

Why Have Educational Evaluators Chosen Not to Do Randomized Experiments?

By
THOMAS D. COOK

This article analyzes the reasons that have been adduced within the community of educational evaluators for not doing randomized experiments. The objections vary in cogency. Those that have most substance are not insurmountable, however, and strategies are mentioned for dealing with them. However, the objections are serious enough, and the remedies partial enough, that it seems hardly warranted to call experiments the “gold standard” of causal inference. Yet even if they are not perfect in research practice, this article shows how they are logically and empirically superior to all currently known alternatives. The article particularly addresses the objection that school personnel will not accept experiments. It shows that hundreds of them have been done there by researchers with backgrounds in psychology and public health who study the prevention of unhealthy behaviors. But experiments are much rarer among researchers trained in education who study changing academic performance. Reasons are adduced for this difference in academic culture within school-based research.

Keywords: experiments; education; research culture

Calls are heard to improve schools through innovations as diverse as school-based management, charter schools, vouchers, more effective teaching practices, higher standards, increased accountability, smaller schools, smaller classes, new technologies, and better trained teachers. Many “brand-name” reforms also exist, including Slavin’s Success for All, Levin’s Accelerated Schools, Comer’s School Development Program, Sizer’s Coalition of Essential Schools, Lezotte’s Effective Schools, and Total Quality Management Schools. Claims also abound about other educational strategies like revisions to special education and bilingual programs; more

Thomas D. Cook is the John Evans Professor of Sociology at Northwestern University. He also has courtesy appointments in psychology and in education and social policy and is a faculty associate of the Institute for Policy Research. His major interest is in the theory and practice of evaluation, but he also publishes on issues of adolescent development, particularly in high-risk settings. He is a fellow of the American Academy of Arts and Sciences and of the American Academy of Political and Social Sci-

DOI: 10.1177/0002716203254764

phonics in early grades; constructivist learning that treats students as active learners; an end to social promotions; and better integration of schools with preschools, families, and after-school activities. And these are only a subset of all the reform proposals recently made.

Most of these ideas have never been seriously evaluated to learn how they affect student performance. This is surprising since there is a profession of educational evaluators, and one might presume its members to be interested in identifying “what works” in schools. Most of these evaluators work within schools of education where their interests overlap with those of substantive researchers who want to learn what will enhance student performance so as to construct better theories and impact on school practice. Some academics outside of education schools also do evaluations in educational settings, as do researchers in private firms that contract with federal, state, and local education authorities.

All these researchers have access to the generally preferred methods for learning what works. These methods depend heavily on the quality of measurement and design, and measuring individual student change and educational practice at both the classroom and school levels are long-recognized strengths of educational research. Design-wise, the randomized experiment is widely known as the best tool for attributing observed student change to whatever classroom or school option is under consideration as a possible cause of the observed change.

Random assignment entails using the equivalent of a fair coin toss to create two or more initially equivalent groups. The option under consideration (“the treatment”) is then assigned to one group, while the other group is exposed to something else—often no explicit treatment but sometimes a qualitatively different one. If an experiment is properly maintained over time, any observed group differences at the end of a study can be reasonably attributed to the intervention. They are certainly not likely to be due to selection—differences in the average person in each group—for the assignment process renders these last differences unlikely. Control groups generated through random assignment provide the best counterfactual to describe what would have happened to students in a treatment group if they had not been exposed to the treatment (Rubin 1974; Holland 1986).

This theoretical rationale for random assignment is complemented by an empirical justification. In the past, results from individual experiments have been contrasted with the results from the major design and statistical alternatives to such experiments. Different effect sizes were found for the experiments versus the nonexperiments (Mosteller, Gilbert, and McPeck 1980; LaLonde 1986; Fraker

ences and has received the Gunnar Myrdal Award of the American Evaluation Association and the Donald T. Campbell Award of the Policy Studies Organization as well as a Distinguished Scientist Award from the American Psychological Association. He is a trustee of the Russell Sage Foundation in New York and of the Textile Museum in Washington, D.C.

NOTE: Thanks for feedback on an earlier draft are due to Tony Byrk, Lee Cronbach, Joseph Durlak, Christopher Jencks, Paul Lingenfelter, and two unnamed reviewers. The financial support of the MacArthur, Rockefeller, and Spencer Foundations is gratefully acknowledged. Parts of this article were read as the Jerry Lee Lecture at the University of Pennsylvania.

and Maynard 1987). This finding also holds in more recent work that used even more sophisticated statistical control alternatives or design alternatives with comparison groups that are presumptively better than the comparisons in the first wave of research that were constructed from national survey data rather than from non-equivalent groups physically close to the treatment groups (e.g., Agodini and Dynarski 2002; Bloom et al. 2002; Friedlander and Robins 1995; Heckman et al. 1997; Wilde and Hollister 2002). The lone exception to this finding about the invalidity of results from nonexperiments comes from Dehejia and Wahba (1999), but Smith and Todd (2002) have shown that this exception has little generalizability. The alternatives to random assignment empirically considered to date do not reliably reproduce the results of experiments, though these alternatives have not yet included the regression discontinuity, interrupted time series, and matched cohort designs that Shadish, Cook, and Campbell (2001) treated as the strongest alternatives to the randomized experiment. Even so, since logic also supports a stronger causal warrant for random assignment than its alternatives, we have to conclude that individual nonexperiments tend to provide more biased causal conclusions than do individual experiments.

However, there have been three instances where a series of experiments on a topic was contrasted with a series of nonexperiments on the same topic (Lipsey and Wilson 1993; Bloom et al. 2002; Glazerman et al. 2002). In each case, the same average effect size was found across the experiments as nonexperiments, implying that the positive and negative biases from nonexperiments have cancelled each other out exactly. Since some scholars believe that the usual unit of progress in science is the review and not the single study, they may place weight on the correspondence of average findings from experiments and nonexperiments more than on the noncorrespondence of findings between single experiments and nonexperiments. But I caution that only three empirical reviews of cumulative bias exist to date and that the Boom et al. (2002) and the Glazerman et al. (2002) work cover many of the same experiments. Given the absence of strong theory about the conditions under which biases sum to zero, we cannot guarantee in any particular instance that the positive and negative biases will exactly counterbalance. Moreover, in the three reviews the standard deviation across all the nonexperimental effect sizes is considerably larger than the standard deviation across all the experimental effect sizes. This suggests that experiments are more efficient than nonexperiments. Since they get to the same answer sooner, they are particularly important in fields where few experiments already exist. As I will demonstrate shortly, education is such a field. So, empirical research on the results of experiments and their alternatives suggests, first, that individual experiments are less biased and, second, that as studies of a topic accumulate, experiments are more efficient.

This implies a pragmatic justification for experiments over and above the logical and empirical warrants detailed above. Experiments are probably less expensive in the long run because, being more efficient about reducing causal uncertainty, fewer of them are needed for the same degree of confidence in the causal conclusion drawn. This can still be true even if individual experiments were to turn out on average to be more expensive than their nonexperimental counterparts.

A second pragmatic justification is even less speculative. What if policy elites incorrectly concluded that Catholic schools are superior to public ones, and did something about this in the policies they created? What if they erroneously concluded that vouchers stimulate academic achievement, and did something about this in terms of funding priorities? What if they falsely concluded that school desegregation does not affect minority achievement when it does, and acted accordingly? Incorrect causal conclusions have costs in terms of dollars, achievements, and dreams.

The superiority of random assignment for drawing inferences about the consequences of planned change attempts is routinely acknowledged in philosophy, medicine, public health, agriculture, statistics, microeconomics, psychology, criminology, prevention research, early childhood education, marketing, and those parts of political science and sociology concerned with improving opinion surveys. It is also acknowledged in all the elementary education method textbooks I have consulted. However, random assignment is relatively rare in educational research, especially for assessing the impact of educational interventions of obvious policy relevance.

Education is not unique in a low base rate of experimentation. Random assignment is also rare in sociology, political science, macroeconomics, and management. Yet causal statements are routinely made in these fields, usually through a process that links substantive theory to various qualitative or quantitative nonexperimental practices. This article does not argue that correct causal conclusions come only from experiments. It does argue, though, that experiments provide a better warrant for such conclusions than any other method. So if experiments can be conducted in schools, they should be. Not to use them requires a very strong justification.

Over the past thirty years, self-ascribed educational evaluators like Alkin, Cronbach, Eisner, Fetterman, Fullan, Guba, House, Huberman, Lincoln, Miles, Provus, Sanders, Schwandt, Stake, Stufflebeam, and Worthen have proposed many justifications for not doing experiments. These theorists want educational evaluation to pursue goals other than describing what works in schools. Most of them want evaluation to improve the organization and management of individual districts or schools, assuming that this will routinely improve student performance. So, they examine ways to provide individual schools or district staff with continuous feedback about strategic planning, program implementation, or student and teacher performance monitoring. The expectations are that local officials will immediately use this feedback and that performance in their schools will consequently improve. This model of particular forms of research and their connection to organizational change is much like what is found in management consulting in the private sector.

Other educational evaluators want evaluation to contribute to developing general theories, especially those that specify the often complex constellation of forces that bring about important school effects. Such evaluators particularly value the identification of generative processes that are effective over a broad set of circumstances, much like engaged time on task. This is a highly general cause of enhanced

academic achievement and can be instantiated as more days of schooling per year, as longer school days, as more time devoted to the “core curriculum,” as textbooks that are engaging, as exposure to teachers who know how to motivate students, and so on. Unfortunately, neither of these two evaluation priorities places the premium where experimentation does—on directly observing student change and unambiguously attributing it to a single policy-related treatment.

This article probes the validity of the intellectual arguments that educational evaluators have adduced for not doing experiments and for taking evaluation in directions other than identifying the effects of circumscribed causal agents of relevance to educational policy. The account I offer places little emphasis on the political and organizational factors within the federal system of support for educational research, factors that Vinovskis (2002) emphasized in his explanation of the paucity of experiments. The emphasis here is on examining the justifications offered by self-ascribed educational evaluators. To the extent that these overlap with the rationales offered by their more substantively oriented colleagues in schools of education, I also address the latter’s objections to doing experiments. Although the specific reasons for downplaying experiments vary across educational evaluators, the total set of arguments can be divided into five types:

1. *Philosophical arguments* designed to show that experiments (a) cannot provide unbiased tests of causal hypotheses and (b) are predicated on a descriptive theory of causation that is less useful than explanatory theories of cause.
2. *Practical arguments* asserting that experiments (a) can rarely be mounted in schools, and when mounted, are often imperfectly realized because of compromises to (b) the planned treatment contrasts and (c) the quality of individual treatment implementation.
3. Arguments about *undesirable trade-offs* because experiments (a) sacrifice external for internal validity and (b) value causal conclusions so highly that a conservative bias results which overlooks useful findings indicated by more liberal criteria.
4. Arguments that *schools will not use experimental results* because (a) experiments meet the interests of federal and state policy makers who are not major actors in educational policy and (b) the experiment’s logic recreates a rational decision-making model that does not describe how schools actually make decisions.
5. Arguments that *experiments are not necessary because better alternatives exist*, including (a) the intensive qualitative case studies preferred by self-styled educational evaluators, (b) the quasi-experiments conducted by substantive researchers who prefer design control over statistical control, and (c) the causal modeling preferred by substantive researchers who do longitudinal work in schools.

Any of the points above casts doubt on the wisdom or practicality of experimenting in schools, and all of these objections have been raised at one time or another by someone in the community of educational evaluation specialists who operate out of schools of education. So the number and variety of the arguments confers a genuine intellectual integrity, making it important to deal with each argument in turn, both to examine its validity and to assess its implications for creating a practical theory of school-based experimentation. Many advocates of experimentation recognize as legitimate only a much smaller array of objections—namely, those that follow only from the process of random assignment. This is because, on close examination, the other objections turn out to refer to the intellectual, political,

social, and economic contexts in which random assignment is embedded, and as such, they refer to other kinds of research as well as experimentation. For instance, some of the philosophical objections pertain to all quantification and not just to experimentation as one type of quantified social science; and questions about treatment-related attrition apply to longitudinal quasi-experiments as well as to experiments. Since my starting point is the objections to random assignment used by a particular group of scholars, I prefer to cast a wider net than the one that advocates of experimentation view as legitimate. The objections I consider are not just those that are truly unique to experimentation; they are those that a particular research community has raised against using random assignment in schools.

Philosophical Beliefs Adduced to Reject Random Assignment

Random assignment is epistemologically discredited

To philosophers of science, positivism connotes a rejection of realism, the formulation of theories in mathematical form, the primacy of prediction over explanation, and the belief that entities do not exist independently of their measurement. Although this epistemology has long been discredited, many educational researchers still use the term “positivism” but connote something less historically precise—namely, quantification and hypothesis testing, both central to experimentation. Kuhn’s (1970) work is at the forefront of their reasons for rejecting positivist science. He argued that theories are “incommensurable”—that is, their postulates cannot be formulated as specifically as philosophical theories of verification or falsification require and are always subject to reinterpretation. He also argued that observations are inevitably “theory-laden”—that is, they are impregnated with researchers’ theories, hopes, wishes, and expectations, thus undermining their neutrality for discriminating between truth claims. In refuting the possibility of totally explicit theories and totally neutral observations, Kuhn’s work seems to undermine science in general and experimentation in particular. Tending in the same direction are the views of other philosophers that educational evaluators like to cite, such as Lakatos, Harre, and Feyerabend. Also relevant is their citation of descriptions showing how often scientists’ behavior in the laboratory deviates from the very scientific norms they espouse (e.g., Latour and Woolgar 1979). All these sources are meant to indicate that science is an emperor without clothes.

However, the critique is overly simple. Even if observations are never theory-neutral, many observations have stubbornly reoccurred whatever the researcher’s predilections. As theories replace each other, most fact-like statements from the older theory are incorporated into the newer one, surviving the change in theoretical superstructure. So, even if there are no “facts” we can independently know to be certain, there are still many propositions with such a high degree of facticity that they can be confidently treated as though they were true. For practicing experi-

menters, the trick is to build multiple theories into how the data are collected, especially the perspectives of theoretical opponents. Independent replications are particularly important, therefore, provided they do not share bias in the same direction (Cook 1985). Kuhn's (1970) work complicates what a "fact" means but does not deny that some claims to fact-like status are strong.

Kuhn (1970) is also correct that theoretical statements are never definitively tested (Quine 1951, 1969). But this does not mean that individual experiments fail to probe theories and the causal hypotheses they generate. When an experiment produces negative results, its advocates are not likely to accept the disappointing result. Instead, they invoke various methodological and substantive contingencies that might have changed the result—perhaps if a different outcome measure had been used or if a different population had been examined. Subsequent studies can then probe these contingency formulations. If the results again prove negative, this might lead to an even more complicated contingency hypothesis designed to explain the latest disconfirmation. A test of this revised hypothesis can then take place, and so on. After a time, this process runs out of steam, so particularistic are the contingencies that remain to be examined. The consensus seems to emerge that

the program could not be shown to be effective under the many different conditions examined. Other conditions could still be probed. But they are so circumscribed that the reform will not be worth much even if it is effective under these conditions.

Kuhn is correct that this process is social and not exclusively logical, and he is further correct that the predicament arises because program theory is not explicit enough to be definitively confirmed or rejected. But the reality of elastic theory does not mean that decisions about causal hypotheses are only social or that they are devoid of all empirical or logical content.

Experiments are predicated on an overly simple theory of causation

Randomized experiments test the impact of only a small subset of potential causes from within the world, often a single one. And at their most elegant, they can responsibly test only a modest number of interactions between treatments. So randomized experiments are best when a causal question is simple, sharply focused, and easily justified. The theory of causation most relevant to this is variously called the manipulability, activity, or recipe theory (Collingwood 1940; Gasking 1955; Whitbeck 1977). It seeks to describe the consequences of a set of activities that can be listed as though they were recipe ingredients and can be actively manipulated as a whole to ascertain what effects the lumped manipulation has. The aim here is to describe the effects of a given cause.

However, the most esteemed theories of cause seek to ascertain the causes of a given effect. They want to explain rather than to describe "if-then" connections. One explanatory theory emphasizes identifying "generative processes" (Bhaskar

1975; Harre 1981). These are forces that bring about effects in a wide variety of circumstances, such as gravity as it affects falling, or a specific genetic defect as it induces phenylketonuria, or time on task as it facilitates learning. However, as simple as these examples seem to be, they are replete with hidden causal contingencies. Thus, the genetic defect does not induce phenylketonuria if an appropriate diet is adopted early in life, and time on task does not induce learning if a student is disengaged or the curriculum meaningless. So, a second understanding of causation requires specifying all the contingencies (co-causes) that impact on an effect, including those that follow from a causal manipulation but are prior to the effect. Yet experiments were not designed for this purpose. They were designed to describe the effects of a multidimensional set of activities deliberately manipulated as a package. Experiments are only explanatory if the manipulations are chosen to help discriminate between competing theories; or if the processes mediating between a cause and effect are specified and measured; or if effect sizes vary in systematic ways across outcomes, populations, or settings.

Cronbach and his colleagues (1980) maintained that explanatory theories of cause are more relevant to schools than the activity theory. They believe that multiple causal factors are implicated in all student or teacher change, and often in complex ways. Their model of real-world causation is system related, more akin to intersecting pretzels than to the experimenter's simple arrow from A to B. They rightly believe that no educational intervention fully explains an outcome; at most, it is just one more determinant of that outcome. Nor are effect sizes constant across student and teacher populations or across school types and times. Causal contingency is the watchword, not simple generalization. Given this priority, the activity theory seems irrelevant since so few variables causally implicated in an effect can be simultaneously manipulated. Believing that experiments cannot faithfully represent a real world of multivariate, nonlinear, and often reciprocal causation, Cronbach and Snow (1976) searched for aptitude-treatment interactions—specifications that a treatment's effect depends on student or teacher characteristics. But they discovered few that were robust. This might reflect an ontological truth—nature is not as contingently ordered as many theorists think. But it might also reflect the methodological problems associated with testing statistical interactions—for example, underspecified theory, partially valid measures, imperfectly implemented treatments, truncated distributions, and noninterval scales. By itself, empirical research cannot inform us whether the real world is more or less contingently ordered than critics of the experiment contend.

Even so, the activity theory is clearly limited as a theory of causation. But to be limited is not to be useless. Notice how many useful conclusions about effective educational practices are today specified without qualification. For instance, small schools are better than large ones, time on task raises achievement, summer school raises test scores, school desegregation hardly affects achievement, and assigning and grading homework raises achievement.

Most educational researchers who espouse a highly contingent theory of causation nonetheless seem willing to accept some noncontingent causal statements.

Critics of the experiment also seem to accept some minimally contingent statements—for example, reducing class size increases achievement, provided that the amount of change is “sizable” and to a level under twenty. Or, Catholic high schools increase graduation rates over public schools, but only in the inner city. Commitment to an explanatory theory of causation has not stopped researchers from acting as though some educational change attempts result in dependable main effects or simple interactions.

Some causal contingencies are irrelevant to educational policy even if they are relevant to full explanation. For policy, the most important contingencies are those that, within normal ranges, modify the sign of a causal relationship and so help

Design-wise, the randomized experiment is widely known as the best tool for attributing observed student change to whatever classroom or school option is under consideration as a possible cause of the observed change.

identify where a treatment is sometimes harmful. Less important are those interactions describing benefits that are generally positive even if they are more positive in some circumstances than others. Policy makers are often constrained in what they can do and cannot assign different treatments to different populations. But they are willing to advocate broad changes with differential effects provided that these effects are rarely negative. When policy concerns are paramount, it is possible to ignore many variables that genuinely contribute to fuller explanation.

I appreciate the arguments of those opponents of experimentation who believe that biased answers to big explanatory questions are more important than unbiased answers to smaller casual-descriptive questions. I also agree with them that random assignment depends on a less comprehensive and less esteemed theory of causation. But acknowledging these limitations does not undermine the justification for experiments. It is still necessary to know about the effects of given causal agents. This is not a trivial knowledge need. And acknowledging the limitations of the activity theory of causation should embolden researchers to design future experiments with a greater explanatory yield than in the past. At a minimum, this means greater sensitivity to identifying moderator and mediating processes and thus building into experiments the sampling and measurement particulars that such sensitivity requires. No more black box experiments.

Practical Reasons for Not Doing Randomized Experiments

Randomized experiments cannot be mounted

Education researchers were at the forefront of the flurry of social experimentation at the end of the 1960s and during the 1970s. Evaluations of Head Start (Cicirelli and Associates 1969), Follow Through (Stebbins et al. 1978), and Title I (Wargo et al. 1972) found few, if any, positive effects. These disappointing results engendered considerable dispute about methods, and many educational evaluators concluded that quantitative evaluation had been tried and failed. So they turned to other methods. Other scholars responded by stressing the need to study school management and program implementation, believing them to be the reasons why results were so disappointing (Berman and McLaughlin 1977; Elmore and McLaughlin 1983; Cohen and Garet 1975).

However, the educational studies most often criticized during this period did not involve random assignment. Indeed, I know of only three randomized experiments on educational topics of policy relevance then available—studies of Sesame Street (Bogatz and Ball 1972), of the Perry Preschool Project (Schweinhart, Barnes, and Weikart 1993), and of only twelve youngsters randomly assigned to a desegregated school (Zdep 1971). Since only the Zdep (1971) study took place in schools, it is not accurate to claim that policy-relevant randomized experiments had been tried in education and had failed. Indeed, to critique randomized experiments, Cronbach et al. (1980) had to reanalyze studies that had nothing to do with schools.

Nonetheless, many district officials do not like the focused inequities in school resources that random assignment generates, fearing negative reactions from parents and staff. They prefer it when individual schools choose the changes they will make, or when changes are districtwide. Principals and other school staff have similar preferences and have additional concerns about disrupting ongoing routines. Also, ethical concerns are often raised about withholding potentially helpful treatments; and some programs are meant to be universal under law, thus precluding the use of no-treatment control groups. Are experiments so unpopular and impractical that they cannot be used to study the effects of school improvement attempts?

It is manifestly false that experiments cannot be done within schools. I shall not consider here the small and highly controlled experiments done by cognitive scientists since these tend to be of little immediate policy relevance. On other topics, there are some experiments, but not many. Looking at the policy areas specified in the first part of this article's first paragraph, I could find no experiments on standards setting. The literature on effective schools reveals no experiments systematically varying the school practices that correlational studies suggest are effective. Recent studies of school-based management reveal only two randomized experiments, both on Comer's School Development Program (Cook et al. 1999; Cook, Hunt, and Murphy 2000), but not on other kinds of whole-school reform. There

seem to be no experiments on Catholic or Accelerated or Total Quality Management schools. On vouchers, there is a study by Witte (1998), sometimes reanalyzed (Greene et al. 1996; Rouse 1998), plus a program of research by Peterson and colleagues in several sites (Howell and Peterson 2002). On charter schools, I know of no relevant experiments. On smaller class sizes, there are six experiments, the best known being the Tennessee class size study sometimes called Project Star (Finn and Achilles 1990; Mosteller, Light, and Sachs 1996, Krueger 1999, Krueger and Whitmore 2001). On smaller schools, I know of only one randomized experiment (Kemple 2001), though there is also an experiment on alternative schools in Stockton, California (Dynarski and Gleason 1998). On teacher training, I know of no relevant experiments, though experiments on dropout prevention do exist for middle and high schools (Dynarski and Wood 1997). The obvious conclusion here is that current knowledge of effectiveness depends heavily on methods less esteemed than random assignment.

It is particularly striking that only two of the policy experiments above were conducted by researchers trained in schools of education or currently so affiliated. The best-known class size experiment was done by educators but was popularized by statisticians and reanalyzed by economists. The Milwaukee voucher experiment was done by political scientists and reanalyzed as a randomized experiment by political scientists and economists. The Comer studies were conducted by sociologists. The research on academies within high schools was done by an education-trained researcher working at Manpower Demonstration Research Corporation, a contract research firm with a strong economics background. The work on school choice was done by political scientists and the work on alternative schools and dropout prevention was done by economists in Mathematica, Inc., a contract research firm. Among those within educational evaluation who do experiments, it is rare to find individuals who call themselves educational evaluation specialists and who operate from schools of education.

To further illustrate the paucity of experiments in educational evaluation, Nave, Miech, and Mosteller (1999) reported that not even 1 percent of the dissertations in education or of the studies archived in ERIC Abstracts involved random assignment. Casual reading of the major journals on school improvement (*American Educational Research Journal* and *Educational Evaluation and Policy Analysis*) tells a similar story.

So does consideration of some reviews on specific substantive topics. The National Institute of Child and Human Development's congressionally mandated National Reading Panel reviewed nearly two thousand published studies on the effects of phonemics (Ehri, Nunes, Willows et al. 2001). Of these two thousand, only fifty-two experiments and quasi-experiments met the criteria for inclusion. If the ratio of experimental to quasi-experimental comparisons reported in the meta-analysis is the same as the ratio of such studies, then about twenty-three of these two thousand studies would have been randomized experiments—slightly more than 1 percent. In a related study of phonics (Ehri, Nunes, Stalh, et al. 2001), thirty-eight experiments and quasi-experiments were found, and of these, fourteen were randomized experiments. But since the number of studies claiming to be

about the causal effects of phonics was not identified in the review, it is not possible to determine what percentage of all studies were experiments. Even on topics like the effects of homework, it is clear that some randomized experiments have been done (Cooper 1989). So while experiments can be done on many pedagogic topics, the reality is that they are relatively rare.

Yet random assignment is very common in schools when the topic is not pedagogic. For instance, school-based prevention studies are designed to improve student health; to prevent school violence; or to reduce teen use of tobacco, drugs, and alcohol (e.g., Cook, Anson, and Walchli 1993; Peters and McMahon 1996; Durlak and Wells 1997a, 1997b, 1998). The reviews by Durlak and Wells of prevention studies prior to 1991 included about 190 randomized experiments and 120 other studies. The number of experiments has certainly increased since then, given the rapid growth of prevention research. So school-based experiments are common when the topic involves preventing negative behaviors.

Experiments are also common in preschool education. They were used with the Ypsilanti Perry Preschool Program (Schweinhart, Barnes, and Weikart 1993), the Abecedarian Project (Ramey and Campbell 1991; Campbell and Ramey 1995), the Comprehensive Child Development Program (Goodson et al. 2000), Olds's home nurse visiting program (Olds et al. 1997), Early Head Start (Raikes and Love 2002), and even the new Head Start study (Cook and Puma 2002). While preschool experiments generally take place in child care centers or homes rather than schools, instruction is usually evaluated, and changes in cognition and social behavior are the major outcomes, just as they are in school-based pedagogic studies.

So, experiments are common in schools if the topic is preventing negative behavior but not if the outcomes are more traditionally educational. And experiments are common at the preschool but not the school level, even though student cognitive performance is routinely assessed in the preschool studies. This raises several interconnected questions. Why are experiments on policy-relevant educational topics so rare? Why have most of the school reform experiments conducted to date been done by contract researchers or academics not in schools of education when one might expect them to be done by the many researchers who call themselves educational evaluators and who operate out of schools of education? Why is there this disciplinary difference in the likelihood of conducting experiments?

One possibility touches on subject matter. The prevention experiments tend to last less than a year, they do not involve changes in major school routines, they do not threaten the performance in math and language arts by which schools are held accountable, and the implementation is usually done by researchers rather than teachers. In contrast, pedagogical interventions are more likely to be multiyear; teachers are more often asked to deliver the treatment that entails changes in their established routines; and the threat exists that performance in core competencies may not rise because of the intervention, compromising accountability goals.

Arguing against this interpretation are two things. First is a personal report to the author by Durlak that some of the 190 prevention experiments he reviewed involve multiyear interventions, whole school changes, and teachers delivering the treatment, though none speak to competencies in core academic areas. (However,

it is not clear how many prevention studies combine all of the features that differentiate educational policy from prevention studies.) Moreover, some pedagogic reforms do not require multiyear and whole school efforts. Yet even under these—the most propitious—circumstances for experimentation, the low base rate of pedagogic experiments documented earlier suggests they are still rare.

The second explanation for the discipline-based difference in the frequency of experiments invokes political will and disciplinary culture. Random assignment is common in the health sciences because it is institutionally supported there by funding agencies, publishing outlets, graduate training programs, the clinical trials tradition, and practices in government health action agencies. Prevention studies in schools tap into a similar research culture, as does preschool education, where congressional mandates play an ancillary role. Moreover, prevention and preschool researchers tend to be trained in psychology, human development, public health, and microeconomics—fields that value experimentation. Gueron (2002) has emphasized how important sponsor and researcher commitment are for getting experiments mounted.

Contrast the norms and structures above with the situation in education. Reports from the Office of Educational Research and Improvement (OERI) are supposed to identify effective school practices. But neither the work of Vinovskis (1998) nor my own haphazard reading of OERI reports suggests any privilege accorded to random assignment. Moreover, one recent report I read on bilingual education repeated old saws about the impossibility of randomizing and claimed that alternatives are just as good—in this case, poorly designed quasi-experiments. And at a recent foundation meeting on teaching and learning, the representative of a reforming state governor spoke about a list of best practices being disseminated to all schools in his state. He did not care, and he believed that no governors cared, about the technical quality of the research. His main concern was that there was consensus among education researchers about each practice. When asked how many of the recommended practices depended on evidence from randomized experiments, he guessed it would be none. Several nationally known educational researchers were also present, and they all agreed that such assignment probably played no role in generating the practices on the list. No one felt any distress at this. As long as such beliefs and feelings are widespread, there will never be the pan-support for experimentation in education that is currently found in health, agriculture, or health in schools.

There has been considerable recent concern in parts of the Office of Education to change this situation and to do more experiments. Several large ones are now on the drawing board. Time will tell how much comes of this new priority and what roles professional educational evaluators will play in carrying out such studies as opposed to the contract research firms and public policy institutes doing them. Government funding priorities are only one source of professional norm setting in education. So, changes in federal funding cannot coerce broad-scale practice changes within a research community that is also supported by the money and prestige associated with teaching and tenure and by research funding from nonfederal sources. Large-scale policy experiments on educational topics will be

more common now that the federal climate has changed, but it is not clear whether the community of educational evaluators will follow suit.

Principal ignorance of random assignment is not a likely cause of the low frequency of pedagogic experiments, though principals usually do prefer other alternatives. Research contexts aside, when new programs are announced in schools, the demand for places sometimes exceeds the supply. In this situation, principals often resort to random assignment to determine who gets a place, being afraid of parental or staff reactions if they allocated slots by merit, need, “first come, first served,” or teacher recommendation. Like other politicians, principals understand the benefits of random assignment when resources are limited, even if they rarely follow up on the situation to study the effects of these resources. Few principals or superintendents are ignorant of the general concept of random assignment, though they may not value it.

*Randomized experiments are best when
a causal question is simple, sharply focused,
and easily justified.*

What is needed to make randomized experiments more common in schools? There is no single road. In Cook et al. (1999), random assignment was sponsored by the school district and all middle schools had to comply. Principals had no choice about participating or about the treatment they eventually received. The district took this step because a foundation-funded network of prestigious scholars—none from education—insisted on random assignment at the school level as a precondition for funding the program and its evaluation. In Cook, Hunt, and Murphy (2000), the principal investigator insisted on random assignment as a precondition for collaborating with the program implementers. In Chicago, participation was restricted to schools where principals applied for the program and agreed in advance to live with the results of the randomization process. No principal had any difficulty appreciating the method’s logic, and most lived with its consequences for up to six years. (But not all. Of the twenty-four Chicago principals, one assigned to the control group dropped out immediately after the coin toss, and three of the treatment principals who retired at the end of the second study year had replacements who abandoned the program.) Also important was honestly acknowledging up front that the program might not be effective and promising the principals that if they were assigned to the control condition, they would be the first to be offered the intervention at the study end, by when it might be improved. Schools were also

paid for participating in the measurement process, and a year was set aside for recruitment.

Contrast this with what happened very recently when the developer of one of the nation's best-known school reform packages tried to evaluate his program using random assignment. He used letters and e-mail to solicit schools to volunteer to be in the experiment. His staff then followed up and found few schools willing to be assigned once the random assignment was explained to them. Nonetheless, the federal funders continued to insist that a randomized experiment be done. So, the program designer developed a *within-school* experiment in which some grades get the intervention but others do not, even though cross-grade contamination is likely in this circumstance.

Given the facts above, can one seriously imagine the developer's staff informing schools they were being recruited because it was not clear that the program would work? Can one imagine them asserting that prior (quasi-experimental) research on the program's effectiveness was not definitive when the developer's staff had done this very research and had used it in many grant applications for program funding? Developers are, and should be, passionate advocates for their programs, not brokers of honest appraisal.

And how much should we expect a program developer in education to know about practical ways to implement random assignment if, as in this case, he has never done such a study before and if education as a field has little recent practical history with such assignment? In the policy realm, random assignment should be in independent hands and carried out by staff with a recent history of successful randomization in complex field settings.

Which conditions are most conducive to random assignment? When the general principles in Shadish, Cook, and Campbell (2001) are applied to schools, such assignment is most feasible when treatments are shorter, teachers training is not required, patterns of coordination among school staff are not modified much, the demand for an educational change outstrips the supply, different treatments with similar goals are compared, the units receiving different treatments cannot communicate with each other, and students are the unit of assignment rather than classrooms or whole schools. Thus, it should be easiest to study different curricula at random, to introduce new technologies at random, to give students tuition rebates for Catholic schools at random, to assign homework variants at random, to assign teachers trained in different ways at random. Even these studies will not be easy. But they should not be very difficult either, provided that researchers have the will to make random assignment happen and have the practical experience to make it happen in real school settings.

*Even when experiments are mounted, many of the planned
between-treatment contrasts become compromised*

Random assignment creates treatment group equivalence at the pretest. But it is at the posttest that groups should be equivalent in everything except treatment

exposure. Sometimes, different kinds of students drop out of the various treatment groups, creating a subsequent group nonequivalence that threatens the integrity of all analyses of outcomes other than dropping out itself. Such differential attrition is most likely when treatments vary in intrinsic value, and since many independent variables involve deliberately created resource differences, in policy studies, differential attrition is not a remote possibility. However, it can be routinely minimized, if not always completely prevented. The keys are eliciting a strong initial staff commitment to staying in the study; providing modest payments to the units experiencing less desirable treatments; and closely monitoring treatment implementation, including monitoring to detect and deal with early dropout trends.

With long-lasting treatments, an additional difficulty arises. Officials leave schools, and their replacements sometimes want to jettison their predecessors' innovations and introduce their own, especially when the predecessors' reforms initially disrupted school routines and so have not yet generated a strong core of supportive teachers. Little can be done in this situation, which, in my experience, often occurs in the first years of whole school reform. If possible, the schools lost to intervention should remain within the measurement framework. But this is not always possible. And such strategies are not perfect. Despite taking the precautions above, Cook, Hunt, and Murphy (2000) still lost four of their twenty-four schools. (Fortunately, they built in strong, fallback, quasi-experimental options, pretest values on the outcome variable being a minimum for this.) But even with some differential attrition, the resulting bias is likely to be less than the bias due to school or teacher self-selection from the start; and statistical selection controls are better the smaller the initial bias and the better selection has been directly observed (Holland 1986).

Experiments can be compromised by treatment crossovers as well as differential attrition. These crossovers occur when units in one treatment condition experience intervention particulars destined for another, thereby reducing the size of the treatment contrast and increasing the chances of falsely concluding there is no treatment effect. One way cross-treatment borrowing occurs is when the units receiving different treatments can communicate with each other. To circumvent this, researchers should work with physically separated units. This is often easy, though not always. However, the more schools are separated, the higher the research budget becomes and the more schools are needed to meet sample size requirements.

Extensive treatment crossovers may be rare, though. Cook et al. (1999) documented that only three of ten control schools borrowed any program elements, and none borrowed the program's most central elements. They did not have access either to School Development Program facilitators or to treatment-specific professional development opportunities, including trips to the program developers at Yale. Of the three documented crossovers, one occurred because a treatment principal was married to a teacher from a control school and they talked about the intervention. Another was caused by the daughter of a Yale program official teaching in a control school and inviting her father to give some lectures at the school. The third involved a control principal becoming interested in the program, reading up

on it, and trying to implement some of its practices without formal program support. Also relevant is that the district program coordinator also did some districtwide professional development and some program details entered into what she taught. So, cross-treatment borrowing occurred, but it was not universal across schools, it involved some but not all program details, and the borrowed particulars varied from school to school. None of this would have been detected, of course, without sensitivity to the possibility of treatment crossovers and without collecting annual data on the matter.

The eternal hope is that treatments will be so innovative that the control units will experience none of the treatment particulars. But even without direct communication control, units often have experiences that overlap with those planned for the treatment group. It is as though program designers are inspired by ideas that appear new but that are, in reality, only a little ahead of the emergent professional mainstream. During the study, they then enter into that mainstream, and so controls pick them up. For evaluation to be maximally useful, it may be that program designers need to be more original. But even when the various treatment groups have become closer in content, this does not mean that the planned contrast is useless for policy purposes. It is still worth learning whether the planned treatment adds something over and above the newly emerged status quo. Program developers dislike this since they believe that their program can be better than the former status quo without being better than the new one. Particularly galling for them is the possibility that their own ideas might have co-caused the new status quo, thus heightening the bar over which their own program now has to jump. Developers want to learn about the effects of their planned treatment at its maximal point of contrast, and so they prefer it be evaluated against the best approximation to the total absence of anything resembling their treatment.

The usual approximation to meeting their needs requires measuring treatment fidelity on each unit in each treatment condition and then using this fidelity measure as the “independent variable” in analyses of the outcome. A selection problem obviously results. However, if such stratification occurs within a randomized experiment, this is one of the few situations where an instrumental variables approach to causal inference credibly deals with selection (Angrist, Imbens, and Rubin 1996), permitting two distinct and important questions to be answered. First, is there an effect of the original treatment assignment, a question that deliberately disregards all variation in treatment exposure that might have occurred either within the planned treatment or within the planned control group? The analysis of this question is usually called the “intent to treat analysis” because the independent variable is the treatment contrast from the original design plan. The second important question targets the treatment variants that students demonstrably received rather than those planned for them. It asks, Is there an effect of the “treatment” actually received, given that some treatment group units might have received little of the planned treatment while some control group units might have received some or all of the treatment? It is ironic that the most defensible current test of the causal impact of self-selected treatments occurs if the self-selection has occurred within the framework of a randomized experiment.

Random assignment assumes fixed program theory and standard implementation, but these treatment-specific assumptions are not valid for school contexts

Experimental results are easier to interpret when the intervention is the product of strong substantive theory, when the achieved implementation faithfully reflects treatment-specific program theory, and when the within-treatment variation in implementation is minimal. These conditions are not often met.

Schools tend to be large, complex social organizations characterized by multiple, simultaneously occurring programs; disputatious building politics; and conflicting stakeholder goals. Management is all too often weak and removed from classroom practice, and day-to-day politics can swamp effective program planning and monitoring. So, many reform initiatives are implemented highly variably. Indeed, when different educational models are contrasted in the same study, the between-model variation is usually small relative to the variation between schools implementing the same model (Rivlin and Timpane 1975; Stebbins et al. 1978). In school research, it is not realistic to assume standard program implementation or total fidelity to program theory (Berman and McLaughlin 1977). To those who assume that schools have severe management and implementation problems, experiments must seem premature.

However, educational research does not need to assume such complex organization. An earlier model treated schools as physical structures with many self-contained classrooms in which teachers tried to deliver effective curricula using instructional practices that had been “shown” to enhance student performance. This approach privileged curriculum design and instructional practice over the schoolwide factors that now dominate—for example, strong leadership, a buildingwide communitarian climate, a focus on learning, undertaking multiple forms of professional development, and creating supportive links to the outside world. Many important consequences follow from how schools were reconceptualized. One is the lesser profile accorded to curriculum and instructional practice and to what happens once the classroom door is closed. Another is the view that random assignment is premature, given the presumption that its implementation depends on positive school management and quality program implementation. And another is the consequence that quantitative techniques are of lesser value, since school management and culture are best understood through ethnographic case studies.

Advocates of random assignment will not be credible if they assume treatment homogeneity or setting invariance in educational contexts. However, random assignment does not require well-specified program theories, good management, standard implementation, or treatments that are totally faithful to program theory. Experiments primarily protect against bias in causal estimates and only secondarily against imprecision in these estimates. So, the complexity and heterogeneity of schools leads to the need for larger school sample sizes and the need to anticipate and measure specific sources of variation to reduce their unwanted influence through statistical control. But just as important, implementation quality should be

studied as a dependent variable to ascertain which types of schools and teachers implement the program better. Variable implementation is important in its own right as well as having implications for budgets and sample sizes. We should also not forget that few educational interventions will be standardized once they are implemented as formal policy. So, why standardize them in an experiment? Treatment standardization is desirable for researchers interested in testing the substantive theory behind a treatment and for those interested in assessing an intervention's potential as policy. But it is not desirable for those seeking to determine a treatment's likely effects in settings where standard implementation cannot be expected.

Random Assignment Entails Undesirable Trade-Offs

Increasing internal validity decreases external validity

Random assignment prioritizes on unbiased answers to descriptive causal questions. But few educational evaluators share this priority, and most believe that it compromises more important research goals. Cronbach (1982) rejected the assertion that internal validity is the sine qua non of experimentation (Campbell and Stanley 1963) because of the neglect this implies for external validity. Experiments are clearly limited in time and space, and nationwide experiments are very rare. Most experiments are limited to the subclass of schools willing to surrender choice over the treatment they will receive and to tolerate the in-school measurement of implementation, mediating processes, and individual outcomes. What kinds of schools are these? Cronbach preferred to sample from a more representative population of schools even if less certain causal inferences result from this.

Science values results that are general. This includes the discovery and explanation of generative causal processes like gravity, relativity, DNA, nuclear fusion, aspirin, personal identity, infant attachment, or engaged time on task. Constructs like these index processes that are capable of bringing about effects across a broad array of contexts. They are more general than the results of a single experiment that typically shows whether one form of instruction affected achievement at a particular time in the particular sample of schools that happened to volunteer for an experiment. Many educational evaluators want their field to identify general causal agents whose operating principles are fully understood. They believe that the path to this is through explaining why programs work and then synthesizing this knowledge into higher-order constructs. Less priority is placed on what an experiment typically achieves—demonstrating that a particular instance of a program is effective in a particular context at a particular time with a particular group of respondents. As a result, scholars espousing this view of the purpose of a practical science are prepared to tolerate more uncertainty than most other scientists about whether a program does in fact work. Educational evaluators espouse the traditional schol-

arly goal of full explanation, but they reject the quantitative methods usually preferred for this. Cronbach (1982) has even asserted that the methods of the historian, journalist, and ethnographer suffice for learning about what happened and why in an educational reform and for combining these results with prior bodies of knowledge to facilitate generalization.

An obvious problem with this high priority assigned to external validity is that it has not led education evaluators to reliably learn what works. It is now thirty years since vouchers were proposed, and we still have no clear answers about them. It is thirty years since Comer began his work on the School Development Program, and almost the same situation holds. It is fifteen years since Levin began accelerated schools, and here too we have no experiments and no answers. The Obie-Porter legislation cites Comer's program as proven effective. But when the legislation passed, the relevant evidence consisted of testimony, studies by the program's own staff that used primitive quasi-experimental designs, and one interrupted time-series study that confounded the court-ordered introduction of the program with a simultaneously ordered reduction of about 40 percent in class sizes (Comer 1988).

Of the other whole school programs, only Success for All has been evaluated moderately well (for a summary, see Herman 1998). But even here, the evaluations have not been independent of the developer, the treatment assignment has never been random, and in no school did fewer than 80 percent of the teachers agree to the program. When trade-offs are made that favor generalization over cause, this risks ending up with the current state of affairs. Many studies exist in many districts at many times, but none is worth much as a study of cause. Experiments are not meant to be representative; they test causal claims.

But causal statements are more useful if they come from experiments where the sampling particulars permit tests of generalization across various types of students, teachers, settings, and times. The hope is either to demonstrate empirical robustness or to identify the boundary conditions under which an effect occurs. The two keys here are research questions that are crystal clear about the populations targeted and then the use of sampling procedures that represent these targets and make heterogeneous all the other irrelevancies that might limit generalization. So, formal sampling is one way to increase external validity within experiments. The ideal is random selection followed by random assignment to achieve an unbiased causal estimate that generalizes without bias to a prespecified population of schools, students, teachers, or sites.

As admirable as this is, there are many reasons why sampling units with known probability has rarely been the path to generalization in the experimental sciences. Volunteering to be in a study is usually required, and this limits generalization. And many of the populations it is practical to sample without volunteering are of only parochial interest. Moreover, random sampling is hardly relevant to estimating causal relationships that might vary by historical period, and it cannot be used to select the outcome measures and treatment variants that are used to represent general cause and effect constructs. So bench scientists use a different generalization model, one that emphasizes how consistently a causal relationship replicates across multiple sources of heterogeneity (Cook 1993). The operative question is

this: Can the same causal relationship be observed across different laboratories, time periods, regions of the country, and ways of operationalizing the cause and effect? This heterogeneity-of-replication model underlies current practice in both clinical trials and meta-analysis and permits purposive rather than random sampling to be used. Vital is only a heterogeneous sampling plan with respect to people, settings, operational definitions, and times—though multisite clinical trials typically sample only one time period.

In deciding whether to adopt a potentially life-saving therapy for a loved one, would experimenters not use the more liberal risk calculus? Why be different in science?

Heterogeneous, purposive sampling is not an easy path to follow for increasing external validity while maintaining high internal validity, especially in education, which has no tradition of multisite clinical trials with national reach. More typical are individual school experiments with unclear reach being done only in Milwaukee or Washington or Chicago or Tennessee. In addition, few reform efforts in education have a fixed protocol. So we can implement vouchers, charter schools, or Total Quality Management in many different ways across districts and even within them. Indeed, the Comer programs in Prince George's County (Maryland), Chicago, and Detroit are different from each other in many ways, given how much latitude districts are supposed to have in how they define and implement that program's specific details. This means that many educational treatments will require even larger samples of settings than do clinical trials in medicine, where the between-site variation in protocols is almost certainly less than in schools.

Very large experiments may not be wise, however. The physical sciences have progressed in generalizing because knowledge claims are routinely and mundanely replicated during the next stage of research on a phenomenon. But in research areas with weaker (and more expensive) traditions of replication—as in education—replication cannot be so haphazard. We need experience-filled theories of replication. Should replication depend on conducting a number of smaller experiments staggered over the years, each adequately statistically powered? Does it make sense to conduct even more experiments, but smaller ones, many of which are inadequately powered—as seems to be the case in school-based prevention studies? Is meta-analysis the only serious answer to the causal generalization problem, so that patience becomes a needed policy virtue because of the time needed to

build up a database? Whatever the merits of particular forms of phased programs of experiments, the point is that individual experiments vary in their sampling reach and in their connections to solid findings from the past. Single experiments rarely produce definitive answers, however large they are. And they certainly do not answer all ancillary questions about the contingencies on which a causal relationship depends.

I have shown that causal generalization can be understood as a single causal estimate for a given population (as in the formal sampling tradition) or as an average effect size derived from heterogeneous studies of the same hypothesis (as in synthesis methods). But causal generalization can also be understood as identifying generative causal processes. For instance, engaged time on task is presumed to stimulate achievement through activities as diverse as more homework, summer classes, longer school days, and more interesting curricula—procedures that can be implemented in any school district in any country at any time. The methods for identifying such explanatory processes place relatively little weight on sampling, instead requiring the collection of data about each of the variables in the presumed generative theory or getting historians and ethnographers to explore what happened and why in each treatment group. Fortunately, it is easier to build these explanatory methods into individual experiments than it is to sample at random or to add populations to the sampling design. In whatever ways are feasible, experiments should be designed to explain the consequences of interventions and not just to describe them. This means adding to an experiment's measurement and sampling plans and abjuring black box experiments.

Prioritizing scientific purity over utility

Critics contend that experimenters value uncertainty reduction about cause so much that conservative criteria are used to protect against wrong inferences, with the result that many effective programs are judged to be ineffective. One example of this is use of the traditional statistical criterion of $p < .05$ rather than, say, $.25$. In deciding whether to adopt a potentially life-saving therapy for a loved one, would experimenters not use the more liberal risk calculus? Why be different in science? Should statistical traditions be so strict that schools not implementing the treatment are included in the analysis as though they had been treated? What about purist experimenters who refuse to explore the data for unplanned treatment interactions with student or teacher characteristics or who view unplanned variation in implementation as a cause for concern rather than an opportunity to explore the origins and consequences of this variation? And why should one persist with the original research question if a more useful question has emerged during a study, even if unbiased answers to this new question are not possible? Better relevant than pure is the implication of these critical questions—all the more so since many experiments take so long to plan, mount, and analyze that the answers they provide turn out to be of more antiquarian than contemporary interest.

Experiments do tend to be so preoccupied with bias protection that other types of knowledge become secondary. But they need not be so secondary. There is no

compelling need for stringent alpha rates; only statistical convention is at play here, not statistical wisdom. Nor need one restrict data analyses to the intent-to-treat group, though such analyses need to be reported. Nor need one ignore all statistical interactions, though probing them should be done with substantive theory and statistical power in mind, and conclusions about substantive interactions should be couched more tentatively than conclusions about main effects. Researchers can also try to replicate experimental results generated from limited samples by reanalyzing data about similar constructs collected from formally representative samples and subjecting them to the best available nonexperimental analyses. Finally, many controlled experiments will be improved by collecting ethnographic data in all the treatment groups. This will help better understand issues of implementation and causal mediation as well as help identify some possible unintended outcomes. And such continuous ethnographic feedback can be provided to the treatment and control schools alike, as long as one is willing to restrict generalization to contexts where such feedback is regularly available. So, experiments need not be as rigid as many clinical trial texts paint them.

Experimental Results Are Not Likely to Be Used in Educational Policy

*Experiments prioritize on the information needs
of central decision makers who are not important in
the decentralized American educational system*

Most of the funds spent on education come from local sources, next most from states, and least from federal sources. Yet experiments are most often designed with the opposite priority ordering. Federal interests come first in how an evaluation's guiding questions are framed, and local interests come last, if they are considered at all. Many opponents of experimentation identify most with the information needs of local school staff, and they contend that such staff rarely want the kind of knowledge for which experiments were designed—knowledge summarizing what a reform has achieved across a sample of schools. They want information about their own school only, and they want it when they need it, not just at the end of a study. Putting the needs of local service deliverers above those of amorphous state and especially federal policy makers prioritizes on utility over truth and on continuous feedback over final reports. A letter to the *New York Times* captured the following:

Alan Krueger . . . claims to eschew value judgments and wants to approach issues (about educational reform) empirically. Yet his insistence on postponing changes in education policy until studies by researchers approach certainty is itself a value judgment in favor of the status quo. In view of the tragic state of affairs in parts of public education, his judgment is a most questionable one. (W. M. Petersen, 20 April 1999)

It is a mistake to believe central decision makers are powerless in education. At both the federal and state levels, their role is steadily increasing. They are especially powerful in inner cities, where the proportion of dollars from local taxes is lower than in the suburbs. They are also especially powerful in certain programs like special education and bilingual education. Moreover, Congress and state government are important sources of political pressure to improve schools in the belief that improved human capital will keep the economy strong in a global context where low-wage industries have moved abroad. It is also a mistake to believe that experiments necessarily preclude continuous feedback to schools. The key requirements are only that such feedback be provided similarly in all treatment conditions and that no premature experimental contrast results be presented. Even so, providing continuous feedback should not be done willy-nilly, for the information provided becomes part of the study context and limits generalization to settings where feedback is part of program design.

It may be true that in the past, local school personnel have used information from site-specific management studies more than from experiments, though this is not clear. However, nothing prevents experiments being done on issues that school staff bring up, and it is probably an accident of history that the causal questions of central decision makers in politics and education have been paramount. But even results from centrally determined experiments have some utility for local school personnel, especially if it is particularly likely that they will get into the corpus of findings in those textbooks and manuals that are used to train the next generation of teachers and principals. The anticipated benefit from experiments is their potential reach across a nation, either through policy decisions or their results getting into textbooks. It is a mistake to limit uses of research to immediate use by the districts or schools in a study.

Experimentation recreates a classical model of rational decision making that has been disproved

Theories of rational decision making require analysts first to lay out the alternatives (the treatments). Then one decides on decision criteria (the outcomes). Next, one collects information on each criterion for each treatment group (data collection). And finally, one uses the observed effect sizes and whatever utilities can be attached to them to make a decision about the merits of the contending alternatives. Empirical work on how social science data are used in policy reveals that such use (termed “instrumental use”) is rare (Weiss and Bucuvalas 1977; Weiss 1988). Instead, use is more diffuse and is better described by an “enlightenment” model. This involves information blended from existing theories, personal testimony, extrapolations from surveys, the consensus of a field, empirical claims from experts who may or may not have interests to defend, and novel concepts that are au courant and broadly applied—like social capital is in sociology and political science today. Describing research utilization in this enlightenment fashion extends no special privilege to science in general or to experiments in particular.

Empirical research also notes that use decisions are multiply rather than singly determined, with central roles being played by politics, personality, windows of opportunity, and values. Also, many decisions are not made in a systematic sense but are rather slipped into or accrete, with earlier, small decisions constraining later, larger ones. In addition, official decision-making bodies change in personnel, with new persons and issues replacing older ones. When studies take time to complete—as with most experiments—the results may not be available until the policy agenda has already changed. Research use is much more complex than simply making an evidence-based rational choice.

Critics also note that experiments rarely provide uncontested verdicts. Disputes typically arise about how well the original causal question was framed, about whether the claimed results are valid, about whether all relevant outcomes were assessed, and about whether the proffered recommendations follow from the results. The logical control over selection that makes experiments so valuable does not mean that all quibbles about causal claims are put to rest. Consider the very different conclusions offered about whether and where effects are warranted in the Milwaukee voucher study (Witte 1998; Greene et al. 1996). Consider, also, the different effect sizes generated from the Tennessee class size experiment and (Finn and Achilles 1990; Mosteller, Light, and Sachs 1996; Hanushek 1999; Krueger 1999; Krueger and Whitmore 2001). Sometimes, real scholarly disagreements are at issue, while in other cases, the disputes reflect stakeholders protecting their interests or confusion between fundamentally different research questions. Policy insiders use multiple criteria for making decisions, and scientific knowledge of causal influences is never uniquely determinative.

Close examination of claims that policy changes were made because of experimental results suggests some oversimplification. The Tennessee class size results are consonant with the results of an earlier meta-analysis (Glass and Smith 1979) and with theories that postulate that children gain more if they are engaged and on-task. The results also conform with teachers' hopes and expectations as well as with parents' commonsense notions. Moreover, when the results were delivered, the Tennessee governor thought he would be able to push through increased investments in education, and he reasoned that lowering class sizes would be popular with both teacher unions and business interests, thus furthering his national political ambitions. So any policy change he might have made could not be attributed to the experimental results alone. In the end, though, Tennessee did not reduce class sizes because of the financial cost and the lack of teachers in this era of a national teacher shortage (Ritter and Boruch 1999).

Other states did make the change to smaller classes, though. Illustrating what then happened provides an important lesson in the limits of generalization of localized experiments like the Tennessee one. To make smaller class sizes possible, California had to recruit more teachers. It did so in part through poaching teachers from other (less wealthy) states, possibly exacerbating state-level inequalities. Teacher transfers also occurred across district lines within California, again probably favoring the wealthier districts. Some new teachers were recruited from the corporate world under the assumption that they would be effective teachers

despite scant formal training and little classroom experience. Smaller classes require more space in a school or new buildings. Given the expense of new buildings, some California students had to be located in suboptimal facilities. Experiments like the Tennessee study exist on a smaller and more local scale than would typically pertain if the services they test were to be implemented nationwide. This scale issue is serious since program dynamics can be different on the larger scale and, thus, entail a different pattern of effects than achieved on the smaller scale (Elmore 1996).

There is some substance to the idea that the theory of use buttressing randomized experiments is at odds with the ways social science data are used. But the objections are exaggerated. Instrumental use does occur (Chelimski 1987), and more often than the very low base rates implied in most research denigrating instrumental usage. Moreover, some results are probably more widely disseminated because random assignment confers credibility on them in many quarters. This happened with the Tennessee class size study and the preschool studies cited earlier. Studies can even be important if political events have rendered the results obsolete when they are announced. This is because many policy initiatives are recycled later—as with vouchers—and because the texts used to train professionals in a field often describe past studies that throw light on particulars of professional practice (Leviton and Cook 1983).

There is also no necessary trade-off between instrumental and enlightenment usage. Experiments also contribute to enlightenment. They teach us about the kinds of interventions that can be better or worse implemented and about how principal turnover seems to affect school management. They inform about the low utility of theories that fail to specify what happens once the teacher closes the classroom door, about the kinds of principals who are most attracted to school-based management, and about the kinds of teachers most open to professional development.

The era of black box experiments is long past. We now want to learn, within experiments, about the determinants and consequences of implementation quality and about the viability of the substantive theory undergirding program design. We also want to collect qualitative and quantitative data as long as the data collection protocol is identical in all treatment groups. And we want to get all stakeholder groups involved in formulating experimental questions and in interpreting the relevance of findings. These steps help generate enlightenment and thus make the experiment more like what its critics claim it is not.

Random Assignment Is Not Needed Because Better Alternatives Already Exist

Intensive case studies

No method will die, whatever its imperfections, unless a demonstrably better or simpler method can replace it. Educational evaluators believe that superior alter-

natives to the experiment already exist. The alternative they generally prefer is the intensive case study. The main reason for their preference lies in its flexibility of purposes. To their way of thinking, the case study can be used to appraise the theory of the program, to assess implementation quality, to record program redesign, to identify whether intended changes have occurred, to identify unplanned effects, to estimate subgroup effect differences, to probe reasons for the effects claimed, and to assess how relevant the findings are for different stakeholder groups. By itself, the experiment cannot match this flexibility, having been designed only to answer one type of causal questions. Few advocates contend that case studies reduce as much uncertainty about cause as an experiment. But they do assert that such studies can reduce the uncertainty to an acceptable level. After all, the general public and political actors often believe the causal claims of investigative journalists, historians, and ethnographers, none of whom do experiments.

When done carefully, case studies constitute a serious form of research whose epistemological premises overlap with those of quantitative science. That is, observations are first used to develop a hypothesis about what works. Researchers then think through other implications of this hypothesis and collect data relevant to these implications. This round of observations is then used to revise the last version of the hypothesis, and so on until closure is reached. This is basically an empirical, falsificationist hypothesis-testing procedure, and theorists of ethnography have long advocated it (e.g., Becker 1958). The results from procedures like the above should be particularly rich for explaining why findings came about because close attention is paid to social processes as they unfold at different stages in a program's progress. There is no doubt that hard thought and nonexperimental empiricism can reduce some uncertainty about cause—sometimes even all the uncertainty. However, it will usually be very difficult to know when this last has happened.

Nonetheless, intensive qualitative case studies do not reduce as much causal uncertainty as well-executed experiments. Case studies rarely involve a convincing causal counterfactual. The absence of control groups makes it difficult to know how a treatment group would have changed in the absence of the reform under analysis. Adding comparison groups helps, but unless they are randomly created, it will not be clear whether the two groups would have changed at similar rates over time. Whether intensive case methods reduce enough uncertainty about cause to be generally useful is such a poorly specified proposition that I cannot answer it. Still, it forces advocates to note yet again that experiments are best justified when a high standard of uncertainty reduction is required about a manifestly important causal claim.

Factors other than flexibility also favor case studies over experiments. First, schools are probably less squeamish about collaborating with ethnographers than experimentalists. Second, case studies produce human stories that can be used to communicate the results. And third, feeding interim results back to the teachers and principals with whom an ethnographer has ongoing relationships may be especially likely to generate local use of the data. This is less grandiose than affecting large numbers of districts and schools, but educational evaluators like Stake (1967) and Guba and Lincoln (1982) consider local use as the ultimate desideratum

because they doubt that schools will comply with policy dictates from outside. Intensive case studies have many advantages, and researchers value them.

However, they are valued most for their role within experiments rather than as alternatives to experiments. They complement an experiment whenever a causal question is central, but it is not clear how successful program implementation will be, why implementation shortfalls may occur, what unexpected effects are likely to

*No method will die, whatever its imperfections,
unless a demonstrably better or simpler
method can replace it.*

emerge, how respondents interpret the questions asked of them, what the casual mediating processes are, and so on. Since these questions are important and not relevant to experimental functions per se, qualitative methods have a central role to play as adjuncts within experimental work on educational interventions. They should not be afterthoughts.

Quasi-experiments

Most researchers who do educational evaluations do not think of themselves as evaluators. They are primarily substantive researchers who want to test the effectiveness of changes within their own subfield. They mostly use quasi-experiments, as in the early evaluations of Comer's program noted earlier and the studies Herman (1998) detailed for Success for All and other whole school reform ideas. According to Herman, qualitative studies are also rare. This last is surprising, given how much emphasis theorists of educational evaluation have placed on qualitative work. Has their advocacy had little influence on their own colleagues? I do not know. But it is not easy to integrate outcome-based qualitative and quantitative case studies into a summary picture of the effects of a school reform.

Quasi-experiments are identical to experiments in purpose and in most structural details, the defining difference being no random assignment. Quasi-experiments use design rather than statistical controls to create the best possible approximation (or approximations) to the missing counterfactual that random assignment would have generated. These design controls include matched comparison groups, age or sibling controls, pretest measures at several times before a treatment begins, interrupted time series, assigning units based solely on a quantitative criterion, assigning the same treatment to different groups at different times, and building multiple outcome variables into studies, some of which should theo-

retically be influenced by a treatment and others not (Corrin and Cook 1998; Shadish, Cook, and Campbell 2001). Quasi-experimental designs are created through a mixing process that tailors the research problem and the resources available to the best design that can be achieved by mixing the design elements above.

In some quarters, “quasi-experiment” has been used promiscuously to connote any study that seeks to test a causal hypothesis and is not a true experiment but that has some form of nonequivalent control group or some pretreatment observation. Yet Campbell and Stanley (1963) and Cook and Campbell (1979) labeled some such studies as “generally causally uninterpretable,” and many of the studies that educational researchers call “quasi-experiments” are of this last kind. They lag far behind the state of the art. Reading quasi-experimental studies of educational reform projects is dispiriting, so weak are the designs and so primitive are the statistical analyses. All quasi-experimental designs and analyses are not equal. Recent advances in the design and analysis of quasi-experiments are not getting into educational research where they are needed. In particular, researchers should lament the dearth of interrupted time series studies, of regression-discontinuity studies, and of nonequivalent control group designs with more than one pretest measurement wave on matched cohort samples.

The best estimate of any quasi-experiment’s internal validity, though, is to compare its results with those from a randomized experiment on the same topic. When this is done systematically and empirically, the results described earlier indicate that quasi-experiments are more likely to be biased and inefficient when compared to experimental results—at least for the classes of quasi-experiment examined to date, which is not all of them and not the strongest of them. Even so, in areas like education where few studies exist on most of the reform ideas being currently debated, randomized experiments are particularly needed. It will take fewer of them to arrive at what might—or might not—be the same answer, and anyway, most scholars trust the answers from experiments more than from quasi-experiments.

Theories of change

It is currently fashionable in many foundations and some scholarly circles to espouse a nonexperimental theory of change for use in evaluations of complex social settings like communities and schools (Connell et al. 1995). The method depends on explicating the substantive theory behind a reform initiative and detailing all the flow-through relationships that should occur if the intended intervention is to impact on a major distal outcome like achievement gains. The method also requires measuring each construct in the substantive theory and then analyzing the data to assess whether the postulated relationships have actually occurred in the predicted time sequence. With shorter time periods, the data analysis will involve only the first part of the postulated causal chain; but over longer periods, the complete model might be testable. This conception of evaluation places a primacy on very specific program theory, on high-quality measurement, and on the valid analysis of multivariate explanatory processes as they unfold in time.

The claim is that such theory-based evaluation can function as an alternative to random assignment that has broader applicability. First, it does not require a causal counterfactual constructed through random assignment or matched comparison groups, a requirement that increases costs and alienates some school officials. Only the group experiencing a treatment is needed. Second, obtaining data patterns congruent with program theory is assumed to validate that theory. This is an epistemology that does not require explicitly rejecting alternative explanations, merely demonstrating a close match between the predicted and obtained data. Finally, the theory of change approach does not depend on attaining the end points typically specified in educational research—usually a cause-effect relationship involving some aspect of student cognitive performance. Instead, when the initial phases of the program theory are corroborated, this can be used to argue for maintaining the program because it might be effective if data on more distal criteria were collected. Initial corroboration also defends against prematurely concluding that a program is ineffective when insufficient time has elapsed for all the intervening processes to occur that might bring about the ultimate change. So, this theory of change approach does not require measuring the end points most often valued in educational studies but that often preclude using the data for early decision making.

Few advocates of experimentation will argue against the greater use of substantive theory to guide measurement and analysis in experimental evaluations. These features improve experimental probes, first, of whether the intervention led to changes in the theoretically specified intervening processes and, second, of whether these processes could then have plausibly caused changes in distal outcomes. The first of these tests will be unbiased because it relates each step in the causal model to the planned treatment contrast. But the second test will entail selection bias if it depends only on stratifying units according to the extent the postulated theoretical processes have occurred. Still, quasi-experimental analyses of the second stage are worth doing, provided that their results are clearly labeled as more tentative than the results of planned experimental contrasts. The utility of analyzing theoretical intervening processes in experiments is beyond dispute; the issue is whether such measurement and analysis can alone provide an alternative to random assignment.

There are reasons for skepticism about the validity of using theories of change to support strong causal conclusions (Cook 2000). First, it has been my experience writing papers on the theory behind a program with its developer (Anson et al. 1991) that the theory is not always very explicit and could be made explicit in several different ways, not just one. Is there a single theory of a program, or several possible versions of it? Second there is the problem that many of these theories seem to be too linear in their flow of influence, rarely incorporating reciprocal feedback loops or external contingencies that might moderate the flow of influence. It is all a little bit too neat for our chaotic world. Third, few theories are specific about timelines, specifying how long it should take for a given process to affect some proximal indicator. Without such specifications, apparently disconfirming results make it difficult to know whether the next step in the model has not yet occurred or whether it will never occur because the theory is wrong in this particular. Fourth, the method places a great premium on knowing how to measure since

failure to corroborate the model could result from partially invalid measures or from an invalid theory. Careful researchers can protect against this by developing more reliable measures. Fifth is the epistemological problem that many different models can usually be fit to any single pattern of data (Glymour, Scheines, Sprites, and Kelly 1987). The implication here is that causal modeling is more valid when multiple competing models are tested against each other rather than when a single model is tested.

The biggest problem though, is the absence of a valid counterfactual to inform us about what would have happened to students or teachers had there been no treatment. As a result, it is impossible to decide whether the observed data result from the intervention or would have occurred anyway. One way to guard against this is with “signed causes” (Scriven 1976) predicting a multivariate pattern among the outcomes that is so unique it could only have occurred because of the reform. But signed causes depend on the availability of well-validated substantive theory and high-quality measurement (Cook and Campbell 1979). So, a better safeguard is to have at least one comparison group, and the best comparison group is a randomly selected one. So, researchers must come back to random assignment and the proposition that theory-based evaluations are useful complements to randomized experiments but not alternatives to them.

Conclusions

1. Random assignment cannot be considered as the “gold standard” for justifying causal inferences in school-based research. It creates only a probabilistic equivalence between the groups being contrasted, and then only at pretest. Moreover, treatment-correlated attrition is likely when treatments differ in intrinsic desirability. Also, treatments are not always independent of each other in practice like they are supposed to be in theory, and many of the ways used to increase internal validity can also reduce external validity.

2. Nonetheless, random assignment is the best mechanism for justifying causal conclusions. More appropriate rationales for it than the “gold standard” are that, even after all the limitations above are taken into account, (1) it still provides the logically most valid causal counterfactual, (2) it almost certainly provides a more efficient counterfactual, and (3) the results it generates are more credible in nearly all academic circles.

3. Despite widespread recognition of the superiority of random assignment, it is still too rare in research on the effectiveness of school-based strategies to improve student performance.

4. Most surprising is that scholars in schools of education who call themselves educational evaluators hardly ever do experiments and usually counsel against doing them. Their professional colleagues who are primarily interested in substantive educational topics do most of the experiments that are done, though experiments constitute only a tiny fraction of all the cause-probing studies they do. There has been a recent flurry of obviously policy-relevant education experiments. But

these have been done by researchers in economics or public policy, often working for contract research firms.

5. Unlike in school-based research on teaching and learning, random assignment is not at all rare in preschool education or in school research on preventing negative behaviors or feelings.

6. Intellectual culture is one possible explanation for this difference. Prevention researchers tend to be trained in public health and psychology, where random assignment is esteemed and where funders and journal editors clearly prefer the technique. Something very similar is true in preschool education. However, the system of training and professional rewards is quite different in schools of education.

7. Capacity may be another explanation for the difference. Most school-based prevention experiments are shorter, researchers do the hands-on implementation rather than school staff, and the research topics probably engage educators less than issues of school governance or teaching practice.

8. Nearly all educational evaluators believe that experiments are of little value. They believe that the theory of causation buttressing experiments is naïve, that experiments are difficult to implement, that they require unacceptable trade-offs, that they deliver a kind of information that is rarely used to change policy, and that the information experiments provide can be gained using simpler and more flexible methods.

9. Some of these beliefs are better justified than others. Beliefs about the viability of alternatives to experiments are particularly poorly warranted since no current alternative provides as convincing a causal counterfactual as the randomized assignment. However, the better criticisms suggest useful additions to school-based experiments, especially as concerns describing program theory, describing implementation quality, relating implementation to outcome changes, measuring whether theoretically specified intervening processes have occurred, and then relating these intervening processes to student outcomes.

10. Educational evaluators will not be persuaded to do experiments simply by outlining their advantages and describing newer methods for implementing randomization. Most educational evaluators share some of the antiexperimental beliefs outlined above. To start a dialog, advocates of experimentation will need to be explicit about the method's limits. They will also have to note the utility of incorporating into experimental practice some of the critics' concerns, especially about program theory, the quality of implementation, the value of qualitative data, the necessity for analyses of causal contingency, and concern to meet the information needs of school personnel as well as other stakeholders.

11. Even so, for the intellectual reasons given above, it will be difficult to enlist the current generation of self-styled educational evaluators behind a banner promoting more experimentation. They are not needed for this task, though. They are not part of the current flurry of controlled experimentation now underway. And while the future demand for experiments cannot be predicted accurately, it may well be possible to meet it just with the staff from contract research firms, university faculty in the policy sciences, and substantive researchers in schools of education.

References

- Agodini, R., and M. Dynarski. 2002. *Nonexperimental replication of social experiments*. Princeton, NJ: Mathematica Policy Research, Inc.
- Angrist, J. D., G. W. Imbens, and D. B. Rubin. 1996. Identification of causal effects using instrumental variables. *Journal of the American Statistical Association* 91:444-62.
- Anson, A., T. D. Cook, F. Habib, M. K. Grady, N. Haynes, and J. P. Comer. 1991. The Comer School Development Program: A theoretical analysis. *Journal of Urban Education* 26:56-82.
- Becker, H. S. 1958. Problems of inference and proof in participant observation. *American Sociological Review* 23:652-60.
- Berman, P., and M. W. McLaughlin. 1977. *Federal programs supporting educational change*. Vol. 8, *Factors affecting implementation and continuation*. Santa Monica, CA: RAND.
- Bhaskar, R. 1975. *A realist theory of science*. Leeds, UK: Leeds University Press.
- Bloom, H. S., C. Michaelopoulos, C. J. Hill, and Y. Lei. 2002. *Can non-experimental comparison group methods match the findings from a random assignment evaluation of mandatory welfare-to-work programs?* New York: Manpower Demonstration Research Corporation.
- Bogatz, G. A., and S. Ball. 1972. *The impact of "Sesame Street" on children's first school experience*. Princeton, NJ: Educational Testing Service.
- Campbell, D. T., and J. C. Stanley. 1963. *Experimental and quasi-experimental designs for research*. Chicago: Rand-McNally.
- Campbell, F. A., and C. T. Ramey. 1995. Cognitive and school outcomes for high-risk African American students at middle adolescence: Positive effects of early intervention. *American Educational Research Journal* 32 (4): 743-72.
- Chelimsky, E. 1987. The politics of program evaluation. In *Evaluation practice in review*, edited by D. S. Cordray, H. S. Bloom, and R. J. Light. San Francisco: Jossey-Bass.
- Cicirelli, V. G., and Associates. 1969. *The impact of Head Start: An evaluation of the effects of Head Start on children's cognitive and affective development*. Vols. 1 and 2, *A report to the Office of Economic Opportunity*. Athens: Ohio University and Westinghouse Learning Corporation.
- Cohen, D. K., and M. S. Garet. 1975. Reforming educational policy with applied social research. *Harvard Educational Review* 45 (1): 17-43.
- Collingwood, R. G. 1940. *An essay on metaphysics*. Oxford: Clarendon.
- Comer, J. P. 1988. Educating poor minority children. *Scientific American* 259 (5): 42-48.
- Connell, J. P., A. C. Kubisch, L. B. Schorr, and C. H. Weiss, eds. 1995. *New approaches to evaluating community initiatives: Concepts, methods and contexts*. Washington, DC: Aspen Institute.
- Cook, R., and M. Puma. 2002. *The national Head Start study*. Princeton, NJ: Mathematica Policy Research.
- Cook, T. D. 1985. Post-positivist critical multiplism. In *Social science and social policy*, edited by R. L. Shotland and M. M. Mark, 21-62. Beverly Hills, CA: Sage.
- . 1993. A quasi-sampling theory of the generalization of causal relationships. In *New directions for program evaluation: Understanding causes and generalizing about them*, vol. 57, edited by L. Sechrest and A. G. Scott. San Francisco: Jossey-Bass.
- . 2000. The false choice between theory-based evaluation and experimentation. *New directions in evaluation: Challenges and opportunities in program theory*. *Theory Evaluation* 87:27-34.
- Cook, T. D., A. Anson, and S. Walchli. 1993. From causal description to causal explanation: Improving three already good evaluations of adolescent health programs. In *Promoting the health of adolescents: New directions for the twenty-first century*, edited by S. G. Millstein, A. C. Petersen, and E. O. Nightingale. New York: Oxford University Press.
- Cook, T. D., and D. T. Campbell. 1979. *Quasi-experimentation: Design and analysis issues for field settings*. Boston: Houghton-Mifflin.
- Cook, T. D., F. Habib, J. Phillips, R. A. Settersten, S. C. Shagle, and S. M. Degirmencioglu. 1999. Comer's School Development Program in Prince George's County, Maryland: A theory-based evaluation. *American Educational Research Journal* 36 (3): 543-97.
- Cook, T. D., H. D. Hunt, and R. F. Murphy. 2000. Comer's School Development Program in Chicago: A theory-based evaluation. *American Educational Research Journal* 37 (2): 535-97.

- Cooper, H. M. 1989. *Homework*. New York: Longman
- Corrin, W. J., and T. D. Cook. 1998. Design elements of quasi-experimentation. *Advances in Educational Productivity* 7:35-57.
- Cronbach, L. J. 1982. *Designing evaluations of educational and social programs*. San Francisco: Jossey-Bass.
- Cronbach, L. J., S. R. Ambron, S. M. Dornbusch, R. D. Hess, R. C. Hornik, D. C. Phillips, D. F. Walker, and S. S. Weiner. 1980. *Toward reform of program evaluation*. San Francisco: Jossey-Bass.
- Cronbach, L. J., and R. E. Snow. 1976. *Aptitudes and instructional methods: A handbook for research on interactions*. New York: Irvington.
- Dehejia, Rajeev, and Sadek Wahba. 1999. Causal effects in nonexperimental studies: Reevaluating the evaluation of training programs. *Journal of the American Statistical Association* 94:1053-62.
- Durlak, J. A., and A. M. Wells. 1997a. Primary prevention mental health programs for children and adolescents: A meta-analytic review. *American Journal of Community Psychology* 25 (2): 115-52.
- . 1997b. Primary prevention mental health programs: The future is exciting. *American Journal of Community Psychology* 25:233-41.
- . 1998. Evaluation of indicated preventive intervention (secondary prevention) mental health programs for children and adolescents. *American Journal of Community Psychology* 26 (5): 775-802.
- Dynarski, M., and P. Gleason. 1998. *What we have learned from evaluations of federal drop-out prevention programs*. Princeton, NJ: Mathematica Policy Research.
- Dynarski, M., and Wood, R. 1997. *Helping high-risk youths: Results from the Alternative Schools Demonstration Program*. Princeton, NJ: Mathematica Policy Research.
- Ehri, L., S. Nunes, S. Stahl, and D. Willows. 2001. Systematic phonics instruction helps children learn to read: Evidence from the National Reading Panel's meta-analysis. *Review of Educational Research* 3:393-447.
- Ehri, L., S. Nunes, D. Willows, B. Schuster, Z. Yaghoub-Zadeh, and T. Shanahan. 2001. Phonemic awareness instruction helps children learn to read: Evidence from the National Reading Panel's meta-analysis. *Reading Research Quarterly* 36:250-87.
- Elmore, R. F. 1996. Getting to scale with good educational practice. *Harvard Educational Review* 66:1-26.
- Elmore, R. F., and M. W. McLaughlin. 1983. The federal role in education: Learning from experience. *Education and Urban Society* 15:309-33.
- Finn, J. D., and C. M. Achilles. 1990. Answers and questions about class size: A statewide experiment. *American Educational Research Journal* 27 (3): 557-77.
- Fraker, T., and R. Maynard. 1987. Evaluating comparison group designs with employment-related programs. *Journal of Human Resources* 22:194-227.
- Friedlander, D., and P. Robins. 1995. Evaluating program evaluations: New evidence on commonly used nonexperimental methods. *American Economic Review* 85 (4): 923-37.
- Gasking, D. 1955. Causation and recipes. *Mind* 64:479-87.
- Glass, G. V., and M. L. Smith. 1979. Meta-analysis of research on the relationship of class size and achievement. *Educational Evaluation and Policy Analysis* 1:2-16.
- Glazerman, S., D. M. Levy, and D. Myers. 2002. *Nonexperimental replications of social experiments: A systematic review*. Princeton, NJ: Mathematica Policy Research.
- Glymour, C., R. Scheines, and P. Spirtes. 1987. *Discovering causal structure: Artificial intelligence, philosophy of science, and statistical modeling*. Orlando, FL: Academic Press.
- Goodson, B. D., J. I. Layzer, R. G. St. Pierre, L. S. Bernstein, and M. Lopez. 2000. Effectiveness of a comprehensive five year family support program on low income children and families: Findings from the Comprehensive Childcare Development Program. *Early Childhood Research Quarterly* 15 (1): 5-39.
- Greene, J. P., P. E. Peterson, J. Du, L. Boeger, and C. L. Frazier. 1996. The effectiveness of school choice in Milwaukee: A secondary analysis of data from the program's evaluation. Mimeograph, University of Houston.
- Guba, E. G., and Y. Lincoln. 1982. *Effective evaluation*. San Francisco: Jossey-Bass.
- Gueron, J. M. 2002. The politics of random assignment: Implementing studies and affecting policy. In *Evidence matters: Randomized trials in education research*, edited by F. Mosteller and R. Boruch. Washington, DC: Brookings Institution.

- Hanushek, E. A. 1999. Evidence on class size. In *When schools make a difference*, edited by S. Mayer and P. E. Peterson. Washington, DC: Brookings Institution.
- Harre, R. 1981. *Great scientific experiments*. Oxford: Phaidon.
- Heckman, J. J., H. Ichimura, J. C. Smith, and P. Todd. 1997. Characterizing selection bias. *Econometrica* 66:1017-98.
- Herman, R. 1998. *Whole school reform: A review*. Washington, DC: American Institutes for Research.
- Holland, P. W. 1986. Statistics and causal inference. *Journal of the American Statistical Association* 81:945-70.
- Howell, W. G., and P. E. Peterson. 2002. *The education gap: Vouchers and urban schools*. Washington, DC: Brookings Institution.
- Kemple, J. J. 2001. *Career academies: Impacts on students' initial transitions to post-secondary education and employment*. New York: Manpower Demonstration Research Corporation.
- Krueger, A. 1999. Experimental estimates of educational production function. *Quarterly Journal of Economics* 114:497-532.
- Krueger, A., and D. Whitmore. 2001. The effects of attending a small class in the early grades on college-test taking and middle school test results. *Economic Journal* 111:1-28.
- Kuhn, T. S. 1970. *The structure of scientific revolutions*. 2d ed. Chicago: University of Chicago Press.
- LaLonde, R. J. 1986. Evaluating the econometric evaluations of training programs with experimental data. *American Economic Review* 76:604-620.
- Latour, B., and S. Woolgar. 1979. *Laboratory life: The construction of scientific facts*. Beverly Hills, CA: Sage.
- Leviton, L. C., and T. D. Cook. 1983. Evaluation findings in education and social work textbooks. *Evaluation Review* 7:497-518.
- Lipsey, M. W., and D. B. Wilson. 1993. The efficacy of psychological, educational, and behavioral treatment: confirmation from meta-analysis. *American Psychologist* 48:1181-209.
- Mosteller, F., R. J. Light, and J. A. Sachs. 1996. Sustained inquiry in education: Lessons from skill grouping and class size. *Harvard Educational Review* 66:797-842.
- Mosteller, F., J. P. Gilbert, and B. McPeck. 1980. Reporting standards and research strategies for controlled trials: Agenda for the editor. *Controlled Clinical Trials* 1:37-58.
- Nave, B., E. J. Miech, and Mosteller, F. 1999. A rare design: The role of field trials in evaluating school practices. Paper presented at the American Academy of Arts and Sciences at Harvard University.
- Olds, D. L., D. Eckenrode, C. R. Henderson, H. Kitzman, J. Powers, R. Cole, K. Sidora, P. Morris, L. M. Pettitt, and D. Luckey. 1997. Long-term effects of home visitation on maternal life course and child abuse and neglect. *Journal of the American Medical Association* 278 (8): 637-43.
- Peters, R. D., and R. J. McMahon, eds. 1996. *Preventing childhood disorders, substance abuse, and delinquency*. Banff International Science Series, vol. 3. Thousand Oaks, CA: Sage.
- Quine, W. V. 1951. Two dogmas of empiricism. *Philosophical Review* 60:20-43.
- . 1969. *Ontological relativity and other essays*. New York: Columbia University Press.
- Raikes, H. H., and John M. Love. 2002. Early Head Start: A dynamic new program for infants and toddlers and their families. *Infant Mental Health Journal* 23:1-13.
- Ramey, C. T., and F. A. Campbell. 1991. Poverty, early childhood education, and academic competence: The Abecedarian experiment. In *Children in poverty: Child development and public policy*, edited by A. C. Huston, 190-221. New York: Cambridge University Press.
- Ritter, G. W., and R. F. Boruch. The political and institutional origins of a randomized controlled trial on elementary school class size: Tennessee's Project STAR. *Educational Evaluation and Policy Analysis* 21:111-126.
- Rivlin, A. M., and M. M. Timpane, eds. 1975. *Planned variation in education*. Washington, DC: Brookings Institution.
- Rouse, C. E. 1998. Private school vouchers and student achievement: An evaluation of the Milwaukee parental choice program. *Quarterly Journal of Economics* 113:553-602.
- Rubin, D. B. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology* 66:688-701.
- Schweinhart, L. J., H. V. Barnes, and D. P. Weikart (with W. S. Barnett and A. S. Epstein). 1993. *Significant benefits: The HighScope Perry Preschool Study through age 27*. Ypsilanti, MI: HighScope Press.

- Scriven, M. 1976. Maximizing the power of causal investigation: The Modus Operandi method. In *Evaluation studies review annual*, vol. 1, edited by G. V. Glass, 101-18. Beverly Hills, CA: Sage.
- Shadish, W. R., and T. D. Cook. 1999. Design rules. *Statistical Science* 14:294-300.
- Shadish, W. R., T. D. Cook, and D. T. Campbell. 2001. *Experimental and quasi-experimental designs for generalized causal inference*. Boston: Houghtin-Mifflin.
- Smith, J. C., and P. Todd. 2002. Does matching overcome LaLonde's critique of nonexperimental estimators? *Journal of Econometrics*. Pier working paper 01-35, Penn Institute for Economic Research, University of Pennsylvania.
- Stake, Robert E. 1967. The countenance of educational evaluation. *Teachers College Record* 68:523-40.
- Stebbins, L. B., R. G. St. Pierre, E. C. Proper, R. B. Anderson, and T. R. Cerba. 1978. An evaluation of Follow Through. In *Evaluation studies review annual*, vol. 3, edited by T. D. Cook, 571-610. Beverly Hills, CA: Sage.
- Vinovskis, M. A. 1998. Changing federal strategies for supporting educational research, development and statistics. Background paper prepared for the National Educational Research Policy and Priorities Board, U.S. Department of Education, Washington, DC.
- . 2002. Missing in practice: Development and evaluation at the U.S. Department of Education In *Evidence matters: Randomized trials in education research*, edited by F. Mosteller and R. Boruch. Washington, DC: Brookings Institution.
- Wargo, M. J., G. K. Tallmadge, D. D. Michaels, D. Lipe, and S. J. Morris. 1972. *ESEA Title I: A reanalysis and synthesis of evaluation data from fiscal year 1965 through 1970*. Final Report, Contract No. OEC-0-71-4766. Palo Alto, CA: American Institutes for Research.
- Weiss, C. H. 1988. Evaluation for decisions: Is anybody there? Does anybody care? *Evaluation Practice* 9:5-20.
- Weiss, C. H., and M. J. Bucuvalas. 1977. The challenge of social research to decision making. In *Using social research in public policy making*, edited by C. H. Weiss, 213-34. Lexington, MA: Lexington.
- Whitbeck, C. 1977. Causation in medicine: The disease entity model. *Philosophy of Science* 44:619-37.
- Wilde, Elizabeth Ty, and Robinson Hollister. 2002. How close is close enough? Testing nonexperimental estimates of impact against experimental estimates of impact with education test scores as outcomes. Institute for Research on Poverty Discussion Paper no. 1242-02, University of Wisconsin-Madison.
- Witte, J. F. 1998. The Milwaukee voucher experiment. *Educational Evaluation and Policy Analysis* 20:229-51.
- Zdep, S. M. 1971. Educating disadvantaged urban children in suburban schools: An evaluation. *Journal of Applied Social Psychology*, pp. 1, 2, 173-86.