

the
Academy
in
Transition

A Brief History of Student Learning Assessment

*How We Got Where We Are and
a Proposal for Where to Go Next*

Richard J. Shavelson

With a Foreword by
Carol Geary Schneider
and Lee S. Shulman



Association
of American
Colleges and
Universities

the
Academy
in
Transition

A Brief History of Student Learning Assessment

*How We Got Where We Are and
a Proposal for Where to Go Next*

Richard J. Shavelson

With a Foreword by
Carol Geary Schneider
and Lee S. Shulman



*Association
of American
Colleges and
Universities*



*Association
of American
Colleges and
Universities*

1818 R Street, NW, Washington, DC 20009

Copyright © 2007 by the Association of American Colleges and Universities.
All rights reserved.

ISBN 978-0-9779210-7-2

The author gratefully acknowledges the support of the Atlantic Philanthropies for the work presented in this paper. An abridgement of this paper appeared as "Assessing student learning responsibly: From history to an audacious proposal," *Change* 2007, 39 (1): 26–33.

Contents

About This Series	x
Foreword	x
I. Introduction	x
II. A Brief History of Learning Assessment	x
III. The Collegiate Learning Assessment	x
IV. A Proposal for Assessing Learning Responsibly	x
V. Concluding Comments	x
Appendix: Summary of Tests and Testing Programs by Era	x
Notes	x
References	x
About the Author	x

the
Academy
in
Transition

Other titles in the series:

Integrative Learning: Mapping the Terrain

By Mary Taylor Huber and Pat Hutchings

General Education and the Assessment Reform Agenda

By Peter Ewell

The Living Arts: Comparative and Historical Reflections on Liberal Education

By Sheldon Rothblatt

College-Level Learning in High School: Purposes, Policies, and Practical Implications

By D. Bruce Johnstone and Beth Del Genio

General Education in an Age of Student Mobility: An Invitation to Discuss Systemic Curricular Planning

Essays by Robert Shoenberg and others

General Education: The Changing Agenda

By Jerry G. Gaff

Globalizing Knowledge: Connecting International and Intercultural Studies

By Grant Cornwell and Eve Stoddard

Mapping Interdisciplinary Studies

By Julie Thompson Klein

Contemporary Understandings of Liberal Education

By Carol Geary Schneider and Robert Shoenberg

For information about these and other AAC&U publications, or to place an order, visit www.aacu.org, e-mail pub_desk@aacu.org, or call 202-381-3760.

About This Series

The Association of American Colleges and Universities (AAC&U) has a long history of working with college leaders across the country to articulate the aims of a liberal education in our time. AAC&U is distinctive as a higher education association. Its mission focuses centrally on the quality of student learning and the changing purpose and nature of undergraduate curricula.

AAC&U has taken the lead in encouraging and facilitating dialogue on issues of importance to the higher education community for many years. Through a series of publications called the Academy in Transition—launched in 1998 with the much-acclaimed *Contemporary Understandings of Liberal Education*—AAC&U has helped fuel dialogue on such issues as the globalization of undergraduate curricula, the growth of interdisciplinary studies, how liberal education has changed historically, and the increase of college-level learning in high school. The purpose of the series, which now includes ten titles, is to analyze changes taking place in key areas of undergraduate education and to provide “road maps” illustrating the directions and destinations of the changing academy.

During transitions, it is important to understand context and history and to retain central values, even as the forms and structures that have supported those values may have to be adapted to new circumstances. For instance, AAC&U is convinced that a practical and engaged liberal education is a sound vision for the new academy, even if the meanings and practices of liberal education are in the process of being altered by changing conditions. As the titles in this series suggest, AAC&U’s vision encompasses a high-quality liberal education for all students that emphasizes connections between academic disciplines and practical and theoretical knowledge, prizes general education as central to an educated person, and includes global and cross-cultural knowledge and perspectives. Collectively, the papers published in the Academy in Transition series point to a more purposeful, robust, and efficient academy that is now in the process of being created. They also encourage thoughtful, historically informed dialogue about the future of the academy.

AAC&U encourages faculty members, academic leaders, and all those who care about the future of our colleges and universities to use these papers as points of departure for their own analyses of the directions of educational change. We hope this series will encourage academics to think broadly and creatively about the educational communities we inherit, and, by our contributions, the educational communities we want to create.

Foreword

Many of us in higher education today are thinking hard about assessment. But often our cogitations tend toward what can only be called wishful thinking.

A wish that has no doubt washed over all of us at one time or another is that if we simply ignore assessment, or hold it off long enough, the issues (like the misguided new administrator) will finally give up and go away. But in our wiser moments, we know that this is not an answer. Indeed, as is clear from Richard Shavelson's lively tour of the twists and turns of assessment over the past century, the names may change, and the technology has evolved, but assessment has stayed with us with great persistence. "Today's demand for a culture of evidence of student learning appears to be new" Shavelson tells us, but it turns out "to be very old," and there's no wishing it away. Moreover, we should not be eager to wish it away.

Nor is there a magic bullet. One of the most dangerous and persistent myths in American education is that the challenges of assessing student learning will be met if only the right instrument can be found—the test with psychometric properties so outstanding that we can base high-stakes decisions on the results of performance on that measure alone.

This wish is not only self-indulgent but self-defeating. Ironically, no test can possess such properties because, to achieve validity, test designers have to narrow the focus on any particular instrument to a sobering degree. Thus, the better the arguments we can make regarding the validity of any given measure—whether of knowledge, skills, or some other virtue—the less appropriate that measure is as the basis for consequential decisions about a student's overall learning gains, much less as the sole determinant of an institution's educational quality.

Thinking of assessment as primarily a technical challenge—though certainly it *is* that, as Shavelson's analysis also makes clear—is another form of wishful thinking. The far-reaching questions raised through assessment cannot be solved through technical ingenuity alone.

What's needed, of course, is *educational* thinking, and happily there has been a good deal of that in the past two decades of assessment activity. With the wave of state mandates for assessment in the mid and later 1980s, and new accreditation requirements in the 1990s, campuses began to organize themselves to respond. Many did so grudgingly, and there were plenty of missteps, misunderstandings, and dead ends. But there were also wonderful examples of what can happen when educators take up the challenge to figure out and clearly articulate what they want their students to know and be able to do: the core task of assessment. Many of these efforts were

funded by the U.S. Department of Education's Fund for the Improvement of Postsecondary Education, and the results, in turn, provided models and momentum for additional campuses that came together—hundreds of them—at the assessment forum led by the American Association for Higher Education until 2005, when that organization closed its doors.

Toward the end of his essay, Shavelson makes a crucial point that campuses committed to assessment know well: *that assessment all by itself is an insufficient condition for powerful learning and improvement.* Of course, more and better evidence of student learning is important, but knowing what to *make* of that evidence, and how to act on it, means getting down to core questions about the character of the educational experience and the goals of liberal learning. These are not questions that higher education can dare leave to the testing companies or to external agencies, no matter how well intentioned and enlightened.

This view of assessment has become central to the work of both the Carnegie Foundation for the Advancement of Teaching and the Association of American Colleges and Universities (AAC&U). Shavelson traces Carnegie's work over the first part of the twentieth century as a story about standards and standardization. But the needs are different today, and Carnegie's more recent work places much greater emphasis on the role of faculty in exploring what our students do—and don't—learn. The foundation's extensive work on the scholarship of teaching and learning, for instance, has helped fuel a movement in which “regular” faculty, across the full spectrum of disciplines and institutional types, are treating their classrooms and programs as laboratories for studying student learning in order to improve it. Seen through the lens of classroom inquiry, assessment is a feature of *the pedagogical imperative* in which faculty see themselves as responsible for the learning of their students and for deepening our collective sense of the conditions in which important forms of learning can occur.

Through its Liberal Education and America's Promise (LEAP) initiative, AAC&U is working with its member campuses to develop assessments that strengthen students' learning and assess their best work rather than just the attainment of a few narrowly defined foundational skills and/or basic knowledge. As AAC&U's board of directors put it in their official statement on assessment (2005, 3), colleges and universities should hold themselves “accountable for assessing [their] students' best work, not generic skills and not introductory levels of learning.”

In its recently released LEAP report, *College Learning for the New Global Century*, AAC&U recommends tying assessment efforts much more closely to the curriculum and to faculty priorities for student learning across the curriculum. The report affirms, as well, that any national assessment measure, however well developed, is only part of the solution to the

problem of underachievement. As the report notes, “standardized tests that stand outside the regular curriculum are, at best, a weak prompt to needed improvement in teaching, learning, and curriculum. Tests can, perhaps, signal a problem, but the test scores themselves do not necessarily point to where or why the problem exists or offer particulars as to solutions” (2007, 40).

A fuller strategy, the LEAP report proposes, would prepare students to produce a substantial body of work—capstone projects and/or portfolios—that require their best efforts. The resulting accomplishments should be assessed for evidence of students’ competence on liberal education outcomes such as analytical reasoning and integrative learning, as well as their achievement in their chosen fields. Standardized assessments can then fill out the emerging picture, providing the ability to benchmark accomplishment against peer institutions, at least on some aspects of student learning.

As the LEAP report notes, “however the assessments are constructed . . . the framework for accountability should be students’ ability to apply their learning to complex problems. Standards for students’ expected level of achievement also will vary by field, but they should all include specific attention to the quality of the students’ knowledge, their mastery of key skills, their attentiveness to issues of ethical and social responsibility, and their facility in integrating different parts of their learning” (2007, 41–2).

Richard Shavelson offers an important historical context to consider as institutions across the country continue to develop new methods of assessment in response to renewed calls for greater accountability and, more importantly, the urgent need to raise levels of student achievement. He helps us better understand the “state-of-the-art” in standardized testing today, and what we should ask from testing agencies in the future. Above all, he helps us understand why psychometricians themselves are so opposed to any efforts at institutional ranking or comparisons based on standardized tests.

We are grateful to Richard Shavelson for taking the time to put current debates in a larger historical and educational forum. Everyone who is thinking today about assessment and public accountability will benefit greatly from the insights this study provides.

Carol Geary Schneider

President, Association of American Colleges and Universities

Lee S. Shulman

President, the Carnegie Foundation for the Advancement of Teaching

I. Introduction

Over the past thirty-five years, state and federal policy makers, as well as the general public, have increasingly been pressuring higher education to account for student learning and to create *a culture of evidence*. While virtually all states already use proxies (e.g., graduation rates) to report on student performance, states are now being pressured to measure learning directly. U.S. Secretary of Education Margaret Spellings' Commission on the Future of Higher Education, for example, has called for standardized tests of students' critical thinking, problem solving, and communication skills.

While the current demand to establish a culture of evidence appears to be new, it has a long lineage. The future development of this culture may very well depend on how well we appreciate the past. *Cultures of evidence will not automatically lead to educational improvement, if what counts as evidence does not count as education*. Narrow definitions and narrow tests of what count as learning outcomes in college may very well distort the culture of evidence we seek to establish. As we shall see from the past, and as we know from current studies (Immerwahr 2000; AAC&U 2007), there is more to be learned and assessed in higher education than the broad abilities singled out by the Spellings Commission for measurement by standardized tests. These additional outcomes include learning to know, understand, and reason in an academic discipline. They also include personal, civic, moral, social, and intercultural knowledge and actions—outcomes the Educational Testing Service has described as “soft.” Some experts say that such “soft” outcomes cannot be measured adequately because “the present state of the art in assessing these skills is not adequate for supporting the institution of a nationwide set of standardized measures” (Dwyer, Millett, and Payne 2006, 20). But this position is unsatisfactory. This set of outcomes—which, following the lead of the Association of American Colleges and Universities, I will call personal and social responsibility (PSR) skills—are every bit as demanding as the academic skills that often get labeled exclusively as *the cognitive skills* and are too important not to be measured. *If we do not measure PSR skills, they will drop from sight as accountability pressures force campuses to focus on a more restricted subset of learning outputs that can be more easily and less expensively measured.*

The outcomes framework depicted in figure 1 demonstrates the importance of extending the range of outcomes we assess beyond broad abilities. Such outcomes could range from the development of specific factual, procedural, and conceptual knowledge and reasoning in a discipline

(such as history) to the development of the skills on which the Spellings Commission focused (critical thinking, problem solving, and communication), to the development of reasoning applicable to a very wide variety of situations, or to the development of intelligence. Moreover, “cognitive” outcomes include PSR skills insofar as reasoning and thinking are involved in personal relations, moral challenges, and civic engagement. The PSR skills are not so soft; they involve cognition and more, *as do academic skills*. Finally, the arrows in figure 1 remind us that general abilities influence the acquisition of knowledge in concrete learning environments, that direct experiences are the stuff on which reasoning and abstract abilities are developed, and that cognitive performance on academic and PSR skills is influenced by the interaction of individuals’ accumulated experience in multiple environments with their inheritance.

Furthermore, the standardized tests that the Spellings Commission and others have in mind for outcomes assessment are not interchangeable. There are different ways to measure student learning; some standardized tests focus only on a narrow slice of achievement, while others focus on broader abilities developed over an extended course of study. Especially for

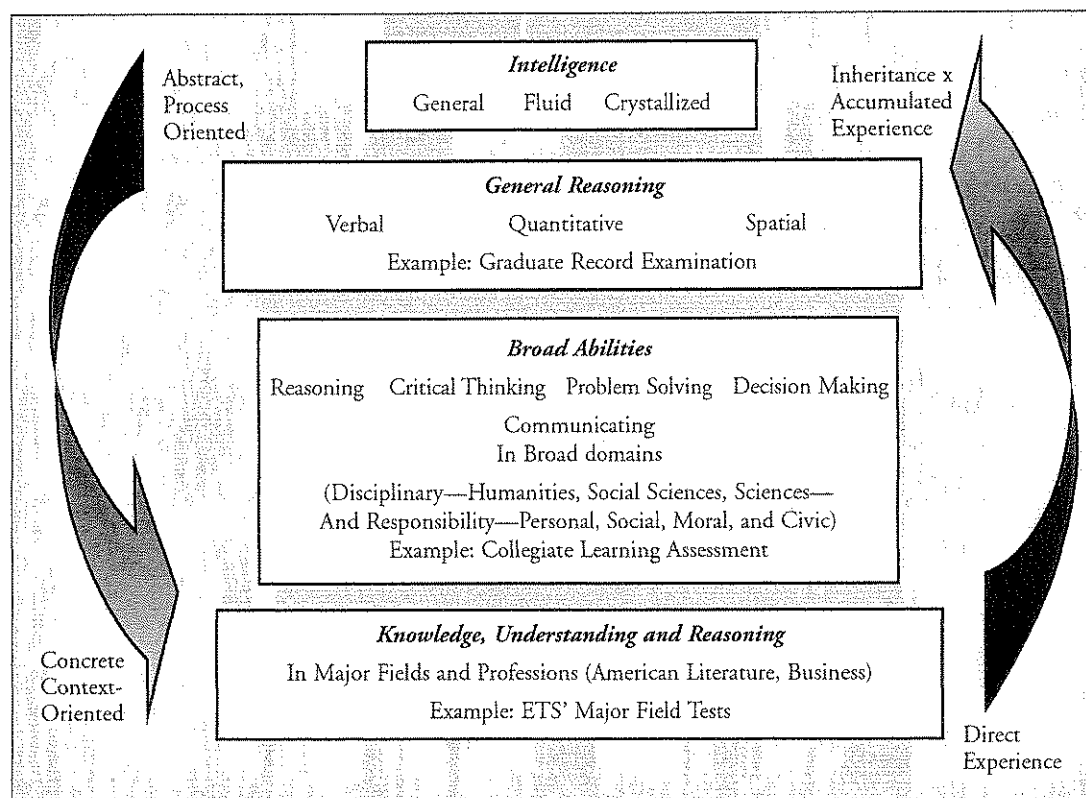


Figure 1. Framework for student learning outcomes. (Adapted from Shavelson and Huang 2003, 14.)

higher education, the different assumptions about what ought to be measured that are embedded in every assessment instrument need to be clarified and carefully considered before specific tests are chosen to assess students' cumulative gains from college study.

The multiple-choice technology developed almost a hundred years ago, for example, is inherently limited when it comes to measuring the full array of student learning outcomes depicted in figure 1. Multiple-choice measures have a long history, as we shall see. They are the basis of the standardized tests that are often used today, including the Collegiate Assessment of Academic Proficiency (CAAP), the Measure of Academic Proficiency and Progress (MAPP), and the College Basic Academic Subjects Examination (CBASE). The MAPP was recommended by the Spellings Commission as ways of assessing student learning in college. But these measures are limited in their ability to get at some of the more complex forms of reasoning and problem solving that are commonly viewed as distinctive strengths of American higher education.

If the learning outcomes of higher education are narrowly measured because cost, capacity, and convenience dictate reductive choices, then we stand the risk of narrowing the mission and diversity of the American system of higher education, as well as the subject matter taught. What we need to do instead is to learn from the rich history of student learning assessment and take responsible steps to develop and measure the learning outcomes our nation values so highly.

II. A Brief History of Learning Assessment

Surprisingly, our journey begins with the Carnegie Foundation for the Advancement of Teaching, which is so well known for the “Carnegie unit” and TIAA. The foundation was arguably the ringleader of student learning assessment. Howard Savage (1953), a staff member and historian of the foundation in its early days, attributes Carnegie’s leadership in college learning assessment to its first president, Henry Pritchett, who was motivated by his concern for the quality of higher education and his recognition of the potential impact that the emergence of “objective testing” might have on monitoring that quality. Walter A. Jessup, the foundation’s third president, later put what had become the foundation’s vision this way:

The central problems [in improving higher education] are three in number: first, the setting up of generally accepted standards of achievement; secondly, the devising of methods of measuring this achievement and holding pupils to performance; and thirdly, the introduction of such flexibility in educational offerings that each individual may receive the education from which he is able to derive the greatest benefit. (Kandell 1936, vii)

Pritchett’s passion was shared by his chief staff member, William S. Learned, “a man who had clear and certain opinions about what education ought to be . . . [with] transmission of knowledge as the *sine qua non*” (Lagemann 1983, 101). Learned became the instrument through which the foundation transformed higher education learning assessment.¹ Together with Columbia College’s Ben D. Wood, who held the view “that thinking was dependent upon knowledge and knowledge dependent upon facts” (Lagemann 1983, 104), Learned led a large-scale assessment of college learning in the state of Pennsylvania. Learned parlayed this experience into the development of the Graduate Record Examination and germinated the idea of a “National Examination Board,” a national testing agency that, twenty years later, became the Educational Testing Service.

The assessment of college learning² evolved through four eras: (1) the origin of standardized tests of learning: 1900–1933; (2) the assessment of learning for general and graduate education: 1933–47; (3) the rise of test providers: 1948–78; and (4) the era of external accountability: 1979–present. For ease of reference, the tests and testing programs discussed below are summarized in table form in the appendix.

The Origin of Standardized Tests of Learning: 1900–1933

The first third of the twentieth century marked the beginning of the use of standardized, objective testing to measure learning in higher education. The Carnegie Foundation led the movement; in 1916, William Learned tested students “in the experimental school at the University of Missouri in arithmetic, spelling, penmanship, reading, and English composition using recognized tests, procedures, and scales, and a statistical treatment that though comparatively crude was indicative” (Savage 1953, 284). E. L. Thorndike’s study of engineering students followed. Thorndike tested students at the Massachusetts Institute of Technology, the University of Cincinnati, and Columbia University on “all or parts of several objective tests in mathematics, English and physics” (Savage 1953, 285). These tests focused on content knowledge, largely tapping facts and concepts (declarative knowledge) and mathematical routines (procedural knowledge). The early tests were “objective”; students responded by selecting an answer where one answer was correct. Compared to the widely used essay examination, these tests gained reliability in scoring and content coverage per unit of time.

The monumental Pennsylvania Study, conducted between 1928 and 1932, emerged from this start. It tested thousands of high school seniors, college students, and even some college faculty members using extensive objective tests of largely declarative and procedural content knowledge. In many ways, the Pennsylvania Study was exemplary; it proceeded with a clear conception of what students should achieve and how learning should be measured. In other ways, however, it reflected its time; the study focused on knowledge and required compliant students to sit for hours of testing.³

In the 1928 pilot study, no less than 70 percent of all Pennsylvania college seniors, or 4,580 students, took the assessment as did about 75 percent of high school seniors, or 26,500 students. Of the high school seniors, 3,859 entered a cooperating Pennsylvania college; 2,355 remained through their sophomore year in college, and 1,187 remained through their senior year (Learned and Wood 1938, 211).

The assessment itself was a whopping twelve hours and 3,200 items long. (The examiners expressed regret at not being more comprehensive in scope!) Comprised of selected-response questions—for example, multiple-choice, matching, and true-false—the assessment covered nearly all areas of the college curriculum. The main study focused on student *learning*, not simply on achievement in the senior year, by testing students during their senior year of high school and then testing them again during their sophomore and senior years in college.

The Pennsylvania Study is noteworthy because it laid out a conception of what was meant by undergraduate achievement *and* learning, assuming that achievement was the result of college learning defined as the accumulation of breadth and depth of content knowledge. It also focused heavily and *comprehensively* at the knowledge level, especially on declarative and procedural knowledge. Nevertheless, because it included an intelligence test, the assessment program tapped the extremes of the outcomes framework: content knowledge and general intelligence. Moreover, the Pennsylvania Study employed technology for assessing learning and achievement—objective testing—that followed directly from the study’s conception of learning. If knowledge were understood as the accumulation of learning content, then objective testing could efficiently verify—literally index—the accumulation of that knowledge (Learned and Wood 1938, 372). Finally, the Pennsylvania Study is also noteworthy because, unlike assessments done today, it collected data in designs that provided evidence of both achievement and learning. In some cases, the comparison was across student cohorts, or “cross-sectional,” including high school seniors, college sophomores, and college seniors. In other cases, it was longitudinal; the same high school seniors tested in 1928 were tested again as college sophomores in 1930 and then as seniors in 1932.

The Assessment of Learning in General and Graduate Education: 1933–47

This era saw the development of both general education and general colleges in universities across the country, as well as the evolution of the Graduate Record Examination (GRE). The Pennsylvania Study had provided an existence proof; comprehensive assessment of student learning was feasible. Individual institutions, as well as consortia, put together test batteries designed primarily to assess cognitive achievement. Perhaps most noteworthy in this progressive-education period, with its focus on the whole student, was the attempt to measure not only cognitive outcomes across the spectrum but also the personal, social, and moral outcomes of general education.

Here, I briefly treat learning assessment in general education because, as it emerged in some important cases, it diverged from rather than adopted the Carnegie Foundation’s view of education and learning assessment. The University of Chicago’s approach presages contemporary developments; the Cooperative Study presages the call for “soft-skills.” I then focus attention on the evolution of the GRE.

General Education and General Colleges. The most notable examples of general education learning assessment in this era are the University of Chicago College program and the Cooperative Study of General Education (for additional programs, see Shavelson and Huang 2003). The former reflected thinking in the progressive era, while the latter had its roots in the Carnegie Foundation's conception of learning but also embraced progressive notions of human development as well.

In the Chicago program, a central university examiner's office, rather than individual faculty in their courses, was responsible for developing, administering, and scoring tests of student achievement in the university's general education program (Frodin 1950). Whereas the Pennsylvania Study assessed declarative and procedural knowledge, the Chicago examinations tested a much broader range of knowledge and abilities: the use of knowledge in a variety of unfamiliar situations; the ability to apply principles to explain phenomenon; and the ability to predict outcomes, determine courses of action, and interpret works of art. The Chicago comprehensive exams were characterized by open-ended essays and multiple-choice questions demanding interpretation, synthesis, and application of new texts (primary sources).⁴

The Cooperative Study of General Education, conducted by a consortium of higher education institutions, stands out from assessment initiatives at individual campuses. The participating institutions believed they would benefit from a cooperative approach to the improvement of general education (Executive Committee of the Cooperative Study in General Education 1947; Dunkel 1947; Levi 1948). To that end, and in order to assess students' achievement *and well-being*, the consortium developed the Inventory of General Goals in Life, the Inventory of Satisfaction Found in Reading Fiction, the Inventory of Social Understanding, and the Health Inventories.

The Evolution of the Graduate Record Examination: From Content to General Reasoning. While learning assessment was in full swing, Learned and Wood parlayed their experience with the Pennsylvania Study into an assessment for graduate education. In proposing the "Co-operative Graduate Testing Program," Learned noted that, with increased demand for graduate education following the Depression, the A.B. degree had "ceased to draw the line between the fit and the unfit" (Savage 1953, 288). Graduate admissions and quality decisions needed to be based on something more than the number of college credits.

In October 1937, Learned's team worked with the graduate schools at Columbia, Harvard, Princeton, and Yale to administer seven tests designed to index the quality of students in graduate education. This was the first administration of what was to be the Graduate

Record Examination (GRE). The program was a success and grew by leaps and bounds (see fig. 2). And at a time when the Carnegie Foundation was struggling to keep its faculty retirement system (TIAA) afloat, it was also a growing financial and logistic burden. Ultimately, the foundation was motivated by these stresses to pursue the establishment of an independent, national testing service.

Like the examinations used in the Pennsylvania Study, the original GRE was a comprehensive and objective test focused largely on students' content knowledge, but it also tapped verbal reasoning and was used to infer students' fitness for graduate study (Savage 1953). In 1936, a set of "profile" tests was developed to cover the *content* areas of a typical undergraduate general education program. To be completed in two half-day sessions totaling six hours, the tests measured knowledge in mathematics, the physical sciences, social studies, literature and fine arts, and one foreign language. The verbal factor was "developed primarily as a measure of ability to discriminate word meanings" (Lannholm and Schrader 1951, 7). In 1939, sixteen Advanced Tests in subject major fields were added to the GRE, and in 1949, a general education section was added to the Profile Tests in order to tap "effectiveness of expression" and to provide a "general education index" (see ETS 1953).

The fall of 1949 saw a landmark in student learning assessment: in a shift from testing content to testing general reasoning, ETS introduced a GRE Aptitude Test with the kind of verbal and quantitative sections we see today. Then, in 1952, it introduced the now standard

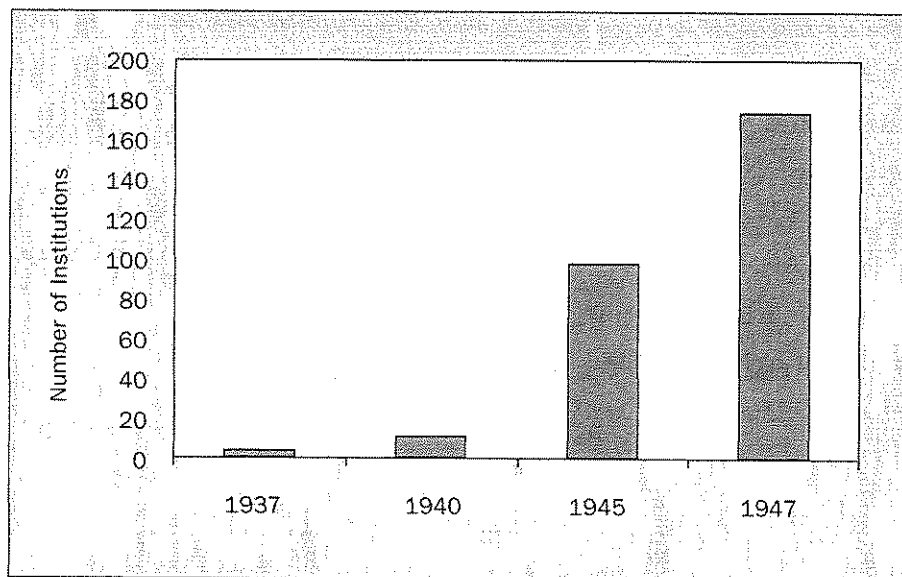


Figure 2. GRE growth over its first ten years.

scale for reporting scores (the normal distribution with mean 500 and standard deviation 100). In 1954, ETS continued the shift away from content and toward general reasoning by replacing both the Profile Tests and the Tests of General Education with the “Area Tests,” which served as a means of assessing broad outcomes of the liberal arts. The Area Tests focused on academic majors in the social and natural sciences and the humanities. They emphasized reading comprehension, understanding, and interpretation, often providing requisite content knowledge “because of the differences among institutions with regard to curriculum and the differences among students with regard to specific course selection” (ETS 1966, 3).

The Rise of the Test Providers: 1948–78

During the period following World War II, with funding from the G.I. Bill of Rights, postsecondary education enrollments mushroomed, as did the number of colleges to accommodate the veterans *and* the number of testing companies to assist colleges in screening them—most notably ETS, created in 1948, and the American College Testing (ACT) program, created in 1959.

Tests Provided by Testing Organizations to Assess Student Learning. By the time the Carnegie Foundation had transferred the GRE to ETS, completing its move out of the testing business, it had left an extraordinarily strong legacy of objective, group-administered, cost-efficient testing using selected response questions—now solely multiple-choice. That legacy has endured into the twenty-first century. The precursors of today’s major learning assessment programs were developed by testing organizations in this era (Shavelson and Huang 2003, 2006). These 1960s and 1970s testing programs included ETS’s Undergraduate Assessment Program, which incorporated the GRE, and ACT’s College Outcomes Measures Project (COMP). The former evolved via the Academic Profile into today’s Measure of Academic Proficiency and Progress (MAPP), and the latter evolved into today’s College Assessment of Academic Proficiency (CAAP).

However, several developments in the late 1970s, reminiscent of the progressive era, augured for a change in the course set by Learned and Wood. Faculty members were not entirely happy with multiple-choice tests. They wanted to get at broader abilities—such as the ability to communicate, think analytically, and solve problems—in a holistic manner. This led to several new developments. ETS studied constructed-response tests that tapped communication skills, analytic thinking, synthesizing ability, and social/cultural awareness (Warren 1978).

ACT experimented with open-ended performance-based assessments that sought to measure skills for effective functioning in adult life in social institutions, in using science and technology, and in using the arts. And the state of New Jersey developed Tasks in Critical Thinking, which sampled real-world tasks in a “performance-based assessment . . . [that measured] the ability to use the skills of inquiry, analysis, and communication” with prompts that “do not assess content or recall knowledge” (ETS 1994, 2). These assessment programs were designed to embrace what college faculty considered to be important learning outcomes.

For a short period, these learning assessments set the mold. But due to time and cost limitations, as well as difficulties in securing and training people to score responses and in achieving adequate reliability, they either faded into distant memory or morphed back into multiple-choice tests. For example, the COMP began as a pathbreaking performance-based assessment. Its content was sampled from materials culled from everyday experience including film excerpts, taped discussions, advertisements, music recordings, stories, and newspaper articles. The test sought to measure three process skills—communicating, solving problems, and clarifying values—in a variety of item formats, including multiple-choice, short answer, essay, and oral response (an atypical format). COMP, then, bucked the trend toward multiple-choice tests of general abilities by directly observing performance sampled from real-world situations.

The test, however, was costly in terms of time and scoring. Students were given six hours to complete it in the 1977 field trials; the testing time was reduced to four and a half hours in the 1989 version. Raters were required to score much of the examination. As both a consequence and a characteristic of trends, a simplified “Overall COMP” was developed as a multiple-choice-only test. In little more than a decade, however, this highly innovative assessment was discontinued altogether due to the costliness of administration and scoring. Roughly the same story describes the fate of Tasks in Critical Thinking (see Erwin and Sebrell 2003).⁵

The influence of the Carnegie Foundation, then, waned in the mid-1970s. However, as we shall see, the foundation’s vision of objective, selected-response testing continued to influence the standardized learning assessment programs of ETS, ACT, and others.

The Era of External Accountability: 1979–Present

By the end of the 1970s, political pressures to assess student learning and hold campuses accountable had coalesced. While in the 1980s only a handful of states had some form of mandatory standardized testing (e.g., Florida, Tennessee), public and political demand for

such testing increased into the new millennium (Ewell 2001). To meet this demand, some states (e.g., Missouri) created incentives for campuses to assess learning, and campuses responded by creating learning assessment programs.

Tests of College Learning. ETS, ACT, and others were there to provide tests. Indeed, a wide array of college learning assessments following in the tradition of the Carnegie Foundation was available. Currently, ETS provides the MAPP, ACT provides CAAP, and the College Resource Center at the University of Missouri, Columbia, offers the College Basic Academic Subjects Examination (CBASE). All are multiple choice test batteries. MAPP measures college-level reading, mathematics, writing, and critical thinking in the context of the humanities, social sciences, and natural sciences to enable colleges and universities to improve their general education outcomes. CAAP measures reading, writing, mathematics, science, and critical thinking to enable postsecondary institutions to evaluate and enhance general education programs. CBASE is a criterion-referenced achievement examination of English, mathematics, science, and social studies that serves both to qualify individuals for entry into teacher education programs and to test general academic knowledge and skills.

Vision for Assessing Student Learning. As we saw, at the end of the 1970s, objective testing was incompatible with the way faculty members either assessed student learning or wanted student learning to be assessed. For them, *life is not a multiple-choice test*. Rather, faculty members like the open-ended, holistic, problem-based assessments exemplified by, for example, Tasks in Critical Thinking. Intuitively, faculty members suspected that the kind of thinking stimulated and performance assessed by multiple-choice and other highly structured tests is different from that stimulated and assessed by more open-ended tasks. And empirical evidence supports their intuition.

While a multiple-choice test and a “constructed-response” test may produce scores that are correlated with each other, this correlation does not mean that the same kind of thinking and reasoning is involved (Martinez 1999; National Research Council 2001). Student performance varies considerably depending upon whether a task is presented as a multiple-choice question, an open-ended question, or a concrete performance task (Baxter and Shavelson 1994). For example, Lythcott (1990, 248) found that “it is possible, though not our intention, for [high school and college] students to produce right answers to chemistry problems without really understanding much of the chemistry involved.” Moreover, Baxter and Shavelson (1994) found that middle school students who solved electric circuit problems hands-on could not solve the same problems represented abstractly in a multiple-choice test; these students did not

make the same assumptions that the test developers made. Finally, using a “think aloud” method to tap into students’ cognitive processing, Ruiz-Primo and colleagues (2001) found that students reasoned differently on highly structured assessments than on loosely structured assessments. In the former case students “strategized” as to what alternative fit best, while in the latter they reasoned through the problem.

To illustrate the difference between multiple-choice and open-ended assessments, consider the following concrete example from the Collegiate Learning Assessment (CLA). College students are asked to assume that they work for “DynaTech”—a company that produces industrial instruments—and that their boss has asked them to evaluate the pros and cons of purchasing a “SwiftAir 235” for the company. Concern about such a purchase has risen with the report of a recent SwiftAir 235 accident. When provided with an “in-basket” of information, some students, quite perceptively, recognize that there might be undesirable fallout if DynaTech’s own airplane crashed while flying with DynaTech instruments. Students are not prompted to discuss such implications; they have to recognize these consequences on their own. There is no way such insights could be picked up by a multiple-choice question.

Finally, consistent with the views of faculty, both the American Association of State Colleges and Universities (AASCU) and members of the Spellings Commission have in mind a particular standardized learning assessment: the CLA.

The best example of direct value-added assessment is the Collegiate Learning Assessment (CLA), an outgrowth of RAND’s Value Added Assessment Initiative (VAAI) that has been available to colleges and universities since spring 2004. The test goes beyond a multiple-choice format and poses real-world performance tasks that require students to analyze complex material and provide written responses (such as preparing a memo or policy recommendation). (AASCU 2006, 4)

This brief history has now arrived at the present. In contrast to the evolution in multiple-choice testing technology, the Council for Aid to Education has taken Tasks in Critical Thinking and COMP to what might be considered the next level by marrying the open-ended assessment of real-world, holistic tasks and the use of computer technology to assess ability and learning.⁶ A closer look at the CLA may provide insight into a next generation of learning assessments.