

The Assessment of Student Achievement: The Hundred Years

War*

Lyle V. Jones

AERA Annual Convention, April 1999

The University of North Carolina at Chapel Hill

*Invited address, Division D of AERA. Montreal, Canada, April 21,

The Assessment of Student Achievement:

The Hundred Years War

I shall begin with some remarks about the use of achievement tests for high-stakes decision making, and then shall review some of Ralph Tyler's contributions to educational assessment and note what has become of them. After referring again to high-stakes testing, I suggest an agenda item for future research.

An elementary school principal comments on high-stakes student testing:

"... a teacher knows that his whole professional status depends on the results he produces and he really is turned into a machine for producing those results; that is, I think, unaccompanied by any substantial gain to the whole cause of education."

This statement, cited by Sutherland (1973, p. 68) is from a schoolmaster in 1887! The reference is to the system of "payment by results" for the publicly supported elementary schools in England and Wales. For 30 years, beginning in 1862, grants to each school were based on annual inspections that entailed the testing of individual students on reading passages and on arithmetic test cards (and, late in this period, on other subjects as well).

Was this a high-stakes testing system? It certainly was for teachers and for school managers, whose salaries were linked to the amount of the grant. And an Inspectors' visit was an anxiety-provoking event. One report reads as follows:

"Two inspectors came once a year and carried out a dramatic examination. The schoolmaster came into school in his best suit; all the pupils and teachers would be listening till at ten o'clock a dog-cart would be heard on the road even though it was eighty yards away. In would come two gentlemen with a deportment of high authority with rich voices. Each would sit at a desk and the children would be called in turn to one or other. The master hovered round, calling children out as they were needed. The children could see him start with vexation as a good pupil stuck at a word in the reading book he had been using all year, or sat motionless with his sum in front of him. The master's anxiety was deep, for his earnings depended on the children's work. One year the atmosphere of anxiety so affected the children that, one after another as they were brought to the Inspector, the boys howled and the girls whimpered. It took hours to get through them." (Cited by Sutherland, 1973, p. 66).

The inspectors, sub-inspectors, and their assistants weren't having much fun, either. In 1873, one of Her Majesty's Inspectors reported:
"I feel as if I were writing a despatch in the midst of a battle, so dire is the din of educational conflict around." (Cited by Sutherland, 1973 p. 69).

Clearly, the assessment war was waging even earlier than 100 years ago. The political demand for school accountability led teachers in England and Wales to concentrate on drilling their pupils to prepare them for the tests, at the expense of broader objectives of instruction. Does that sound familiar now, over a century later?

One commentator, after evaluating "payment by results," concluded: "if indeed it is possible to point to a simple moral from this dismal episode in England's educational history, perhaps it is that true accountability in education should not be facilely linked to mechanical examination results, for there is a very distinct danger that the pedagogical methods employed to attain those results will themselves be mechanical and the education of children will

be so much the worse." (Rapple, 1994). (By 1897, "payment by results" in England and Wales was replaced by block grants to schools, with sharply higher levels of funding.)

In mid-20th century, Ralph Tyler did his best to effect a truce in an ongoing war. He invented the term assessment to distinguish it from three other forms of educational appraisal: first, testing achievements of individual students to assign grades or to select students for further opportunities; second, diagnosing learning difficulties of a student (or of a class) to plan subsequent teaching; and third, evaluating the effectiveness of a curriculum or a set of teaching methods. In contrast, Tyler proposed that the focus of assessment be not on individual students, classrooms, schools or school systems. Assessments furnish information about the educational attainments of large numbers of people, perhaps of different ages, different demographic groupings, and different geographic regions. The purpose of assessment is analogous to that of estimating the Gross National Product, or the Consumer Price Index, or health and mortality indices, to provide dependable information about population and sub-population change over time, in this case about the progress of education.

As the National Assessment of Educational Progress (NAEP) was designed (see Jones, 1996; Tyler, 1966), only a sample of students would be assessed, no student would take more than a fraction of the exercises, and no score would be obtained from any student's performance. Exercises -- many of them in the form of hands-on problems to be solved, and some entailing discussion by a group of children -- would represent a broad range of difficulty and a full array of educational objectives in ten different subject areas. Professional administrators were trained to provide highly controlled assessment conditions. Exercises were read aloud, so that deficiencies in reading would not prevent good performance in math, say, or in citizenship. An "I don't know" alternative was offered to discourage guessing and to reduce non-response. Results of periodic assessments would be reported, exercise by exercise, so that for four

different age groups, 9, 13, 17, and young adult, the public would have concrete evidence about what respondents know and can do.

Over the past 30 years, some of Tyler's objectives for NAEP have survived, but just barely. First, the desired rich variety of exercises was compromised, in favor of more traditional multiple-choice and short-answer items, the kinds of items with which testing companies were familiar. Exercises became quite homogeneous in difficulty, with fewer very easy or very difficult ones. The young-adult sample was eliminated, and school grade has replaced age as the primary unit of assessment. The ten subject areas have received uneven attention, with math, reading, science, and writing assessed far more often than literature, social studies, art, music, citizenship, and career development. No longer are exercises read aloud, nor has an "I don't know" alternative been retained. For state assessments, local school personnel now administer the exercises, which raises questions about the uniformity of administration.

Instead of reporting a percent-correct score for each exercise, scale scores were developed, for large clusters of exercises. More recently, reporting has been by "achievement levels," so as to compare actual performance with how good performance "should be".

Using IRT technology, scores now are imputed for each child in the sample, even though different children take different sets of exercises. Imputed scores then are averaged for any specified subgroup of children.

Many of these changes were well-intentioned, and some clearly are supported by psychometric considerations and by the need to better communicate results to the public. Nonetheless, some changes have compromised Tyler's vision of assessment. While Tyler's vision has materialized fully in New Zealand's promising National Education Monitoring Project (e.g., Flockton & Crooks, 1997), many features have disappeared from NAEP.

Along with the changes in NAEP have come pressures to report scores not just for large subpopulations, but by school district, by school, by classroom, and for each individual child. Indeed, President Clinton and the U.S. Department of Education promote "a personalized version of NAEP" (Riley, 1997) for every child, a voluntary national test derived from the NAEP model. And most states have adopted or are developing high-stakes tests in reading and math, often to be used as a basis for promotion from grade-to-grade, for receipt of a high school diploma, and for salary supplements to teachers -- shades of 1887!

Numerous unresolved problems attend these procedures. What accommodations should be provided for students who are visually handicapped or in special education classes? What special efforts will be dedicated to helping students who fail to pass the tests? Will the testing programs lead to increased levels of school dropout thereby helping some students at the expense of others? Will teachers focus on preparing students for the tests at the expense of failing to meet other important educational objectives?

One unanticipated side effect recently came to my attention. In North Carolina, teacher bonuses are awarded by the State to all teachers in schools for which test scores rise substantially above "an expected amount". That now is influencing teacher recruitment. Promising new teachers ask about the bonus and are attracted to schools that have earned the bonus. These prospective teachers then are less likely to accept positions at schools in which they are even more needed.

Illustrative of problems created by an over-emphasis on test scores is a school principal's recent announcement to teachers at a middle school in North Carolina. The school goal is to elevate end-of-grade test scores in reading and math so that by 2003, 95% of the students will score at or above grade level. Between now and then, students scoring below grade level will not be permitted

to enroll in elective courses, but will be assigned to remedial math and reading classes instead.

As one might expect, the teachers of drama, art, dance, and music are appalled. They cite many instances of students whose engagement in their subjects result in a heightened commitment to school and increased performance in core subjects as well.

Attributed to William Butler Yeats is the statement:
"Education is not the filling of a pail, but the lighting of a fire."

Quite obviously, good education can be both, instilling knowledge, but also lighting fires, encouraging engagement. And neither goal is sufficient without the other. For accountability, an imbalance has been created because we have learned to measure achievement much better than we can measure engagement. Isn't it time to develop measures of the latter?

A recent study by Mahoney & Cairns (1997) supports the importance of engagement for reducing high school dropout, especially for at-risk students. In a longitudinal study, the authors followed about 400 students through middle school and high school. Students were classified as highly competent, marginally competent, or at risk, based upon annual ratings by teachers of each student's academic and personal competence. Engagement was measured by the extent of extracurricular involvement in athletics, student government, student clubs, music or drama, journalism, and school assistantships. For all three competence groups, but most dramatically for at-risk students, the proportion of dropouts in high school declines for students with higher levels of extracurricular involvement.

Such promising results suggest that devising systematic measures of student engagement may have a considerable payoff. Adding measures of school

engagement to measures of school achievement could help to restore a balance between two important components of effective schooling.

Perhaps the most serious impediment to educational reform is confusion about the criterion. For many states -- and perhaps for the nation as well -- test scores seem to be accepted as the criterion for effective education. So the message to teachers is to train students to elevate their test scores. How shortsighted it is to accept test scores as an ultimate criterion of the benefits of education. At best, tests can serve to predict one component of successful schooling. As emphasized by Haertel (1999) in his Presidential address for the National Council on Measurement in Education, we must question the validity of test scores not only for that component, but for others as well. By establishing a richer and more appropriate criterion, we would be challenged to develop richer and more appropriate predictors, certainly to include indices of commitment to learning -- that is, the lighting of fires -- as well as indicators of what students have learned.

References

- Flockton, L. & Crooks, T. (1997). Reading and speaking assessment results, National Assessment Monitoring Report 6. Dunedin, New Zealand: Educational Assessment Research Unit, University of Otago.
- Haertel, E. H. (1999). Validity arguments for high-stakes testing: In search of the evidence. Presidential address presented at the annual meeting of the National Council on Measurement in Education, Montreal, Canada, April 21.
- Jones, L. V. (1996). A history of the National Assessment of Educational Progress and some questions about its future. *Educational Researcher*, 25, no. 7, 15-22.
- Mahoney, J. L. & Cairns, R. B. (1997). Do extracurricular activities protect against early school dropout? *Developmental Psychology*, 33, 241-253.
- Rapple, B. A. (1994). Payment by results: An example of assessment in elementary education from nineteenth century Britain. *Educational Policy Analysis Archives*, 2, no.1, 1-22.

Riley, R. W. (1997). Testimony of Secretary Richard W. Riley. Statement to the Senate Labor, Health & Human Services & Education Subcommittee of the Senate Appropriations Committee, September 4.

Sutherland, G. (1973). Policy-making in elementary education 1870-1897. London: Oxford University Press.

Tyler, R. W. (1966). The objectives and plans for a National Assessment of Educational Progress. *Journal of Educational Measurement*, 3, no. 1, 1-4.

14

1