma sanction for test failure would punish black students for deficiencies created by an illegally segregated school system which had provided them with inferior physical structures, course offerings, instructional materials, and equipment).

¹⁴⁹ See *Larry P. v. Riles*, 793 F.2d at 980-981, 983 (finding that IQ tests the state used had not been validated for use as the sole means for determining that black children should be placed in classes for educable mentally retarded students); *Sharif*, 709 F. Supp. at 354 (observing that the SAT under-predicts success for female college freshmen as compared with males). See also *Parents in Action on Special Educ. v. Hannon*, 506 F. Supp. 831, 836-837 (N.D. Ill. 1980) (court's analysis of items on I.Q. test found only minimal amount of cultural bias not resulting in erroneous mental retardation diagnoses given other information considered in process).

¹⁵⁰ See *Groves v. Alabama State Bd. of Educ.*, 776 F. Supp. 1518, 1530-1531 (M.D. Ala. 1991) (finding test required for admission to undergraduate teacher training program would not be educationally justified if the passing score is not itself a valid measure of the minimal ability necessary to become a teacher); *Richardson v. Lamar County Bd. of Educ.*, 729 F. Supp. 806, 823-825 (M.D. Ala. 1989) (evidence revealed that cut off scores had not been set through a well-conceived, systematic process nor could the scores be characterized as reflecting the good faith exercise of professional judgment), *aff'd sub nom.*, *Richardson v. Alabama State Bd. of Educ.*, 935 F.2d 1240 (11th Cir. 1991).

¹⁵¹ OCR enforces five nondiscrimination statutes, Title VI of the Civil Rights Act of 1964, 42 U.S.C. §§ 2000d, *et seq.* (2000); Title IX of the Education Amendments of 1972, 20 U.S.C. §§ 1681 *et seq.* (1999); Section 504 of the Rehabilitation Act of 1973, as amended, 29 U.S.C. §§ 794 (1999); Title II of the Americans with Disabilities Act of 1990, 42 U.S.C. §§, 12131, *et seq.* (1995 and Supp. 1999); and the Age Discrimination Act of 1975, *as amended*, 42 U.S.C. §§ 6101, *et seq.* (1995 and Supp. 1999). Regulations issued by the United States Department of Education implementing Title VI, Title IX, and Section 504, respectively, can be found at 34 C.F.R. Part 100, 34 C.F.R. Part 106, and 34 C.F.R. Part 104. These regulations can be found on OCR's web-site at www.ed.gov/offices/OCR. For regulations implementing Title II of the ADA, see 28 C.F.R. Part 35. Title III of the ADA, which is enforced by the U.S. Department of Justice, prohibits discrimination in public accommodations by private entities, including schools. Religious entities operated by religious organizations are exempt from Title III.

discriminatory disparate impact on students based on their race, national origin, sex, or disability.¹⁵²

This section describes two central analytical frameworks for examining allegations of discrimination as set forth in federal nondiscrimination regulations: different treatment and disparate impact.¹⁵³ It also includes a further discussion of legal principles that apply specifically to students with limited English proficiency and to students with disabilities.

A. Different Treatment

Under federal law, policies and practices generally must be applied consistently to similarly situated individuals or groups, regardless of their race, national origin, sex, or disability.¹⁵⁴ For example, a federal court concluded that a school district had intentionally treated students differently on the basis of race where minority students whose test scores qualified them for two or more ability levels were more likely to be assigned to the lower level class than similarly situated white students, and no explanatory reason was evident.¹⁵⁵

In addition, educational systems that were previously segregated by race in violation of the Fourteenth Amendment and have not achieved unitary status have an obligation to dismantle their prior *de jure* segregation. In such instances, when a school district or other educational system uses a test or assessment procedure for a high-stakes purpose that has racially disproportionate effects, the school district or other educational system must show that the disparity is not traceable to prior intentional segregation or that the test or assessment procedure does not perpetuate the adverse effects of such

¹⁵² 34 C.F.R. § 100.3(b)(2); 34 C.F.R. §§ 106.21(b)(2), 106.36(b), 106.52; 34 C.F.R. § 104.4(b)(4)(i); and 28 C.F.R. § 35.130(b)(3).

The authority of federal agencies to issue regulations with an "effects" standard has been consistently acknowledged by U.S. Supreme Court decisions and applied by lower federal courts addressing claims of discrimination in education. See, e.g., *Lau v. Nichols*, 414 U.S. 563, 568 (1974); *Guardians Ass'n v. City Service Comm'n. of City of N.Y.*, 463 U.S. 582, 584-593 (1983); *Alexander v. Choate*, 469 U.S. 287, 289-300 (1985). See also Memorandum from the Attorney General for Heads of Departments and Agencies that Provide Federal Financial Assistance, "Use of the Disparate Impact Standard in Administrative Regulations under Title VI of the Civil Rights Act of 1964," July 14, 1994.

¹⁵³ Intentional racial discrimination is a violation of both the Fourteenth Amendment to the United States Constitution and federal civil rights statutes in cases where evidence demonstrates that an action such as the use of a test for high-stakes purposes is motivated by an intent to discriminate. See *Elston v. Talladega County Bd. of Educ.*, 997 F.2d 1394, 1406 (11th Cir. 1993). As explained further in this section, the regulations promulgated under the federal civil rights statutes prohibit the use of neutral criteria having disparate effects unless the criteria are educationally justified. See *Guardians Ass'n v. Civil Service Comm'n.*, 463 U.S. at 598.

¹⁵⁴ For example, under the Fourteenth Amendment and Title VI, different treatment based on race is permitted only when such action is narrowly tailored to further a compelling state interest. See *Regents of the Univ. of Cal. v. Bakke*, 438 U.S. 265 (1978); *Adarand Constructors, Inc. v. Peña*, 515 U.S. 200 (1995).

¹⁵⁵ See *People Who Care v. Rockford Bd. of Educ.*, 851 F. Supp. 905, 958-1001 (N.D. Ill. 1994), *remedial order rev'd, in part*, 111 F.3d 528 (7th Cir. 1997). On appeal, the Seventh Circuit Court of Appeals stated that the appropriate remedy in this case was to require the district to use objective, non-racial criteria to assign students to classes, rather than abolishing the district's tracking system. 111 F.3d at 536.

segregation.¹⁵⁶ The school district is under "a 'heavy burden' of showing that actions that increase [] or continue [] the effects of the dual system serve important and legitimate ends."¹⁵⁷

B. Disparate Impact

Discrimination under federal law may also occur where the application of neutral criteria has discriminatory effects and those criteria are not educationally justified. The federal nondiscrimination regulations provide that a recipient of federal funds may not "utilize criteria or methods of administration which have the effect of subjecting individuals to discrimination."¹⁵⁸ It is important to understand that disparities in student performance based on race, national origin, sex, or disability, alone, do not constitute disparate impact discrimination under federal law. Furthermore, nothing in federal law guarantees equal results. (For a further discussion of issues related to testing of students with disabilities, see pp. 56 - 60.)

Courts applying the disparate impact test have examined three questions to determine if the practices at issue are discriminatory: (1) Does the practice or procedure in question result in substantial differences in the award of benefits or services based on race, national origin, or sex? (2) Is the practice or procedure educationally justified? and (3) Is there an equally effective alternative that can accomplish the institution's educational goal with less disparity?¹⁵⁹

¹⁵⁶ See *United States v. Fordice*, 505 U.S. at 731-732 (finding state's requirement that students have higher ACT scores for admission to historically white colleges than historically black colleges to be constitutionally suspect where the requirement was enacted for discriminatory purposes, emanated from the prior *de jure* system that continue to have segregative effects and was not shown to be justified in educational terms); *Debra P. v. Turlington*, 644 F.2d at 407 ("[defendants] failed to demonstrate either that the disproportionate failure [rate] of blacks was not due to the present effects of past intentional segregation or, that as presently used, the diploma sanction was necessary [in order] to remedy those effects"); *McNeal v. Tate County Sch. Dist.*, 508 F.2d 1017, 1020-1021 (5th Cir. 1975) (since ability grouped classroom assignments preserved effects of past intentional discrimination, defendants were required to show educational benefits of assignment practice on remand or propose an educationally sound alternative); *GI Forum v. Texas Educ. Agency, No. SA-97-CA-1278-EP*, 2000 U.S. Dist. LEXIS 153, slip op. at 56-57 (W.D. Tex. 2000) (upholding use of graduation test where the test is used to identify educational inequalities and attempt to address them).

¹⁵⁷ *Dayton Bd. of Educ. v. Brinkman*, 443 U.S. 526, 538 (1979) (quoting *Green v. County School Board*, 391 U.S. 430, 439 (1968)).

¹⁵⁸ See 34 C.F.R. § 100.3(b)(2) (Title VI); 34 C.F.R. § 104.4(b)(4)(i) (Section 504); and 28 C.F.R. § 35.130(b)(3)(i) (Title II). See also 34 C.F.R. § 106.31 (Title IX). In *Guardians*, 463 U.S. at 589-590, the U.S. Supreme Court upheld the use of the effects test, stating that the Title VI regulation forbids the use of federal funds, "not only in programs that intentionally discriminate on racial grounds but also in those endeavors that have a[n] [unjustified racially disproportionate] impact on racial minorities."

¹⁵⁹ See *Georgia State Conf.*, 775 F.2d at 1417. See also *Elston*, 997 F.2d at 1407 & n.14; *Larry P.*, 793 F.2d at 982 & n. 9; *Groves*, 776 F. Supp. at 1523-1524, 1529-1532; *Sharif*, 709 F. Supp. at 361. Many courts use the term "equally effective" when discussing whether the alternative offered by the party challenging the test is feasible and would effectively meet the institution's goals. See, e.g., *Georgia State Conf.*, 775 F.2d at 1417; *Sharif*, 709 F. Supp. at 361. Other courts use the term "comparably effective" in evaluating proposed alternatives. See, e.g., *Sandoval*, 7 F. Supp. 2d at 1278; *Elston*, 997 F.2d at 1407; *Fitzpatrick v. City of Atlanta*, 2 F.3d 1112, 1118 (11th Cir. 1993). Review of the decisions in these cases indicate that the courts appear to be using the terms synonymously.

The party challenging the test has the burden of establishing disparate impact. If disparate impact is established, the educational institution must provide sufficient evidence of an educational justification. If an educational justification is established, then the party challenging the test must demonstrate that an alternative with less disparate impact is equally effective in meeting the institution's educational goals or needs in order to prevail.¹⁶⁰

1. Determining disproportionate impact

The first question in the disparate impact analysis is whether there is information indicating a significant disparity in the award of benefits or services to students based on race, national origin, or sex.¹⁶¹ To determine if a significant disparate impact exists, courts have focused on evidence of statistical disparities.¹⁶² Generally, a test has a disproportionate adverse impact if a statistical analysis shows a significant difference from the expected random distribution.¹⁶³ There is no rigid mathematical threshold regarding the degree of disproportionality required; however, courts have used various statistical methods to identify disparities

Generally, if a statistical analysis shows that the success rate for a particular group of students is significantly lower (or the failure rate is significantly higher) than what would be expected from a random distribution, then the test has disproportionate adverse impact.

National Research Council, High Stakes: Testing for Tracking, Promotion, and Graduation, 1999: 59

that are sufficiently substantial to raise an inference that the challenged practice caused the disparate results.¹⁶⁴ To establish disparate impact in the context of a selection system, the comparison must be made between those selected for the educational benefit or service and a relevant pool of applicants or test-takers.¹⁶⁵

¹⁶⁰ See *Georgia State Conf.*, 775 F.2d at 1417. See also the Department of Justice's Title VI Legal Manual at p. 2.

¹⁶¹ For a further discussion of the legal principles regarding students with disabilities under the IDEA, Section 504 and Title II of the ADA, see pp. 38-40.

¹⁶² See *Watson v. Fort Worth Bank & Trust*, 487 U.S. 977, 994-997 (1988) (O'Connor, J., plurality opinion).

¹⁶³ See *Watson*, 487 U.S. at 995; *Groves*, 776 F. Supp. at 1526-1528.

¹⁶⁴ See *Watson*, 487 U.S. at 994-995; *Groves*, 776 F. Supp. at 1526-1527. A variety of methods are commonly used by courts to distinguish differences between outcomes that are statistically and practically significant from those that are random. Some have used an 80% rule whereby disparate impact is shown when the rate of selection for the less successful group is less than 80% of the rate of selection for the most successful group. Another type of statistical analysis considers the difference between the expected and observed rates in terms of standard deviations, with the difference generally expected to be more than two or three standard deviations. Another test is known as the "Shoben formula" in which the difference or Z-value in the groups' success rates must be statistically significant. *Groves*, 776 F. Supp. at 1526-1528 (discussing these methods and the cases in which they were used).

¹⁶⁵ When determining disparate impact in the context of a selection system, the comparison pool generally consists of all minimally qualified test-takers or applicants. When tests are used to determine placement or some other type of educational treatment, the comparison is between those identified by the test for the placement or educational treatment and the relevant pool of test takers. The precise composition of the comparison pool is determined on a case-by-case basis. See *Wards Cove Packing Co. v. Atonio*, 490 U.S. 642, 650-651 (1989); *Watson*, 487 U.S. at 995-997; *Groves*, 776 F. Supp. at 1525-1526.

In general, a specific policy, practice or procedure must be identified as causing the disproportionate adverse effect on the basis of race, national origin, or sex.¹⁶⁶ For example, when a particular use of a test is being challenged, the evidence should show that the test use, rather than other selection factors, accounts for the disparity.¹⁶⁷

2. Determining educational necessity

Where the use of a test results in decisions that have a disparate impact on the basis of race, national origin, or sex, the test use causing the disparity must significantly serve the legitimate educational goals of the institution.¹⁶⁸ This inquiry is usually referred to as determining the "educational necessity" of the test use *or* determining whether the test is "educationally justified."¹⁶⁹ The test need not be "essential" or "indispensable" to achieving the institution's educational goal;¹⁷⁰ rather, the educational institution must show a manifest relationship between use of the test and the institution's educational purposes.¹⁷¹

In evaluating educational necessity, both the legitimacy of the educational goal asserted by the institution and the use of the test as a valid means to advance this goal may be at issue. Courts generally allow educational institutions to define their own educational goals and focus on whether the challenged test serves the institution's articulated objectives.¹⁷²

¹⁶⁶ Elements of a decision-making process that cannot be separated for purposes of analysis may be analyzed as one selection practice. See Title VII of the Civil Rights Act of 1964, 42 U.S.C. § 2000e-2[k][1][B][i]. This is necessary because limiting the disparate impact analysis to a discrete component of a selection process would not allow for situations "where the adverse impact is caused by the interaction of two or more components of the process." See *Graffam v. Scott Paper Co.*, 870 F. Supp. 389, 395 (D. Me. 1994), *aff'd*, 60 F.3d 809 (1995).

¹⁶⁷ As noted in *Watson*, 487 U.S. at 994, courts have found it "relatively easy," when appropriate statistical proof is presented, to identify a standardized test as causing the racial, national origin, or sex related disparity at issue. See also *GI Forum v. Texas Educ. Agency*, No. SA-97-CA-1278-EP, 2000 U.S. Dist. LEXIS 153, slip op. at 35-40 (W.D. Tex. 2000) (given legally meaningful differences in the pass rates of minority and majority students, plaintiffs made a prima facie showing of disparate impact resulting from a minimum competency test).

¹⁶⁸ See *Wards Cove*, 490 U.S. at 659.

¹⁶⁹ See *Board of Educ. v. Harris*, 444 U.S. 130, 151 (1979); *Elston*, 997 F.2d at 1412.

¹⁷⁰ See *Wards Cove*, 490 U.S. at 659; *Elston*, 997 F.2d at 1412 (citing *Georgia State Conf.*, 775 F.2d at 1417-1418).

¹⁷¹ See *Georgia State Conf.*, 775 F.2d at 1418 (showing required that "achievement grouping practices bear a manifest demonstrable relationship to classroom education"); *Sharif*, 709 F. Supp. at 362 (defendants must show a manifest relationship between use of the SAT and recognition of academic achievement in high school). As explained in *Elston*, 997 F.2d at 1412, "from consulting the way in which . . . [courts] analyze the 'educational necessity' issue, it becomes clear that... [they] are essentially requiring . . . [the educational institution to] show that the challenged course of action is demonstrably necessary to meeting an important educational goal." In other words, the institution can defend the challenged practice on the grounds that it is "supported by a 'substantial legitimate justification.'" See *Elston*, 997 F.2d at 1412 (quoting *Georgia State Conf.*, 775 F.2d at 1417); see also *Georgia State Conf.*, 775 F.2d at 1417-1418; *Groves*, 776 F. Supp. at 1529-1532.

¹⁷² See, e.g., *Debra P.*, 644 F.2d at 402 (indicating that the court is not in a position to determine education policy and; state's efforts to establish minimum standards and improve educational quality are praiseworthy).

In conducting this analysis, courts have generally considered relevant evidence of validity, reliability, and fairness¹⁷³ provided by the test developer and test user to determine the acceptability of the test for the purpose used, giving appropriate deference to the expertise and experience of educators and testing professionals.¹⁷⁴ The educational justification inquiry thus generally looks at technical questions regarding the test's accuracy in relation to the nature and importance of the educational institution's goals, the educational consequences to students, and the relationship of the educational

¹⁷³ In general, courts have said that validity refers to the accuracy of conclusions drawn from test results. See *Allen v. Alabama State Bd. of Educ.*, 976 F. Supp. 1410, 1420-1421 (M.D. Ala. 1997) ("Generally, validity is defined as the degree to which a certain inference from a test is appropriate and meaningful", quoting *Richardson v. Lamar County Bd. of Educ.*, 729 F. Supp. 809, 820 (M.D. Ala. 1989), *aff'd*, 164 F.3d 1347 (1999), *injunction granted*, 2000 U.S. Dist. LEXIS 123 (2000).) See also *Richardson*, 729 F. Supp. at 820-821 ("[A] test will be valid so long as it is built to yield its intended inference and the design and execution of the test are within the bounds of professional standards accepted by the testing industry."); *Anderson*, 520 F. Supp. at 489 ("Validity in the testing field indicates whether a test measures what it is supposed to measure.").

¹⁷⁴ See, e.g., *United States v. LULAC*, 793 F.2d 636, 640, 649 (5th Cir. 1986) (pointing to substantial expert evidence in the record, including validity studies, indicating that the tests involved were valid measures of the basic skills that teachers should have). The sponsors of the newly revised *Joint Standards* advise that the *Joint Standards* are intended to provide guidance to testing professionals in making such judgments. See *Joint Standards*, Introduction, p. 4. The *Joint Standards* are discussed more fully in Chapter One of this guide.

Where the evidence indicates that the educational institution is using a test in a manner that does not lead to valid inferences, educational justification may be found lacking. See *United States v. Fordice*, 505 U.S. at 736-737 (ruling that Mississippi's exclusive use of ACT scores in making college admissions decisions was not educationally justified, since, among other factors, the ACT's administering organization discouraged this practice); *Groves*, 776 F. Supp. at 1530 (requiring minimum ACT score for admission to undergraduate teacher education programs violated the Title VI regulations since ACT scores had not been validated for this purpose); *Sharif*, 709 F. Supp. at 361-363 (in ruling on a motion for preliminary injunction, court found that the state's use of SAT scores as the sole basis for decisions awarding college scholarships intended to reward high school achievement was not educationally justified for this purpose in that the SAT had been designed as an aptitude test to predict college success and was not designed or validated to measure past high school achievement).

Psychometric or scientific evidence is not the only way that validity can be demonstrated, however. Courts can draw inferences of validity from a wide range of data points. See *Watson v. Fort Worth Bank & Trust*, 487 U.S. 977, 998 (1988) (referring to procedures used to evaluate personal qualities of candidates for managerial jobs).

institution to the student.¹⁷⁵ Where a test is used for promotion or graduation purposes, courts may also consider whether the skills tested have been taught in the program.¹⁷⁶

3. Determining whether there are equally effective alternatives that serve the institution's educational goal with less disparity

If the educational institution provides sufficient evidence that the test use in question is justified educationally, the party challenging the test has the opportunity to show that there exists an equally effective alternative practice that meets the institution's goals with less disparity.¹⁷⁷ The feasibility of an alternative, including costs and administrative burdens, is a relevant consideration.¹⁷⁸

II. Testing Of Students With Limited English Proficiency

Testing of students with limited English proficiency in the elementary and secondary education context raises a set of unique issues. To understand the obligations of states and school districts with regard to high-stakes testing of such students, it is important to understand the basic obligations of school districts and states under Title VI and related federal law that relate to language minority students who are learning English.

¹⁷⁵ See, e.g., *Georgia State Conf.*, 775 F.2d at 1417-1420; *Groves*, 776 F. Supp. at 1530-1531; *Larry P.*, 793 F.2d at 980. In the educational context, tests play a complex role that bears on evaluation of educational justification. As noted by the court in *Larry P.*,

[I]f tests can predict that a person is going to be a poor employee, the employer can legitimately deny that person a job, but if tests suggest that a young child is probably going to be a poor student, the school cannot on that basis alone deny that child the opportunity to improve and develop the academic skills necessary to success in our society.

793 F.2d at 980 (quoting *Larry P.*, 495 F. Supp. at 969). Because determining whether a test is a valid basis for classifying students and placing them in different educational programs may be even more complex and difficult than determining if a test validly predicts job performance, particular sensitivity is needed to all of the interests involved. The question may be not only whether a test provides valid information about a student's ability and achievement, but whether the educational services provided to the student as a consequence of the test serve the student's needs. Inequality in the services provided to students prior to the test, as well as in the services provided as a consequence of the test, may also be a factor considered as part of the educational justification for using a test in a particular way. See *Debra P.*, 644 F.2d at 407-408 (agreeing with the statement that Title VI would not be violated if the test were a fair test of what students were taught); *Debra P.*, 730 F.2d 1405, 1407, 1410-1411, 1416 (1984) (affirming that the extent of remedial efforts to address test failure is relevant to evaluation of test use).

¹⁷⁶ See *Debra P.*, 644 F.2d at 408.

¹⁷⁷ See *New York Urban League v. New York*, 71 F.3d 1031, 1036 (2d Cir. 1995) (stating "... the plaintiff may still prove his case by demonstrating that other less discriminatory means would serve the same objective"). See also *Albemarle Paper Co. v. Moody*, 422 U.S. 405, 425 (1975); *Richardson v. Lamar County Bd. of Educ.*, 729 F. Supp. at 815.

¹⁷⁸ See *Wards Cove*, 490 U.S. at 661 (indicating that factors such as costs or other burdens are relevant in determining whether the alternative is equally effective in serving employer's legitimate goals); *Sharif*, 709 F. Supp. at 363-364 (finding defendant's claim that proposed alternative was not feasible and excessively burdensome not persuasive since most other states used proposed alternative); *MacPherson v. University of Montevallo*, 922 F.2d 766, 773 (11th Cir. 1991) (holding that plaintiff must show that the alternative is economically feasible).

Title VI prohibits discrimination based on race, color, or national origin. On May 25, 1970, the United States Department of Health, Education, and Welfare's Office for Civil Rights issued a policy memorandum entitled "Identification of Discrimination and Denial of Services on the Basis of National Origin." The May 25th memorandum clarified the responsibility of school districts, under Title VI, to provide equal educational opportunity to national origin minority group students whose inability to speak and understand the English language excludes them from effective participation in the education program offered by the school district.¹⁷⁹ This memorandum was cited with approval by the Supreme Court in its decision in *Lau v. Nichols*, which held that the district's policy of teaching national origin minority group children only in English, without any special assistance, deprived them of the opportunity to benefit from the district's education program, including meeting the English language proficiency standards required by the state for a high school diploma.¹⁸⁰ The *Lau* case held that such policies are barred when they have the effect of denying such benefits, even though no purposeful design is present.¹⁸¹

Subsequently, *Castaneda v. Pickard*,¹⁸² relying on the language of the Equal Educational Opportunities Act (EEOA), explained the steps school districts must take to help students with limited English proficiency overcome language barriers to ensure that they can participate meaningfully in the district's educational programs.¹⁸³ The court stated that school districts have an obligation to provide services that enable students to acquire English language proficiency. A school system that chooses to temporarily emphasize English over other subjects retains an obligation to provide assistance necessary to remedy academic deficits that may have occurred in other subjects while the student was focusing on learning English.

Under the *Castaneda* standards, school districts have broad discretion in choosing a program of instruction for limited English proficient students. However, the program must be based on sound educational theory, must be adequately supported so that the program has a realistic chance of success, and must be periodically evaluated and revised, if necessary, to achieve its goals.

The disparate impact framework discussed above may also be used to examine whether tests used for high-stakes purposes result in a discriminatory impact upon students with limited English proficiency. As part of this analysis, questions may arise regarding the

¹⁷⁹ See *Identification of Discrimination and Denial of Services on the Basis of National Origin*, 35 Fed. Reg. 11595 (1970). The Department of Health, Education and Welfare was the predecessor of the U.S. Department of Education.

¹⁸⁰ See *Lau*, 414 U. S. at 566-568.

¹⁸¹ *Id.* at 568, citing, among other legal authority, the predecessor of 34 C.F.R. § 100.3 (b)(2).

¹⁸² See *Castaneda*, 648 F. 2d at 1005-1006, 1009-1012. The analytical framework in *Castaneda* which was decided under the Equal Educational Opportunities Act (EEOA), 20 U.S.C. §§ 1701 *et seq.*, has been applied to OCR's Title VI analysis. See *Williams Memorandum*, *supra* note 39. The EEOA contains standards related to limited English proficient students similar to the Title VI regulations.

¹⁸³ See *Castaneda*, 648 F.2d at 1011.

validity and reliability of the test for these students.¹⁸⁴ Depending upon the purpose of the test and the characteristics of the populations being tested, in some situations, accommodations or other forms of assessment of the same construct may be necessary. In short, the obligation is to ensure that the same constructs are being measured for all students.

There are three particularly important areas involving high-stakes testing of students with limited English proficiency: (1) tests used to determine a student's proficiency in the areas of speaking, listening, reading, or writing English for the purpose of determining whether the student should be provided with a program to enable the student to acquire English language skills (and, later, for the purpose of determining whether the student is ready to exit the program); (2) tests used to determine if the student meets the criteria for other specialized instructional programs, such as gifted and talented or vocational education programs; and (3) system-wide tests administered to determine if students have met performance standards.

Tests used to determine a student's initial and continuing need for special language programs should be appropriate in light of the district's own performance expectations and otherwise valid and reliable for the purpose used. Tests used by schools to help select students for specialized instructional programs, including programs for gifted and talented students, should not screen out limited English proficient students unless the program itself requires proficiency in English for meaningful participation.¹⁸⁵ When a state or school district adopts content and performance standards, and uses high-stakes tests to measure whether students have mastered these standards, a critical factor is whether the overall educational program provided to students with limited English proficiency is reasonably calculated to enable the students to master the knowledge and skills that all students are expected to master. When education agencies institute standards based testing, it is important for them to examine their programs for students with limited English proficiency to determine when and how these students will be provided with the instruction needed to prepare them to pass the test in question.

In addition, students with limited English proficiency may not be categorically excluded from standardized testing designed to increase accountability of educational programs for effective instruction and student performance. If these students are not included, the test data will not fairly reflect the performance of all students for whom the education agency is responsible.¹⁸⁶ Such test data can also help a district to assess the effectiveness of its content and English language acquisition programs.

¹⁸⁴ See pages 38-42 for a discussion of the psychometric principles involved in determining the reliability and validity of tests used with limited English proficient students.

¹⁸⁵ See *Williams Memorandum, supra*, note 39.

¹⁸⁶ Indeed, Title I of the Elementary and Secondary Education Act explicitly requires States to include limited English proficient students in the statewide assessments used to hold schools and school districts accountable for student performance. Title I of the Elementary and Secondary Education Act, 20 U.S.C. § 6311(b)(3)(F)(iii). If a school district uses the results of a test given for program accountability purposes to make educational decisions about individual students, the high-stakes use of the test must also be valid and reliable for this purpose.

For information on the factors that help ensure accuracy of tests for limited English proficient students, see pages 38 - 40 above. In making decisions about testing limited English proficient students, factors such as the student's level of English proficiency, the primary language of instruction, the level of literacy in the native language, and the number of years of instruction in English may all be pertinent.¹⁸⁷ When students participate in assessments designed to meet the requirements of Title I of the Elementary and Secondary Education Act, as amended, those assessments must be implemented in a manner that is consistent with both the requirements of Title VI and Title I.

III. Testing Of Students With Disabilities

Three federal statutes provide basic protections for students with disabilities. Section 504 of the Rehabilitation Act of 1973 (Section 504) and Title II of the Americans with Disabilities Act of 1990 (Title II) prohibit discrimination against persons with disabilities by public schools.¹⁸⁸ The Individuals with Disabilities Education Act (IDEA) establishes rights and protections for students with disabilities and their families. It also provides federal funds to state education agencies and school districts to assist in educating students with disabilities.¹⁸⁹ Under Section 504, Title II, and the IDEA,¹⁹⁰ school districts have a responsibility to provide students with disabilities with a free appropriate public education. Providing effective instruction in the general curriculum for students with disabilities is an important aspect of providing a free appropriate public education.

The regulations implementing Section 504 and Title II specifically provide that a recipient of federal funds may not "utilize criteria or methods of administration which have the effect of subjecting individuals to discrimination."¹⁹¹ Under Section 504, Title II, and the IDEA, tests given to students with disabilities must be selected and administered so that the test accurately reflects what the student knows or is able to do, rather than the student's disability (except when the test is designed to measure disability-related skills). This means that students with disabilities must be given appropriate accommodations and modifications in the administration of the tests. Examples include

¹⁸⁷ For more information on appropriate ways of testing students who are learning English, see *Ensuring Accuracy in Testing for English Language Learners*, (CCSSO, 2000).

¹⁸⁸ Although this part of the chapter deals only with students with disabilities attending public elementary and secondary schools, private schools that are not religious schools operated by religious organizations are covered by Title III of the ADA. Title II of the Americans with Disabilities Act of 1990, 42 U.S.C. §§ 12181 *et seq.* In addition, Title I of the Elementary and Secondary Education Act of 1965, as amended, contains important provisions regarding students with disabilities in the Title I program and their participation in assessments of Title I programs. 20 U.S.C. § 6311(b)(3)(F).

¹⁸⁹ The Individuals with Disabilities Education Act, 20 U.S.C. § 1400(d)(1)(c).

¹⁹⁰ The Section 504 regulation is found at 34 C.F.R. Part 104 (1999). The Title II regulation is found at 28 C.F.R. Part 35. The IDEA regulation is found at 34 C.F.R. Part 300.

¹⁹¹ See 34 C.F.R. § 100.3(b)(2) and similar provisions under Title IX, Section 504, and the ADA. In *Guardians*, 463 U.S. at 589, the United States Supreme Court upheld the use of the effects test, stating that the Title VI regulation forbids the use of federal funds, "not only in programs that intentionally discriminate on racial grounds but also in those endeavors that have a [racially disproportionate] impact on racial minorities."

oral testing, large print tests, Braille versions of tests, individual testing, and separate group testing.

Generally, there are three critical areas in which high-stakes testing issues arise for students with disabilities: (1) tests used to determine whether a student has a disability and, if so, the nature of the disability; (2) tests used to determine if the student meets the criteria for other specialized instructional programs, such as gifted and talented or vocational education programs; and (3) system-wide tests administered to determine if students have met performance standards.¹⁹²

Under Section 504, Title II, and the IDEA, before a student can be classified as having a disability, the responsible education agency must individually evaluate the student in accordance with specific statutory and regulatory requirements, including requirements regarding the validity of tests and the provision of appropriate accommodations.¹⁹³ These requirements prohibit the use of a single test score as the sole criterion for determining whether a student has a disability and for determining an appropriate educational placement for the student.¹⁹⁴

When tests are used for other purposes, such as in making decisions about placement in gifted and talented programs, it is important that tests measure the skills and abilities needed in the program, rather than the disability, unless the test purports to measure skills or functions which are impaired by the disability and such functions are necessary for participation in the program.¹⁹⁵ For this reason, appropriate accommodations may need to be provided to students with disabilities in order to measure accurately their performance in the skills and abilities required in the program.

Furthermore, federal law requires the inclusion of students with disabilities in state- and district-wide assessment programs, including high-stakes tests, except as participation in such tests is individually determined to be inappropriate for a particular student. Such assessments provide valuable information which benefits students, either directly, such as in the measurement of individual progress against standards, or indirectly, such as in evaluating programs. Given these benefits, exclusion from assessment programs based on disability generally would violate Section 504 and Title II. If a student with a disability will take the system-wide assessment test, including a high-stakes test, the student must be provided appropriate instruction and appropriate test accommodations.¹⁹⁶

¹⁹² Tests used for college admission are discussed on pp. 4-5.

¹⁹³ See 34 C.F.R. § 104.35(b) for specific provisions covering the use of tests for evaluation purposes.

¹⁹⁴ See 34 C.F.R. § 104.35(c), requiring placement decisions to consider information from a variety of sources.

¹⁹⁵ See 34 C.F.R. § 104.35(b)(3) and 34 C.F.R. § 300.532.

¹⁹⁶ See *Brookhart*, 697 F.2d at 183-184. Some courts have held that a student with a disability may be denied a diploma if, despite receiving appropriate services and testing accommodations, the student, because of the disability, is unable to pass the required test or meet other graduation requirements. *Id.* at 183; *Anderson*, 520 F. Supp. at 509-511; *Board of Educ. v. Ambach*, 458 N.Y.S.2d 680, 684-685, 689 (N.Y. App. Div. 1982), *aff'd*, 469 N.Y.S.2d 669 (1983).

In addition, the Individuals with Disabilities Education Amendments of 1997 specifically require states, as a condition of receiving IDEA funds, to include students with disabilities in the regular state- and district-wide assessment programs, with appropriate accommodations, where necessary.¹⁹⁷ The IDEA requirements cover tests with high-stakes consequences given to measure individual achievement as well as tests given for program accountability purposes. The IDEA also requires state or local educational agencies to develop guidelines for the relatively small number of students with disabilities who cannot take part in state- and district-wide tests to participate in alternate assessments.¹⁹⁸

For children with disabilities, school personnel knowledgeable about the student, the nature of the disability, and the testing program, in conjunction with the student's parent or guardian, determine whether the student will participate in all or part of the state- or district wide assessment of student achievement.¹⁹⁹ The decision must be documented in the student's individualized education program (IEP), or a similar record such as a Section 504 plan. These records must also state any individual accommodations in the administration of the state- or district-wide assessments of student achievement that are needed to enable the student to participate in such assessment. An IEP, developed under the IDEA, must also explain how the student will be assessed if it is inappropriate for the student to participate in the testing program even with accommodations.²⁰⁰

Section 504 and Title II also prohibit discrimination in virtually all public and private post-secondary institutions. The regulatory requirements related to disability discrimination are different in post-secondary education than in elementary and secondary education. Post-secondary institutions are not required to evaluate students or to provide them with a free appropriate education.

High-stakes testing issues at the post-secondary level generally relate to tests used in admissions, including tests given by an educational institution or other covered entities as prerequisites for entering a career or career path, and tests of academic competency required by the institution to complete a program. This guide is not intended to offer a complete or detailed explanation of each of these testing situations, but only a brief synopsis.²⁰¹

¹⁹⁷ See 34 C.F.R. § 300.138(a).

¹⁹⁸ See 34 C.F.R. § 300.138(b). The IDEA Final Regulations, Attachment I--Analysis of Comments and Changes, 64 Fed. Reg. 12406, 12564 (1999) projects that there will be a relatively small number of students who will not be able to participate in the district or state assessment program with accommodations and modifications, and will therefore need to be assessed through alternate means. These alternate assessments must be developed and conducted beginning not later than July 1, 2000.

¹⁹⁹ See 34 C.F.R. § 300.347(a)(5) for the IEP requirements applicable to assessment of students with disabilities under IDEA and 34 C.F.R. § 104.33 for the more general evaluation requirements under Section 504.

²⁰⁰ See 34 C.F.R. § 300.347(a)(5).

²⁰¹ Test providers that are not higher education institutions may be covered by Section 504 if they receive federal funds; by Title II if they are parts of governmental units; or by Title III if they are private entities. Each of these laws has its

The Section 504 regulation specifically provides that higher education institutions' admissions procedures may not make use of any test or criterion for admission that has a disproportionate, adverse impact on individuals with disabilities unless (1) the test or criterion, as used by the institution, has been validated as a predictor of success in the education program or activity and (2) alternative tests or criteria that have a less disproportionate, adverse impact are not shown to be available.²⁰² In administering tests, appropriate accommodations must be provided so that the person can demonstrate his or her aptitude and achievement, not the effect of the disability (except where the functions impaired by the disability are the factors the test purports to measure).²⁰³

For other high-stakes tests that an institution might administer, such as rising junior tests, similar requirements apply.²⁰⁴ The institution must provide adjustments or accommodations and auxiliary aids and services that enable the student to demonstrate the knowledge and skills being tested.²⁰⁵

Students are required to notify the educational institution when accommodations are needed and supply adequate documentation of a current disability and the need for accommodation. The student's preferred accommodation does not have to be provided as long as an effective accommodation is provided.

Test accommodations are intended to provide the person with disabilities the means by which to demonstrate the skills and knowledge being tested. Although Section 504 and Title II require a college or university to make reasonable modifications, neither Section 504 nor Title II requires a college or university to change, lower, waive, or eliminate academic requirements or technical standards, including admissions requirements, that can be demonstrated by the college or university to be essential to its program of instruction or to any directly related licensing requirement.²⁰⁶ Accommodations requested by students need not be provided if they would result in a fundamental alteration to the institution's program.²⁰⁷

own requirements. For more information regarding testing under Title III of the ADA, consult the U.S. Department of Justice.

²⁰² 34 C.F.R. § 104.42(b)(2). Appendix A to the Section 504 regulation, Subpart E-Post-secondary Education, No. 29, notes that the party challenging the test would have the burden of showing that alternate tests with less disparate impact are available.

²⁰³ See 34 C.F.R. § 104.42(b)(2). Appendix A to the Section 504 regulation, Subpart E-Post-secondary Education, No. 29, notes that the party challenging the test would have the burden of showing that alternate tests with less disparate impact are available.

²⁰⁴ Some undergraduate college programs require students to pass a rising junior examination to determine whether students have met the college's standards in writing or other academic skills as a prerequisite for advancement to junior year status.

²⁰⁵ See 34 C.F.R. § 104.44(a) & (d).

²⁰⁶ See 34 C.F.R. § 104.44 (a).

²⁰⁷ See *Southeastern Community College v. Davis*, 442 U.S. 397, 413 (1979); *Wynne v. Tufts Univ. Sch. of Med.*, 976 F.2d 791, 794-796 (1st Cir. 1992), cert. denied, 507 U.S. 1030 (1993).

IV. Constitutional Protections

In addition to applying federal nondiscrimination statutes, courts have also considered constitutional issues that may arise when public school districts or state education agencies require students to pass certain tests that are intended to certify that students have attained a level of competency in skills or knowledge taught in the program.²⁰⁸ Constitutional challenges to testing programs under the Fourteenth Amendment have raised both equal protection and due process claims. The equal protection principles involved in discrimination cases are, generally speaking, the same as the standards applied to intentional discrimination claims under the applicable federal nondiscrimination statutes.²⁰⁹

The due process clause of the Fourteenth Amendment is particularly associated with cases challenging the adequacy of the notice provided to students prior to this type of test and the students' opportunity to learn the required content.²¹⁰ In analyzing such due process claims, courts have generally considered three issues:

²⁰⁸ The U.S. Department of Education, Office for Civil Rights, does not have jurisdiction to resolve constitutional cases. However, some cases involve constitutional issues that overlap with discrimination issues arising under federal civil rights laws.

²⁰⁹ Federal cases may involve equal protection challenges to a jurisdiction's use of tests in which the claim is not based on intentional race or sex discrimination, but, instead, on the alleged impropriety of the jurisdiction's use of tests to separate out those students who should not be allowed to graduate. As a general matter, courts express reluctance to second guess a state's educational policy choices when faced with such challenges, although recognize that a state cannot "exercise that [plenary] power without reason and without regard to the United States Constitution." *Debra P.*, 644 F.2d at 403. When there is no claim of discrimination based on membership in a suspect class, the equal protection claim is reviewed under the rational basis standard. In these cases, the jurisdiction need show only that the use of the tests has a rational relationship to a valid state interest. *Id.* at 406. See also *Erik V.*, 977 F. Supp. at 389.

²¹⁰ A review of relevant cases reveals the highly fact and context-specific nature of the conclusions reached by federal courts considering alleged violations of the due process clause. In *Debra P.*, 644 F.2d at 404, the Fifth Circuit held that students' due process rights were violated when a newly imposed minimum competency test required for high school graduation was instituted without adequate notice and an opportunity for students to learn the material covered by the test. Three years later, in *Debra P. v. Turlington*, 730 F.2d at 1416-1417, the court held that students who now had six years notice of the exam were afforded the opportunity to learn the relevant material, given the state's remedial programs. For additional courts identifying due process violations in the way in which a competency test was instituted, see *Brookhart*, 697 F.2d at 186-187 (holding that district-required minimum competency test for graduation denied due process to students with disabilities where notice was inadequate and students had not been exposed to 90% of the material covered by the test); *Crump v. Gilmer Indep. Sch. Dist.*, 797 F. Supp. 552, 556-557 (E.D. Tex. 1992) (granting temporary restraining order where district had not demonstrated validity of graduation examination in light of actual instructional content); *Anderson*, 520 F. Supp. at 508-509 (finding that school district failed to show that minimum competency test required for high school graduation covered material actually taught at school). Other cases have concluded that adequate notice was provided, the test or criterion at issue was closely related to the instructional program, or the promotion decision was not shown to be outside the discretion of school authorities. See *Erik V.*, 977 F. Supp. at 389-390 (finding that promotion decision was within proper purview of school authorities); *Williams v. Austin Indep. Sch. Dist.*, 796 F. Supp. 251, 253-254 (W.D. Tex. 1992) (considering students to have had seven years advance notice of high school competency exam although standards of performance were recently raised). See also promotion cases in which students were required to demonstrate adequate reading skills, although a separate test was not apparently involved. *Bester v. Tuscaloosa City Bd. of Educ.*, 722 F.2d 1514, 1516 (11th Cir. 1984) (finding reading standards required for promotion to merely reinforce district policy of retention for substandard work); *Sandlin v. Johnson*, 643 F.2d 1027, 1029 (4th Cir. 1981) (finding denial of second grade promotion for failing to attain required level in reading series within discretion of school district). For a testing case raising similar due process issues at the post-secondary level, see *Mahavongsanan v. Hall*, 529 F.2d 448, 450 (5th Cir. 1976) (finding no violation of due

(1) Is the purpose of the testing program legitimate?

Federal courts typically defer to educators' policy judgments regarding the value of the educational benefits sought from testing programs, as long as these judgments are not arbitrary or capricious.²¹¹ Improving the quality of elementary and secondary education through the establishment of academic standards has been seen as a reasonable goal of a testing program, and colleges and universities are generally given wide latitude in framing degree requirements and making academic decisions.²¹²

(2) Have students received adequate notice of the test and its consequences?

In the elementary and secondary context, courts have required sufficient advance notice of tests required for graduation to give students a reasonable chance to learn the material presented on the test.²¹³ A particularly important concern in some of these decisions is the adequacy of notice provided to students. This issue has arisen in cases where racial minority students and students with disabilities received inadequate notice and did not receive a program of instruction that prepared them to pass the test.²¹⁴ In looking at the length of the transition period needed between announcement of a new requirement and its full implementation, the kind of test and the context in which it is administered are central factors to be considered. Specific circumstances taken into account include the nature of instructional supports, including remediation, that accompany the test,²¹⁵

process where the university's decision to require a comprehensive examination for receipt of a graduate degree was a reasonable academic regulation, plaintiff received timely notice that she would be required to take the examination, she was allowed to retake the test, and the university afforded her an opportunity to complete additional course work in lieu of the examination).

²¹¹ The determination as to whether a testing program is rationally related to a legitimate educational goal is technically considered as one of substantive due process under the Fourteenth Amendment. Courts have approved testing as a rational means of improving educational outcomes. See *Debra P.*, 644 F.2d at 406; *Anderson*, 520 F. Supp. at 506. Insofar as due process cases may involve additional questions of the validity of the test used to address institution's goal, these issues are discussed in the portions of the guide addressing discrimination under federal civil rights laws.

²¹² See *Ewing*, 474 U.S. at 222, 226-227 (acknowledging that courts will not review academic decisions of colleges and universities unless the decision is such a substantial departure from accepted academic norms as to demonstrate that professional judgment was not actually exercised or where discrimination is claimed); *Debra P.*, 644 F.2d at 402 (finding praiseworthy a state's effort to set standards to improve public education).

²¹³ Although there are important exceptions, see *United States v. LULAC*, 793 F.2d at 648, and *Anderson*, 520 F. Supp. at 505, courts have often considered the issue of adequate notice to be one of procedural due process. For procedural due process to apply, a protected property or liberty interest must be identified. See *Debra P.*, 644 F.2d at 404 (finding sufficient to trigger due process protection a state-created mutual expectation that students who successfully complete required courses would receive diploma); *Brookhart*, 697 F. 2d at 185 (identifying a liberty interest, based on stigma of diploma denial, that disastrously affected plaintiffs' future employment and educational opportunities); *Erik V.*, 977 F. Supp. at 389-390 (finding no property interest in grade level promotion warranting preliminary injunction).

²¹⁴ See *Brookhart*, 697 F. 2d at 186-188; *Debra P.*, 644 F.2d at 404.

²¹⁵ See *Debra P.*, 730 F.2d at 1407, 1410-12, 1415-1416; *Anderson*, 520 F. Supp. at 505.

whether re-testing is permitted,²¹⁶ and whether the decision to promote or graduate the student considers other information about the student's performance.²¹⁷

- (3) Are students actually taught the knowledge and skills measured by the test?

Several courts have found that "fundamental fairness" requires that students be taught the material covered by the test where passing the test is a condition for receipt of a high school diploma.²¹⁸ In analyzing this issue in a case involving a state where there had been past intentional segregation in elementary and secondary schools before a statewide diploma test was required, and where racial minority students had a disproportionate failure rate on the test, the courts took the state's past intentional segregation into account in determining whether racial minority students had had adequate opportunities to learn the material covered by the test.²¹⁹ For the test to meaningfully measure student achievement, the test, the curriculum, and classroom instruction should be aligned. In cases examining system-wide administration of a test, courts require evidence that the content covered by the test is actually taught, but may not expect proof that every student has received the relevant instruction.²²⁰

²¹⁶ Re-testing was available in *Erik V.*, 977 F. Supp. at 388-389, and in *Anderson*, 520 F. Supp. at 505.

²¹⁷ See *Erik V.*, 977 F. Supp. at 387 (reading performance of students with grades of A, B, or C on grade level work was further reviewed by teacher and principal to determine if student should be promoted notwithstanding the failing test score).

²¹⁸ The question of instructional or curricular validity is usually posed as one of substantive due process. See *Brookhart*, 697 F.2d at 184-187; *Debra P.*, 644 F.2d at 406; *Anderson*, 520 F. Supp. at 509.

²¹⁹ *Debra P.*, 644 F.2d at 407 (where black students disproportionately failed a statewide test necessary to obtain a high school diploma, and, due to the prior dual school system, black students received a portion of their education in unequal, inferior segregated schools, and where the state was unable to show that the diploma sanction did not perpetuate the effects of that past intentional discrimination, the court found that immediate use of the diploma sanction punished the black students for deficiencies created by the dual school system in violation of their constitutional right to equal protection); *Debra P.*, 474 F. Supp. at 257 ("punishing the victims of past discrimination for deficits created by an inferior educational environment neither constitutes a remedy nor created better educational opportunities").

²²⁰ See *Anderson*, 540 F. Supp. at 765.

APPENDIX A: Glossary of Legal Terms

This glossary is provided as a plain language reference to assist non-lawyers in understanding commonly used legal terms that are either used in this guide or are important to know in understanding the terms in the guide. Legal terms are often "terms of art." In other words, they mean something slightly different or more specific in the legal context than they do in ordinary conversation.

Burden of proof—the duty of a party to substantiate its claim or defense against the other party. In civil actions, the weight of this proof is usually described as a preponderance of the evidence. See BLACK'S LAW DICTIONARY 196-197 (6th ed. 1990). See *Disparate impact*.

Constitutional rights—the rights of each American citizen that are guaranteed by the United States Constitution. See *Brown v. Board of Education*, 347 U.S. 483 (1954); *Bolling v. Sharpe*, 347 U.S. 497 (1954); BLACK'S LAW DICTIONARY 312 (6th ed. 1990).

De jure segregation or discrimination— term applied to systemic school segregation that was mandated by statute or that was accomplished through the intentionally segregative actions of local school districts or state agencies.

Different Treatment—a claim that similarly situated persons are treated differently because of their race, color, national origin, sex or disability. Under federal non-discrimination laws, policies and practices must be applied consistently to an individual or group of students regardless of their race, national origin, sex, or disability, unless there is a lawful reason for not doing so. To prove different treatment, one must show that "a challenged action was motivated by an intent to discriminate." *Elston v. Talladega County Bd. of Educ.*, 997 F.2d 1394, 1406 (11th Cir. 1993). This requires a showing that the decision-maker was not only aware of the person's race, national origin, sex, or disability, but that the recipient acted, at least in part, because of the person's race, national origin, sex or disability. However, the record need not contain "direct evidence of bad faith, ill will or any evil motive," on the part of the recipient. *Elston*, 997 F.2d at 1406, (quoting *Williams v. City of Dotham*, 745 F.2d 1406, 1414 (11th Cir. 1984)). Evidence of discriminatory intent may be direct or circumstantial. Different treatment may be justified by a lawful reason, for example, to remedy prior discrimination. See generally *Wygant v. Jackson Bd. of Educ.*, 476 U.S. 267, 290-291 (1986); *United States v. Fordice*, 505 U.S. 717, 728-730 (1992); *Regents of the Univ. of Cal. v. Bakke*, 438 U.S. 265, 305-320 (1978), *Hopwood v. Texas*, 78 F.3d 932, 948-950 (5th Cir. 1996), *cert. denied*, 518 U.S. 1033 (1996); BLACK'S LAW DICTIONARY 470 (6th ed. 1990).

Disparate impact—disparate impact analysis applies when the application of a neutral criterion or a facially neutral practice has discriminatory effects and the criterion or practice is not determined to be "educationally justified" or "educationally necessary." In contrast to intentional discrimination, the disparate impact analysis does not require proof of discriminatory motive. Under the disparate impact analysis, the party

challenging the criterion or practice has the burden of establishing disparate impact. If disparate impact is established, the party defending the practice must establish an "educational justification." If the educational institution provides sufficient evidence that the test use in question is justified educationally, the party challenging the test has the opportunity to show that there exists an alternative practice that meets the institution's goals as well as the challenged test use and that would eliminate or reduce the adverse impact. See *Board of Educ. v. Harris*, 444 U.S. 130, 143 (1979); *Groves v. Alabama State Bd. of Educ.*, 776 F. Supp. 1518 (M.D. Ala. 1991); *Georgia State Conf. of Branches of NAACP v. Georgia*, 775 F.2d 1403, 1412 (11th Cir. 1985).

Dual system—a previously segregated educational system in which black and white schools, ostensibly similar, existed side-by-side. See *Brown v. Board of Educ.*, 347 U.S. 483 (1954); *Anderson v. Banks*, 520 F. Supp. 472, 499-501 (S.D. Ga. 1981).

Due process—a constitutionally guaranteed right. The Fifth Amendment states that no citizen shall "be deprived of life, liberty, or property, without due process of law." The Fourteenth Amendment applied this passage to the states as well. Today it is used by the judiciary to define the scope of fundamental fairness due to each citizen in his or her interactions with the government and its agencies. Some courts have held that a student's expectation in receiving a high school diploma in return for meeting certain attendance and academic criteria is a form of a property right or liberty interest. See *Debra P. v. Turlington*, 644 F.2d 397 (5th Cir. 1981); *Crump v. Gilmer Indep. Sch. Dist.*, 797 F. Supp. 552, 555-556 (E.D. Tex. 1992); But see *Board of Educ. v. Ambach*, 458 N.Y.S.2d 680, (N.Y. App. Div. 3d Dep't 1982), *aff'd*, 457 N.E.2d 775 (1983); BLACK'S LAW DICTIONARY 500-501 (6th ed. 1990). See also *Procedural Due Process, Substantive Due Process*.

Educational necessity—once the party challenging the practice has shown a significant disparate impact, the educational institution using the challenged practice must present sufficient evidence that it is justified by educational necessity. Educational necessity generally refers to a showing that practices or procedures are necessary to meeting an important educational goal. See *Elston v. Talladega County Bd. of Educ.*, 997 F.2d 1394, 1412 (11th Cir. 1993) (citing *Georgia State Conf. of Branches of NAACP v. Georgia*, 775 F.2d 1403, 1412, 1417 (11th Cir. 1985)). In the context of testing this means the test or assessment procedure must serve a legitimate educational goal and be valid and reliable for the purpose used.

Equal protection—classifications based on race, sex or other grounds may be challenged under the equal protection clause of the Fourteenth Amendment to the U.S. Constitution when imposed by state or local government agencies. Distinctions explicitly based on race or ethnicity, neutral criteria having a discriminatory purpose or other intentionally discriminatory conduct based on race or ethnicity will violate the Fourteenth Amendment, unless the action is narrowly tailored to serve a compelling purpose. Intentional sex discrimination will violate the Fourteenth Amendment unless there is an exceedingly persuasive justification. *United States v. Virginia*, 518 U.S. 515 (1996).

Distinctions based on other grounds will not violate the equal protection clause unless they are not rationally related to a legitimate governmental objective.

Facially neutral—a regulation, rule, practice or other activity that does not appear to be discriminatory. Facially neutral practices may be found to violate regulations implementing federal civil rights laws if they adversely impact a group based on race, national origin, sex or disability without a legitimate educational justification. See *Larry P. v. Riles*, 793 F.2d 969 (9th Cir. 1984); *Lau v. Nichols*, 414 U.S. 563 (1974).

High-stakes educational decisions for students—decisions that have significant impact or consequences for individual students. These decisions may involve student placement in gifted and talented programs; decisions concerning whether a student has a disability; the appropriate educational program for a student with a disability; promotion or graduation decisions; and higher education admissions decisions and scholarship awards. See Jay P. Heubert & Robert Hauser, eds., *HIGH STAKES: TESTING FOR TRACKING, PROMOTION, AND GRADUATION* 1-2 (1999); *Larry P. v. Riles*, 793 F.2d 969 (9th Cir. 1984); *Sharif v. New York State Educ. Dep't*, 709 F. Supp. 345 (S.D.N.Y. 1989).

Less discriminatory alternative—if the education institution presents sufficient evidence that the test use or educational practice in question is justified educationally, the party challenging the test has the opportunity to show that there exists an equally or comparably effective alternative practice that meets the institution's goals and that would eliminate or reduce the adverse impact. *Elston v. Talladega County Bd. of Educ.*, 997 F.2d 1394, 1407 (11th Cir. 1993); *Georgia State Conference of NAACP Branches v. State of Georgia*, 775 F.2d 1403 (11th Cir. 1985). Costs and administrative burdens are among the factors considered in assessing whether the alternative practice is equally effective in fulfilling the institution's goals. *Ward's Cove Packing Co. v. Atonio*, 490 U.S. 642, 661 (1989); *Sharif v. New York State Educ. Dep't*, 709 F. Supp. 345, 363-364 (S.D.N.Y. 1989) (defendant's claim that proposed alternative was not feasible and excessively burdensome not persuasive since most other states used proposed alternative).

Procedural due process—the right each American citizen has under the Constitution to a fair process in actions that affect an individual's life, liberty or property. Procedural due process includes notice and the right to be heard. Some courts have found that procedural due process applies to the implementation of minimum competency examinations required for high school graduation. See *Debra P. v. Turlington*, 474 F. Supp. 244, 263-64 (M.D. Fla. 1979), *aff'd in part and vacated in part*, 644 F.2d 397 (5th Cir. 1981); *Erik V. v. Causby*, 977 F. Supp. 384, 389-90 (E.D.N.C. 1997); *Crumpp v. Gilmer Indep. Sch. Dist.*, 797 F. Supp. 552, 555-56 (E.D. Tex. 1992); BLACK'S LAW DICTIONARY 1203 (6th ed. 1990).

Significantly disproportionate—when statistical analysis shows that the success rate of members of an identified group is significantly lower than would be expected from random distribution within the appropriate qualified pool, the test in question is said to

have a disproportionate adverse impact. There is no set formula to determine when a sufficient level of adverse impact has been reached; the Supreme Court has stated that statistical disparities must be sufficiently substantial that they raise an inference of causation. Courts have advanced percentage disparities, standard deviations or other statistical formulae to address this component. Disparate impact itself does not necessarily mean that discrimination has taken place, but it does trigger an inquiry regarding the educational justification of the challenged practice. See *Watson v. Fort Worth Bank & Trust*, 487 U.S. 977, 994-95 (1988); *Richardson v. Lamar County Bd. of Educ.*, 729 F. Supp. 806, 815-16 (M.D. Ala. 1989), *aff'd*, 935 F.2d 1240 (11th Cir. 1991); *Groves v. Alabama State Bd. of Educ.*, 776 F. Supp. 1518, 1529-32 (M.D. Ala. 1991).

Statutory rights—rights protected by statute, as opposed to constitutional rights, which are protected by the Constitution.

Substantive due process—often stated as "fundamental fairness." In an education context, proof that students had not been taught the material on which they were tested might be a substantive due process violation. Some courts have held that students have the equivalent of a property or liberty interest in graduating or being promoted according to the expectations given them. See *Debra P. v. Turlington*, 644 F.2d 397 (5th Cir. 1981); *Crump v. Gilmer Indep. Sch. Dist.*, 797 F. Supp. 552, 555-56 (E.D. Tex. 1992). BLACK'S LAW DICTIONARY 1429 (6th ed. 1990).

Unitary system—a desegregated school system. The Supreme Court has held that all previously intentionally segregated school systems are required to become unitary systems. Although the term has been interpreted in different ways by different courts, a "unitary system" is typically one in which all vestiges of past discrimination and segregated practices have been eliminated. See *Freeman v. Pitts*, 506 U.S. 467, 486-489 (1992); *Board of Educ. v. Dowell*, 498 U.S. 237, 243-246, 249-251 (1991); *Keyes v. School Dist. No. 1*, 413 U.S. 189, 208, 257-258 (1973); *Debra P. v. Turlington*, 474 F. Supp. 244, 249-257 (M.D. Fla. 1979) *aff'd* in part and vacated in part, 644 F.2d 397 (5th Cir. 1981); *Bester v. Tuscaloosa City Bd. of Educ.*, 722 F.2d 1514, 1517 (11th Cir. 1984); *Georgia State Conference of Branches of NAACP v. Georgia*, 775 F.2d 1403, 1413-1416 (11th Cir. 1985).

APPENDIX B: Glossary of Test Measurement Terms

This glossary is provided as a plain language reference to assist readers in understanding commonly used test measurement terms used in this guide or terms relevant to issues discussed in the guide. For additional relevant information, readers are encouraged to review the Glossary in the *Joint Standards*, as well as the appropriate chapters in the *Joint Standards*.

Achievement level/ proficiency levels—Descriptions of a test taker's competency in a particular area of knowledge or skill, usually defined as ordered categories on a continuum, often labeled from "basic" to "advanced," that constitute broad ranges for classifying performance.

Accommodation—A change in how a test is presented, in how a test is administered, or in how the test taker is allowed to respond. This term generally refers to changes that do not substantially alter what the test measures. The proper use of accommodations does not substantially change academic level or performance criteria. Appropriate accommodations are made in order to level the playing field, i.e., to provide equal opportunity to demonstrate knowledge.

Alternate Assessment—An assessment designed for those students with disabilities who are unable to participate in general large-scale assessments used by a school district or state, even when accommodations or modifications are provided. The alternate assessment provides a mechanism for students with even the most significant disabilities to be included in the assessment system.

Assessment—Any systematic method of obtaining information from tests and other sources, used to draw inferences about characteristics of people, objects, or programs.

Bias—In a statistical context, a systematic error in a test score. In discussing test fairness, bias may refer to construct underrepresentation or construct irrelevant components of test scores. Bias usually favors one group of test takers over another.

Bilingual—The characteristic of being relatively proficient in two languages.

Classification accuracy—The degree to which neither false positive nor false negative categorizations and diagnoses occurs when a test is used to classify an individual or event.

Composite score—A score that combines several scores according to a specified formula.

Content areas—Specified subjects in education, e.g. language arts, science, mathematics, or history.

Content domain—The set of behaviors, knowledge, skills, abilities, attitudes or other characteristics to be measured by a test, represented in a detailed specification, and often organized into categories by which items are classified.

Content validity—Validity evidence which analyzes the relationship between a test's content and the construct it is intended to measure. Evidence based on test content includes logical and empirical analyses of the relevance and representativeness of the test content to the defined domain of the test and the proposed interpretations of test scores.

Content standard—Statements which describe expectations for students in a subject matter at a particular grade or at the completion of a level of schooling.

Construct—The concept or the characteristic that a test is designed to measure.

Construct equivalence—1. The extent to which the construct measured by one test is essentially the same as the construct measured by another test. 2. The degree to which a construct measured by a test in one cultural or linguistic group is comparable to the construct measured by the same test in a different cultural or linguistic group.

Construct irrelevance—The extent to which test scores are influenced by factors that are irrelevant to the construct that the test is intended to measure. Such extraneous factors distort the meaning of test scores from what is implied in the proposed interpretation.

Construct underrepresentation—The extent to which a test fails to capture important aspects of the construct that the test is intended to measure. In this situation, the meaning of test scores is narrower than the proposed interpretation implies.

Constructed response item—An exercise for which examinees must create their own responses or products rather than choose a response from an enumerated set. Short-answer items require a few words or a number as an answer, whereas extended-response items require at least a few sentences.

Criterion validity—Validity evidence which analyzes the relationship of test scores to variables external to the test. External variables may include criteria that the test is expected to be associated with, as well as relationships to other tests hypothesized to measure the same constructs and tests measuring related constructs. Evidence based on relationships with other variables addresses questions about the degree to which these relationships are consistent with the construct underlying the proposed test interpretations. *See* predictive validity.

Criterion-referenced—Scores of students referenced to a criterion. For instance, a criterion may be specific, identified knowledge and skills which students are expected to master. Academic content standards in various subject areas are examples of this type of criterion.

Criterion-referenced test—A test that allows its users to make score interpretations in relation to a functional performance level, as distinguished from those interpretations that are made in relation to the performance of others. Examples of criterion-referenced interpretation include comparison to cut scores, interpretations based on expectancy tables, and domain-referenced score interpretations.

Cutscore—A specified point on a score scale, such that scores at or above that point are interpreted or acted upon differently from scores below that point. *See* performance standard.

Discriminant validity—Validity evidence based on the relationship between test scores and measures of different constructs.

Error of measurement—The difference between an observed score and the corresponding true score or proficiency. This unintended variation in scores is assumed to be random and unpredictable and impacts the estimate of reliability of a test.

False negative—In classification, diagnosis, or selection, an error in which an individual is assessed or predicted not to meet the criteria for inclusion in a particular group but in truth does (or would) meet these criteria.

False positive—In classification, diagnosis, or selection, an error in which an individual is assessed or predicted to meet the criteria for inclusion in a particular group but in truth does not (or would not) meet these criteria.

Field test—A test administration used to check the adequacy of testing procedures, generally including test administration, test responding, test scoring, and test reporting. A field test is generally more extensive than a pilot test. *See* pilot test.

High-stakes decision for students—A decision whose result has important, direct consequences for examinees.

Internal consistency estimate of reliability—An index of the reliability of test scores derived from the statistical interrelationships of responses among item responses or scores on separate parts of a test.

Inter-rater agreement—The consistency with which two or more judges rate the work or performance of test takers; sometimes referred to as *inter-rater reliability*.

Local evidence—Evidence (usually related to reliability or validity) collected for a specific and particular set of test takers in a single institution, district, or state, or at a specific location.

Local norms—Norms by which test scores are referred to a specific, limited *reference population* of particular interest to the test user (e.g., institution, district, or state); local norms are not intended as representative of populations beyond that setting.

Norm-referenced—Scores of students compared to a specified reference population.

Norm-referenced test—A test that allows its users to make score interpretations of a test taker's performance in relation to the performance of other people in a specified reference population.

Norms—Statistics or tabular data that summarize the distribution of test performance for one or more specified groups, such as test takers of various ages or grades. The group of examinees represented by the norms is referred to as the reference population. Norm reference populations can be a local population of test takers, e.g. from a school, district or state, or it can represent a larger population, such as test takers from several states or throughout the country.

Percentile rank—Most commonly, the percentage of scores in a specified distribution that fall below the point at which a given score lies. Sometimes the percentage is defined to include scores that fall at the point; sometimes the percentage is defined to include half of the scores at the point.

Performance assessments—Product- and behavior-based measurements based on settings designed to emulate real-life contexts or conditions in which specific knowledge or skills are actually applied.

Performance standard—1. An objective definition of a certain level of performance in some domain in terms of a cut score or a range of scores on the score scale of a test measuring proficiency in that domain. 2. A statement or description of a set of operational tasks exemplifying a level of performance associated with a more general content standard; the statement may be used to guide judgements about the location of a cut score on a score scale. The term often implies a desired level of performance. *See* cutscore.

Pilot test—A test administered to a representative sample of test takers to try out some aspects of the test or test items, such as instructions, time limits, item response formats, or item response options. *See* field test.

Portfolio assessments—A systematic collection of educational or work products that have been compiled or accumulated over time, according to a specific set of principles.

Precision of measurement—A general term that refers to a measure's sensitivity to error of measurement.

Predictive validity—Validity evidence that analyzes the relationship of test scores to variables external to the test that the test is expected to predict. Predictive evidence indicates how accurately test data can predict criterion scores that are obtained or outcomes that occur at a later time. When test scores are used to predict a dichotomous criterion, such as a diagnosis, false positive and false negative errors can occur. *See* criterion evidence of validity; false positive error and false negative error.

Random error—An unsystematic error; a quantity (often observed indirectly) that appears to have no relationship to any other variable.

Reference population—The population of test takers represented by test norms. The sample on which the test norms are based must permit accurate estimation of the test score distribution for the reference population. The reference population may be defined in terms of size of the population (local or larger), examinee age, grade, or clinical status at time of testing, or other characteristics.

Reliability—The degree to which test scores for a group of test takers are consistent over repeated applications of a measurement procedure and hence are inferred to be dependable, and repeatable for an individual test taker; the degree to which scores are free of errors of measurement for a given group.

Sample—A selection of a specified number of entities called sampling units (test takers, items, schools, etc.) from a large specified set of possible entities, called the population. A random sample is a selection according to a random process, with the selection of each entity in no way dependent on the selection of other entities. A stratified random sample is a set of random samples, each of a specified size, from several different sets, which are viewed as strata of the population.

Sampling from a domain—The process of selecting test items to represent a specified universe of performance.

Score—Any specific number resulting from the assessment of an individual; a generic term applied for convenience to such diverse measures as test scores, absence records, course grades, ratings, and so forth.

Scoring rubric—The established criteria, including rules, principles, and illustrations, used in scoring responses to individual items and clusters of items. The term usually refers to the scoring procedures for assessment tasks that do not provide enumerated responses from which test takers make a choice. Scoring rubrics vary in the degree of judgement entailed, in the number of distinct score levels defined, in the latitude given scorers for assigning intermediate or fractional score values, and in other ways.

Selection—A purpose for testing that results in the acceptance or rejection of applicants for a particular educational opportunity.

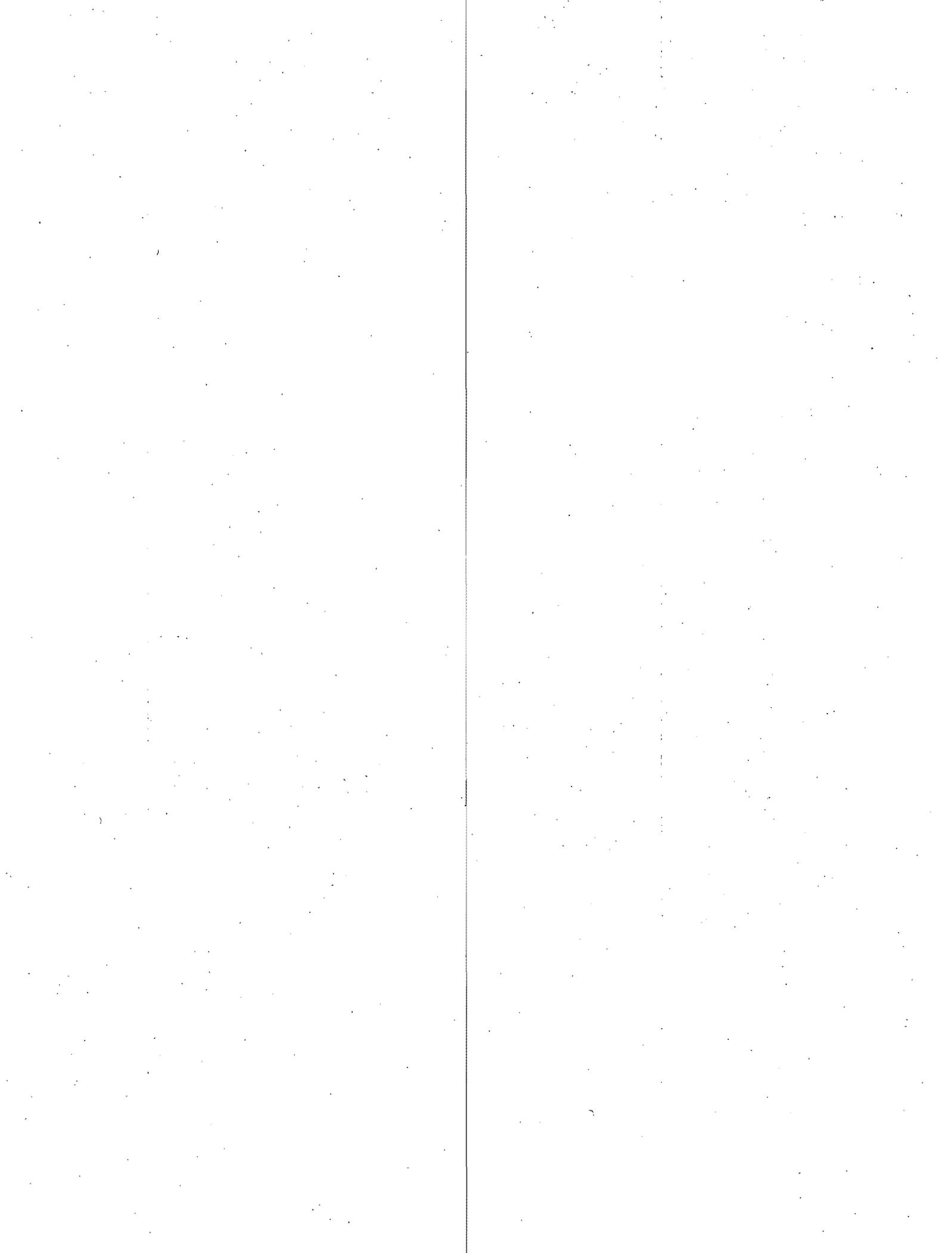
Sole criterion—When only one standard (such as a test score) is used to make a judgement or a decision. This can include a step-wise decision making procedure where students must reach or exceed one criterion (such as a cutscore of a test) before other criteria can be considered.

Speed test—A test in which performance is measured primarily or exclusively by the time to perform a specified task, or the number of tasks performed in a given time, such as tests of typing speed and reading speed.

Standards-based assessment—Assessments intended to represent systematically described content and performance standards.

Systematic error—A score component (often observed indirectly), not related to the test performance, that appears to be related to some salient variable or sub-grouping of cases in empirical analyses. This type of error tends to increase or decrease observed scores consistently in members of the subgroup or levels of the salient variable. *See* bias.

Technical manual—A publication prepared by test authors and publishers to provide technical and psychometric information on a test.



Test developer—The person(s) or agency responsible for the construction of a test and for the documentation regarding its technical quality for an intended purpose.

Test development—The process through which a test is planned, constructed, evaluated and modified, including consideration of content, format, administration, scoring, item properties, scaling, and technical quality for its intended purpose.

Test documents—Publications such as test manuals, technical manuals, user's guides, specimen sets, and directions for test administrators and scorers that provide information for evaluating the appropriateness and technical adequacy of a test for its intended purpose.

Test manual—A publication prepared by test developers and publishers to provide information on test administration, scoring, and interpretation and to provide technical data on test characteristics.

Validation—The process through which the validity of the proposed interpretation of test scores is evaluated.

Validity—The degree to which accumulated evidence and theory support specific interpretations of test scores entailed by proposed uses of a test.

Validity Evidence—Systematic documentation which empirically demonstrates, under the specific conditions of the individual analysis, to which extent, for whom, and in which situations test score inferences are valid. No single piece of evidence is sufficient to document validity of test scores; rather, aspects of validity evidence must be accumulated to support specific interpretations of scores.

Validity Evidence for Relevant Subgroups—In order to support that proposed interpretations of test scores are valid for subgroups who take the test, separate validity evidence is collected for subgroups when a prior probability suggests that interpretations may differ. For instance, if a test will be used to predict future performance, validity evidence should document that the scores are as valid a predictor of the intended performance for one subgroup as for another.

Validity argument—An explicit scientific justification of the degree to which accumulated evidence and theory supports the proposed interpretation(s) of test scores.

APPENDIX C: Accommodations Used by States

Table 1
Accommodations for Limited English Proficient Students

PRESENTATION FORMAT

Translation of directions into native language
Translation of test into native language
Bilingual version of test (English and native language)
Further explanation of directions
Plain language editing
Use of word lists/ dictionaries
Bilingual dictionary
Large print

ADMINISTRATION FORMAT

Oral reading in English
Oral reading in native language
Person familiar to students administers test
Clarification of directions
Use of technology
Alone, in study carrel
Separate room
With small group
Extended testing time
More breaks
Extending sessions over multiple days

RESPONSE FORMAT

Allow student to respond in writing in native language
Allow student to orally respond in native language
Allow student to orally respond in English
Use of technology

OTHER

Out-of-level testing
Alternate scoring of writing test

Adapted from: Council of Chief State School Officers, *Annual Survey: State Student Assessment Programs*,
Washington D.C., 1999

Table 2
Accommodations for Students with Disabilities

PRESENTATION FORMAT

Braille edition
Large-print editions
Templates to reduce visual field
Short segment testing booklets
Key words highlighted in directions
Reordering of items
Use of spell checker
Use of word lists/dictionaries
Translated into sign language

ADMINISTRATION FORMAT

Oral reading of questions
Use of magnifying glass
Explanation of directions
Audiotape directions or test items
Repeating of directions
Interpretation of directions
Videotape in American Sign Language
Interpreter signs test in front of classroom/student
Signing of directions
Amplification equipment
Enhanced lighting
Special acoustics
Alone in study carrel
Individual administration
In small groups
At home with appropriate supervision
In special education classes separate room
Off campus
Interpreter with teacher facing student; student in front of classroom
Adaptive furniture
Use place marker
Hearing aids
Student wears noise buffers
Administrator faces student
Specialized table
Auditory trainers
Read questions aloud to self
Colored transparency
Assist student in tracking by placing students finger on item

Typewriter device to screen out sounds
Extended testing time
More breaks
Extending sessions over multiple days
Altered time of day that test is administered

RESPONSE FORMAT

Mark responses in booklet
Use template for recording
Point to response
Lined paper
Use sign language
Use typewriter/computer/ word processor
Use Braille writer
Oral response, use of scribe
Alternative response methods, use of scribe
Answers recorded on audiotape
Administrator checks to ensure that student is placing responses in correct area
Lined paper for large script printing
Communication board

OTHER

Out-of level testing

Adapted from: Council of Chief State School Officers, *Annual Survey: State Student Assessment Programs*,
Washington D.C., 1999

APPENDIX D: Compendium of Federal Statutes and Regulations

This compendium provides a description of the federal nondiscrimination statutes and regulations that are relevant to testing issues and constitute the primary sources of legal authority in the guide. Specifically, this appendix primarily provides information on federal civil rights laws, including Title VI, Title IX, Section 504, and Title II of the Americans with Disabilities Act.

A. Title VI and Title IX

Title VI of the Civil Rights Act of 1964, 42 U.S.C. 2000d, prohibits race and national origin discrimination in programs and activities that receive Federal financial assistance. Title IX of the Education Amendments of 1972, 20 U.S.C. 1681 *et seq.*, prohibits sex discrimination in education programs that receive Federal financial assistance. For the regulations issued by the Department of Education implementing these statutes, see 34 C.F.R. Part 100 (Title VI) and 34 C.F.R. Part 106 (Title IX). Under the Civil Rights Restoration Act of 1987, OCR generally has institution-wide jurisdiction over the recipient of Federal funds. *See* 42 U.S.C. § 2000d-4 (1989).

The Title VI and Title IX statutes bar only intentionally discriminatory conduct. However, the regulations promulgated under these statutes prohibit the use of neutral criteria having disparate effects unless the criteria are educationally justified. *Guardian's Association v. Civil Service Commission*, 463 U.S. 582 (1983).

The regulations implementing Title VI do not specifically address the use of tests and assessment procedures, but bar discrimination based on race, color or national origin in any service, financial aid or other benefit provided by the recipient. 34 C.F.R. 104.3(b)(2), which prohibits criteria or methods of administration having an unjustified discriminatory effect, is often applied in testing cases.

In addition to general prohibitions against discrimination, the regulations implementing Title IX specifically prohibit the discriminatory use of tests or assessment procedures in admissions, 34 C.F.R. § 106.21, employment, 34 C.F.R. § 106.52, and counseling 34 C.F.R. § 106.36.

See also 34 C.F.R. § 100, Appendix B, part K (Guidelines for Eliminating Discrimination and Denial of Services on the Basis of Race, Color, National Origin, Sex, and Handicap in Vocational Education Programs) ("if a recipient can demonstrate that criteria [that disproportionately exclude persons of a particular race, color, national origin, sex, or disability] have been validated as essential to participation in a given program and that alternative equally valid criteria that do not have such a disproportionate adverse effect are unavailable, the criteria will be judged nondiscriminatory. Examples of admission

criteria that must meet this test or assessment procedure are ... interest inventories ... and standardized test or assessment procedures").

B. Section 504 of the Rehabilitation Act of 1973

Section 504 prohibits discrimination based on disability in programs and activities receiving federal financial assistance. OCR enforces Section 504 and its regulations in education programs. The regulations implementing Section 504 contain certain sections that are particularly relevant to testing situations:

34 C.F.R. 104.4(b)(4) prohibits criteria or methods of administration that have the effect of discriminating against qualified persons with disabilities.

34 C.F.R. 104.42(b)(2) prohibits admissions procedures by higher educational institutions that make use of any test or criterion for admission that has a disproportionate, adverse impact on qualified individuals with disabilities unless (1) the test or criterion, as used by the institution, has been validated as a predictor of success in the education program or activity and (2) alternate tests or criteria that have a less disproportionate, adverse impact are not shown to be available. 34 C.F.R. 104.42(b)(3) requires admissions tests used by post-secondary institutions to be selected and administered so as best to ensure that, when a test is administered to an applicant with a disability, the test results accurately reflect the applicant's aptitude or achievement, rather than reflecting the student's disability (except where disability-related skills are the factors the test purports to measure). 34 C.F.R. 104.44(a) and (d) require higher education institutions to provide adjustments or accommodations and auxiliary aids and services that enable the student to demonstrate the knowledge and skills being tested.

34 C.F.R. 104.44(a) states that academic requirements that the institution can demonstrate are essential to the program of instruction or to any directly related licensing requirement will not be regarded as discriminatory.

34 C.F.R. 104.35 (b) requires public elementary and secondary education programs to individually evaluate a student before classifying the student as having a disability or placing the student in a special education program; tests used for this purpose must be selected and administered so as best to ensure that the test results accurately reflect the student's aptitude or achievement or other factor being measured rather than reflecting the student's disability, except where those are the factors being measured. These provisions also require that tests and other evaluation materials include those tailored to evaluate the specific areas of educational need and not merely those designed to provide a single intelligence quotient.

C. Title II of the Americans with Disabilities Act (ADA)

Title II of the Americans with Disabilities Act of 1990 (ADA), 42 U.S.C. §12134, prohibits discrimination on the basis of disability by public entities. Regulations implementing Title II, issued by the U.S. Department of Justice, can be found at 28

C.F.R. Part 35. OCR enforces Title II as to public schools and colleges. Like the Section 504 regulations, the regulations implementing Title II prohibit "criteria and methods of administration which have the effect of discriminating" against qualified persons with disabilities. 28 C.F.R. 35.130(b)(3). The regulations also require public entities to make reasonable accommodations to policies, procedures, and practices when the modifications are necessary to avoid discrimination unless the public entity can demonstrate that the modification would fundamentally alter the nature of the service, program, or activity. 28 C.F.R. 35.130(b)(7).

D. Individuals with Disabilities Education Act (IDEA)

Although not a discrimination law per se, IDEA contains important provisions related to testing students with disabilities in elementary and secondary schools. IDEA is enforced by the Office of Special Education Programs in the U.S. Department of Education. As amended in 1997, IDEA requires inclusion of students with disabilities in state and district-wide assessment programs, with appropriate accommodations, if necessary, unless the student's individual education team decides that participation in all or part of the testing program is not appropriate. The student's individualized education program (IEP) should also state any individual modifications in the administration of State or district-wide assessments of student achievement that are needed in order for the student to participate in such assessment. If the IEP team determines that the student will not participate in a particular State or district-wide assessment of student achievement (or part of such an assessment), the student's IEP must include statements of why that assessment is not appropriate for the student and how the student will be assessed. IDEA also requires state or local educational agencies to develop guidelines for the alternate assessment of the relatively small number of students with disabilities who cannot take part in state and district-wide tests to participate in alternate assessments. These alternate assessments must be developed and conducted not later than July 1, 2000. See 20 U.S.C. 1412(a) (16) and (17), 1413 (a)(6), and 1414(d)(1)(A) and (d)(6)(A)(ii), and regulations at 34 C.F.R. 300.138, 300.139, 300.240, and 300.347.

APPENDIX E: Resources and References

Office for Civil Rights U.S. Department of Education

Minority Students and Special Education: Legal Approaches for Investigation, 1995.
Provides an overview of the legal theories and approaches employed in OCR investigations examining disproportionate representation of minority students in special education.

Policy Update On Schools' Obligations Toward National-Origin-Minority Students With Limited-English Proficiency, 1991.

Used by OCR staff to determine schools' compliance with their Title VI obligation to provide any alternative language programs necessary to ensure that national-origin-minority students with limited English proficiency have meaningful access to programs. Provides additional guidance for the December 1985 and May 1970 memoranda.

The Office for Civil Rights' Title VI Language-Minority Compliance Procedures, 1985.
Focuses on the treatment of limited English proficient students in programs that received funds from the Department.

Identification of Discrimination and Denial of Services on the Basis of National Origin, May 1970, 35 Fed. Reg. 11595.

Clarifies school district responsibilities to limited English proficient students. Memo was the foundation for the U.S. Supreme Court decision *Lau v. Nichols* and was affirmed in that decision.

Office of Elementary and Secondary Education U.S. Department of Education

Peer Reviewer Guidance for Evaluating Evidence of final assessments Under Title 1 of the Elementary and Secondary Education Act (ESEA), 1999.

Informs the states about types of evidence that would be useful in determining the evaluation of assessments under Title 1.

Taking Responsibility for Ending Social Promotion, 1999.

Provides strategies for preventing academic failure and give information about how these strategies can be sustained through ongoing support for improvement.

Handbook for the Development of Performance Standards: Meeting the Requirements of Title 1 (with Chief State School Officers, 1998).

Describes the best practices and current research on the development of academic performance standards for K-12.

Standards, Assessments and Accountability, 1997.

Overview of the major provisions under Title 1 of the Elementary and Secondary Education Act.

**National Research Council
National Academy Press, Washington D.C.**

Heubert, Jay P. and Hauser, Robert M., ed., *High Stakes: Testing for Tracking, Promotion and Graduation*, 1999.

Discusses how tests should be planned, designed, implemented, reported and used for a variety of educational policy goals. Focuses on the uses of tests that make high-stake decisions about individuals and on how to ensure appropriate test use.

Beatty, Alexandra; Greenwood, M. R. C. and Linn, Robert L., ed., *Myths and Tradeoffs: The Role of Tests in Undergraduate Admissions*, 1999.

Four recommendations regarding test use for admission are made to colleges and universities, including a warning to schools to avoid using scores as more precise and accurate measures of college readiness than they are. One recommendation is made to test producers, which is to make clear the limitations of the information that the scores provide.

Elmore, Richard F. and Rothman, Robert, ed., *Testing, Teaching and Learning: A Guide for States and School Districts*, 1999.

Practical guide to assist states and school districts in developing challenging standards for student performance and assessment as specified by Title I. Discusses standards-based reform and specifies components of an education improvement system, which are standards, assessments, accountability and monitoring the conditions of instruction.

August, Diane and Hakuta, Kenji, ed., *Improving America's Schooling for Language Minority Children: A Research Agenda*, 1997.

Summarization of extensive study of limited English proficient students. Gives state of knowledge review and identifies research agenda for future study. Includes discussion of student assessment and program evaluation.

Morison, Patricia; White, S.H. and Feuer, Michael J., ed., *The Use of I.Q. Tests in Special Education Decision-Making and Planning: Summary of Two Workshops*, 1996.

Report provides a synthesis of the key themes and ideas discussed at workshops, including: an overview of legal, policy and measurement issues in use of IQ tests in special education; validity and fairness of IQ testing for student classification and placement; alternative assessment methods used in combination with or as substitutes for IQ tests.

McDonnell, Lorraine M.; McLaughlin, Margaret J. and Morison, Patricia, ed. *Educating One & All: Students with Disabilities and Standards-Based Reform*, 1997.

Twelve recommendations are given regarding how to integrate students with disabilities in standards-based reform, including: participation of students with disabilities should be maximized; that any test alterations must be individualized and have a compelling educational justification; include these students' test results in any accountability system; ensure opportunity for students with disabilities to learn the material tested; and use the IEP process for decision-making on the participation of individual students.

Recommendations for policy-makers include: revising policies that discourage the inclusion of students with disabilities in high-stake tests; giving parents enough information to make informed choices about participation; monitoring possible unanticipated consequences of participation, both for standardized testing and for students with disabilities; designing realistic standards; and designing a long-term research agenda.

Hyde, Lorraine D.; Robertson, Gary J. and Krug, Samuel E., *et al.*, *Responsible Test Use: Case Studies for Assessing Human Behavior*, 1993.

Casebook for professionals using educational and psychological test data, which was developed to apply principles to proper test interpretation and actual test use. Cases are organized under eight sections: general training, professional responsibility training, test selection, test administration, test scoring and norms, test interpretation, reporting to clients and administrative or organization policy issues.

Test Measurement Standards

Joint Committee on the Standards for Educational and Psychological Testing, *Standards of Educational and Psychological Testing*, 1999.

Provides criteria for the evaluation of tests, testing practices, and the effects of test use. Begins with discussion of the test development process, which focuses on test developers, and moves to specific test uses and applications, which focus on test users. One chapter centers on test takers.

National Council on Measurement in Education, *Code of Professional Responsibilities in Educational Measurement*, 1995.

Association for Measurement and Evaluation in Counseling and Development,
Responsibilities of Users of Standardized Tests, 1992.

Joint Committee on Testing Practices, *Code of Fair Testing Practices in Education*,
American Psychological Association, Washington, D.C., 1988.

Measurement Texts

Linn, Robert L., ed., *Educational Measurement*, 3rd edition, American Council on
Education, New York: Macmillan Publishing Company, 1989.

Includes 11 chapters, including Messick's classic chapter on validity, and organizes them
in two parts: theory and general principles; and construction, administration and scoring.

Messick, Samuel, Validity of psychological assessment: Validation of inferences from
persons' responses and performances as scientific inquiry into score meaning,
September 1995, *American Psychologist*. Gives a new cohesive definition of validity that
looks at score meaning and social values. Six perspectives of construct validity are
defined: content, substantive, structural, generalizability, external and consequential.

Thurlow, Martha; Elliott, Judy and Ysseldyke, Jim, *Testing Students With Disabilities*,
Thousand Oaks, CA: Corwin Press, 1998.

This document provides guidance about how students with disabilities should be included
in large scale tests, considerations about how to select the appropriate accommodations
for which students, and discussions about the role of state and local educators in ensuring
proper test use, the use of alternate tests, and appropriate reporting considerations.

Kopriva, Rebecca J., *Ensuring Accuracy in Testing for English Language Learners*,
Washington, D.C.: Council of Chief State School Officers, 2000.

This resource provides guidance to states, districts, and test publishers about developing,
selecting, or adapting large-scale, standardized assessments of educational achievement
that are appropriate and valid for English language learners. The guide's practical
recommendations identify the "who, what, when, why and how" associated with
developing, selecting, or adapting tests for institution use, including how to select the
appropriate accommodations for which students, how to collect appropriate validity
evidence, and a discussion of salient reporting considerations.

Test Publisher Materials

Most test publishers produce materials that explain the appropriate use of their tests. We
encourage interested readers to obtain these materials from the publishers of the tests they
administer or from publishers of tests in which they are interested. Readers can also
contact the Association of Test Publishers, 655 15th St. NW, Washington, D.C., 20005,
telephone 202-857-8444 for more information.

Other Resources

There are many books and other materials that might be helpful to educators and policymakers as they develop policies, and design and implement programs which include the use of tests in making high-stakes decisions for students. The following web sites will provide additional information and links to some of these resources.

Council for Chief State School Officers

<http://www.CCSSO.org>

The National Center on Education Outcomes

<http://www.coled.umn.edu/NCEO>

Center for Evaluation, Research, Standards and Student Testing

<http://cresst96.cse.ucla.edu>

National Clearinghouse for Bilingual Education

<http://www.ncbe.gwu.edu>