

DOCUMENT RESUME

ED 279 680

TM 870 067

AUTHOR Haney, Walt; Madaus, George
TITLE Effects of Standardized Testing and the Future of the National Assessment of Educational Progress. Working Paper for the NAEP Study Group.

PUB DATE Sep 86
NOTE 37p.; One of 46 papers commissioned by the Study Group on the National Assessment of Student Achievement and cited in Appendix B to their final report "The Nation's Report Card" (TM 870 049). For other papers in this group, see TM 870 050-094.

PUB TYPE Viewpoints (120)

EDRS PRICE MF01/PC02 Plus Postage.
DESCRIPTORS *Achievement Tests; *Educational Assessment; Educational Policy; Educational Testing; *Educational Trends; Elementary Secondary Education; Evaluation Utilization; *Futures (of Society); Higher Education; National Surveys; Standardized Tests; *Testing Problems; Testing Programs; *Test Use

IDENTIFIERS *National Assessment of Educational Progress

ABSTRACT

Topics related to the National Assessment of Educational Progress (NAEP) are discussed. Issues affecting the impact of testing on education include: what is tested; how scores are referenced; internal versus external sources of testing; and the rewards, sanctions, or stakes associated with test results. Testing has both intentional (direct) and indirect effects. Seven principles regarding the impact of testing are especially true when the tests are used for important social decisions: (1) measuring any social indicator results in its distortion or change in status; (2) the power of testing is determined by the perception of its effects; (3) teachers will teach to the test; (4) previous tests influence curriculum content; (5) instruction and learning are adjusted to test format; (6) test results become a major goal of education; and (7) control of the curriculum is transferred to the testing agency. NAEP can use test information: (1) to inform policy makers about the current state of education, or (2) as administrative devices in policy implementation. NAEP is the most valid source of national achievement data and has been designed to inform educational policy making, rather than to direct it. Increasing its policy-making ability might have distorting effects or threaten test validity. It has been, and will continue to be useful as a model for other testing activities. (A list of ERIC references concerning assessment is included.) (GDC)

* Reproductions supplied by EDRS are the best that can be made *
* from the original document. *

ED279680

Effects of Standardized Testing and the Future of the National Assessment of Educational Progress

Working Paper for the NAEP Study Group

by

Walt Haney and George Madaus
Center for the Study of Testing, Evaluation and Educational Policy
Boston College
Chestnut Hill, MA 02167

ph 617/552-4521

September 1986

"PERMISSION TO REPRODUCE THIS
MATERIAL HAS BEEN GRANTED BY

W. Haney

TO THE EDUCATIONAL RESOURCES
INFORMATION CENTER (ERIC)."

U.S. DEPARTMENT OF EDUCATION
Office of Educational Research and Improvement
EDUCATIONAL RESOURCES INFORMATION
CENTER (ERIC)

This document has been reproduced as
received from the person or organization
originating it.
 Minor changes have been made to improve
reproduction quality.

• Points of view or opinions stated in this docu-
ment do not necessarily represent official
OERI position or policy.

Paper commissioned by
THE STUDY GROUP ON THE NATIONAL ASSESSMENT OF STUDENT ACHIEVEMENT
1986

TM 870 067



I. Introduction

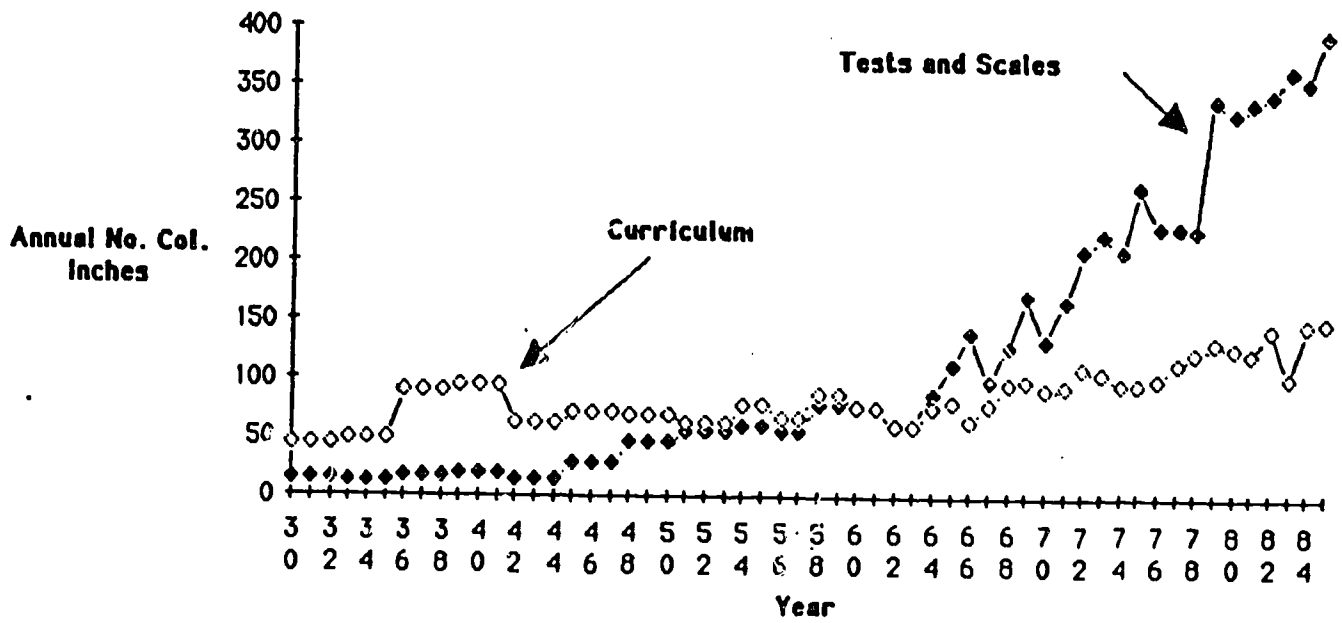
Since tests have a powerful directive influence on teaching and the study of pupils, a major policy to follow is to establish a testing program that faithfully reflects the objectives sought by the school. In this way the influence of testing is to reinforce the objectives sought by the school.

Ralph Tyler, "What Testing Does to Teachers and Students," 1959, In Anastasi, 1966, p. 49.

Standardized testing in the schools is on the increase. As Phipo (1985, p.19) observed "nearly every large education reform effort of the past few years has either mandated a new form of testing or expanded uses of existing tests." The increasing prominence of testing over the last five years is linked directly to efforts to reform education, particularly at the state level. For example, a 50-state survey of reform measures conducted by Education Week found that: 29 states required competency tests for students, and 10 other states had such a requirement under consideration; 15 states required an exit test for graduation, 4 additional states had such a measure under consideration; 8 states employed a promotional "gates" test, while 3 others were considering such a mandate; finally, 37 states had some sort of state assessment program, and 6 additional states had such a program under consideration. This growing use of tests in the policy sphere by agencies external to the local education agencies (LEAs), we will argue, is having increasing impact on what is taught and learned in schools.

As a means of documenting this increasing attention to testing, and contrasting it with curriculum concerns, we charted the amount of space in Education Index over the last 50 years devoted to citations concerning testing and curriculum. As shown in Figure 1, the average annual number column inches devoted to citations concerning curriculum has increased only modestly over the last half-century -- from 50 - 100 inches per year in the

Figure 1: Education Index Listings Under Testing and Curriculum 1930 -1985



Haney and Medaus, paper for NAEP Study Group, 1986, p. 2.

1930s and 40s to only 100 - 150 in recent years. In contrast, attention devoted to testing in Education Index has increased dramatically, from only 10 - 30 column inches in the 1930s and 1940s to well over 300 inches in the 1980s.

Measuring column inches in Education Index is, of course, a fairly crude way of charting what is happening in the world of education, but these data certainly suggest a trend that we suspected even before looking at the Education Index. It is that increasingly, standardized testing seems to have become the coin of the educational realm. In recent years, it seems that the aims of education and the business of our schools are addressed not so much in terms of curriculum - the courses of study that are or should be followed - - as in terms of what gets tested. The data from the Education Index showing that the relative attention to curriculum and testing issues has undergone a ten-fold change in the last 50 years, clearly suggests this.

Before reviewing arguments and evidence bearing on the impact of testing, we need to comment briefly on what is meant by the term "standardized tests." Essentially by this term we refer to standardized achievement tests of reading and math skills, including state- and local-education-agency sponsored standardized tests, commercially developed nationally normed tests, and tests that have sometimes been called basic skills or minimum competency tests. Under the rubric of "standardized achievement tests" one might even include the Scholastic Aptitude Test (SAT) which is, we think, essentially a general achievement test of verbal and math skills rather than a test of aptitude (though we will not in this brief presentation attempt to delve into the controversy over the aptitude/achievement distinction, either in general or with respect to the SAT).

By way of introduction it should be noted that such tests may be used for a variety of purposes, including general evaluation of schools and programs, diagnosis of student strengths and weaknesses, student grade promotion, high school graduation, guidance, college admissions and even teacher evaluation. However, despite the initial purposes of either test developers or sponsors for achievement testing, once a test is administered the results often get used for other quite different purposes. Perhaps the most prominent example of this is the way in which college admissions tests, despite occasional protestations by test sponsors, have come to be used as general indicators of the educational health of states, and the nation as a whole (in this regard see, for example, the United States Office of Education's annual Wall Chart, or the recent study, Trends in Achievement, by the Congressional Budget Office of the U.S. Congress, 1986. Both use SAT and ACT test data as two sources of evidence regarding national trends in achievement).

Thus, in this brief paper we will not attempt to distinguish between the different, ostensible, commonly recognized, purposes of standardized tests. Instead we will attempt to set out some fairly broad ideas about the effects of testing on education, and on curriculum in particular; offering some ideas on the past and future of the National Assessment of Educational Progress (NAEP). Specifically the following four sections of this paper are organized around these topics:

II. Conditions of Testing Affecting Its Impact

III. Seven Principles Regarding the Impact of Testing

IV. NAEP as an Instrument for Informing Educational Policy; and

V. The Future of NAEP

II. Conditions of Testing Affecting Its Impact

Rather than commenting on the myriad specific ostensible uses of testing, instead, in this section, we offer some general observations on four broad issues that condition the impact of testing: (1) What is tested, (2) How scores are referenced; (3) The source of testing; and (4) The rewards or sanctions associated with test results

2.1 What is tested

A wide range of variables has been the subject of measurement, though the main emphasis has been on the measurement of cognitive rather than affective characteristics. The cognitive variables which have attracted most attention are "intelligence", and achievement in basic skill areas of the curriculum (reading, arithmetic). As one proceeds up the educational ladder into secondary schools, where instruction is organized around subject matter/content areas, rather than around specific skills, commercially available test batteries become less specific and less related to what is taught. High school test batteries closely resemble elementary school batteries in that they are more oriented to the basic skills of numeracy and literacy than to what is taught in specific subject fields like math, physics, history, English literature, etc. As a result, such test scores are less relevant to the work of the high school teacher. However, some of the recent reform reports call for the development of exams for specific secondary school curricula areas. This has profound implications for curriculum and instruction at that level. This is because in general, it seems that other

things being equal, the impact of testing is greater when tests are keyed to specific courses.

2. How scores are referenced.

The main distinction here is between norm-referenced tests, on which performance is assessed by reference to the performance of other students, and criterion-referenced tests, on which performance is assessed by reference to the mastery of specific content domains. It should be noted that norm-referenced information can also be structured to provide criterion-referenced interpretations and vice versa. While criterion-referenced tests are increasingly hailed as superior to norm-referenced ones in terms of information provided to teachers, norm-referenced information is valuable for comparative purposes. Further, the specificity of criterion-referenced information from commercially available tests, relative to what is actually taught at the local level can often be dubious.

The general point to be made here is that however tests are referenced, if they are poorly matched to what is taught in schools, and if they are linked to important decisions, they can have great impact, for good or ill.

3. Internal vs. external testing programs.

An important factor regarding the impact of testing is its source. The main distinction here is between internal testing, and external testing. An internal testing program is one which is carried out within a school at the initiative and under the control of the school superintendent, principal or teacher. In this category we include the traditional norm-referenced standardized achievement testing programs that have been used by school

Honey and Madans, paper:
systems since the 19
available criterion re
limited to off-the-shel
made tests which we
be built by school sys
system.

External testing
external authority, su
state legislature or a
Examination Board (CI
Aptitude Test (SAT).
a student, teacher or s
requirements. State n
promotion or remedial
Leaving Certificate Exa
all examples of the lat
criterion-referenced te
instruments, tests built
a contractor.

External tests ha
Until the advent of mir
the United States had t
programs were relative
examinations being not
testing is being mandai
of their efforts at educa

ly Group, 1986, p. 6.

as the use of the newer commercially
movement tests. Internal testing is not
1 tests; it also includes traditional teacher
it in this paper. In larger systems tests can
nel or customized by a contractor for the

ire those controlled and/or mandated by an
education education agency (SEA), the
cy such as the College Entrance
be a semi-voluntary test like the Scholastic
only, it is a test which the state mandates if
it are to fulfill certain mandated
npetency exams linked to graduation,
is, the English 'O' and 'A' Levels, the Irish
and the New York State Regents Exams are
1 testing programs can use norm- or
the-shelf commercially available
, or those developed for such a sponsor by

n a part of European educational systems.
 competency testing, most educational testing in
1 in origin and control. External testing
New York Regents Exams and the CEEB
lons. However, as we have noted, external
igly by state boards or legislatures as part
m.

There are two reasons the distinction between internal and external testing programs is important. First, external tests tend to have a greater mismatch with what is taught at the local level than do internal tests; for the simple reason that curriculum and teaching methods are determined, either explicitly or implicitly at the "grass-roots level" -- the LEA, the school and ultimately behind the classroom door. Second, external testing programs tend to have clearer consequences associated with them than do locally initiated and controlled tests. And as we argue in the next section the size of the stakes associated with testing programs is a key determinant of the educational impact.

4. High Stakes vs Low Stakes Testing Programs

A test whose results are seen - rightly or wrongly - by students teachers, administrators, parents, or the general public, to make important decisions that immediately and directly impact on them are what we shall term high-stakes tests. High-stakes student tests can be norm- or criterion-referenced, internal or external in origin. Examples include tests directly linked to such important decisions as: (1) graduation, promotion or placement of students; (2) the evaluation or rewarding of teachers or administrators; (3) the allocation of resources to schools or school districts; and (4) school or school system certification. In all of these examples, the perception of people that test results are linked to a high-stakes decision is in fact accurate. Policy makers have mandated that the results be used automatically to make such decisions.

However, there are other uses of test results that do not actually impact immediately and directly on students but nonetheless, are generally perceived by people as involving high-stakes. For example, SAT and ACT

Haney and Madaus, paper for NAEP Study Group, 1986, p. 8.

results are of secondary importance in admission decisions in the vast bulk of colleges trying to fill vacant seats in the face of adverse demographics. Individuals and school systems, nevertheless, very often act on the perception that these college admissions tests are of crucial and singular importance, if not to individual students, at least to the public's perception of the quality of schooling. Thus, we find high schools are increasingly offering courses to prepare students to take these tests, and commercial coaching schools are doing a land office business.

In contrast to a high stakes test, a low stakes test is one which is perceived as having no important rewards or sanctions tied directly to test performance. Traditional school district standardized norm referenced testing programs, where results are reported to teachers, but there is no immediate, automatic decision linked to performance, are examples of this sort of testing. Teachers are free to ignore any results that they feel are discrepant from their own perceptions of students, and the results are not perceived by them as being used to evaluate their performance. This does not mean that test results from such programs do not affect teachers' perceptions of students, nor does it mean that student placement decisions are not related to test performance. The important distinction is that teachers, students, and parents do not perceive test performance as a direct vehicle of reward or sanction.

In short, it is clear that "high stakes" testing has the greatest impact on schooling, regardless of whether the stakes are associated with specific decisions made on the basis of the results, or with the perceived importance of the tests.

Arguments Concerning the Effects of Standardized Testing

If these are conditions which can affect the degree of impact of testing, what of the nature of the impact? Is it for good or ill? Over the last decade a wide variety of arguments have been advanced concerning the positive and negative effects of standardized testing. These arguments have been advanced in a variety of forums, including professional journals and meetings, the popular press, legislative debates, meetings of school governance agencies, and in both federal and state courts.

Without trying to document such myriad sources, it is useful to consider the nature of the arguments advanced regarding the positive and negative effects of standardized achievement testing. Among the most common affirmative arguments have been that such testing:

- helps focus instruction on skills (e.g. basic skills);
- motivates students;
- provides teachers with diagnostic information to improve instruction;
- identifies curriculum areas in need of improvement; and
- helps hold teachers and schools accountable for the learning of their charges.

Among the most common charges levelled against standardized achievement tests are that they:

- Are biased against certain kinds of students;
- Do not match what students have been taught;
- Constrain teachers' initiative and creativity in teaching;
- Promote teaching to the tests;
- Narrow the curriculum.

Haney and Madans, paper for NAEAP Study Group

There are of course many variations in arguments and their relative merits from program to program.

One point worth noting about the literature suggests that standardized tests work at different levels: by motivating individual instructional efforts of teachers, and by guiding improvement efforts at school. This is important, we think, because discussions focus on the instrumental role of tests to help make specific decisions about student arguments just cited suggest that indirect and we suspect, especially over time, rather than direct effects related to specific communication to students, teachers, and parents. In large, what is valued in an educational system is hardly surprising that with the increasing emphasis on the skill of simply taking tests is emerging independent of what is ostensibly tested. One of the test taking skills is that the phrase "test taking skills" is a separate indexing category in volume 1 of the July 1982 - June 1983 period; in the number of items under this rubric has

III. PRINCIPLES ABOUT THE II

Perhaps the most convenient work in the literature on the impact of testing on

Sp, 1986, p. 10.

1 such arguments. And different
subtleties vary from place to place and

ch arguments, however, is that they
their effects, for good or ill, on several
ial students, by informing the
r providing information which may
ate or even national levels. This point
ions of testing programs often seem to
g in which results may be used to
lents, programs or schools. Yet the
ct effects of testing are important,
ay have a larger cumulative impact
cisions. What gets tested, after all,
ministrators and the community at
endeavor in myriad ways. Thus, it is
r in testing over the last decade or so,
ing as a focus of concern, almost
d. A minor sign of the concern over
t taking skills" first emerged as a
3 of the Education Index, covering
e subsequent two volumes the
ontinued to grow.

EFFECT OF TESTS

r to summarize the considerable
ucation is to list major conclusions

from our review of this literature in the form of seven principles regarding the impact of testing.

Principle 1:

The more any quantitative social indicator is used for social decision-making, the more subject it will be to distort and corrupt the social processes it is intended to monitor.

This Principle comes directly from Donald Campbell's work on social indicators. It is not limited to testing per se; being much more general in scope, extending to any social indicator that is used to describe, make decisions about, or influence an important social process. This Principle reminds us that educational testing is a form of measurement subject to a social version of Heisenberg's uncertainty principle. Any measurement of the status an educational institution, no matter how well contrived, inevitably changes its status. And when the testing is used for important social decisions, the change tends to be both large and corrupting.

Principle 2

The power of tests and examinations to affect individuals, institutions, curriculum or instruction is a perceptual phenomenon: if students, teachers, or administrators believe that the results of an examination are important, it matters very little whether this is really true or false -- the effect is produced by what individuals perceive to be the case.

Ben Bloom writing in the 63th NSSE Yearbook coined this second Principle. Its importance lies in the fact that when people perceive a phenomeon to be true, their actions are guided by the importance perceived to be associated with it. The greater the stakes perceived to be linked to test results the greater the impact on instruction and learning. A high-stakes test is one where the results are seen by students, teachers, administrators, parents, or the general public, to be used to make important decisions that

immediately and directly impact on individuals with whom they are concerned. Examples of high-stakes student tests include those perceived to be directly linked to such important decisions as: (1) graduation, promotion or placement of students; (2) the evaluation or rewarding of teachers or administrators; (3) the allocation of resources to schools or school districts; and (4) school or school system certification.

PRINCIPLE 3

If important decisions are presumed to be related to test results, then teachers will teach to the test.

This sort of accommodation by teachers to a high-stakes test is seen as having both positive and negative consequences. High-stakes tests, it is often argued, can focus instruction, giving students and teachers specific goals to attain. If the test is measuring basic skills, preparing students for the skills measured by the test could, proponents argue, serve as a powerful lever to improve basic skills. Unfortunately, however, the only sort of evidence generally available to bolster this proposition is that scores on high stakes tests do tend to increase over time. But standardized tests are indirect measures of the real skills of interest, and what repeated experience shows is that there are many, many ways of raising test scores without changing the levels of the skills the tests are intended to measure. People too often fail to distinguish between the skill and the indirect indicator of it.

If the test is specific to a specialized curriculum area, e.g. college preparatory physics, then the examination will eventually narrow instruction and learning, focusing only on those things measured by the tests. Indeed, this narrowing of the curriculum has been one of the enduring complaints leveled at external certification examinations used for the

important functions of certifying the successful completion of elementary or secondary education, and admission to third level education and/or certain jobs.

A review of the effects of such exams on the curriculum over many years, and in several countries, indicates that faced with a choice between objectives which are explicit in the curriculum or course outline, and a different set which are implicit in the certifying examination, students and teachers generally choose to focus on the latter. Spaulding's 1938 report pointed out that teachers in New York disregarded the objectives in local curriculum guides in favor of those tested in the Regents' examinations. Morris found the rigidity of the exams was the principal reason that the chemistry curriculum in Australia remained almost unchanged from 1891 to 1959. He concluded that the the proportion of instructional time spent on various aspects of the syllabus was "seldom higher than the predictive likelihood of its occurrence on the examination paper." Similar observations about the influence of the exams on the curriculum have been made in India, Japan, Ireland, and in England. Turner sums up the English experience: "One only has to look at the timetable of the typical comprehensive school to see that the curriculum consists almost entirely of subjects which can be taken in public examinations."

Why does this happen? First, there is tremendous social pressure on teachers to see to it that their students acquit themselves well on the certifying examinations. Second, the results of the examination are so important to students, teachers, and parents that their own self interest dictates that instructional time focus on test preparation.

As we have indicated, however, a high-stakes test can serve as a lever for the introduction of new curricular material. New curricula in physics, chemistry, and mathematics made an immediate impact in Irish schools when in the early 1960's they were prescribed and examined for the Leaving Certificate Examination. Primary teachers in Belgium accepted curricular reforms only when, in 1936, the external exams given at the end of primary school were modified to incorporate the ideas of the new curriculum. In New York State, curriculum specialists from the State Department of Education had little success in moving the emphasis in modern language teaching from grammar and translation to conversation and reading skills, until the corresponding changes had been incorporated in the content of the Regents' examination. Revisions of the College Entrance Examination Board (CEEB) math achievement tests to include modern math played an important part in the introduction of "modern" mathematics curricula in the 1960's.

Despite the ability of high stakes examinations to help introduce new material, the weight of examination precedent strongly influences not just what is taught and learned but also how. The question that educators must ask themselves is whether the positive aspects associated with this phenomenon outweigh the disadvantages. The answer is a value judgment and depends on one's view of education, the learner, teaching, curriculum development and testing. Our view of the "driving" of instruction via external examinations is that in the long term, given the nature of most commonly used tests, the narrowing of instruction and learning associated with this phenomenon far outweigh any advantages.

Principle 4:

In every setting where a high-stakes test operates, a tradition of past exams has developed and it eventually comes to define the curriculum de facto.

Given Principle 3, the question still remains: "How do teachers cope with the pressure of the examination?" The answer is relatively simple. Teachers see the kind of intellectual activity required by previous test questions and prepare the students to meet these demands. Some have argued strongly that if the skills are well chosen, and if the tests truly measure them, then coaching is perfectly acceptable. This argument sounds reasonable, and in the short term, it may even work. But we need to take a longer view, because the argument ignores a fundamental fact of life: When the teacher's professional worth is estimated in terms of exam success, teachers have great incentive to produce gains; and they can do more simply by teaching test-taking strategies based on previous exam questions.

Further, the expectations and deep-seated primary agenda of students and their parents for exam success will further corrupt the process. The view that we can coach for the skills apart from the tradition of test questions is a staggeringly optimistic view of human nature which ignores the powerful pull of self-interest. It simply doesn't consider the long term effects of the examination sanctions.

An interesting aspect of Principle 4 is that if the examination is perceived as important enough, a commercial industry develops, outside the schools, to prepare students for it. In this country this phenomenon can be seen in the rise of commercial firms in virtually every major city selling coaching services to students for the SAT college admissions test. Another sign of the increasing prominence of test coaching in the United States is the fact that the phrase "test taking skills" first appeared as a separate indexing

Haney and Madans, paper for NAEP Study Group, 1986, p. 16.

category in volume 33 of the Education Index covering education literature for the July 1982 - June 1983 period. Most of the articles referenced under this new category dealt with improving admissions test scores through coaching provided by commercial or computer tutorial programs. In Japan it is common for parents to enroll their children in special extra-study schools known as *juku*. Beginning in the 19th century a whole industry of private coaching schools called "crammers" developed in Europe to prepare students, for a fee, for high stakes examinations. The important point about these coaching schools is not whether or not they are successful in preparing students for the exam, it is instead, that the public perceives them as helpful and willing to pay for their services.

Principle 5:

Teachers and students pay particular attention to the form of the questions on a high-stakes test, e.g. short answer, essay, multiple choice, and adjust their instruction accordingly.

The problem here is that the form of the test question can narrow instruction, study and learning to the detriment of other broader learning. Rentz recounts an example of this phenomenon which occurred as a result of the Georgia Regents' Testing Program, a program designed to assess minimum competencies in reading and writing on the part of college students in that state. The head of an English department lamented:

Because we now are devoting our best efforts to getting the largest number of students past the essay exam ..., we are teaching to the exam, with an entire course, English III, given over to developing one type of essay writing, the writing of a five-- paragraph argumentative essay written under a time limit on a topic about which the author may or may not have knowledge, ideas, or personal opinions. Teaching this one useful writing skill has the beneficial effect of bringing large numbers of weak students to a minimal level of literacy, but at the same time, it devastates the content of the composition program that

should be offering the better students challenges to produce writing of high quality. Because the Regents Test is primarily designed to establish a minimal level of literacy, our teaching to this test, which its importance forces us to do, tends to make the minimal acceptable competency the goal of our institution, a circumstance that guarantees mediocrity.

Principle 5 has profound implications for the curriculum specialists. Given our free enterprise system, publishers have begun to look at state mandated minimum competency or basic skills tests in order to design materials to better train pupils to take them. Children are apt therefore to find themselves spending more and more time filling out ditto answer sheets or work books. Deborah Meier, a successful principal of a public school in Manhattan, testified at the 1981 NIE sponsored hearings on Minimum Competency Testing (MCT) that in New York City, reading instruction has come to closely resemble the practice of taking reading tests. In reading class, students, using commercial materials, read dozens of little paragraphs about which they then answer multiple choice questions. Meier described the materials as evolving to resemble more and more the tests students will take in the spring. She went on to point out that when synonyms and antonyms were dropped from the New York City test of word meaning, teachers promptly dropped commercial material that stressed them. It is also interesting to note that in 1983, sales of ditto paper were way up nationally while sales of lined theme paper were down.

Principle 6:

When test results are the sole or even partial arbiter of future educational or life choices, society tends to treat test results as the major goal of schooling rather than as a useful but fallible indicator of achievement.

Of all of the effects attributed to tests, this may well be the most damaging to education. It is illustrated in the following observation from a

nineteenth century British school inspector who observed first hand the negative effects of linking teacher salaries in England and Ireland to pupil examination results:

Whenever the outward standard of reality (examination results) has established itself at the expense of the inward, the ease with which worth (or what passes for such) can be measured is ever tending to become in itself the chief, if not sole, measure of worth. And in proportion as we tend to value the results of education for their measurableness, so we tend to undervalue and at last to ignore those results which are too intrinsically valuable to be measured.

Sixty years later, Ralph Tyler echoed the same message in the 62nd NSSE Yearbook when he warned readers that society conspires to treat marks in certifying examinations as the major end of secondary schooling, rather than as a useful but not infallible indicator of student achievement.

We see the importance society places on test scores to the exclusion of other indicators in such things as: the media attention to declines in SAT scores; reports that our schools score lower than those of other countries in math and science; the Education Department's wall chart that ranks states by their performance on the SAT or ACT; newspapers ranking school districts and/or schools within districts by their performance on standardized tests; the use of test results by real estate agents in selling homes; the money spent by parents on coaching schools for the SATs; the list could go on and on.

Principle 7

A high-stakes test transfers control over the curriculum to the agency which sets or controls the exam.

The agency responsible for a high-stakes test assumes a great deal of power or control over what is taught, how it is taught, what is learned and how it is learned. This phenomenon is well understood in Europe where a

system of external certification examinations, controlled by the central government or by independent examination boards, operates at the secondary level. And while this shift in power is also understood in this country by policy makers who are mandating graduation and promotion tests, the implications of the shift from the local educational authority (LEA) to the state department of education (SEA) has not received sufficient attention and discussion. Further, since most state level testing programs are developed and validated for SEAs by outside contractors, it is important to realize that the state may be effectively delegating this very real power to a commercial company whose interest is primarily financial and secondarily educational.

Despite the negative consequences historically associated with Principles 1-7, the use of tests perceived as involving high-stakes, is growing in the United States. Policy makers are well aware of the high symbolic value of tests. By mandating a test they are seen to be making things happen. The test becomes a synecdoche for standards. The general public is gratified to find the tests restoring confidence in the schools. The numerical scores from high-stakes tests have an objective, scientific, almost magical persuasiveness about them, that the general public and policy makers are quick to accept. Curriculum specialists, teachers, and administrators therefore, increasingly are forced to deal with the consequences of mandated high-stakes tests. Until recently, these tests have generally involved basic skills or minimum competencies. However, recently various education reform reports have recommended a move to secondary school certification tests in individual subjects that closely resemble the British 'O' and 'A' examinations.

In order to minimize the negative side effects associated with the impact of such tests, and to maximize possible benefits, it is important that professionals understand the value and limitations that historically have been associated with their use. They also need to understand how testing and the use of tests is changing in the United States: we are moving from LEA control to SEA control. Further, there are proposals that would alter the gathering and reporting of NAEP data to permit inter and intra state comparisons. The full philosophical, political and educational issues associated with such shifts in the use of tests need to be better understood and more fully debated.

Certainly, testing has an important function to play in American education. Test results when used in conjunction with other data about pupil performance, can help teachers to improve their instruction and to make educationally sound decisions. Test results can also provide important independent information to parents, the public and school administrators about the schools. Much of this information to which they are entitled. What we need right now is an awareness of the limitations and fallibility of tests. Test scores need to be regarded as one important element in decision making. The scores should be used by teachers and administrators in conjunction with other indicators of students' progress when making important decisions.

IV. NAEP as an Instrument for Informing Educational Policy

There are two basic ways in which test may affect educational policies. The first is the use of test information to inform policy makers about the current state of education. The second is the use of tests as administrative devices in the implementation of policy. In the former case test results are

used to describe the present state of education or some aspect of it, or in lobbying efforts for new programs or for reform proposals. The effects of this informational, descriptive use of test results on the educational process are indirect. This is in sharp contrast to the administrative use of test results where by results automatically trigger a direct reward or sanction being applied to an individual, or institution.

The Use of Tests Results to Inform Policy

The 1867 Act establishing the Department of Education recognized the need for gathering descriptive information about "the condition and progress of education in the several states and territories." Of course, at that time testing as we now know it did not exist. From the 1920's to the 1960's, standardized tests had little or nothing to do with state or federal policy. It was not until the early 1960's with the passage of the Elementary and Secondary Education Act of 1965, and the establishment of The National Assessment of Educational Progress (NAEP), that the Department began systematically to gather test data as part of its original mandate. Further, state departments of education have only recently begun to systematically collect test data to describe the status of education..

There were several reasons for this shift. First, the concept of equality of educational opportunity evolved from a concern about equality of inputs, resources, and access to programs into a preoccupation with achieving equality of outcomes. As a result test scores began to be used as a primary indicator of educational outcomes. Second, advocates for minority groups began to point to the large discrepancies between the test scores of middle class students and their constituents to lobby successfully for compensatory funds for programs to reduce these disparities. Third, the large expenditures

Haney and Madaus, paper for NAEP Study Group, 1986, p. 22.

in the 50's & 60's for curriculum development and compensatory programs led policy makers to ask for student test data as an indicator of the effectiveness of these programs. Fourth, as noted above NAEP was designed to provide a basis for public discussion and broader understanding of educational progress and problems of the nation.

More recently, the numerous educational reform reports have used test results, including SAT and NAEP data, to bring to the attention of the country what they conclude to be the mediocre state of American education, as well as to lobby for improvement programs to redress these weaknesses. Test data clearly form an important basis for the current negative descriptions of the status of American academic education. The question is how valid are these inferences? Is the academic performance of our students as poor as it is painted in the various reports? While there are weak spots - particularly at the secondary level with higher order skills - one could look at the same data and conclude that our schools are doing quite a creditable job; that declines and weak spots may be due in large part to non-school factors. In general the reports use test results in ways that accentuates the negative, eliminates the positive and leaves no room for Mr. In-Between.

An illuminating example of how these indicators are actually used to inform such reports is provided by the Twentieth Century Fund Report. It opens with the following gloomy assertion, "The nation's public schools are in trouble. By almost every measure, the performance of our schools falls far short of expectations". However, in a commissioned background paper, published as an appendix to the Report itself, Peterson examines all of the available indicators, including test scores, and concludes that, " Nothing in

these data permits the conclusion that educational institutions have deteriorated badly." It would seem that the Task Force did not take cognizance of its own commissioned paper.

Stedman and Smith in their excellent review of these reform proposals point out that they are quintessentially political documents. Testing evidence was used selectively to buttress arguments and evidence was often ignored that might have lessened the impact of the message. Stedman and Smith examined critically the way test score indicators were interpreted and concluded not only was interpretation often sloppy, but also that we have little in the way of valid, longitudinal national indicators of the academic performance of students nationwide. NAEP, however, was cited as the single exception to the latter indictment.

V. The Future of NAEP

In considering the future of NAEP, we agree with Stedman and Smith and other observers of the currently available indicators of the nation's educational health; NAEP is the most valid source of national data on our students general achievement level over time. In light of our previous discussion, we note too that NAEP has clearly been designed to inform educational policy-making, not drive it. This is in sharp contrast to many new testing programs implemented over the last decade at the state level which have been designed not just as means for informing policy, but instead as vehicles to implement educational policies. For instance, tests that are directly linked to decision-making about individual students (e.g. regarding grade to grade promotion, remediation or graduation), teachers (e.g. for teacher evaluation), schools or school systems (e.g. used allocate funding or other resources, or to rank publicly schools, school systems or

Haney and Madaus, paper for NAEP Study Group, 1986, p. 24.

even states) all are, at least in some degree, instruments that not only inform policy, but also implement it.

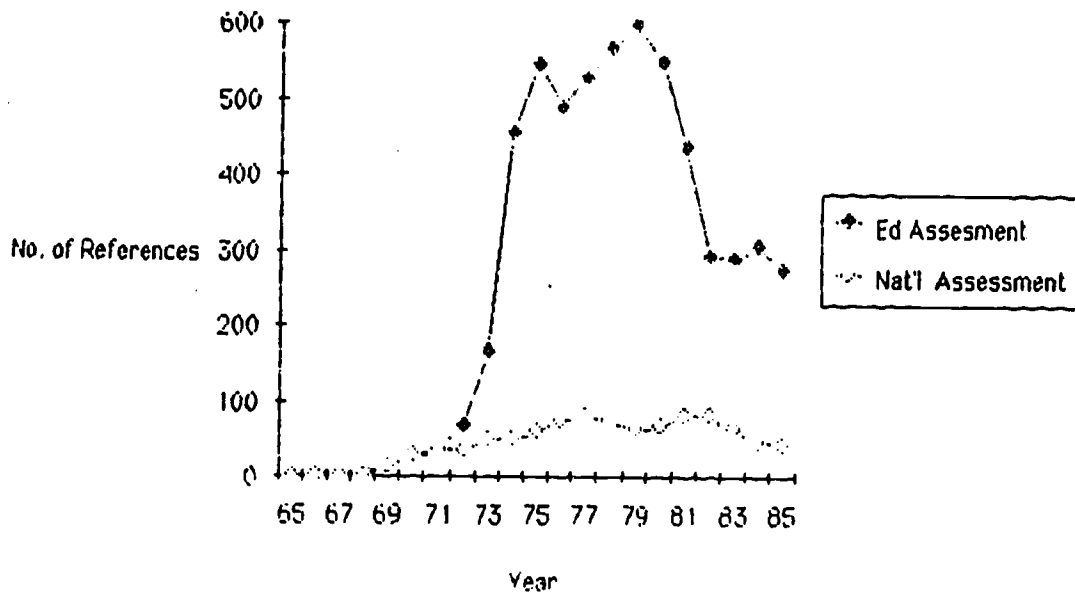
Despite the relatively high praise that NAEP has received as a source of valid data on the status of our nation's students, it also has been subjected to a variety of criticisms. For example, the NAEP move from reporting results at only the item level to reporting results aggregated across items, it was criticized for a serious failure to collect and analyze appropriate construct validity evidence (e.g. Haney, 1982). It is our understanding, however, that this problem has been -- or at least is being -- remedied via new latent trait scaling procedures designed by the Educational Testing Service since it has taken over the administration of NAEP from the Education Commission of the States .

A more common criticism of NAEP by almost every review of it over the last two decades has been that its results have not been terribly useful (e.g. Greenbaum, Garet and Solomon, 1977, Hazlett, 1974, , GAO, 197, Wirtz and Lapointe , Haney 1982). In order to provide a rough check on the utility of NAEP over time we performed a search of the ERIC data base to see the number of occurrences of the expression "national assessment" by year. To provide some perspective we also ascertained the total number of citations in the ERIC data base by year, and as a points of contrast with the search for "national assessment" items, the number of occurrences, also by year of items containing the phrase "educational assessment." Details of how we conducted this search and its limitations are described in Appendix 2, and the results are depicted in Figure 2.

Our findings suggest two major points. One is that in the citations from around 1970 (specifically 1969, 1970, and 1971), the number of occurrences of the phrase "national assessment" exceeded the number of

ERIC Database References

Total No. of References by Year



occurrences of the phrase "educational assessment." However, beginning in 1973 the number of times the phrase "educational assessment" is found in the ERIC databases increases dramatically. This is evidence, we suspect of one of the indirect affects of NAEP, namely that of promoting broader attention in the educational research community to the idea of educational assessment. More concretely, of course, the NAEP model also contributed in the 1970s to the founding of a number of statewide assessment programs. (Given the scale of the vertical axis of Figure 2, it is hard to discern the pattern of citations, for the pre-1973 years, but readers interested in details may consult the data table in the appendix.)

More recently, in the mid-1980s, the number of entries under both terms "educational assessment" and "national assessment" have decreased sharply, roughly one-third to one-half from what they were at their peaks (in 1982 the peak number for "national assessment" was 83, and for "educational assessment" the peak occurred in 1979 with 598). In part these changes may simply represent the trend of enthusiasms and jargon in the broad field of educational reasearch and testing. For it seems apparent that the enthusiasm for "assessment" in the 1970s gave way to interest in "competency testing" and "basic skills testing" in the late 1970s and early 1980s.

Nevertheless, the figures represented in Figure 2 give cause for concern about the continuing utility of NAEP. And it is against this background that we wish to consider ideas about the future of NAEP. Specifically, we wish to consider the idea of extending NAEP to be used in statewide testing programs, and the possibility of NAEP assessment exercises being used not just in sampling assessments but also in programs in which

Haney and Mc

whole pop
district).

The
assessment
testing pr
makes em
developm
population
associated
(particular
be tested)

The
instrumen
local progr
threaten ti
of NAEP m
been assoc

1. N
distorting
more usefu
decisions t
comparisor
mainly a s
into an ins

In st
administrat
effects des

aus, paper for NAEP Study Group, 1986, p. 26.

lations of students are tested (e.g. within a state or school

otion of using NAEP exercises (or items) not just in national
s based on stratified random samples, but also in statewide
grams in which all students at a particular grade might be tested
ent economic sense. It would allow the considerable costs of
it of NAEP exercises to be amortized over much larger

. However, what we wish to warn against are three dangers
with the possibility of using NAEP in statewide assessments
r where whole populations of students rather than samples might

angers we wish to warn against are : 1) If NAEP becomes an
of educational policy it likely would have distorting effects on
ms; 2) If NAEP becomes an instrument of policy, it would
validity of NAEP findings, and 3) Such wider instrumental use
ht threaten some of the indirect benefits which we think have
ted with NAEP in its first 20 years.

EP as an instrument of educational policy would have
effects. Here our worry is that in the efforts to make NAEP
if it develops into a program which would allow specific
be made on the basis of NAEP results (or even state by state
, it would change from what it has been over the last 20 years,
rce of high quality information for informing educational policy,
ment of policy.

rt our concern is that to the extent that NAEP becomes an
re device for implementing educational policy, the negative
bed in the seven principles above would come into play. We

Honey and M

should ha

NAEP for

or gradua

or district

2.

Our secon

actually t

note that

high quali

notable ex

gather cor

of NAEP e

largely so

NAEP mor

instrumen

effects on

assessment

concern is

measurme

because it

obtrusive

of course,

best and h

through th

practitione

3. Indirect

er for NAEP Study Group, 1996, p. 27.

id that we know of no specific plans by anyone to use
pecific decisions about students (such as grade promotion
wever, as we have argued, even comparisons such as state
on the basis of test results can lead to distortions.

Instrument of policy, might have less validity.

is that efforts to improve the utility of NAEP might
e validity of assessment data derived from it. Here, we
as been one feature of NAEP universally praised, it is the
rmation derived from NAEP assessments (with the
prior to the ETS administration of NAEP, of the failure to
lity evidence supporting the reporting of data on groups
- but as we said, we think that this problem has been
new ETS scaling procedures). However, if efforts to make
ave the effect of making NAEP an administrative
tional policy, this would likely not only have distorting
ial practices, but also would threaten the validity of
rived from NAEP. Another way of communicating this
e distinction between obstrusive and unobtrusive
the real virtues of unobyrusive measurement is that
usive it tends to have higher measurement validity than
ent. Unobstrusive measurement need not be unuseful,
ur view that educational testing and assessment works
untoward consequences when their effects are mediated
d judgments of educational policymakers and

ts

As we have already noted, one of the relatively early indirect benefits of NAEP seems to have been that in the mid-1970s, it helped promote the development of broader interest in educational assessment and specifically several statewide assessment programs. However, in the late 1970s as we also already noted, many state programs moved in the direction of "minimum competency" or "basic skills" testing. In this regard we feel that one of the real indirect benefits of NAEP is that it has provided a model of broader assessment, broader both with regard to the range of knowledge and skills tested, and with regard to the methods of assessment used.

In this regard and in thinking about the possible indirect effects of standardized testing on education in general and on curriculum specifically, two general points seem apparent. One is that though most standardized tests may focus on reading and math skills, the general aims of education as represented in curriculum are considerably broader. Take for example the "New Basics" advocated by the National Commission on Excellence in Education (1983). This group, though only one of many groups recently advocating reforms of education, was perhaps the most prominent. It recommended that:

State and local high school graduation requirements be strengthened and that at a minimum, all students seeking a high school diploma be required to lay the foundations in the Five New Basics by taking the following curriculum during the 4 years of high school: (a) 4 years of English; (b) 3 years of math; (c) 3 years of science; (d) 3 years of social studies; and (e) one-half year of computer science. For the college-bound, 2 years of foreign language are strongly recommended in addition to those taken earlier. (National Commission on Excellence in Education, 1983, p. 24.)

Honey and Medina, p.

This is of course
recommendations
striking point about
"standardized" tests
completely, neglecting
the case with NAEP

A second general
with NAEP again
employ the multiple
striking contrast that
that students should
identify and work
this regard we are
the predominance
bias". Frederiksen
order problem solving
problems, multiple
domain of skills and
the preponderance
surely has contributed
communicating exams
schools to be teaching
article on the real tests
from Ralph Tyler first

Frederiksen put it:
The "real tests"
tests on teaching
efficient tests
untaught. A

see for NLEED Study Group, 1996, p. 20.

rise only one of many recent reform
for improving education in the United States. But a
t this recommendation, in contrast to the coverage of most
s, is that three of the five "new basics" are badly, if not
ed in most statewide testing programs. This has not been.

eral feature of almost all large scale testing programs,
ing perhaps the most notable exception, is that they
e-choice format almost exclusively. This format is in
the calls in many of the recent education reform reports
l be taught not just to solve pre-set problems, but to
ut solutions to problems they find for themselves. In
eminded of Norman Frederiksen's article suggesting that
f the multiple-choice format represents "the real test
reviews evidence showing that particularly for higher
ng skills, involving solving open-ended or ill-structured
choice items may not be good measures of the broad
abilities that we might hope are taught in schools. Thus,
f tests employing the multiple-choice format, which
ed to recent interest in "test-taking skills," may be
ly the wrong message about what it is that we want our
g. In this regard it is notable that Frederiksen closed his
st bias with a message remarkably similar to the one
een year earlier, quoted at the start of this article. As

t bias" in my title has to do with the influence of
ing and learning. Efficient tests drive out less
s, leaving many important abilities untested and
important task for those involved in testing is to

Haney and Madaus, paper for NAEP Study Group, 1986, p. 30.

develop instruments that will better reflect the whole domain of educational goals and to find ways to use them in improving the educational process. (p.19)

In the future, we therefore hope that there will be maintained the past tradition of NAEP, of employing diverse methods of assessment and helping to do what both Tyler and Frederiksen advocate -- namely making sure that our assessments span as broad a range of the curriulums and goals of our schools as possible.

Appendix 1

The Education Index, published by the H. W. Wilson Company since the 1932 is an author and subject index to educational material in the English language. Though primarily a periodical index it also covers proceedings, yearbooks, monographs and material printed by the U.S. Government. In terms of periodical coverage, the index covers both professional and more popular journals. Actual selection of periodicals for indexing is accomplished by subscriber vote, at least since 1970. In voting, subscribers are asked to place primary emphasis on the reference value of periodicals under consideration.

Figure 1 was constructed simply by measuring the number of column inches devoted to listings concerning testing and curriculum in every volume of the Index from 1932 (volume 1, covered material published from January 1929 - June 1932) through 1985 (volume 35; July 1984 - June 1985). Over these volumes there were occasionally some changes in the index rubrics concerning testing and curriculum. However the primary rubrics considered as to pertaining to the two topics of interest were:

Testing:

- Tests and Scales
- Testing programs (introduced in vol. 6)
- Tests of general educational development (vol. 23)
- Testing instruments (vol. 25)

Curriculum:

- Curriculum
- Curriculum Making
- Curriculum laboratories (introduced in vol. 4)
- Curriculum satisfaction (vol. 8)
- Curriculum selection (vol. 8)
- Curriculum development (replacing Curriculum Making in vol. 11)
- Curriculum studies (vol. 16)

It should be noted too that only since 1964 (volume 14 covering 7/63-7/64) has the Index been produced annually. Before that time it was issued on either biennial or triennial bases. Thus for years prior to 1964 we have shown the average annual numbers of column inches listed under the relevant testing and curriculum rubrics. Finally it is worth mentioning that over all 35 volumes of the Index, the type size has remained the same and that 1 column inch, including headings and subheadings, amounts on average to roughly 2 to 3 references.

Haney and Medaus, paper for NAEP Study Group, 1986, p. 32.

References

Anastasi, A. (Ed.) 1966. Testing problems in perspective. Washington, D.C.: American Council on Education.

Congressional Budget Office, 1986. Trends in Educational Achievement. Washington, D.C.: author.

Education Index. 1932 -1985. Volumes 1 - 35. New York: H. W. Wilson.

Frederiksen, N. 1981. The real test bias. (ETS Research report 81-40). Princeton, NJ: Educational Testing Service.

Pipho, C. 1985. Tracking the reforms: Part 5: Testing -- Can it measure the success of the reform movement? Education Week, May 22, 1985, p. 19.

National Commission on Excellence in Education. A Nation at Risk: The Imperative for Educational Reform. Washington, D.C.: U.S. Department of Education.

Tyler, R. "What Testing Does to Teachers and Students," 1959, In Anastasi, 1966, pp. 46-52. .

ERIC REFERENCES CONCERNING ASSESSMENT

ERIC Citations. Ed'l Ass't.	Total. Regarding Nat'l Ass't.	By Year			
		Year	Total Ref's	Ed'l Ass't	Nat'l Ass't
	2959	65	1	0	0.00
	4856	66	5	1	0.20
	7115	67	6	0	0.00
	9824	68	9	2	0.22
	25317	69	15	16	1.07
	28954	70	26	30	1.15
	32499	71	36	37	1.03
	34063	72	68	35	0.51
	35208	73	164	50	0.30
	35274	74	453	50	0.11
	38168	75	545	59	0.11
	37417	76	489	76	0.16
	36914	77	528	82	0.16
	38471	78	564	75	0.13
	35984	79	598	64	0.11
	34591	80	548	66	0.12
	31516	81	436	80	0.18
	30661	82	293	83	0.28
	30506	83	290	64	0.22
	29591	84	306	45	0.15
	24275	85	275	44	0.16