

NBER WORKING PAPER SERIES

SOCIAL EXPERIMENTS IN THE LABOR MARKET

Jesse Rothstein
Till von Wachter

Working Paper 22585
<http://www.nber.org/papers/w22585>

NATIONAL BUREAU OF ECONOMIC RESEARCH
1050 Massachusetts Avenue
Cambridge, MA 02138
September 2016

We thank Ben Smith and Audrey Tiew for sterling research assistance, and Angus Deaton, Larry Katz, Jeff Smith, and conference participants for helpful comments. The views expressed herein are those of the authors and do not necessarily reflect the views of the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2016 by Jesse Rothstein and Till von Wachter. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Social Experiments in the Labor Market
Jesse Rothstein and Till von Wachter
NBER Working Paper No. 22585
September 2016
JEL No. H53,I38,J22,J24,J31,J65

ABSTRACT

Large-scale social experiments were pioneered in labor economics, and are the basis for much of what we know about topics ranging from the effect of job training to incentives for job search to labor supply responses to taxation. Random assignment has provided a powerful solution to selection problems that bedevil non-experimental research. Nevertheless, many important questions about these topics require going beyond random assignment. This applies to questions pertaining to both internal and external validity, and includes effects on endogenously observed outcomes, such as wages and hours; spillover effects; site effects; heterogeneity in treatment effects; multiple and hidden treatments; and the mechanisms producing treatment effects. In this Chapter, we review the value and limitations of randomized social experiments in the labor market, with an emphasis on these design issues and approaches to addressing them. These approaches expand the range of questions that can be answered using experiments by combining experimental variation with econometric or theoretical assumptions. We also discuss efforts to build the means of answering these types of questions into the ex ante design of experiments. Our discussion yields an overview of the expanding toolkit available to experimental researchers.

Jesse Rothstein
Goldman School of Public Policy
and Department of Economics
University of California, Berkeley
2607 Hearst Avenue
Berkeley, CA 94720-7320
and NBER
rothstein@berkeley.edu

Till von Wachter
Department of Economics
University of California, Los Angeles
8283 Bunche Hall
MC 147703
Los Angeles, CA 90095
and NBER
tvwachter@econ.ucla.edu

I.	Introduction.....	3
II.	What are Social Experiments? Historical and Econometric Background	10
a.	A Primer on the History and Topics of Social Experiments in the Labor Market.....	10
b.	Social experiments as a tool for program evaluation.....	15
i.	The benchmark case: Experiments with perfect compliance.....	16
ii.	Imperfect compliance and the local average treatment effect.....	19
c.	Limitations of the experimental paradigm	21
i.	Spillover Effects and the Stable Unit Treatment Value Assumption	22
ii.	Endogenously observed outcomes	22
iii.	Site and Group Effects.....	23
iv.	Treatment Effect Heterogeneity and External Validity.....	23
v.	Hidden Treatments	24
vi.	Mechanisms and Multiple Treatments.....	25
d.	Quasi-experimental and Structural Research Designs	25
III.	A more thorough overview of labor market social experiments	26
a.	Labor Supply Experiments	27
b.	Training experiments.....	34
c.	Job Search Assistance	44
d.	Practical Aspects of Implementing Social Experiments.....	51
IV.	Going Beyond Treatment-Control Comparisons to Resolve Additional Design Issues.....	54
a.	Spillover effects and SUTVA	56
i.	Addressing the issue ex post.....	57
ii.	Addressing the issue ex ante through the design of the experiment	60
b.	Endogenously observed outcomes	62
i.	Addressing the issue ex post.....	64
Parametric selection corrections.....	65	
Non- and semi-parametric selection corrections	66	
ii.	Addressing the issue ex ante through the design of the experiment	72
c.	Site and group effects	74
i.	Addressing the issue ex post.....	76
ii.	Addressing the issue ex ante through the design of the experiment	83
d.	Treatment effect heterogeneity and external validity	84
i.	Addressing the issue ex post.....	85
ii.	Addressing the issue ex ante through the design of the experiment	90
e.	Hidden treatments	94
i.	Addressing the issue ex post.....	95
ii.	Addressing the issue ex ante through the design of the experiment	97
f.	Mechanisms and multiple treatments	98
i.	Addressing the issue ex post.....	99
ii.	Addressing the issue ex ante through the design of the experiment	110
V.	Conclusion	112

I. Introduction

There is a very long history of social experimentation in labor markets. Experiments have addressed core labor market topics such as labor supply, job search, and human capital accumulation, and have been central to the academic literature and policy discussion, particularly in the United States, for many decades.

By many accounts, the first large-scale social experiment was the New Jersey Income Maintenance Experiment, initiated in 1968 by the U.S. Office of Economic Opportunity to test the effect of income transfers and income tax rates on labor supply. Where many subsequent experiments have been designed to evaluate a single program or treatment each, the Income Maintenance Experiment was intended instead to map out a response surface. Participants were assigned to a control group or to one of eight treatment arms that varied in the income guarantee to a family that did not work and the rate at which this was taxed away as earnings rose. Three follow-up experiments – in rural North Carolina and Iowa; in Gary, Indiana; and in Seattle and Denver – with varying benefit levels and tax rates (and, in Seattle and Denver, a cross-cutting set of counseling and training treatments) were begun before data collection for the New Jersey experiment was complete.

Other early labor market experiments examined the effects of job search encouragement for Unemployment Insurance recipients; job training and job search programs; subsidized jobs for the hard-to-employ; and programs designed to push welfare recipients into work (Greenberg and Robins, 1986; Gueron, this volume). These topics have been returned to repeatedly in the years since as researchers

have sought to test new program designs or to build on the limitations of earlier research. There have also been many smaller-scale experiments, on bonus pay schemes, management structure, and other firm-level policies.¹

From the beginning, the use of random assignment experiments (also known as randomized controlled trials, or RCTs) has been controversial in labor economics.² The primary, powerful appeal of RCTs is that they solve the assignment, or selection, problem in program evaluation. In non-experimental studies (also known as “observational” studies), program participants may differ in observed and unobserved ways from those who do not participate, and econometric adjustments for this selection rely on unverifiable, often implausible assumptions (Lalonde 1986; Fraker and Maynard 1987; though see also Heckman and Hotz, 1989). With a well-executed randomization study, however, the treatment and control groups are comparable by design, making it straightforward to identify the effect of the treatment under study.

But set against this very important advantage are a number of drawbacks to experimentation. Early on, it was recognized that RCTs can be very expensive and hard to implement successfully. For example, it is not always possible to ensure that everyone assigned to receive a treatment receives a full dose, while those assigned to the control group receive none, though this is the experimental ideal. Sometimes it is not feasible to control participants’ behavior, and many participants deviate

¹ We omit here audit studies aimed at uncovering discrimination in the labor market and elsewhere (e.g., Bertrand and Mullainathan 2004; Kroft, Lange, and Notowidigdo 2013; Farber, Silverman, and von Wachter 2015). These are covered by Bertrand and Duflo, elsewhere in this volume.

² For recent criticisms of reliance on RCTs with particular relevance to labor market studies, see Deaton (2010) and Heckman (2010). See also Heckman and Smith (1995).

from their intended treatment assignments. In other cases, ethical, political, or operational considerations make it undesirable to limit access to alternative treatments. Although this can be partly addressed within the basic experimental paradigm, it does limit what can be learned.

More generally, while random assignment solves the assignment problem, it alone is not sufficient to resolve other problems that researchers often face. Many questions of interest can be answered only with something more than the familiar two-armed randomized control trial – a more complex experimental design, the augmentation of experimental data with additional, non-experimental data, theoretically grounded assumptions, or a combination of these. We consider a number of such questions in this chapter. These include:

- *Questions about impacts on endogenously observed outcomes.* Consider the effect of job training on wages. Because wages are observed only for those who have jobs, and because training may affect the likelihood of working, the contrast in mean wages between randomly assigned treatment and control groups does not compare like to like and thus does not solve the assignment problem for this outcome.
- *Questions about spillovers and market-level impacts.* When one individual's outcome depends on others' treatment assignments, experimental estimates of treatment effects can be misleading about a program's overall effect. In the context of labor market programs, an increase in job search effort by a treatment group may lower the control group's job-finding chances, leading to an overstatement of the program's total effect (which will itself depend

- importantly on the scale at which the program is implemented). Similar issues can arise if subjects communicate with each other, leading to a dilution in treatment contrasts when access to information is part of the treatment.
- *Questions about heterogeneity of treatment effects.* Experiments have limited ability to identify heterogeneity of treatment effects, especially if heterogeneity is not fully characterized by well-defined observable characteristics. This is often of first-order importance, as in many cases the relevant question is not *whether* to offer a program (e.g., job training) but *for whom* to make it available, or *which* versions of the program are most effective (and why).
 - *Questions about generalizability.* While in ideal cases experiments have high internal validity for the effect of the specific program under study on the specific experimental population, in the setting in which it is studied, they may have limited external validity for generalizations to other locations, to other programs (or even to other implementations of the same program), or to other populations. For example, a reemployment bonus program may have a very different effect in a full-employment local economy than when the local area is in a recession, or the same program offered in different sites may have dramatically different effects due to variation in local program administration or context.
 - *Questions about mechanisms.* Many questions of interest in labor market research do not reduce to the effects of specific “treatments” on observed outcomes, but relate, at best, to the mechanisms by which those effects arise.

For example, an important question for the analysis of unemployment insurance programs is whether the unemployed are liquidity constrained or whether they can borrow or save to smooth consumption optimally across periods of employment and unemployment. And important questions about the design of welfare and disability policy turn on whether observed non-employment is due to high disutility of work or to moral hazard. In each case, we want to distinguish income and substitution effects, a distinction that is in general not identified from the simple effect of a treatment on an observed outcome. Carefully designed experiments can shed light on the phenomena of interest, but may not be able answer them directly.

To be clear, all of these questions are thorny under any methodological approach, and are generally no easier to answer in quasi-experimental studies than in randomized experiments. One vocal group of critics of experimentation points to the importance of identifying the “structural” parameters – a full characterization of program enrollment decisions and the behavioral processes that lead to the observed outcomes – that determine program selection and impacts (see, e.g., Keane 2010). In principle, many of the design issues above could indeed be avoided or addressed with estimates of the underlying structural parameters. But these structural parameters are difficult to measure. So-called structural methods generally trade off internal validity in pursuit of more external validity, but a study that fails to solve the assignment problem is unlikely to be any more generalizable than it is internally valid.

Unfortunately, while experiments can sometimes be designed to identify a few key structural parameters, or at least important combinations of them, it is rarely possible to design an experiment that directly identifies all of the structural parameters of interest. Thus, there can be value in combining the two paradigms. This involves imposing untestable assumptions about the processes of interest, while still resting on experimentation (or other empirical methods that offer high internal validity) where possible. The additional assumptions can dramatically enhance external validity if they are correct, though if they are incorrect – and this is generally untestable – both internal and external validity suffer.

The current frontier for labor market research – as in other fields – thus involves combining the best features of the two approaches to permit answers to more questions than are addressed by simple experiments, while retaining at least some of the credibility that these experiments can provide.

In this chapter, we discuss a variety of questions common in labor market research that require this sort of approach. We distinguish two broad strategies for answering these questions using experimental data. First, one can augment traditional randomized experiments by imposing additional structure, either economic or econometric, after the fact. In many cases, the amount of structure required, and the strength of the additional assumptions that are necessary, is small relative to the value of the results that can be obtained. Our review gives a snapshot of an expanding toolkit with which researchers can address a wider range of

questions based on variation from RCTs.³

The second broad strategy is to address the limitations of traditional experiments *ex ante*, via design of the experimental intervention or evaluation itself. In many cases, clever design choices – multiple treatment arms, carefully designed stratification, or randomization both across and within groups, for example – can allow for richer conclusions than would be possible via traditional experiments. This sort of approach has a long history – indeed, the very first large-scale social experiments, the income maintenance experiments of the late 1960s and early 1970s, can be seen as a version of it. But the pendulum swung away for a long time, and researchers have only recently begun to return to experimental designs that synthesize random experimental variation with more structural modeling. Recent examples of this approach include Kling, Liebman, and Katz (2007) who use it to address potential biases from endogenous attrition, and Crepon, Duflo, Gurgand, Rathelot, and Zamora (2013), who quantify the importance of spillovers. In our view, approaches like these represent the current research frontier.

The rest of this chapter proceeds as follows. In Section II, we give brief overviews of the history of social experiments in the labor market and of the value of RCTs for solving selection problems, and summarize potential design issues that remain even with random assignment. In Section III, we review the types of programs and questions that have been analyzed, their main findings, and practical

³ This includes analyses of issues such as endogenously observed outcomes (e.g., Ahn and Powell 1993, Grogger 2005, Lee 2009); hidden treatments (e.g., Kline and Walters 2014, Feller, Grindal, Miratrix, and Page 2014, Pinto 2015); heterogenous treatment effects (e.g., Kline and Walters 2014, Heckman and Vytlačil 2005); and multiple treatments and mechanisms (e.g., Card and Hyslop 2005, Schmieder, von Wachter, and Bender 2016, Della Vigna, Lindner, Reizer and Schmieder 2016).

challenges that labor market experiments often confront. Section IV discusses approaches to addressing the design challenges from Section II and thereby expanding the range of questions that can be answered. We discuss both ex ante and ex post approaches to resolving (or at least ameliorating) the issues. Section V offers some concluding comments.

II. What are Social Experiments? Historical and Econometric Background

a. A Primer on the History and Topics of Social Experiments in the Labor Market

As the so-called “credibility revolution” has swept over empirical economics in the last generation, the role and status of experimental evidence has grown. Over the same period, the field of experimental economics has segmented – List and Rasul (2011) and Harrison and List (2004), for example, draw careful distinctions between social experiments and artefactual, natural, and framed field experiments. Briefly, social experiments tend to be conducted at a large scale and to focus on the overall evaluation of policies or programs, often already in place. By contrast, the various types of field experiments are typically smaller in scale and are more likely to use artificial treatments (e.g., behavioral games) that would not correspond directly to any specific policy but are designed primarily to uncover particular behavioral tendencies or parameters.

Although all of the many varieties of experiments have been used to study topics related to the labor market, this chapter focuses on large-scale social experiments, which in our view have had the largest impact on policy.

The social experiment/field experiment distinction corresponds roughly to the distinction drawn above between program evaluation and the identification of structural parameters – social experiments are, at root, evaluations of programs or policies, where field experiments are designed primarily to uncover one or more specific structural parameters.⁴ As we discussed above, this distinction is less clear than it once was – scholars are increasingly drawing on program evaluation samples to understand structural relationships, and using structural parameters to inform the design and interpretation of program evaluations. But while the distinction has been blurred, it has not been obliterated, and nearly all of the social experiments that we discuss in this chapter are designed, at least in part, to evaluate programs that either have been or might plausibly be implemented in roughly the form used in the experiment.

Another, related distinction has to do with the communities that conduct the different types of experiments. Social experiments are typically conducted at a large scale by an organization that specializes in this – historically, the “Big Three” players (Greenberg and Shroder 2004) have been Mathematica, the Manpower Demonstration Research Corporation (MDRC), and Abt Associates – and has been hired by a government agency (most notably OPDR, the Office for Policy Development and Research within the Department of Labor’s Employment and Training Administration, and ASPE, the Assistant Secretary for Policy and Evaluation within the Department of Health and Human Services) or a large

⁴ Kling et al. (forthcoming) refer to experiments aimed at understanding mechanisms rather than at evaluating programs as “mechanism experiments.” Gueron (this volume) discusses the tension between program evaluation and understanding mechanisms in early social experiments.

foundation (e.g., the Ford Foundation) for a specific study. By contrast, field experiments are more often overseen by individual scholars and their students, perhaps with the cooperation of a company or government agency that is not otherwise closely involved in the design.

The differences in the composition and organizational structure of social experimental and field experimental research teams relate to the scope of the work being carried out. A research team implementing a social experiment faces a number of practical and implementation challenges that are largely absent from laboratory experiments and closely related types of field experiments. Researchers rarely have access to a sampling frame corresponding to the population of interest; face practical, ethical, and political difficulties in randomly assigning access to treatment; have limited or no control over treatment alternatives that control participants may obtain or over the specific implementation of the treatment, which is often under the control of an agency rather than the experimenter; and lack ready access to outcome measures for use in assessing the program's impact (or even to a well-defined set of outcomes of interest). Addressing these challenges often requires a large staff to collect pre- and post-treatment data, to minimize attrition between survey waves, and to monitor both the randomization of treatment and the fidelity of treatment delivery to the program model. The required scale is often out of the reach of individual researchers.

Most authors agree that the first large-scale social experiment in the labor market was the New Jersey Income Maintenance Experiment (hereafter, *IME*; this is also known as the New Jersey Negative Income Tax experiment), first initiated in

1968 and extended in various ways in other locations over the next several years. Consistent with the above dichotomy, this was a large-scale experiment that was initiated by the Office of Economic Opportunity (OEO), then an independent agency within the Federal government that played a lead role in the War on Poverty. But in other ways it more closely resembles what would now be called a field experiment, albeit at a massive scale: It was first conceptualized by an individual researcher, Heather Ross, who proposed it to OEO, and it was designed not to evaluate a specific, well developed program but to map out the surface of labor supply responses to a range of tax parameters and thereby to uncover semi-structural economic parameters, the income and substitution effects of changes in tax rates.

Nearly all analyses of IME data went beyond simple treatment-control contrasts, using the data to estimate parametric or semi-parametric labor supply models.⁵ These models often incorporated corrections for the selection introduced by nonparticipation that relied on strong functional form assumptions (e.g., Tobits) and in some cases also rested on structural specifications of the response to nonlinear tax schedules. In many of these studies, the treatment and control groups were effectively pooled and it can be difficult to identify the extent to which the parameters are identified from experimental vs. non-experimental variation.

Another sense in which the IME diverged from much modern social experimental practice was in the source of outcome measures. The main outcome

⁵ Indeed, in 1990 – seven years after the final experimental report from the follow-up Seattle-Denver Income Maintenance Experiment, and after many published analyses of the data – Ashenfelter and Plant (1990) are apparently the first to report the results as simple means by randomly assigned treatment group.

measures for the IME analyses were payments under the IME and labor supply measures drawn from participants' self-reports as part of the program's administration. But as in other experiments, many subjects failed to complete the follow-up surveys. Unfortunately, the design of the IME program meant that the private returns to continued reporting varied dramatically with both treatment status and endogenous outcomes, as the income maintenance payments were made on the basis of these reports. Differential attrition made the results quite difficult to interpret (Ashenfelter and Plant 1990).

In the wake of the Income Maintenance Experiments, the field exploded. Greenberg, Shroder, and Onstott (1999; see also Greenberg and Shroder 2004) identified 21 social experiments between 1962 and 1974, largely in education and health. By contrast, they identify 52 between 1975 and 1982 and 70 between 1983 and 1996, and most of these are directly related to the labor market. (There has not been as systematic a census of post-1996 experiments, but the pace of large scale labor market experiments seems to have dropped off since then, at least in the United States. There has been rapid growth of social experiments in education over this period, however.) Greenberg et al. (1999; hereafter GSO) highlight important changes in the post-1975 experiments. In contrast to the IME, most involved only one or two treatment arms plus a control, and were designed more as "black box" evaluations of the programs encapsulated in the treatments – often modifications on existing programs (Gueron, this volume) – than as efforts to map out a response surface.

GSO emphasize that the vast majority of the experiments they identified focused on low-income populations, a fact that does not seem to have changed since their survey. Several topics stand out as central:

- *Human capital development.* Over one-third of the studies in GSO's sample include at least one treatment arm involving a supported work experience, on-the-job training, vocational education or training, or basic education (including GED programs).
- *Labor supply.* A number of experiments have involved interventions aimed at increasing labor supply, including the income maintenance experiments, studies of re-employment bonuses for unemployment insurance recipients, and a broad group of welfare-to-work experiments conducted as part of the mid-1990s welfare reform movement.
- *Job search assistance.* Another common category of experiments examines interventions aimed at making disadvantaged workers' job search efforts more effective, through counseling, job clubs, or job placement services.

These are not mutually exclusive. In particular, a number of programs and experiments combined job search assistance with either job training or incentives to find work.

b. Social experiments as a tool for program evaluation

Random assignment solves the selection problem that often plagues non-experimental program evaluations, and makes it possible to generate uniquely credible evidence on the effects of well-defined, successfully implemented

programs. In the absence of random assignment, people who participate in a program (those who are “treated”) are likely to differ in observed and unobserved ways from those who do not participate, and the effect of this selection can be distinguished from the causal effect of the program only via the imposition of unverifiable assumptions about the selection process. This is a very important advantage of the experimental paradigm over other research methodologies (so-called “observational” comparisons), and we do not intend to minimize its contributions to the field of economics, public policy, and beyond.

But experiments have limitations as well – while they can have very high internal validity, at closer inspection this is true only for certain types of programs and certain types of outcomes; and even then there can be other challenges, such as difficulties in generalizing from the experimental results to a broader setting.

In this subsection, we discuss the value of experiments as a means of solving the selection problem. We then discuss some of the limitations of the experimental paradigm for program evaluation and policy analysis. Our discussion draws heavily on the Angrist-Imbens-Rubin (1996) “potential outcomes” framework. Some of the limitations we discuss can be addressed via careful design of the experimental study, while others require augmenting experimental methods with other tools. We take up these topics in Section IV.

i. The benchmark case: Experiments with perfect compliance

The appeal of randomized experiments is that they make transparent the assumptions that permit causal inference and create a direct link between the implementation of the experiment and the key selection assumption. The simple

contrast between those randomly assigned to participate in the program and those randomly excluded identifies the effect of being assigned to participate, subject only to the assumption that the randomization was conducted correctly. Moreover, in many cases this effect is identical to the effect of the program on its participants (known as the “effect of the treatment on the treated,” or TOT), which is often the main parameter of interest; in other cases, it is straightforward to convert the effect of assignment to participate (often known as the “intention to treat, or ITT, effect) into an estimate of the program treatment effect for a subpopulation of interest.

These results are well known (see, e.g., Athey and Imbens, this volume), and we do not review them at length here. But it will be useful to have notation later. We use Donald Rubin’s potential outcomes framework for causal inference as set forth in Holland (1986). We consider the evaluation of a simple, well-defined program, such as an in-class job training course or a bonus scheme to encourage rapid return to work after a job displacement, where it is possible to assign individuals separately to participate or to be excluded from participation in the program.⁶ For each individual i , one can imagine two possible outcomes: One that would obtain if i participated in the program, y_{i1} , and one that would obtain if he or she did not participate, y_{i0} .⁷ The program’s causal effect on person i is simply the difference between the outcome which would obtain if he/she participated and that which

⁶ In the case of the bonus scheme, the “treatment” is eligibility for the bonus, not actual receipt.

⁷ This notation rests on an assumption about the mechanisms by which the program operates, known as the “stable unit treatment value assumption,” or “SUTVA.” We discuss SUTVA at greater length below.

would obtain if she did not, $\tau_i = y_{i1} - y_{i0}$. When $\tau_i > 0$, i would have a higher outcome if he/she participated than if he/she did not; when $\tau_i < 0$, the opposite is true.

Let D_i be an indicator for participation, with $D_i = 1$ if i actually participates in the program and $D_i = 0$ if i does not. The simplest estimator of the program's effect is the contrast between the average outcomes of those who participate and those who do not. This can be written as:

$$\begin{aligned} E[y_i | D_i = 1] - E[y_i | D_i = 0] &= E[y_{i1} | D_i = 1] - E[y_{i0} | D_i = 0] \\ &= E[\tau_i | D_i = 1] + (E[y_{i0} | D_i = 1] - E[y_{i0} | D_i = 0]). \end{aligned}$$

Thus, the simple participant-nonparticipant contrast combines two distinct components: The effect of the treatment on the treated, $\tau^{TOT} = E[\tau_i | D_i = 1]$, and a selection term, $E[y_{i0} | D_i = 1] - E[y_{i0} | D_i = 0]$, that captures the difference in outcomes that would have been observed between those who participated in the program and those who did not, had neither group participated (for example, had the program not existed). This second term arises because the process by which people select (or are selected) into program participation may generate differences between participants and non-participants other than their participation statuses. If so, the treatment-control difference cannot be interpreted as an estimate of the effect of the program.

In a simple social experiment, D_i is randomly assigned. This ensures that the distributions of y_{i0} and τ_i are each the same for those with $D_i = 0$ as for those with $D_i = 1$. The first implies that the selection term is zero; the second, that the TOT effect equals the average treatment effect (ATE), $E[\tau_i]$, in the population represented by

the study sample. Thus, the average causal effect is identified, not just in the treated subgroup but in the larger population.⁸

This, in a nutshell, is the value of randomization in program evaluation. In a simple randomized control trial, the identification assumption that justifies causal inference is simply that the randomization was correctly executed. Of course, in any finite sample there may be differences in the sample averages of y_{0i} or τ_i between treatment and control groups. But this variation is captured by the standard error of the experimental estimate. The estimate is unbiased, with measurable uncertainty, so long as the groups are the same in expectation.

ii. Imperfect compliance and the local average treatment effect

A complication that often arises, and that will be central to some of our discussion below, is that it is not always possible to control subjects' program participation. Some subjects who are assigned to receive job training may not show up to their course, while others who are assigned to the control group, and thus not to receive training, may find another way into the program. This can be formalized by introducing an additional variable, Z_i , representing the experimenter's intention for individual i : An individual with $Z_i = 1$ is intended to be served, and one with $Z_i = 0$ is not to be. Z_i is related to D_i , but imperfectly: Some (non-randomly selected) individuals who are assigned $Z_i = 1$ will wind up with $D_i = 0$ (e.g., those who fail to

⁸ This holds if the entire population of interest is part of the experiment. If the study sample is not representative of the broader population, the ATE identified will be local to the subpopulation represented by the sample.

arrive for their assigned training course), and others who are assigned $Z_i=0$ will wind up with $D_i=1$, (e.g., those who talk their way past the program screener).

With partial compliance, the experiment identifies neither the average treatment effect (ATE) nor the average effect of the treatment on the treated (TOT). Rather, the best that can be identified is the *local* average treatment effect, or LATE, for the subgroup of experimental subjects who comply with their experimental assignment. Specifically, let D_{i0} represent the individual's treatment status if assigned $Z_i = 0$ and D_{i1} represent the treatment status if assigned $Z_i = 1$. The “complier” subpopulation is defined as those with $D_{i0} = 0$ and $D_{i1} = 1$ – those who receive the treatment if and only if they are assigned to receive it. The contrast between the average outcomes of those assigned to receive and not to receive treatment is then:

$$E[y_i | Z_i = 1] - E[y_i | Z_i = 0] = Pr\{D_{i0}=0, D_{i1}=1\} * E[\tau_i | D_{i0}=0, D_{i1}=1].^9$$

This is known as the “intention to treat” (ITT) effect. The first term is the complier share of the experimental population; the second is the local average treatment effect (LATE) for compliers.

In many cases, the ITT is the effect of primary interest. It represents the actual effect of offering access to the program in the setting in which the experiment takes place. Often, it is only possible to manipulate the option to participate (consider, for example, the offer of job training – one can never force individuals to

⁹ We assume here, as in nearly all analyses of experiments with partial compliance, that there are no “defiers” who receive the treatment if and only if they are assigned *not* to receive it ($D_{i0} = 1$ and $D_{i1} = 0$).

participate in a training program), so the effect of manipulating this offer is the key parameter for evaluation of the programs under consideration.

In other cases, however, one might want to identify the effect of program participation (as distinct from the offer to participate). One can recover the LATE for compliers by dividing the ITT by the complier share, which can be identified as $E[D_i | Z_i = 1] - E[D_i | Z_i = 0]$; equivalently, the LATE can be recovered from an instrumental variables regression using Z_i as an instrument for D_i .

The LATE may differ from the ATE or even from the TOT. For example, in many settings one would expect that people who will receive the largest benefits from treatment to make disproportionate efforts to obtain it, even if assigned to the control group; in this case, the TOT will exceed the LATE. Unfortunately, the compliers are not always the population of primary interest. Further structure, or successful randomization of D_i itself, is required to identify the ATE or TOT.

c. Limitations of the experimental paradigm

The basic experimental paradigm is invaluable for its ability to resolve the fundamental problem of causal inference, by ensuring that estimated program effects are not confounded by selection into treatment. But it cannot solve all identification problems faced by program evaluators, nor answer all questions posed by labor economists seeking to understand the workings of the labor market. In the remainder of this section, we will briefly introduce six (partially overlapping) design issues that commonly arise in labor market experiments. In each case, identifying the effects of interest may require moving beyond the treatment-control

contrast in outcomes from a simple randomized experiment. We discuss each in more detail in Section IV, where we also discuss potential solutions to each.

i. Spillover Effects and the Stable Unit Treatment Value Assumption

The above brief overview of the econometrics of experiments glosses over an important assumption, known as the “stable unit treatment value assumption,” or SUTVA (Angrist, Imbens, and Rubin 1996; Athey and Imbens, this volume). Intuitively, this assumption states that the outcome of individual i is unaffected by the treatment status of each of the other study participants. Without this assumption, each individual has not two but 2^N potential outcomes, making analysis intractable. For many program evaluations, SUTVA is innocuous. But in other cases it can be quite restrictive. For example, the provision of job search assistance to some individuals may create “congestion” in the labor market, reducing the job-finding rates of others participating in that market. This is a violation of SUTVA, and will lead a simple randomized trial to overstate the total effect of job search assistance. Another potential violation of SUTVA occurs if members of the treatment group interact with each other or with the control group in a way that dilutes the treatment difference between them – for example, if the treatment involves information provision but treated individuals pass that information on to the controls.

ii. Endogenously observed outcomes

In many labor market experiments, some outcomes of interest are observed only for a subset of individuals. For example, weekly hours of work (labor supply),

hourly wages, job characteristics, career advancement, and retention on the job are observed only for those who are able to find jobs, not for those who are unemployed. Even ideal experiments with perfect compliance may not identify the causal effects of interest on these outcomes.

iii. Site and Group Effects

Another large class of limitations in experiments has to do with generalizing beyond the experimental sample. Extrapolations to other programs, other samples, or other treatment regimes can be hazardous. We will discuss in this paper three broad classes of external validity issues.

One class has to do with variations in the treatment on offer across program locations. In many programs, the treatment is not homogeneous across locations; in other cases, the treatment may be homogeneous but outcome distributions vary. In either case, one might be interested in identifying how treatment effects vary across locations.

The second class derives from observed differences between the population of interest and that included in the experimental sample – one might want to understand a program’s effect on a population that differs in observable ways from that represented in the experimental sample, or on a subpopulation other than the experimental compliers.

iv. Treatment Effect Heterogeneity and External Validity

The third class of external validity issues arises from *unobserved* differences in individual treatment effects – when the effect of the treatment varies across

individuals in ways that are not captured by observed participant characteristics, and when the parameters of interest extend beyond the average treatment effect in the population from which the experimental sample is drawn. This can occur when, for example, the experimental complier share is not expected to match the take-up rate when the program is offered more generally, or when one expects to offer the program to a population that may differ in its treatment effect distribution from the experimental population. While conceptually similar to differences along observed characteristics, the econometrics behind addressing unobserved differences in treatment effects is sufficiently complex and self-contained that we discuss it separately.

v. *Hidden Treatments*

Interpreting estimated program effects and extrapolating to other settings can be complex even in the case of uniform treatments and uniform populations. For example, if non-compliers have access to alternatives to the program under study (e.g., to courses offered by alternative job training providers), this will lead to variation in treatment effects even without treatment effect heterogeneity or non-compliance in treatment assignment in the standard sense. The alternative treatments are often “hidden,” as administrative data on the program under study will not reveal whether participants have received alternatives elsewhere. In this case, the experimental impact identifies the treatment’s effect relative to a poorly specified alternative that may not differ dramatically, and may be a poor guide to the program’s value relative to no treatment. In multi-site studies, differential take-up of such hidden treatments by the control group may create the appearance of

treatment effect heterogeneity across sites and hinder extrapolation to other settings.

vi. *Mechanisms and Multiple Treatments*

In many instances, we are interested in understanding the mechanism generating a particular treatment effect. In some cases, the effects of separate mechanisms are of inherent interest. In complex experiments with multiple treatments, it is important to understand which treatments were particularly effective, and why. For example, many job training programs include job search assistance, and vice versa. In other cases, understanding the mechanisms is crucial in extrapolating from the particular experimental setting to other situations. For example, in the Canadian Self-Sufficiency Program (SSP) workers have to first establish eligibility to then participate a wage subsidy program, creating endogenous selection that makes it difficult to interpret how the subsidy program affects labor supply (Card and Hyslop 2005). Without additional information or additional structure, multiple mechanisms are not separately identified, leading to potential serious limitations in understanding of the program and in external validity.

d. **Quasi-experimental and Structural Research Designs**

It is not always possible to use a true randomized experiment to evaluate a program or mechanism of interest, due to operational, financial, or ethical constraints. Quasi-experimental studies rely on aspects of the program or policy variation as a source of plausibly as-good-as-random variation in treatment

assignment – examples include regression discontinuity designs, regression kink designs, and difference-in-differences (see Angrist and Krueger 1999). These can be useful alternatives when true experiments are infeasible or simply not available. When the quasi-experimental variation is as good as randomly assigned, the various quasi-experimental designs can recover treatment effects just as can experiments.

But even if the assumptions governing assignment are correct, quasi-experimental designs generally solve only the assignment problem, and do not necessarily address the additional issues discussed above. The same is true for selection-on-observables estimators (e.g., matching estimators): The “unconfoundedness” assumption eliminates the selection problem, if it holds, but does nothing to address other design issues.

In contrast, structural approaches that explicitly specify all aspects of the choice problem and resulting outcomes can in principle resolve both assignment and other design issues simultaneously. However, this approach hinges on the model being correctly specified, and hence may come at a substantial cost to internal validity.

III. A more thorough overview of labor market social experiments

It is no accident that we discuss design issues of RCTs in the context of social experiments in the labor market, since many of the major design issues discussed in Section II arise in the evaluation of important labor market programs. In this section we review some of the main characteristics of existing social experiments in labor economics in light of these design issues. We distinguish three broad substantive

topics that have been studied extensively via social experiments: Labor supply, particularly of low-income families, welfare recipients, and unemployment insurance recipients; job training and skill development; and job search. In this Section, we discuss each in turn. For a more detailed discussion of the experiments we mention here, we refer the reader to our summary tables, and excellent overviews provided elsewhere.¹⁰

a. Labor Supply Experiments

One can broadly categorize social experiments providing incentives to increase labor supply into three groups, following their program structure, target group, and time period: The Income Maintenance Experiments in the late 1960s and early 1970s; welfare reform experiments in the late 1980s through the mid-1990s; and reemployment subsidy experiments, which span a longer time period.

The Income Maintenance Experiments

A first wave of experiments were the Income Maintenance Experiments (IME) already discussed in Section II, which treated low-income households with various combinations of lump-sum transfers and taxes on earnings. By randomly assigning treatment and control groups to multiple treatment arms with varying combination of tax rates and subsidies, and by separately targeting groups of different income levels, the experiments allowed tracing out labor supply responses

¹⁰ See among others Greenberg and Shroder (2004), Heckman, Lalonde, and Smith (1999), Meyer (1995). Our overview focuses almost exclusively on U.S. experiments. For an overview of active labor market policy evaluations, drawing largely on European evidence, see Card, Kluve, and Weber (2010).

in different parts of the budget constraint and under varying financial conditions. There were four such experiments, initiated between 1968 and 1971, in New Jersey, Seattle-Denver, Gary (IN), and in rural areas. Table 1 provides detailed information about these experiments. While the sample sizes were moderate by later standards, the total cost was substantial compared to most randomized evaluation of labor supply incentives that would follow. This is in important part because the program – the payments themselves – was expensive on a per-participant basis. Complex, stratified experimental designs were used in efforts to minimize these costs, but even with these the studies were major investments.

Across each of the income maintenance studies and various comparison groups (e.g., husbands, wives, and single female household heads), labor supply results were fairly consistent: The combination of a lump-sum transfer and a positive tax rate reduced participants' earnings (i.e., labor supply), by more so when the transfer and tax rate were larger. This reflects a combination of income and substitution effects; Robins (1985) combines the various studies and uses contrasts among the different treatment arms to separately identify the income and substitution elasticities of labor supply. He concludes that these elasticities were fairly stable across studies, but fairly small: The substitution elasticity was under 0.1 for husbands, just above 0.1 for single female heads, and more variable but averaging 0.17 for wives. Income elasticities were less consistent, but centered around -0.1.

In retrospect, these experiments encountered a number of the design issues that we identified in Section II and discuss at greater length below. For example,

because of the high attrition rates, which as Ashenfelter and Plant (1990) note were differential across treatment groups, they also can be seen as an example of the endogenously observed outcomes problem. Similarly, without additional assumptions it is impossible to estimate the effect of these programs on hours worked or wages. Interestingly, in contrast to most randomized evaluations that followed, they were primarily focused on identifying the mechanisms – income vs. substitution effects – behind any labor supply responses, rather than the simple treatment effect of an existing program. This motivated the use of a large number of treatment arms, an option we discuss below as one way of addressing questions about mechanisms.

Welfare Reform Experiments

A second wave of social experiments related to labor supply was initiated between the late-1980s and the mid-1990s, and evaluated the effect of employment incentives for welfare recipients. While the IME experiments were funded almost exclusively by the federal government, these later evaluations concerned state-level programs and were funded mostly at the state level.¹¹ In contrast to the relatively straightforward structure of the negative income tax treatments, these were usually randomized evaluations of entire, complex programs, often designed as replacements for traditional AFDC, that included components designed to strengthen work incentives along with others (e.g., child care or job search assistance) designed to reduce barriers to work.

¹¹ For a detailed historical account, see the chapter by Judith Gueron in this volume.

We have identified welfare RCTs in at least 13 states. Table 1 includes a selection of four social experiments on this topic, implemented in California, Connecticut, Florida, and Minnesota, though there were many more not listed here. A common component to most new programs (experimental treatments) was the introduction of lifetime time-limits of welfare receipt and increases in earnings disregards, both eventual components of the 1996 federal welfare reform – prior to this reform, implementation of such changes required a waiver from the U.S. Department of Health and Human Services, and this was often conditioned on an experimental evaluation. The exact nature of both the new programs and the traditional welfare benefit varied by state. Other program features varied widely as well, including job search assistance, access to child care, changes in case management, and provision of job training.

Two examples to which we will refer to later are Connecticut’s Jobs First and Florida’s Family Transition Program. In both cases, control group members faced a welfare benefit schedule that had no time limits and high implicit taxes on working.¹² Jobs First and the Family Transition Program each introduced time limits for welfare receipt and benefit schedules with lower implicit tax rates. Under Jobs First, eligible welfare recipients saw no reduction in their benefits while working until earnings hit the federal poverty line. Under the Family Transition Program, a working welfare recipient could keep \$200 a month, plus 50% of all earnings above

¹² In Connecticut, welfare recipients were eligible for a fixed earnings disregard of \$120 for the twelve months following the first month of employment while on assistance and \$90 afterwards. Recipients were also eligible for a proportional disregard of earnings above \$120 (\$90): 51% for the four months following the first month of employment and 27% afterwards. In Florida, after the first four months of work, the marginal tax rate on earnings for AFDC recipients was 100% if they earned over \$90 per month.

\$200. Both programs also modified other welfare program features, including enhanced enforcement of work requirements, changing the duration of access to Medicaid benefits, setting asset limits for welfare receipt, and providing child care assistance, among others.

The randomized evaluation of the two programs captured the combined effects of all of these changes on employment and earnings. Each program led to higher earnings and higher total incomes, inclusive of welfare payments, in the treatment group, though in each case this effect diminished over time. Total governmental costs were higher for the Connecticut treatment group than for controls, but the reverse was true in Florida. An important caveat is that these results largely reflect the period before time limits bound.

In many of the welfare-to-work experiments, key outcomes of interest included hours of work among those who are employed and wages or earnings. Neither of these is observed for those who are not employed. Thus, although many studies report experimental effects on endogenously observed outcomes, these are understood to suffer from serious selection problems. Another issue to take into account in interpreting these experiments is the possibility of spillover effects. These were typically not small pilot studies but involved broad changes to welfare rules, sometimes applied to all program participants except for a hold-out control group.

Another major question regarding welfare-to-work programs concerns heterogeneity in treatment effects. One might imagine that there is a subpopulation of recipients who are responsive to work incentives and another group of hard cases

who are much less responsive. The average treatment effects that can be estimated from these experiments might substantially overstate the employability of the latter participants.

Reemployment Subsidy Experiments

A third broad group of labor supply-related experiments evaluated direct reemployment subsidies. One set of such programs had incentives structured like a negative income tax and were targeted to welfare recipients or low-income individuals, sometimes as part of the same AFDC reforms discussed above. These took place mostly in the mid- to late-1990s, and included the Canadian Self-Sufficiency Program (SSP), Minnesota's Family Investment Program (FIP), and Wisconsin's New Hope Project. These RCTs can be seen as evaluations of welfare-like programs, but included subsidies that were conditional on sustaining a certain amount of employment. Not surprisingly, these programs generally led to increased earnings among treatment group participants (though FIP was an exception); different studies varied in whether the additional income of participants was larger or smaller than the extra welfare costs borne by the government.

Another set of such programs were schemes that paid lump-sum subsidies conditional on employment – effectively, bonuses for finding work. These include the well-known reemployment bonus experiments targeted at unemployed workers receiving unemployment insurance in Illinois, Pennsylvania, and Washington State in the mid-1980s. These studies found that eligibility for a relatively large reemployment bonus led to shorter unemployment insurance spells, with no

detectable impact on the quality of the job obtained, but that the effects were relatively small and thus the programs were not cost effective.

More recently, a bonus for welfare recipients who found a job and who remained reemployed for a certain time was evaluated in the context of Texas' Employment Retention and Advancement (ERA) project in the early 2000s (Dorsett et al. 2013). The Texas evaluation was part of a large-scale randomized evaluation of 12 different service combinations in different U.S. cities from 2000 to 2004 under the ERA project umbrella (Hamilton and Scrivener 2012). The main focus of ERA was to expand workforce services to recently reemployed welfare recipients or low-wage workers to maintain successful labor force attachment (though three sites, including Texas, combined pre- and post-employment assistance). The evaluation tested a broad range of services, with at best mixed results regarding the effect of post-employment services tested.

An important feature of several of these employment subsidy programs was that potential recipients had to become eligible for the subsidy, usually by working a minimum amount of hours. Hence, while the main goal of the programs was to help workers build attachment to the labor force, effects of the subsidy (as distinct from the subsidy offer) on the duration of employment could be estimated only for those who found jobs in the first place, a subsample that was differentially selected in the treatment and control groups. Card and Hyslop (2005)

refer to this as an 'eligibility effect'; in our earlier taxonomy of design challenges, this can be seen as a case where the mechanisms underlying the

treatment effect are of primary interest. Under any name, it complicates the interpretation of the outcomes of a simple RCT.

Overall, randomized studies of a range of labor supply incentive programs have found labor supply responses to changes in implicit or explicit financial incentives as predicted by theory. However, a broad theme emerges that employment effects have mostly been short-lived, and effects on total participant income inconsistent. A challenge in interpreting these studies has been that typically a number of treatments were varied simultaneously, including implicit tax rates and lump-sum transfers, training programs, job search assistance, enforcement and/or time limits. Hence, extrapolating from these findings to new programs providing different combinations of treatments is difficult without understanding the underlying behavioral responses, which typically requires additional assumptions.

b. Training experiments

From 1964 to today, we count over 50 RCTs that evaluate job training programs of various forms. These include large-scale evaluations conducted at the national level, state-level evaluations, and evaluations of programs at the local level. The programs evaluated varied substantially in the type of training, which ranged from vocational and general classroom based training of different durations to on-the-job training by actual employers. Most training programs were complemented by some kind of job search assistance, but in the studies we review here this was not the emphasis. Table 2 provides an overview of a selected group of these RCTs.

Training programs are less easily classified than labor supply programs. While the first job training social experiment of which we are aware focused on laid off workers (the General Education in Manpower Training experiment, begun in 1964), the vast majority of training programs are targeted to welfare recipients, to low-income individuals generally, or to low-income youth. Moreover, while one can broadly distinguish phases of experimental evaluation parallel to the patterns in the evaluation of welfare programs outlined above, randomized evaluations of training programs occurred more evenly from the 1980s to today. It is also harder to discern common patterns in the types of training provided or programs evaluated.

The first large-scale evaluation of a mix of on-the-job experience and supervision for hard-to-employ individuals was the National Supported Work Demonstration (NSWD), which ran from 1975 to 1980. The NSWD was a large and expensive social experiment implemented by the U.S. at the national level, but did not evaluate an established training program. Rather, the NSWD relied on local non-profits to organize a program in which treatment participants were placed in teams of up to 10 participants working under a foreman, who also served as a counselor and later provided job search assistance, on small-scale projects, typically in construction, light manufacturing, or social service provision. Participants received as much as one year of work experience, under conditions of increasing demands, close supervision, and work in association with a crew of peers. The study targeted four groups of workers: women that had been on AFDC for at least 30 months; ex-addicts; ex offenders; and young high-school dropouts. It took place at 10 sites, and

at each sites enrollees were selected randomly from a group of volunteers.¹³

Participation had large positive effects on AFDC recipients and smaller positive effects on ex-addicts, but benefits for other groups were smaller and generally statistically insignificant.

The data used to evaluate NSWDC came from a series of follow-up surveys.¹⁴ Attrition was an issue here: After 27 months, only 72% (68%) of the treatment (control) groups of the NSWDC completed interviews. As in the NIT studies, this can be seen as a variant of the endogenously observed outcomes problem.

The NSWDC study was followed by a range of evaluations of state-level programs in the early- to mid-1980s. These were targeted almost exclusively at welfare recipients, and largely financed by the federal government. These evaluations continued, with greater involvement of state governments, through the late 1980s and mid-1990s. While many of these RCTs were relatively small, some were substantial. Examples include the California GAIN and Ohio JOBS program evaluations, beginning in 1988 and 1989, respectively. Detailed characteristics of some of these evaluations are shown in Table 2. The California program, which was mandatory for welfare recipients, included job search assistance, basic education, and skills training. It had large positive effects on earnings and negative effects on welfare receipt, particularly for single parents. Effects were largest in Riverside County, where administrators emphasized job placement as the central goal.

¹³ The Manpower Demonstration Research Corporation (MDRC) was founded in 1974 to manage the NSWDC study. For a detailed summary of the program and findings, see Manpower Demonstration Research Corporation Board of Directors (1980).

¹⁴ The NSWDC has been examined by an extensive literature, including Lalonde (1986), Dehejia and Wahba (2002), and Smith and Todd (2005).

However, a reanalysis of the long-term effects of GAIN by Hotz et al. (2006) found that the effects in Riverside County were short-lived relative to those in Los Angeles County, which focused more on human capital development and where effects were initially smaller but rose over time.¹⁵ The Ohio program was similar in design but encountered more problems in implementation, and yielded smaller effects.

An exception to the trend towards evaluation of state-level or local training programs was the large-scale, national evaluation of the main federal training program aimed at low-income adults and disadvantaged youth – the National Job Training Partnership Act (JTPA) Study. The JTPA was a federal program enacted in 1982, and was administered at the state and local level. JTPA training programs provided employment training for specific occupations and services, such as job search assistance and remedial education, to roughly one million economically disadvantaged individuals per year. While the program and some services were administered directly by JTPA staff, training was provided through local service providers, such as vocational-technical high schools, community colleges, proprietary schools, and community-based organizations. Training lasted three to four months, on average, but duration varied widely across individuals and program sites.

Congress, in part responding to limitations of non-experimental evaluations of the predecessor program to JTPA, the Comprehensive Employment and Training Act, mandated a randomized evaluation of JTPA in 1986. Control subjects were

¹⁵ Hotz et al. (2006) also point out that the treatment group was selected differently between the four GAIN sites, possibly contributing to the estimated ‘site’ effects. For example, the Riverside County RCT sample included a smaller fraction of the more disadvantaged welfare recipients.

excluded from obtaining JTPA services for 18 months. To assess short- and medium-term program impacts on employment and earnings, the evaluation both collected survey data and drew from administrative state-level records.¹⁶ The evaluation took place at 16 JTPA program sites (so called Service Delivery Areas, SDAs).

Participation by SDAs in the evaluation was voluntary, and some SDAs objected to randomly excluding eligible applicants. The participating SDAs did not differ from others in observable characteristics (e.g., Bloom et al. 1997), but may have differed in unobserved ways that would be relevant to an extrapolation to the overall effect of the national program.

An explicit goal of the JTPA evaluation was to obtain differential impacts for a wide range of target groups, including adult women, adult men, female youths, and male youth with and without an arrest record. Adult women saw the largest earnings gains, followed by adult men; effects on youth were smaller and generally not significant (though there were significant effects on attainment of high school diplomas for both adult women and female youth). In addition to demographic subgroup analyses, heterogeneity in program impacts was estimated along several other dimensions, including JTPA services recommended by program intake staff, ethnicity and prior labor market experience. While the subgroup effects of interest were largely pre-specified, this does not fully eliminate multiple-comparisons problems, particularly when the number of pre-specified comparisons is so large, and thus there is an enhanced risk of a false positive.

¹⁶ See Bell et al. (1994) and Bloom et al. (1997) for descriptions of the JTPA evaluation. There is a substantial literature on the evaluation of the JTPA program. See Heckman, Lalonde, and Smith (1999) for a summary.

Job training evaluations slowed after welfare reform in the mid-1990s, then began to pick up again in the early 2000s. Some evaluations in this period focused on sector-specific employment, such as the Sectoral Employment Impact Study (e.g., Maguire et al. (2010) and evaluations of similar smaller, local programs.¹⁷ There was also a randomized evaluation of combined training and job placement services under the Workforce Investment Act (WIA) from 2005 to 2015 (the Work Advancement and Support Center Demonstration), and more recently a study of the return from community college attendance under the Trade Adjustment Assistance Community College and Career Training (TAACCCT) Grants Program.

A distinct broad strand of randomized evaluations of training programs focuses on low-income youths. Again, these programs offer a broad range of different types of training augmented by varying combinations of support services. Social experiments in this area have included a range of federally and nationally funded evaluations ranging from the early 1980s to the mid-1990s that culminated in the National Jobs Corps Study, described below. As in other job training studies, the pace of experimentation slowed in the mid-1990s, but several new studies were undertaken in the mid-2000s. Some randomized evaluations, such as New York City's Summer Youth Employment Program (strictly, a natural experiment, as randomization is part of the rationing process and not a decision made in order to facilitate an evaluation), are ongoing. Again, the broad trend was from a federal

¹⁷ These include, among others, the Georgia Works programs, Project Quest in San Antonio, the Wisconsin Regional Training Partnership in Milwaukee, Per Scholas in New York City, and the Jewish Vocational Service in Boston.

monopoly on funding towards a greater involvement of local and private funding sources.

The largest and perhaps best known study of a training program for disadvantaged youths is the National Jobs Corps Study. The Job Corps was created in 1964 as part of the War on Poverty, and currently operates under the provisions of the Workforce Innovation and Opportunity Act of 2013, which consolidated programs authorized under the Workforce Investment Act of 1998. Job Corps services are geared towards economically disadvantaged youths aged 16 to 24. Core services are delivered by a Job Corps center, usually residential, and include vocational training, academic education, residential living, health care, and a wide range of other services, including counseling, social skills training, health education, and recreation.¹⁸ About a quarter of the over 100 centers are operated directly by the U.S. government, with the remainder operated by private contractors. The average duration of the program is eight months, though by its philosophy the duration responds to the participant's needs and actual duration varies widely. For six months after the youths leave the program, placement agencies help participants find jobs or pursue additional training.

The Job Corps evaluation was based on an experimental design in which, with a few exceptions, all youths nationwide who applied to Job Corps in the 48

¹⁸ The majority of training is vocational, and curricula were developed with input from business and labor organizations and emphasize the achievement of specific competencies necessary to work in a trade. Academic education aims to alleviate deficits in reading, math, and writing skills and to provide a GED certificate. Although most Job Corps services are residential, there have been nonresidential participants (mostly women with children). There have been efforts to evaluate non-residential Job Corps services (e.g., Greenberg and Shroder 2004, Schochet et al. 2008).

contiguous states between November 1994 and December 1996 and were found to be eligible were randomly assigned to either a program group or a control group. Program group members were allowed to enroll in Job Corps; control group members were excluded for three years after random assignment. The comparisons of program and control group outcomes represent the effects of Job Corps relative to other available programs that the study population would enroll in if Job Corps were not an option.¹⁹ The control and treatment groups were tracked with a series of interviews immediately after randomization and continuing 12, 30, and 48 months after randomization.

The evaluation of Job Corps followed the outcomes of over 15,000 experimental subjects for up to eight years using survey and administrative data. The effect of training on earnings became gradually positive as individuals graduated from the program, and then remained statistically significantly different from the control group for up to four years afterwards. At the same time, government transfers and crime rates fell (e.g., Schochet et al. 2008). There was substantial heterogeneity in outcomes – the effects were strongest for those 20-24 year old at the time of training, and weakest for Hispanics.

A concern with these findings was that the overall level of earnings and the size of the treatment effects were quite different in the administrative data than in the survey data. While survey data are more to be affected by endogenous attrition, administrative data are not a panacea: They exclude under-the-table employment,

¹⁹ Of course, if Job Corps did not exist, the ecosystem of other available programs would presumably change. This is formally a SUTVA violation, and implies that control group mean outcomes may not equal what would be seen in the absence of the program.

which may be common in the Job Corps population.²⁰ They also cannot address the problem that wages are observed only for those who are employed, itself an intermediate outcome of the program (e.g., Lee 2009)

An important question regarding Job Corps is the relative performance of the different Job Corps centers, which operate in different labor markets and are (sometimes) run by contractors rather than directly by the government. Schochet and Burghardt (2008) use the Job Corps evaluation data to estimate separate treatment effects by site, finding that these are not strongly correlated with the non-experimental measures that have been used to assess site performance.

A final issue in the Job Corps evaluation, not to our knowledge addressed in the literature, is that the program may be large relative to the relevant labor markets, creating the possibility of important spillovers from treated to control study participants.

A final, smaller category of large-scale social experiments of training programs focused specifically on unemployed (displaced) workers. As we will discuss below, some of these RCTs evaluated programs providing a broad array of reemployment services that also included some degree of training. This raises a similar issue to what we highlighted above with welfare experiments – experimental evaluations generally identify the “black box” effect of the overall programs, but not the components or mechanisms responsible for those effects.

²⁰ Kornfeld and Bloom (1999) show that this is the case for participants in the Job Training Partnership Act (JTPA) evaluation.

The Individual Training Account (ITA) Experiment running from 2001 to 2005 directly evaluated different modes of training provision prescribed by the 1998 Workforce Investment Act. WIA allowed local agencies to impose different degrees of counseling and supervision of workers' training choices, and the ITA experiment evaluated the effect of these choices on actual training received and labor market outcomes. Effectively, the ITA experiment compared three service models. Guided Choice and Maximum Choice had standardized subsidies for training, but the former required counseling by a case worker while the latter had no counseling requirement. A third model, Structured Choice, was effectively like Guided Choice but offered individualized, and typically more generous, training awards.²¹

The findings indicated that either more generous awards (Structured Choice) or less counseling (Maximum Choice) led to a higher incidence of training (Perez-Johnson et al. 2011). Earnings increased for workers in Structured Choice relative to Guided Choice five years after the treatment. (Earnings effects were higher but not statistically different for Maximum Choice relative to Guided Choice or to a control group.) While Structured Choice was estimated to be cost efficient to society, it was more expensive for the workforce system, and most agencies adopted Guided Choice as the leading model. More recently, an ongoing experiment (the WIA Adult and Dislocated Worker Programs Gold Standard Evaluation, discussed below) evaluates directly the intensive and training services provided under WIA.

²¹ Originally, under Structured Choice case workers were supposed to play a more active role in training choice. However, most case workers did not feel they had enough knowledge of local labor markets or the worker's skills to take on such an active role.

An issue that is common to all of the job training experiments is the possibility that individuals assigned to the control group may have received training through other channels that would not necessarily have been tracked in the experimental data. These hidden treatments are likely to attenuate the estimated training effects – insofar as control participants are receiving substitute treatments, the evaluations identify only the *differential* effect of the public training program, rather than the overall effect of training relative to none. While this could partly explain low estimated treatment effects, this has not been examined carefully in the literature (though, as we discuss below, it has received substantial attention in some other domains, most notably the evaluation of early childhood education).

Although a broad range of findings from different treatments makes it hard to generalize, two themes have emerged from training program social experiments. First, while training for less advantaged adults and the unemployed can have beneficial effects, most training programs for disadvantaged youths fail to achieve strong results. An important exception is Job Corps, which has shown short- and medium-term positive effects for at least some of its participants. Second, the effects of training tend to accrue gradually over time, making them hard to detect in research designs that combine multiple treatments or that do not have sufficient data or samples to precisely estimate medium- to long-term effects.

c. Job Search Assistance

From the inception of welfare programs in the U.S. it was suspected that neither better work incentives nor better human capital would be sufficient to place

hard-to-employ welfare recipients or disadvantaged youth into lasting employment, and that part of the challenge derived from disconnection from the world of work. At the same time, it was not clear which of a range of support services aiding job placement would be effective. Hence, a large number of RCTs have evaluated a range of job search assistance (JSA) programs for low-income workers and youth. Other studies have focused on unemployment insurance recipients and other unemployed workers, who have traditionally been eligible for search assistance from the U.S. government. Hence, while training evaluations have mostly concerned programs aimed at low-income workers, job search assistance experiments have evaluated programs geared towards a wider range of unemployed workers from the mid-1970s to today. As in training evaluations, however, an important challenge in studies of job search assistance is measuring the counterfactual: What sort of assistance, if any, was received by those excluded from the program under study?

An early wave of JSA program experiments geared towards welfare recipients occurred from the early 1970s to the mid-1980s, alongside similar studies of labor supply and training programs aimed at the same population. These were mostly evaluations of local programs funded by the federal government. There is a long history of programs providing placement and training services for welfare recipients in the United States, going back at least to the Work Incentive Program (WIN) initiated in 1967. WIN was criticized on a range of fronts (e.g., Gold 1971). The first wave of federally-funded evaluations tested services provided by the WIN program and alternative programs for WIN-eligible welfare recipients (e.g., Grossman and Roberts 1989). These culminated in the National Evaluation of

Welfare-to-Work Strategies (NEWWS) in 1990, which was a large-scale evaluation of 11 programs combining JSA, training, and enforcement of job search requirements in 7 different sites in the U.S.

The results from randomized evaluation of different WIN services were mixed (e.g., Greenberg and Shroder 2004). The evaluation of so-called “job clubs” in 1976-1979 showed substantial increases in employment and reduction in welfare receipt. As result, job clubs became an integral part of services received by welfare recipients. However, the evaluation was based on a relatively small sample, follow-up was limited to one year, and the results indicated substantial, hard-to-explain heterogeneity in the findings across subgroups and treatment sites. In contrast, the evaluations discussed in Grossman and Roberts (1989) show less consistent effects of JSA under the WIN program.

The much larger evaluation of NEWWS found short-term increases in employment and reductions in welfare receipt. These effects dissipated during the five year follow-up period. As in other evaluations occurring in the early to mid-1990s, such as GAIN discussed above, this may be due in part to the high-pressure labor market of the 1990s. The presence of such cyclical effects is a potentially important confounder limiting the interpretation of the effects of labor market program studies.

A second wave of experiments occurred in the run-up to welfare reform in the mid-1990s, and again saw substantial state-level involvement. As with labor supply and training studies in this period, these studies tended to study contemplated changes to existing programs and to involve large samples. These

included Project Independence in Florida in 1990 (over 13,000 treatment and 4,000 control subjects), the Indiana Welfare Reform Evaluation in 1995 (over 67,000 treatment and 4,000 control subjects), and the LA Jobs First GAIN evaluation in 1995 (over 15,000 treatment and 5,000 control subjects).²² Among these, only the GAIN evaluation discussed above allows inference about the role of JSA alone. The findings confirms that JSA can yield substantial gains in employment, at least in the short term.

In parallel, another group of experiments evaluated JSA services provided to recipients of unemployment insurance. Most of these included a combination of direct job search assistance, instructions on how to search for a job, and verification of job search. These experiments, to a large extent discussed in Meyer (1995), included Nevada (1977, 1988), Charleston (1983), Texas (1984), New Jersey (1986), and Washington State (1986). Another set of experiments during same period, assessed only the effect of verification of job search requirements. Ashenfelter, Ashmore, and Deschenes (2005) discuss experiments in Connecticut, Massachusetts, Tennessee, and Virginia.²³

As summarized by Meyer (1995), a core finding of these studies is that JSA reduces unemployment insurance (UI) receipt, at least in the short run. The effects are small, but cost effective from the point of view of the UI agency. The effects on earnings tend to be imprecise, consistent with the possibility that the program

²² There also have been evaluations of JSA services explicitly directed at low-income youth, but most such RCTs that we found were relatively small. The evidence on this subject quoted most frequently is related to the job search component provided in the JTPA and Jobs Corps programs.

²³ Other such experiments include Minnesota (1988), Maryland , (1994), and Washington D.C./Florida (1995-1996), see Greenberg and Schroder (2004).

impacts derive from workers who leave the UI system without finding jobs. Little is known about which components of JSA matter. Experiments in Nevada and Minnesota suggest that intensive JSA has much stronger effects than do more limited treatments. There is mixed evidence as to whether the verification requirement alone matters: The experiments discussed in Ashenfelter et al. (2004) indicate no effects, while a Maryland study summarized in Klepinger, Johnson, and Joesch (2002) did. This question is a key aspect of ongoing evaluations of the Reemployment and Eligibility Assessment system, discussed below.

Since this early wave of UI experiments, the component of the UI system offering job search assistance and training has been repeatedly reformed, with several evaluations along the way. The Worker Profiling and Reemployment Services (WPRS) program was instituted in 1993. Under the WPRS states are required to profile their UI claimants in order to identify those most likely to exhaust UI benefits and refer them to employment-related services.²⁴ This program was evaluated via a natural experiment in Kentucky beginning in 1994 (Black, Smith, Berger, and Noel 2003, Black, Galdo, and Smith 2007). The findings from the WPRS study suggest that receiving a letter asking individuals to come into the office for JSA services alone reduces UI receipt and raises earnings. An important open

²⁴ The services include (1) an orientation session to explain what reemployment services are available; (2) an assessment of the claimant's specific needs; and (3) development of an individual plan for services based on the assessment. Claimants referred to reemployment services must participate in them as a condition of continuing eligibility. Allowable services include job search assistance and job placement services, such as counseling, testing, and providing occupational and labor market information; job search workshops; job clubs and referrals to employers; and other similar services

question is whether this influential finding is replicated in a true RCT and in less favorable labor market conditions.

The Workforce Investment Act (WIA) of 1998 combined most job placement services and training services provided under the auspices of the federal government under one roof, the so-called one-stop centers (e.g., Jacobson 2009). These centers, renamed America's Jobs Centers in 2012, provide both "core" employment services (e.g., job search assistance) and "intensive" WIA services (e.g., career counseling and training) to the three core constituencies – unemployed worker, welfare recipients, and hard-to-employ young workers.

As the structure of service provision has evolved, additional RCTs have evaluated the system's effectiveness at placing workers. For example, in 2005 the Department of Labor's Employment and Training Administration launched a program called Reemployment and Eligibility Assessment (REA), mandatory in-person visits aimed at speeding the reconnection of UI claimants to the workforce.²⁵ The REA meeting includes an eligibility review, provision of labor market information, development of a reemployment plan and referral to more specific reemployment services. The first wave of randomized evaluation of the effectiveness of the REA counseling process took place in nine states beginning in 2005; a second wave of evaluations took place in four states in 2009. In both cases, the evaluations found that the REA requirement and services reduce UI benefit

²⁵ The REA program was instituted to counteract the trend towards processing of UI claims by telephone and the internet. The concern was that the net effect of these changes was to reduce in-person contact and hence the opportunity to monitor job search activity and orient UI claimants to services available to speed their reemployment (e.g., O'Leary 2006)

receipt (Benus et al. 2008, Poe-Yamagata et al. 2011). Earnings outcomes were studied in only one state (Florida), and were positive. An ongoing REA evaluation examines the difference in the effect of enforcing the interview requirement alone relative to the combined effect of the interview plus services (Klerman et al. 2013). A simultaneous evaluation begun in 2011, the WIA Adult and Dislocated Worker Programs Gold Standard Evaluation²⁶ complements the evaluations of REA, WPRS and earlier JSA programs by focusing on the effectiveness of WIA's intensive *and* training services geared to unemployed adults not covered by the earlier evaluations.

Summarizing the wide range of studies of JSA indicates important heterogeneity of effects by the population targeted. For welfare recipients, a difficulty in assessing the effect of JSA is that many experiments tested JSA in conjunction with other programs. Those studies that focus mainly on the effects of JSA, such as the randomized evaluations of WIN, NEWWS or GAIN, often find positive effects on employment and earnings and negative effects on welfare receipt (but mixed effects at best on total income). These effects tend to be short-run lived, and less is known about the longer-term outcomes. There is also little known about the potentially important role played by context, such as local labor market conditions.

In studies of JSA for UI recipients, a common result is a precisely estimated but rather small effect – e.g., a reduction of about one week of UI benefits, with no

²⁶ See <http://www.mathematica-mpr.com/our-publications-and-findings/projects/wia-gold-standard-evaluation>.

corresponding positive effect on earnings –unless the services provided are very intensive. The frontier in this area is assessing to what extent these effects arise from the threat of enforcement of service requirements spelled out by law, basic JSA themselves, or more intensive services.

d. Practical Aspects of Implementing Social Experiments

Clearly, the implementation of large-scale social experiments is complex and faces a range of practical hurdles that can affect the quality of the results. Sections II.c and IV of this paper focus on a number of design issues that can limit the ability of even an ideal experiment to provide answers to the questions of interest.

Beyond these conceptual design issues, there are some common challenges and practical considerations that have come up over and over in the conduct of social experiments in the labor market. These play important roles in influencing the topics and questions that are studied via social experiments and in informing the study designs.

One set of challenges derives from the fact that, as noted above, one of the defining characteristics of social experiments is that they intend to examine programs that are already in place or might be put in place in essentially the same form that was used in the experiment. For this purpose, the experimental samples and hence the sampling frame need to be representative of the population that the program serves. This is a challenge in the case of many labor market programs, in

part because the sampling frame is often available only to program operators or the government, and may be difficult to access due to formal approval processes.

Once the sampling frame is obtained, it is necessary to randomly assign some members of the sample to the program of interest and others to a control condition, which might be exclusion from the program or an alternative program design. This, too, can be difficult when the program is already in place. For example, if the program in question exists within an ecosystem of other programs, services, and service providers, it may be hard to exclude participants from the program or, if this is done, to avoid also excluding them from other programs that are administratively integrated. For example, excluding a participant from job search assistance offered under the Workforce Investment Act (WIA) might also in practice exclude him or her from job training and other programs, as the same offices that provide job search assistance also do screening and referrals for other services. While some of these problems might be reduced by studying programs *not* already in place, as in the case of the Negative Income Tax experiments or the National Supported Work Demonstration, this can be quite costly, as the sorts of programs typically studied involve substantial program costs – commonly in the thousands of dollars per participant.

A second group of challenges has to do with the difficulty of enforcing compliance with randomization after it is conducted. Again, the use of actual programs tested in real-world settings limits the options. A common challenge in early experiments was that service delivery was delegated to individual case-workers or sites that were both widely dispersed and not closely involved with the

experimental design. This raises the possibility that caseworkers may deviate from random assignment, for example ensuring that a potential participant viewed as especially needy is not assigned to the control group. For example, a key concern in the National Job Corps Study was to ensure that local program operators properly implemented the randomization. Modern practice centralizes the random assignment process, carefully tracking participants' initial assignments to ensure that participants assigned to undesirable treatment conditions do not re-enter the randomization to obtain a better assignment.²⁷

A third set of challenges has to do with the measurement of participant outcomes. Once again, this challenge derives, in large part, from the use of real-world populations as experimental subjects and from the large and heterogeneous subject pools common in social experiments. These make it more expensive to ensure high response rates than in smaller and more targeted field experiments.

In many cases this challenge can be addressed by using administrative data to measure some outcomes. Administrative records may come either from the program under study – for example, unemployment insurance payment records for studies of job search incentives for unemployment insurance recipients – or from other records from other government programs (e.g., tax records). While this can resolve the attrition problem at low cost, it is often contingent on government cooperation or approval. Such cooperation is more likely in large-scale social experimental evaluations of existing programs than in other types of studies.

²⁷ For a discussion of approaches to address this problem, including related software, see, e.g., Crepon et al (2013).

Administrative data can also limit the set of impacts that can be studied, potentially creating important ambiguities in the interpretation of estimated treatment effects. In the unemployment insurance case, for example, it is not clear whether a negative effect of increased job search enforcement on unemployment benefit payments indicates that people are finding jobs faster, or just that many people are leaving the program before finding jobs as a way of avoiding onerous enforcement procedures.

IV. Going Beyond Treatment-Control Comparisons to Resolve Additional Design Issues

Whether one is interested in structural parameters or program evaluation, many questions of policy or scientific interest in labor and public economics require going beyond the basic RCT design described in Section II.a. We discussed a number of these questions in Section II.c. Here, we discuss ways to extend the basic RCT design to provide answers to these questions.

We organize our discussion around the major potential design issues we mentioned in Section II.c. For each, we discuss proposed solutions and, where relevant, point out potential extensions and limitations. We begin by discussing studies that address aspects relating to *internal validity*, including SUTVA violations (e.g., potential general equilibrium effects) and endogenously observed outcomes. We then discuss studies that address *external validity* concerns, including site and sub-group effects; effects on subpopulations other than experimental compliers; hidden or multiple treatments; mechanisms for treatment effects; and studies of optimal or simply alternative policies.

In some cases, the identified issues can be addressed *ex post* (after an experiment is complete), generally by imposing additional structure. In many of these examples the additional structure imposed is justified by appeal to theoretical considerations and is just sufficient to extend the RCT to address a specific question and the design issue it raises. In that sense, the studies can be viewed as an effort to bridge pure experimental or quasi-experimental approaches, credibly identifying a limited number of (potentially composite) causal parameters, with more traditional structural estimation that obtains a fuller characterization of the economic problem via the imposition of substantial additional assumptions. In the ideal case, they maintain the best of both worlds, though they also share some of the limitations of each.

Another possibility is to build the structural questions of interest into the design of the experiment *ex ante*. This can provide credible identification with even fewer structural assumptions than are required for after-the-fact analyses, though can sometimes require a quite complex – and potentially difficult to administer – experimental design. There are fewer existing examples of this, but we discuss them where appropriate.

We discuss each of the design issues identified earlier in turn. Our discussion is meant to highlight the different approaches, as well as to clarify the scope, potential, and difficulties that arise when extending inference from standard RCTs to a broader range of questions.

a. Spillover effects and SUTVA

Social experiments in labor economics typically occur in the context of the local or regional labor market. If the number of workers participating in the program is large relative to the relevant segment of the labor market, the program could have an effect on the labor market outcomes of the control group. This would be a violation of SUTVA – the difference in outcomes between treated and control individuals would differ from the overall effect of the program on the entire population relative to not implementing the program, which is often the effect of primary interest.

Many social experiments in the United States have not raised serious spillover issues, as the treated populations have been small relative to the local labor market. However, this may not be true for large experiments, such as the National Jobs Corps Study. Welfare experiments may also create spillover effects if labor markets for former welfare recipients are sufficiently segmented.

A related issue is that comprehensive program evaluations in many cases *should* include spillover effects that are not captured by small-scale pilot studies. If the pilot programs are eventually scaled to broader populations of low-income workers – which has happened, among others, in the case of welfare reform, of training provided through WIA, or job search assistance services provided by WPRS or REA – then the potential extent of spillover effects would nevertheless matter, since any spillover effect would have to be included in a welfare assessment of the program. This would create systematic differences between the outcomes of the pilot study and the program effects of interest.

i. Addressing the issue ex post

Despite its potential prevalence in social experiments in the labor market, relatively few studies have dealt directly with the issue of spillovers or other failures of SUTVA. A handful of studies have tried to estimate spillover effects directly using inter-regional comparisons (e.g., Blundell, Dias, Meghir, and Van Reenen 2004; Ferracci, Jolivet, and van den Berg 2010; Gautier, Muller, Rosholm, Svarer, and van der Klaauw 2012). There are roughly two approaches, neither of which is able to fully identify the spillover effect. One approach is to compare control group outcomes to those of observably similar individuals in areas where no one is treated. Of course, there may be other explanations for differences seen in this observational comparison. Another approach is to compare the effect of treatment across sites with different treatment intensity or labor market conditions. This is again typically an observational comparison, as in most cases neither the treatment site nor the size of the treatment group (and hence the amount of potential spillover) is randomly assigned. For example, Hotz (1992) discusses the non-random selection of sites for the JTPA evaluations. Alcott (2015) studies the sources of observed bias from site-selection in a large electricity conservation experiment. A recent paper by Crepon, Duflo, Gurgand, Rathelot, and Zamora (2013; see also Baird et al. 2015), discussed further below, resolves this problem in the context of a job search assistance program by randomly assigning both the treatment and the number of workers treated.

Absent such a multi-stage experimental design, relatively few options are available to researchers to assess the degree of the actual or potential spillover

effects present in the context of their evaluation. An area of research where spillover effects have received substantial recent attention is the analysis of the employment and welfare impacts of extensions in unemployment insurance benefits. Here, spillover effects arise because treated and untreated individuals compete for the same positions; the degree of the spillover effect therefore depends on the job creation response to the treated group's labor supply change. To assess the potential degree of spillovers, one can in principle use estimates of the matching function to adjust micro-econometric estimates of the effect of policy-induced changes in unemployment insurance durations on unemployment duration or exit hazards for the presence of crowding.²⁸ Such ad-hoc simulations are partial-equilibrium in nature, and could be interpreted as a short-run effect, when vacancies have not yet adjusted. Landais, Michaillat, and Saez (2015) specify a general equilibrium model of the labor market that incorporates both crowding and vacancy responses. In a standard, competitive search-matching model, the vacancy response to changes in labor supply is sufficiently strong to offset the crowding effect completely.

In the spirit of using random variation in the treatment across localities to assess the presence of spillover effects, a couple of recent papers have tried to exploit region-specific changes in policy-induced UI variation in the U.S. to assess the full effect of the policy on the entire labor market (Hagedorn, Karahan, Manovskii, and Mitman 2015, Hagedorn, Manovskii, and Mitman 2015). Since UI

²⁸ One added difficulty in the case of UI is that in most cases in the U.S. the policy-induced changes in the level or duration of UI benefits are a function of labor market conditions – making it crucial to properly control for the direct effect of local labor market conditions.

variations usually depend on economic conditions at the state level, these studies use border communities unaffected by the policy change as counterfactuals.²⁹ A concern with this approach is that the presence of spatial spillovers between adjacent or related labor market areas would again constitute a failure of SUTVA.³⁰

Another source of SUTVA failures are interactions between treatment and control participants. Such ‘dilution’ effects can lead to an underestimation of the treatment effect. If possible, a typical approach to circumvent such interactions is to raise the level of randomization (say, from a sub-group within a site to a whole site). This approach can help to avoid interactions between individuals in the treatment and control groups. It does not resolve potential interactions between treated participants. This may be part of the mechanism of the treatment; it may also be a potentially unintended source of variation in treatment intensity that we discuss under site effects. In either case, when designing an evaluation, it would be valuable to consider ways of keeping track of social interactions, perhaps by asking about friends in a baseline survey, or monitoring (or manipulating) the use of certain kinds of social media. Another valuable target for data collection is factors relating to how treatment was obtained or take up was decided. Such information may be used to stratify the analysis by the predicted degree of SUTVA violations or at least assess the potential for significant departures from SUTVA.

²⁹ A key practical difficulty there is that measures of unemployment rates at the sub-state level is often very noisy. Estimates using administrative employment data based on the universe of private employees show little sign of spillover effects (Johnston and Mas 2015).

³⁰ Cerqua and Pellegrini (2014) develop alternative estimates to the TOT that take into account the degree of spatial spillover effects. The Hagedorn et al. papers have been quite controversial; see, for example, responses from Chodorow-Reich and Karabarbounis (2016) and Coglianesi (2015)

ii. *Addressing the issue ex ante through the design of the experiment*

In some circumstances it may be possible to avoid, or study, spillover effects by appropriately structuring a randomized experiment. For example, in the spirit of the non-experimental studies cited above, treatment and control groups could be chosen to be sufficiently distant to avoid spillover effects. Alternatively, the treatment group could be chosen to be sufficiently small that spillover effects are unlikely to be a problem. If the spillover effects themselves are of direct interest, the experimental manipulation could be combined with pre-existing variation in the strength of potential spillover effects (e.g., across submarkets), if available. The risk of such ad hoc or hybrid approaches is to potentially lose comparability of the control group, or to confound spillover with other variation in treatment effects.

A preferable approach if spillover effects are potentially present is to manipulate both the treatment and the size of the treatment group (and hence the amount of spillover) experimentally. Baird et al. (2015) develop this strategy formally. Crepon, Duflo, Gurgand, Rathelot, and Zamora (2013) implement it in the context of a public program assisting unemployed workers in their search for a job in France. The researchers manipulate both who gets assigned into the job search assistance program *within* a region (the classic experimental design), as well as randomly vary *between* regions the share of individuals assigned to the treatment group. The manipulation of both regional treatment share and individual treatment status allows separate experimental identification of the effect of the program holding the spillover effect constant and the combined program and spillover effects at various treatment intensities. The latter parameters are ultimately relevant for a

cost-benefit or welfare analysis of the program and for extrapolation to alternative policy settings.

Similar strategies are available for other SUTVA failures, arising for example if some individuals in the control group get accidentally treated, or if treatment compliance depends on the take up rate among peers. In some cases, one may choose the experimental setting to try to minimize SUTVA problems. For example, one can devise strategies to limit the potential for non-compliance (e.g., in case of web-based information treatments, access could be based on hardware address rather than passwords).

Another potentially interesting strategy is to make the degree and structure of SUTVA violations part of the analysis, as in the discussion of spillovers above. This may provide insights into the “black box” of how a program might work in a real life setting and hence enhance external validity.³¹ For example, one could experimentally vary the number of treated units in a reference group or network (e.g., classrooms, friends, etc.), examining interactions among individual treatment status, group treatment share, and perhaps also predetermined factors (such as the tightness of the group) that determine the degree of departure from SUTVA. Depending on the context, it may be possible to more explicitly manipulate interactions between individuals by introducing an additional treatment to the experimental design – for example, a forum in which interactions are facilitated.

³¹ Note that there is a parallel here with the issue of treatment compliance and heterogeneous treatment effects. Here, the compliance function is assumed to depend on treatment status of other individuals, and hence experimentally manipulating compliance probabilities is presumably more complex. Yet, as in the standard case of heterogeneous treatment effects, for external validity it is important to trace out the potential compliance-related interactions as fully as possible.

b. Endogenously observed outcomes

In many labor market experiments, key outcomes include measures observed only for individuals who are employed, such as hours worked and wages. Hence, the impact of, say, welfare-to-work programs or job training programs can only partially be assessed based on simple RCTs alone. Although many studies report experimental impacts on the endogenously observed outcomes, these are understood to suffer from serious selection problems. In the same way, non-random attrition in follow-up data collection can bias the results of nearly any evaluation.

To illustrate, consider a program aimed at unemployed workers that includes skill development and job search assistance modules. We are interested in whether the program raises the probability that a participant is employed one year after participation and whether it makes them more productive when employed. For simplicity, we assume that participation is randomly assigned and compliance is perfect.

We have two outcomes here. We denote employment status by $y_i = D_i y_{1i} + (1 - D_i) y_{0i}$. For those who are employed at the follow-up survey, we observe the wage $w_i = D_i w_{1i} + (1 - D_i) w_{0i}$. Treatment effects of the program on the two outcomes are τ^y_i and τ^w_i . (We can imagine that w_{di} is well defined for an individual with $y_{di} = 0$, $d \in \{0, 1\}$, but simply not observed. It can be thought of as the individual's *latent* productivity, that which he/she would be paid if a job were found.)

Estimation of $E[\tau^y_i]$ is straightforward, as discussed above. But the impact on wages is much harder. In general, it is not possible to identify the average treatment effect $E[\tau^w_i]$; the treatment-on-the-treated effect $E[\tau^w_i | D_i = 1]$; or even the average

treatment effect for the subpopulation that would have been employed with or without the program (for whom τ^{w_i} is least problematic), $E[\tau^{w_i} | y_{0i} = y_{1i} = 1]$.

The problem here is that it is impossible to distinguish, within each D_i group, between those workers who would also have worked in the counterfactual and those who would not have. Consider the treatment-control difference in mean observed wages:

$$\begin{aligned}
 E[w_i | y_{1i} = 1, D_i = 1] - E[w_i | y_{0i} = 1, D_i = 0] &= \\
 &= E[w_{0i} + \tau^{w_i} | y_{1i} = 1, D_i = 1] - E[w_{0i} | y_{0i} = 1, D_i = 0] \\
 &= E[\tau^{w_i} | y_{1i} = 1, D_i = 1] + (E[w_{0i} | y_{1i} = 1, D_i = 1] - E[w_{0i} | y_{0i} = 1, D_i = 0]) \\
 &= E[\tau^{w_i} | y_{1i} = 1, D_i = 1] + \\
 &\quad + (E[w_{0i} | y_{0i} = 1, y_{1i} = 1, D_i = 1] - E[w_{0i} | y_{0i} = 1, y_{1i} = 1, D_i = 0]) \\
 &\quad + (E[w_{0i} | y_{1i} = 1, D_i = 1] - E[w_{0i} | y_{0i} = 1, y_{1i} = 1, D_i = 1]) \\
 &\quad - (E[w_{0i} | y_{0i} = 1, D_i = 0] - E[w_{0i} | y_{0i} = 1, y_{1i} = 1, D_i = 0]).
 \end{aligned}$$

The first term here is the average treatment effect in the subpopulation that works under treatment. It may not equal the overall average treatment effect, but insofar as the potential wages of those who do not work are not relevant to social welfare, it is arguably the parameter of interest. The second term solely reflects selection into treatment, and is zero under random assignment. But the third and fourth terms have to do with selection into employment, not selection into treatment. Random assignment does not ensure that they are zero, and the treatment-control contrast

among workers may therefore be badly biased relative to the impact on wages for any fixed group of workers.³²

One fallback approach is to examine only the program's effect on the share of participants earning high wages, treating low-wage-workers and non-workers the same. This effect can be estimated without bias. Another fallback is to include the non-employed in the wage analysis, with wages set to zero. This in some cases is the impact of interest in any case, and is correctly identified by the experiment.

However, it is quite misleading if interpreted as the magnitude of the effect on productivity, either for the full population or for the subgroup that would have been employed with or without treatment. Without an ability to measure *counterfactual* employment status at the individual level, the latter effects are not identified.

i. Addressing the issue ex post

Non-random attrition in particular has been a long-standing concern in the experimental literature in labor economics (e.g., Hausman and Wise 1979). A classic experimental design would be deemed successful if attrition is low and balanced in terms of magnitude and observable characteristics between the treatment and control groups. If this is the case, reweighting the samples may still recover the

³² Consider a training and job-search assistance program. Suppose 60% of workers will be always low productivity ($w_{1i} = w_{0i} = w^L$), 20% will be always high productivity ($w_{1i} = w_{0i} = w^H$), and 20% will become high productivity if exposed to the training sequence ($w_{0i} = w^L$, $w_{1i} = w^H$). All of the second and third groups will find jobs, with or without search assistance ($y_{0i} = y_{1i} = 1$), but those in the first group of low-skill, impossible-to-train workers will find work if and only if they receive search assistance ($y_{0i} = 0$, $y_{1i} = 1$). In this setting, the program's average treatment effect on employment is 0.6; the average effect on latent productivity is $0.2*(w^H - w^L)$; and the average effect on wages of those who would work with or without the program is $0.5*(w^H - w^L)$. The estimated treatment effect on wages conditional on employment is $-0.1*(w^H - w^L) < 0$. Selection has led to a perverse estimate here: The training program has a positive effect on 20% of participants and a negative effect for no one, but the experiment appears to indicate that it reduces earnings.

effect of the TOT or LATE among the original set of compliers (e.g., Ham and Li 2011). Yet, there are relatively few explicit attempts in the literature to address selection bias in other contexts.

A large literature in labor economics has dealt with sample selection problems, especially in the analysis of wages and hours in the context of the classic human capital and labor supply models. Largely based on that literature, here we will review several approaches to deal with selection bias: the use of control functions to address selection; estimation of percentile effects instead of mean impacts; use of additional data to control for selection; construction of bounds based on selection probabilities; and construction of bounds using theory.

Parametric selection corrections

The 'classic' approach to control for selection bias in estimating the effects of treatment effects on wages or hours worked is based on control functions. Labor supply theory, along with parametric assumptions, is used to derive an explicit expression for the selection bias in terms of the participation probability, which under monotonicity determines the amount of sample selection. This is then accounted for directly in the outcome equation (e.g., Gronau 1973, Heckman 1979).

Early on it was recognized that absent experimental variation in participation (e.g., an exogenous instrument affecting only participation and not the outcome equation), identification is only based on functional form assumptions, and results can be quite misleading if these assumptions are even slightly incorrect. By contrast, a substantial literature has shown that once an instrument for participation is available, treatment effects in the outcome equation can be

identified under quite general functional form and distributional assumptions (e.g., Newey, Powell, and Walker 1990). For example, Ahn and Powell (1993) show that under assumptions of a single, strictly monotonic index for selection, variation in the probability of participation independent from the variables in the outcome equation suffices to control for selection. The difficulty is, of course, that often such independent source of variation is not available.

Card and Hyslop (2005) consider a special case in which an RCT does generate exogenous variation in participation: An employment subsidy program. They show that if the program only has positive effects on labor supply and does not affect the wages for those who would have worked without it, then the experimental effect on the hourly wage can be consistently estimated by the ratio of the treatment effect on total earnings divided by the treatment on total hours worked.

Card and Hyslop's assumptions are inappropriate for any program designed to affect wages and not just participation. Below we discuss how the experimental design itself may be modified to obtain exogenous variation in participation, even in programs with effects on multiple margins.

Non- and semi-parametric selection corrections

Absent an instrument for participation, in the presence of selection the treatment effect on mean wages is not identified. However, several studies have exploited the fact that under certain assumptions quantile-treatment effects (QTEs) may be consistently estimated even in presence of selection. A QTE for the q -th quantile is defined as the difference in the q -th quantile of the outcome distribution

in the treatment and control groups, respectively.³³ It is not necessary to observe each individual's outcome to compute the q -th quantile; it suffices to know that someone is above or below that quantile. Thus, if one can assume that all those who are not employed have potential wages in the bottom q percent of the distribution, one can estimate the treatment effect on the q th quantile of potential wages by merely assigning all non-workers the minimum observed value (e.g., Powell 1984, Buchinsky 1994). Hence, under this assumption all quantiles above the value of the rate of nonemployment of the respective group can be identified. The lower value of nonemployment of the treatment and control group determines which QTE can be identified.

A variant of this approach is to examine the simple treatment-control difference in the probability of being observed in employment with a wage greater than some relatively high threshold w^* . For many program evaluations, understanding the impact on this outcome may be sufficient – it may not matter greatly whether the impact derives from moving some people from non-employment into high-wage employment or from simply lifting those who would have worked anyway into higher-wage jobs. And even when the latter component is the one of interest, this would be identified so long as those pulled into employment by the treatment have wages that are uniformly below w^* .

³³ For any random variable Y having cumulative density function $F(y) = \Pr[Y < y]$, the q th quantile of F is defined as the smallest value, such that $F(y_q) = q$. If we consider two distributions F_0 and F_1 , then $\text{QTE}(q) = y_q(1) - y_q(0)$, where $y_q(g)$ is the q th quantile of distribution F_g .

It is not clear, however, that the required assumption holds – as pointed out by Altonji and Blank (1999), among others, at any given time, some high-wage individuals may be nonemployed. Moreover, this strategy is only useful in so far as differences in quantiles of the outcome are deemed sufficient for evaluating the effect of the program.

Another approach uses reservation wages to measure selection into the subsample of observed wages. This works because – if correctly measured – the reservation wage captures the lowest wage for which an individual is willing to work. Hence, the reservation wage provides the censoring point for an individual's wage-offer distribution, allowing one to make inferences about potential wages for those individuals not working in the treatment and control group. Johnson, Kitamura, and Neal (2000) use the minimum of all observed wages for an individual in longitudinal data to bound the reservation wage, under the assumption that it is stable over time. Grogger (2005) uses directly reported reservation wage information from a randomized evaluation of Florida's Family Transition Program, a welfare-to-work program with emphasis on work incentives and time limits. With this information, he estimates the treatment effect of the program on wages using a bivariate, censored regression model that allows for classical measurement error in both observed wages and reservation wages. Once Grogger (2005) controls for selection, he finds the program had statistically significantly positive effects on wages.

Addressing the selection problem using direct measures of reservation wages makes intuitive use of the reservation wage concept. Moreover, often

information on reservation wages is already being collected in the context of programs providing job search assistance, or if not they are at least in principle relatively easy to elicit if the experimental design includes a survey component. However, recent research suggests that in practice reported reservation wages appear to only partly reflect the properties of the theoretical concept (e.g. Krueger and Mueller 2016), casting some doubt on the robustness of this approach. In particular, Krueger and Mueller report that a substantial number of workers accept (reject) jobs offering wages below (above) their reservation wage, implying that care should be taken in using reservation wages of the nonemployed to make inferences about unobserved wage offers.

Yet another approach is to attempt to derive bounds for the treatment effect under conditions more general than the monotonicity assumption inherent in the Ahn and Powell (1993) and similar estimators. This allows researchers to investigate how severe the bias from selection could possibly be and what can be learned under general assumptions rather than to try and to obtain a point estimate under more restrictive assumptions.

One bounding approach is proposed by Horowitz and Manski (2000). This strategy asks how much the estimated treatment effect would be inflated if all missing treatment observations were assumed to have the highest possible outcomes and all missing control observations the lowest; then it asks how much it would be depressed if the opposite assumptions were made. Unfortunately, these bounds are typically not very tight, particularly when the outcome variable's support is potentially unbounded as for example in the case of wages.

Lee (2009) proposes a strategy for obtaining tighter bounds, via stronger assumptions: He assumes that anyone not employed in the control group would also have been non-employed had they been in the treatment group, so that selection bias arises solely from participants in the treatment group who are employed but would not have been had they been assigned to be controls.³⁴ He can then bound the treatment effect by making extreme assumptions about this latter group. Denote the excess fraction employed in treatment group by p . The upper (lower) bound is constructed by removing the lowest (highest) fraction p observations from the treated subsample and recomputing the mean outcome for the treatment group – effectively making the worst-case assumption that selection was fully responsible for the entire upper or lower tail of values. Lee (2009) shows that the resulting bounds are sharp and provides formulas for the standard errors. In the case of Job Corps, the procedure results in informative bounds suggesting positive wage effects from training – albeit a zero effect is contained in the confidence interval.

Lee's (2009) approach based on trimming requires relatively weak assumptions. It presumes only that selection is monotonic in the treatment – that treatment either only increases, or only reduces, selection into employment. Monotonicity is implied by standard empirical binary choice models typically used to model participation choices (e.g., Vytlačil 2002), and hence bounds based on trimming are applicable to a wide range of problems, including selective employment, survey non-responses, or sample attrition.

³⁴ The role of treatment and control groups are reversed if the treatment reduces employment.

If one is willing to impose further structure from theory, one may obtain tighter bounds more specific to a particular problem. This is especially useful if the theory has explicit predictions about how the endogenous outcome responds to incentives.³⁵ This is pursued by Kline and Tartari (2016), who analyze the randomized evaluation of Connecticut's Jobs First welfare-to-work program. While previous analyses had found only small responses in hours (the intensive margin), absent an instrument for participation (the extensive margin) sample selection makes such estimates hard to interpret. Kline and Tartari (2016) use revealed preference arguments in the context of a canonical but non-parametric static labor supply model to describe which observed responses to the treatment at the intensive and extensive margin are consistent with the theory. Given the nature of the program studied, the result is a mapping of discrete counterfactual outcomes (including non-participation as well as participation at different intensities) under treatment and non-treatment, with restrictions on the allowable counterfactuals. The question then is how likely certain transitions are, and in particular whether changes at the intensive and extensive margin occur with positive probabilities. Since Kline and Tartari can only observe the marginal distribution across states for the treatment and control groups, they cannot point-identify the transition probabilities. Instead, they construct bounds for transition probabilities among the entire (discretized) distribution of states, including the probability of changes in the

³⁵ This may be more easily done for hours, which is typically assumed to be a choice variable, than for wages. Yet, to some degree wage may be a choice variable as well, for example if jobs offer wage and effort combinations among which workers choose. This is the approach taken in some modern public finance, which often substitutes hours worked with taxable earnings as the choice variable in analyses of intensive-margin labor supply.

intensive margin due to the treatment. Their approach also allows them to test the restrictions from the model.

This approach is useful, since it allows Kline and Tartari (2016) to learn about intensive margin responses to the Jobs First program in the presence of selection. Their results could also be used to think about the likelihood of intensive margin responses for similar programs in similar populations. Alternatively, the estimated bounds from the matrix of transition probabilities could be used, along with the marginal distribution of labor supply under an existing program (AFDC, the program of the control group), to construct bounds for the intensive and extensive labor supply responses that could arise if Jobs First was implemented at another site. A potential issue is that the procedure is complex and the analysis is specific to the Jobs First program. Hence, while the general approach may be applicable to a range of problems, this would require careful specification of the decision problem, of the restrictions imposed by revealed preference theory, and of counterfactuals for each case. Nevertheless, since many social experiments are concerned with welfare and other programs that provide explicit variation in employment incentives and hence useful information on the likelihood of counterfactual outcomes, it is useful to consider the role that theory can play in providing bounds on treatment effects on endogenous outcomes.³⁶

ii. *Addressing the issue ex ante through the design of the experiment*

³⁶ Similar approaches have been pursued in Blundell, Bozio, and Laroque (2011).

The endogenous outcome problem is often easily anticipated when designing an experiment, as it arises whenever outcomes like wages or hours are of interest and non-employment is a realistic possibility. There are various ways to adjust the experimental design to facilitate analysis of potential sample selection bias. For example, suppose in the case of the effect of a training program on wages the researcher believes that there are exogenous factors determining a worker's labor supply decision. If these factors can be measured ex ante, the randomization could be stratified by the likelihood of employment as predicted by the exogenous instruments. Stratification would ensure sufficient sample sizes in each exogenous labor supply tier. (If only available ex post, say, in a follow-up survey, even absent stratification such variables can be still used as instruments for participation if sample sizes are sufficiently large.)

However, as it is usually difficult to come by good instrumental variables, the real power of a well-designed RCT would be to manipulate sample selection directly. In the training example, this would entail adding a second source of randomization that explicitly modifies the incentive to work (or the likelihood of finding a job) but does not otherwise affect the endogenous outcome. Whether this is feasible depends on the context. However, sample size considerations need not be a hurdle to adding a second treatment, since with cross-classified treatments the addition of a second treatment has little effect on the power for analyzing the effects of the first in isolation. This approach is particularly useful if one is interested in external validity, since the two-dimensional experimental variation may allow one to trace out the

treatment effect of training for sub-populations with different employment probabilities.

In the case of non-random attrition, a version of this approach would be to randomly select a group of participants to follow up more intensively, perhaps stratified within groups with different ex-ante attrition probabilities. The contrast between mean outcomes in this subgroup and for other participants (again, perhaps within strata) identifies the selectivity of attrition, and can be used to adjust the full-sample estimated treatment effects. This is the approach pursued in the follow-up waves of the Moving To Opportunity experiment (e.g., Kling, Liebman, and Katz 2007). Another solution worth pursuing is to obtain administrative data for the universe of initial participants, including those who have failed to respond to follow-up surveys. Although these data can also be selected – they typically do not include earnings from informal jobs – the selection is different from that created by survey attrition, so the combination of sources can be valuable (though sometimes confusing, as in the Job Corps evaluation discussed above). Since merges to administrative data can usually only be conducted only with identifying information from the survey and permission from participants, it is a good idea to factor the need for additional data into the initial research design.

c. Site and group effects

In many cases an essential problem is to identify the subpopulations that benefit most from a program, so as to target them for treatment. However, there are often many possible subgroups to examine. When many comparisons are estimated,

the chance of a false discovery – a treatment-control contrast that is statistically significant, even though the true treatment effect is zero – rises toward one.

Avoiding incorrect inferences in such a setting requires care.

A version of the subgroup effects problem is to identify variation in treatment effects across program locations or sites. Such variation might arise from observed local characteristics – e.g., treatment effects of training or job search experiments may depend on the tightness of the local labor market. Where the relevant characteristics of the labor market are clear *ex ante* and their dimension is limited, this is relatively straightforward. But if the relevant dimensions are not clear or the number of potential contrasts is large, the multiple comparisons problem becomes relevant. Alternatively, there might be unintended variation in treatment intensity or in the fidelity or effectiveness of treatment delivery among treatment sites. Such site effects render the interpretation of the estimated treatment effect of the overall treatment difficult and limit external validity. If they are potentially important, we need estimates of each site's separate effect. This implies that there are as many treatment effects to be estimated as there are sites at which the experiment is implemented.

A conceptual issue in evaluating the success of social experiments with site variation is to decide whether the parameter of interest is the effect of the program in its most successful variants, with strong local partners and appropriate local conditions, or the average effect across a range of local circumstances. When the latter is of interest, the ideal experimental design would involve drawing participants from all sites. But this is often impractical. More commonly, social

experiments have been carried out at one or a few sites. These are often chosen because the local management is willing to participate, or because they are seen as exemplars of the program. This makes it difficult to interpret the experimental results as representative of the program as a whole (see, e.g., Hotz 1992 and Alcott 2015), but may come closer to identifying the program effect under close-to-ideal circumstances.³⁷

i. Addressing the issue ex post

On its face, it is straightforward to estimate heterogeneity of treatment effects along observed dimensions (e.g., race, gender, or past work experience) using data from an already-completed randomized trial: One simply constructs treatment-control contrasts separately for each subgroup. Many authors emphasize the importance of conducting the randomization separately for each subgroup of interest. This is not in principle necessary – unconditional random assignment ensures that assignment is random conditional on predetermined characteristics as well – but can add power for subgroup comparisons, especially in smaller samples.

A more important issue is the potential number of comparisons to be estimated. If enough subgroup estimates are computed, even a program that has no effect on anyone will be likely to show a statistically significant effect for some subgroup. (A similar problem arises when considering effects on multiple

³⁷ A related but distinct problem is the question of ensuring “fidelity of implementation” in an RCT – a close alignment between the program’s intended design and the services that are actually delivered. While this is important for maximizing the statistical power of the experiment and for testing whether the program’s theory of action is correct, it limits the external validity for use in making judgments about the likely overall impact of real-world programs, which may not be implemented with high fidelity.

outcomes.) Researchers have taken a number of approaches to this multiple comparisons problem. One is to specify the subgroups that will be considered, and the hypotheses of interest, before analyzing the data. This can limit the scope for unconscious data mining. It also ensures that the number of comparisons that were considered is known, so that the p-values of simple treatment-control contrasts can be adjusted for the multiplicity of the comparisons being estimated. An appropriate adjustment makes it possible to obtain accurate p-values for the test of whether the program had any effect on any subgroup. But two issues remain: These tests typically have *very* low power. In addition, even when they do reject they are often not able to identify *which* subgroups have non-zero treatment effects. A full discussion of adjustment for multiple comparisons is beyond the scope of this chapter, but Anderson (2008) is a useful reference.

Multiple comparisons approaches can be useful as well for the analysis of treatment effects by site and/or provider. But the questions of interest regarding site effects are not generally whether each site's effect is or is not different from zero, which is what multiple comparisons adjustments are designed to answer, but rather the magnitude and correlates of variation in treatment effects across sites. Moreover, the fact that the site-specific treatment effects can in some sense be seen as draws from a larger distribution opens up new options for analysis that are not available in traditional studies of subgroup treatment effects.

The mid-1990s National Job Corps Study, discussed above, illustrates some of the issues involved.³⁸ As mentioned previously, the random-assignment study indicated that the program has a positive average effect on earnings four years after participation, of a magnitude roughly comparable to the return to a full year of education (Schochet, Burghardt, and McConnell 2008). (At the time of the evaluation, the average participant was enrolled for about eight months.)

But like other job training programs, the specific “treatment” provided to Job Corps participants varies substantially across individuals, according to perceived needs. Moreover, Job Corps services are delivered at 110 mostly residential centers, the majority of which are operated by private contractors. Some providers may be better at delivering an effective program (or at guiding participants to the types of services that they need) than are others. The center-specific treatment effects are thus of great interest.

The Department of Labor (DOL) has long used a performance measurement system to track performance of the different centers and inform decisions about contract renewal. Performance measures are non-experimental, and include statistics like the GED attainment rate or average full-time employment rate of program participants at each center. But it is not clear that these performance indicators successfully distinguish center impacts from differences in the populations served by the various centers.

³⁸ Other studies that examine similar questions are Bloom, Hill, and Riccio (2005) and Barnow (2000). See also our discussion of treatment spillovers above.

Schochet and Burghardt (2008; hereafter “SB”) attempt to use the random-assignment Job Corps Study to validate DOL’s performance indicators (see also Barnow, 2000, who carries out a similar exercise for JTPA). In principle, estimation of site-level causal effects using the experiment is straightforward: One simply compares mean outcomes of the treatment and control groups at each site, relying on the overall random assignment to ensure balance of each site-level comparison. But a few challenges arise.

First, in the Job Corps Study randomization took place before applicants were assigned to centers. Thus, treated individuals are associated with centers, but control individuals are not. SB address this by using intake counselors’ assessments of the center that the applicant would most likely attend, collected prior to randomization. To ensure that treatment and control individuals are treated comparably, they use this prediction for both groups, even when it differs from the actual treatment assignment. Differences occurred for only 7 percent of treatment group enrollees, largely because participants tend to enroll in the closest center or in one that offers a particular vocational program.

Second, even a large RCT sample – the Job Corps Study included over 15,000 participants – can have very small sample sizes at the individual site level. Rather than estimate center-specific treatment effects, SB divide centers into three groups based on their non-experimental performance measures and estimate mean treatment effects for each group. Interestingly, they find that mean program impacts do not differ significantly across groups, suggesting that the performance measurement system is not successfully identifying variation in centers’ causal

impacts. A related exercise is carried out by Bloom, Hill, and Riccio (2005), who first estimate statistically significant variation in treatment effects across 59 local offices that participated in three welfare-to-work experiments, then use a multi-level model to estimate the relationship between office characteristics – mostly having to do with the way that the treatment was implemented in each site, though they also include the local unemployment rate – and office-level treatment effects. In contrast to the Job Corps study, they do find significant associations of the treatment effect with both their implementation measures and the local unemployment rate.

Bloom, Hill, and Riccio's (2005) interest is in identifying which program features are most effective. It is important to emphasize, however, that the association between site-level characteristics X_j and the site-specific treatment effect τ_j is observational, not experimental, and does not bear a strong causal interpretation. It is quite possible that what appears, for example, to be a strong association between the emphasis that sites place on quick job placement and the site-level treatment effect instead reflects a non-random distribution of this emphasis across sites that vary in other important ways.

Like the Job Corps study, Bloom et al. (2005) do not investigate variation in site impacts conditional on X_j . In many settings, that variation might be of substantial interest. One might like, for example, to estimate effects of individual sites, or to ask which of a number of available performance measures do the best job of predicting experimental impacts. The latter question is a natural one to ask regarding the Job Corps Study, but to our knowledge it has not been pursued with

experimental data (though see Barnes et al. 2014 for a related investigation using non-experimental data).

Much work on the estimation of site effects themselves comes out of efforts to measure hospital, school, or teacher performance (see, e.g., Jackson, Rockoff, and Staiger 2014 and Rothstein 2010). These studies are program evaluations, treating each site or teacher as a distinct “program,” but cannot rely on random assignment to identify program effects. As in the Job Corps Study, there are many sites but samples are frequently small at the site level, so – even if selection biases are set aside – site-specific treatment effect estimates are quite noisy. One consequence is that actual treatment effects will typically be closer to the average than are estimated effects, even when the research design permits unbiased estimation of each effect. Thus, it is common in these literatures to “shrink” the estimated treatment effects toward the mean. The procedure goes by many different names – e.g., shrinkage, Empirical Bayes, regularization, partial pooling, multi-level modeling – but the basic idea is that the posterior estimate of a site’s effect equals a weighted average of the unbiased estimate of that site’s effect and the mean site effect, with weights that depend on the precision of the site estimate.

Let τ_j represent the impact of the program at site j , and suppose that across sites, $\tau_j \sim N(\bar{\tau}, \omega^2)$. Suppose that we have a noisy but unbiased estimate of the site j effect: $t_j | \tau_j \sim N(\tau_j, \sigma^2)$. Then the former can be treated as a prior distribution for τ_j . By Bayes’ Rule, the posterior mean of τ_j given the observed estimate is

$$E[\tau_j | t_j] = \bar{\tau} + f (t_j - \bar{\tau}),$$

where

$$f = \omega^2 / (\omega^2 + \sigma^2)$$

is the reliability ratio of the site-specific effect estimate.³⁹

When the treatment effect varies systematically with site-level covariates – characteristics either of the treatment or of the counterfactual – this can be used to improve precision. If the site effects are modeled as a function of site characteristics, $\tau_j = X_j \beta + v_j$, with $v_j \sim N(0, \sigma_v^2)$, then the noisy site-level estimate t_j should be shrunk toward the conditional mean rather than to the grand mean:

$$E[\tau_j | t_j, X_j] = X_j \beta + f (t_j - X_j \beta),$$

where f is the conditional reliability ratio, $f = \omega^2 / (\omega^2 + \sigma_v^2)$. This is sometimes known in the statistics literature as “partial pooling.”

One use of the shrinkage approach is by Kane and Staiger (2008), who use a random-assignment experiment to validate non-experimental estimates of teachers’ treatment effects on their students. They shrink the non-experimental estimates, under the assumption that these estimates are valid, and ask whether the result is an unbiased predictor of a teacher’s treatment effects under random assignment.

Kane and Staiger focus on “value-added” scores, estimates of teachers’ effects on their students’ test scores from observational regressions, as the sole non-experimental estimate. They fail to reject the hypothesis that these scores are unbiased predictors of the experimental effects, consistent with the view that they are unconfounded by student sorting. But the experiment has quite low power to

³⁹ The posterior mean is also known as an Empirical Bayes estimate. It is an unbiased predictor of the true site-level treatment effect τ_j if the site-specific estimates t_j are unbiased estimates (Rothstein 2016).

distinguish alternative explanations, and Rothstein (2016) argues that the question remains unresolved.⁴⁰

Angrist et al. (2015) explore the optimal combination of experimental estimates with potentially biased but more precise non-experimental estimates to obtain minimum mean-squared-error predictions of schools' treatment effects. A related question is whether non-experimental measures of other parameters (e.g., classroom observations) can improve the prediction of experimental effects. If so, one might want to use a weighted average of the available measures, weighted to best predict the experimental treatment effect, for performance measurement purposes. To our knowledge, no study has attempted to estimate these weights in an experimental setting (though see Mihaly et al., 2013, for a non-experimental analysis).

ii. Addressing the issue ex ante through the design of the experiment

Ultimately, small sample sizes have limited analysts' ability to identify site- or group-level variation in treatment effects. But there may be ways to design experiments to better support these investigations. Most obviously, resources can be put into collecting data on variation in the quantity and types of treatments delivered, to support analyses (like that of Schochet and Burghardt 2008 or Bloom et al. 2005) of how site treatment effects vary with observable measures of site treatment variation. Large-scale program evaluations often include implementation analyses alongside randomized impact evaluations, and if these two portions were

⁴⁰ For more on the topic of teacher value-added, see Chetty, Friedman, and Rockoff (2014) and Rothstein (2016).

closely integrated the results of the implementation study could be used to inform an analysis of site effects in the impact evaluation sample. Power can also be improved by conducting randomization within site-level strata and by minimizing non-compliance rates (and carefully measuring treatments actually received).

d. Treatment effect heterogeneity and external validity

The empirical literature on program evaluation has been increasingly aware of the importance of potential heterogeneity in treatment effects for interpreting estimates of program impacts and assessing their external validity. Many evaluation samples are drawn from specific populations – individuals in particular regions or cities, individuals entering a program in a certain way, or individuals thought suitable for a proposed alternative program. If treatment effects vary, generalizing from these samples to a broader population is hazardous. Another variant of the external validity problem arises when the compliance rate in the experimental sample differs from what would be expected outside the experiment, as the experimental LATE may not correspond to an appropriate complier population for the program evaluation of interest.

There are several potential sources of heterogeneity. In the previous section, we have discussed differences in characteristics of the environment (such as state of the labor market, including business cycle and industry or occupation structure, population density, or labor market discrimination), differences in aspects of the program (such as unintended differences in the intensity of treatment, something we address under site effects). In this section, we focus on the case where treatment

effects vary because of differences in characteristics at the individual level (such as preferences, abilities, health, beliefs, resources, family environment, or access to networks). Below and in Section IV.f, we also discuss variation treatment effects arising because of variation in structural aspects of the program, such as differences in work incentives.

i. Addressing the issue ex post

The literature is broadly in agreement on how to deal with heterogeneity in treatment effects by *observable* characteristics of study participants. As discussed in Section IV.c, the experimental design implies that one can obtain consistent estimates of the treatment impact for each subgroup, subject to having sufficiently large sample sizes. One can then extrapolate the TOT and ATE to settings with other distributions of observable characteristics by constructing appropriately weighted averages of subgroup effects and corresponding standard errors. As a more common alternative, one can directly estimate TOT and ATE for another population by reweighting the original sample to match the distribution of observable characteristics of the target population (e.g., DiNardo, Fortin, and Lemieux 1996). If multiple treatment sites are available, in principle a similar approach can be used to assess the effect of environmental characteristics, such as labor market conditions or industrial structure.

The case of heterogeneity by *unobserved* characteristics has presented greater challenges. Unfortunately, the individual-level treatment effect is generally not identified either by experimental nor non-experimental methods. Even with perfect compliance, an experiment identifies only the average treatment effect

conditional on observed characteristics.

Some argue that average treatment effects are sufficient for most purposes, as we care only about the distributions of outcomes under alternative policies and not about the positions of particular individuals within those distributions. This is a controversial claim, however – in many contexts, a program that helped some individuals but hurt others by an equal amount, with zero average effect, would be judged worse than nothing.

Moreover, average effects may not be generalizable beyond the population (with perfect compliance, experimental participants, or with imperfect compliance, the subgroup of compliers) identified by an experiment. With heterogeneous treatment effects, neither the TOT nor the complier LATE may be relevant for other populations of interest. A key question then is how representative the experimental compliers are of the group of people that would be potentially affected by the program in question. In many cases the program compliers are likely to be similar to the population of interest, in which case the complier LATE is likely to approximate the relevant parameter. In other cases – for example when compliance is likely to differ between the study and the program at scale – the estimated LATE from one program evaluation may be less useful.

Heckman and Vytlacil (2005) propose a conceptual framework to analyze heterogeneity in treatment effects that relies on the concept of the marginal treatment effect (MTE). If τ_i denotes the individual treatment effect, X_i is a vector of observed individual characteristics, and v_i is the error in the equation determining take up of treatment, then the marginal treatment effect is defined as $E[\tau_i | X_i = x,$

$v_i=v$]; of interest is how this varies with v . This structure provides a framework for considering external validity. The traditional LATE obtained from analyses of experiments with noncompliance can be seen as the integral of the MTE over a particular range of v , but proposals to expand or roll back programs may implicate MTEs at other v values.

To move beyond the LATE, we require a multi-valued instrument that can map out the full distribution of v (or, equivalently, the full range of $\Pr(T = 1 | X)$). If such an instrument is available, the MTE can be obtained by a non-parametric regression of the outcome on the fitted probability of program participation resulting from the first stage equation.⁴¹

This is not possible in the case of a simple RCT. However, when the RCT implemented at multiple sites, and if one is willing to assume that heterogeneity of site effects is limited to compliance rates with no variation in effects on the outcome, one can examine the relationship between the site-specific compliance rate and the site-specific estimated treatment effect (i.e., the site-specific LATE).⁴² (Alternatively, one could directly regress the site-specific treatment effect on the estimated probability of take up and obtain the MTE for different compliance rates.) This relationship could in principle be used to forecast the local average treatment effect

⁴¹ Many other relevant parameters, including LATE and ATE, can be expressed as functions of the MTE. However, to estimate the ATE or the TOT, say, one needs to obtain the MTE for each value of X for the full range of complier probabilities, i.e., from 0 to 1. While in many cases this may be infeasible due to data limitations, if available this could be used to extrapolate the ATE or TOT for populations with different compliance rates and distribution of characteristics.

⁴² Note that the weighting function of the LATE estimator for multi-valued instruments in Angrist and Imbens (1995) is proportional to the differences in take up probabilities between different values of the instrument (ordered by the values' impact on take up). This difference can be interpreted as the difference in compliance between instrument values.

at a potential alternative treatment site (possibly reweighting to adjust for differences in observable characteristics), given a forecast of the new site's compliance rate. More generally, this approach would allow inferring the effect of any intervention affecting the cost of compliance and hence the compliance rate itself.

At times it is useful to go further, to estimating the full distribution of treatment effects. The above method will not accomplish this. Heckman, Smith, and Clements (1997) show that without additional assumptions, experimental data is essentially uninformative about the treatment effects distribution. Moreover, they demonstrate that quite strong assumptions on the dependence of counterfactual outcomes in the control and treatment states are needed to obtain plausible estimates of the distribution of the effect of training in the context of the National Job Training Partnership Act (JTPA) study. Nevertheless, as mentioned at the outset, knowledge of the distribution of heterogeneous treatment effects is undoubtedly important in assessing the impact of a particular program. (though it is less straightforward how such information can be used to address the issue of external validity if treatment effects vary purely with unobserved characteristics).

One approach that has been used to make inferences about heterogeneity in treatment effects is estimation of quantile treatment effects (QTE). As discussed in Section IV.a, the QTE for the q -th quantile is defined as the difference in the q -th quantile of the outcome distribution in the treatment and control groups, respectively. It is clear that absent strong assumptions, such as rank stability, QTEs do not recover the distribution of treatment effects (though they do recover the

effect of the treatment on the outcome distribution, which may be sufficient for many purposes; see Athey and Imbens, this volume). Yet, it can be a helpful and easy-to-implement diagnostic device in at least two senses. First, a QTE analysis can be used to test the assumption of constant treatment effects, which would imply that the QTE is equal at all quantiles. Second, in some cases particular features of a program allow one to derive predictions as to responses in different quantiles of the outcome distribution (see below). More generally, QTE may provide a broad descriptive sense of potential treatment responses.

One source of treatment effect heterogeneity is differences in the structure of the program to be evaluated. In this case, theory may provide weak assumptions that allow making inference on the distribution of treatment effects. Welfare programs represent a good example, since they usually combine a range of different labor supply incentives arising among others from welfare payments, earnings disregards, implicit tax rates, or phase-out regions. Clearly, these incentives interact locally with individual heterogeneity in preferences or ability, something we will return to below. But the additional structure can make for more natural identifying restrictions than in the case of a program that is at least intended to be uniform, such as a training course. A series of papers has addressed this question in the context of evaluation of Connecticut's welfare-to-work program, Jobs First, against the then-prevailing alternative welfare program. For example, to assess the degree of heterogeneity in treatment responses Bitler, Gelbach, and Hoynes (2006) implement a QTE analysis as described above, and relate the resulting estimates to prediction from a standard labor supply model. The Kline and Tartari (2016) study

discussed above, aimed at bounding transition probabilities between counterfactual states, takes advantage of across-participant observable differences in the nature of the decision problem faced to construct revealed-preference restrictions on the set of potential transitions. This is an important diagnostic device for assessing the range of counterfactual treatment responses to the program itself. As discussed above, a potential drawback is that the procedure is rather complex and only applies to the particular program studied. One also has to contend with possibly wide bounds.

In principle, Kline and Tartari's approach can also be used for predicting the effect on the distribution of marginal outcomes of moving from traditional welfare to a welfare-to-work program of the same structure at another site (see Section III.a). Yet, it is worth keeping in mind that the estimated bounds have the LATE property, i.e., they may depend on the particular distribution of individual characteristics and the local environment. Extrapolating to different populations or environments in their context would require imposing additional assumptions on the underlying static labor supply model, and thus trade off additional predictions with robustness.

ii. Addressing the issue ex ante through the design of the experiment

There may be an opportunity to make more progress on this type of treatment effect heterogeneity by building it into the randomization design. Cross-classified and multiple treatment group experiments can be quite helpful for identifying variation in treatment effects.

In some cases, we are directly interested in understanding the distribution of

treatment effects. When a plausible structural model (perhaps something as simple as a Heckman-Vytlacil (2005) Roy model) is available, one might use the structural model to predict individual treatment effects, then stratify the experiment based on these predictions. The NIT studies can be seen as a version of this, as these were stratified based on prior earnings, a potentially strong predictor of the treatment effect.

In other cases, concerns about heterogeneity are driven by potential differences between the complier LATE and the population ATE. Rather than simply assigning participants to be offered or not offered the treatment, one might also vary the extent of efforts to enforce compliance with the experimental assignment. When the relevant selection is thought to be based in part on the anticipated individual treatment effect, as in Heckman and Vytlacil (2005), one can identify the MTE curve directly by randomly assigning participants to multiple values of the incentive (or cost) to obtain the treatment.

Which of these is appropriate depends on the nature of the selection into compliance in the experiment, and how it relates to what would be observed in a non-experimental setting. To make things concrete, we will consider a study in which applicants are randomly assigned to be eligible or ineligible to receive training offered at a particular job-training center. One might expect that non-compliance rates will be low for those assigned to the treatment group for whom it is inconvenient to travel to the program site. One might then expect the LATE to vary with travel costs, but in a simple experiment there is no way to estimate how much of this is due to differences in average treatment effects between those who

live close to and far away from the program site and how much to differences in selection into the complier group.

One way to learn about this would be to implement a more complex, multiple treatment arm experiment in which a subset of individuals offered access to the training are also offered transportation to the training site. If the distance-treatment effect curves differ between the two treatment arms, one can conclude that selection into participation is important, and this can then be used (with a parametric selection model) to estimate how the LATE for a similarly-selected complier population varies with distance. This may be important if the goal is to generalize from the experiment to a scaled-up program that would offer training at a wider number of sites.

One can also use the three-arm experiment to identify the MTE curve, but only with strong restrictions on the shape of this curve (which correspond to strong parametric assumptions about the selection process; see Brinch, Mogstad, and Wiswall forthcoming). These restrictions may be unattractive. If an important goal of the study is to understand how treatment effects vary with the costs of participation, an even more complex experimental design might be called for. Rather than assigning individuals to a treatment group that receives training at zero cost or a control group that is denied access to training at any price, one might use multiple groups that are offered training at different price points (including potentially negative prices). Variation in outcomes across these groups will trace out several points on the MTE curve and can be used to identify a more flexibly shaped curve under weaker assumptions.

Cross-classified and multiple treatment arm experiments raise a number of practical issues that are not confronted in classical treatment/control studies. First, allocating observations across many arms reduces power to detect differences in outcomes between any pair of treatments. Researchers designing experiments must therefore trade off the benefits of a multiple-treatment-arm experiment against reduced ability to detect particular pairwise contrasts. This issue can sometimes be addressed, however, when the alternative arms can be seen as varying the dosage of a single well-defined treatment. An experiment where all treated individuals are assigned a treatment dose of 1 gives *less* power for identifying a linear dose-response relationship than one where the same individuals are assigned varying doses with a mean of 1 (for example, when half are assigned a dose of 0.5 and half are assigned 1.5); moreover, the latter design provides at least the chance of detecting nonlinear effects.

Cross-classified experiments, with a fraction p assigned to treatment A and a fraction q independently assigned to treatment B, can also be seen as sacrificing power, though again the reality is more complex. Let y_{abi} represent the potential outcome for individual i when the program A assignment is a ($a=0$ or 1) and the program B assignment is b . The traditional estimand for evaluation of program A is $E[y_{10i} - y_{00i}]$. Only $(1-q)N$ of the N observations in the cross-classified experiment can be used for estimating this quantity, as the other qN observations are assigned to receive treatment B. But the experiment has full power for estimating the alternative treatment effect $E[((1-q)y_{10i} + qy_{11i}) - ((1-q)y_{00i} + qy_{01i})]$. This can be seen as a weighted average of two treatment effects of program A, one that applies

to individuals who also receive program B and one for those who do not. In some cases, this may be of more interest than the traditional estimand – e.g., when the scaled-up version of program A will coexist with program B.

e. Hidden treatments

A long-standing issue in the interpretation of job training program evaluations is that these evaluations commonly have substantial rates of non-compliance and crossovers. Many people assigned to receive training do not complete their courses, and it has been operationally and politically difficult to exclude people assigned to the control group from receiving treatment, either from the same provider that serves the treatment group or from an alternative provider. Indeed, in some cases, ethical concerns led to decisions to actively inform control group individuals about alternative sources of training.

Much of the literature treats this as non-compliance of the type discussed in Section II.b.ii, so estimates the training effect by dividing the ITT effect by an estimate of the complier share (see, e.g., Heckman, Hohmann, Smith, and Koo, 2000). But this is unsatisfactory when the control group non-compliers receive a different treatment – e.g., training from a different provider – from that given to the treatment group. In technical terms, this is a violation of SUTVA; practically, it means that assignment to treatment may affect outcomes even for the always-takers who receive (some type of) training in any case. To our knowledge, this issue has not been addressed in the enormous literature on job training experiments. (Heckman et al., 2000, note the issue, but their analyses focus on non-random

selection into training and heterogeneity of training effects, which are related but distinct issues.)

Even the IV approach, unsatisfactory as it is, is often not feasible: It requires measuring the share of the control group that crosses over. In many cases, this is not available: The experimental data includes information on the receipt of services from the program under study but not on services obtained from other sources. In this case, only intention-to-treat (ITT) estimates can be computed. But these are attenuated by the failure to measure the “hidden” alternative treatments.

i. Addressing the issue ex post

A very recent literature takes up this topic in the context of the Head Start pre-school program. The Head Start Impact Study randomly assigned Head Start applicants to be offered care or turned away. Many of the control group applicants (and a smaller share of the treatment group) wound up receiving alternative center-based childcare that is thought to be less effective but may be a partial substitute. Where traditional IV estimators treat this as equivalent either to the Head Start treatment or to the receipt of no services, it might be more appropriate to treat it as a distinct, “hidden” treatment.

Walters (2014) estimates heterogeneity in the Head Start effect across centers (sites), finding (among other results) that the LATE of Head Start participation is smaller when more of the complier group is drawn from other centers rather than home-based care. This is suggestive that other center-based care is distinct from home-based care.

Kline and Walters (2014) explicitly model the hidden alternative center treatment, using variation in the compliance patterns across participants' observable characteristics (e.g., parental education) to identify a multinomial variant of a Heckman (1979) parametric selection correction and thus obtain partially experimental estimates of the separate effects of the two types of child care. Their approach leverages variation across observable characteristics (X) in the share of experimental compliers who are drawn from alternative center care, together with a utility-maximizing choice model that constrains how selection on *unobservables* varies with X . With the restrictions imposed by this model, they find large effects of Head Start relative to home-based care. As the Head Start experiment did not directly manipulate the choice between home-based and other center care, they are not able to estimate the relative effect of these with any precision in their least restrictive model, though point estimates are consistent with an effect of other centers comparable to that of Head Start. When Kline and Walters impose stronger restrictions on the selection process, they obtain similar point estimates but with more precision.

Feller et al. (2014) also examine the hidden treatments issue in the Head Start Impact Study sample. They use a principal post-stratification approach that, like Kline and Walters, exploits variation across observables in selection into the two treatments. They couple this to a finite mixture modeling strategy that treats the separation of the two complier subgroup distributions as a deconvolution exercise. Parametric assumptions about these distributions are used to identify the local average treatment effects of the two treatments. Results are similar to Kline

and Walters: Head Start has positive effects on those who would otherwise be at home, but little effect on those who would otherwise receive alternative center-based care.

Another example of the analysis of hidden treatments is Pinto's (2015) analysis of the Moving to Opportunity experiment. In one view, the MTO study involved two treatment arms: One offered a housing voucher that could be used anywhere, and the other restricted the voucher to a low-poverty neighborhood. Straightforward experimental comparisons identify the ITT and LATE of usage of each type of voucher. In another view, however, the relevant treatment is the type of neighborhood in which the participant lives. Kling, Liebman, and Katz (2007) use variation across the two treatment arms and across sites to identify effects of neighborhood poverty (under restrictions on treatment effect heterogeneity). Pinto (2015) adds more structure, using revealed preference restrictions – anyone offered an unrestricted voucher who moves to a low-poverty neighborhood can be assumed to choose the same type of neighborhood in the counterfactual where she receives a restricted voucher – to identify parameters of interest concerning the distribution of neighborhood-type treatment effects.⁴³

ii. Addressing the issue ex ante through the design of the experiment

The Pinto (2015) study takes advantage of the multiple-treatment arms in the MTO experiment, while the Head Start papers discussed above exploit, in

⁴³ Pinto's analysis assumes that the set of neighborhoods in which a voucher can be used is the only relevant difference between the two treatment arms. But in MTO low-poverty voucher recipients were also offered counseling that may have had independent impacts on neighborhood choice or even on outcomes.

various ways, the use of centers as strata in that experiment. This suggests, correctly, that complex experimental designs may be useful in resolving hidden treatment problems, and that a researcher interested in these problems might be able to design an experiment with them in mind. In the neighborhood effects example, one might want to have several treatment arms that vary in the restrictions they place on neighborhood choice; for Head Start, one might explore a third treatment arm that provides a voucher usable either at a Head Start center or at an alternative center. This design might also be useful for a job training evaluation.

In each of these cases, it is *crucial* to collect information about the type and amount of treatment that each participant actually receives; without this, the complex experimental designs are of little value.

f. Mechanisms and multiple treatments

The history in Section III makes clear that many labor market experiments involve variation in more than one aspect of a given program. This is clearly the case when programs consisting of suites of services and incentives are evaluated, such as in randomized evaluations of welfare-to-work programs or of large-scale training programs with a range of integrated services such as JTPA or Job Corps. Yet, even the interpretation of many RCTs of smaller training programs is made difficult by the fact that some form of job search assistance is provided. Simple RCTs do not identify which of the components of the treatment are responsible for the impact. Learning about such mechanisms, besides being of interest in its own right, is

particularly desirable if one wishes to extrapolate to new programs or learn about underlying behavioral parameters. This is for example recognized explicitly in the ongoing evaluation of the REA program discussed in Section III, which aims explicitly at distinguishing the effect of a ‘hassle’ due to being summoned to appear from the actual job search assistance provided.

Even when the treatment has only one component, in many cases that component is sufficiently complex that the average treatment effect is not enough – we want to understand the underlying mechanism. The simplest example of this is labor supply experiments, for which it is often important to distinguish income and substitution effects. It also arises in many of the welfare reform programs, which can create complex changes in intertemporal budget constraints due to time limits or eligibility effects.

i. Addressing the issue ex post

Researchers have used a number of strategies to extract from experimental data evidence on the mechanisms underlying the treatment effects identified by the experiment. In the simplest case, it is sometimes possible to use experimental variation to distinguish the relevant mechanisms, with only minimal restrictions derived from theory. This is most feasible when the experiment involves more than two groups. The first large-scale social experiments, the Negative Income Tax studies, were used in this way. The “treatment” here was a tax schedule described by two parameters: The transfer received if earnings were zero and the tax rate applied to any earnings. The main outcome was labor supply, and a key concern of these studies was to distinguish income from substitution effects.

With a single treatment arm and a single control group, this would not be possible: The net effect of the treatment would be identified, but there would be no way of distinguishing substitution from income effects. (One exception would be if the treatment were designed to be a fully compensated change in the marginal tax rate – this would have no income effect, so the treatment effect would equal the substitution effect. But the NIT treatments were not designed this way.) With multiple treatments that vary both the base transfer and the marginal tax rate, and with an assumption that both income and substitution effects are linear in the relevant tax variable, the two effects can be estimated separately.

To see this, suppose a labor supply function that relates hours of work (H) to the wage rate (w), non-labor income (N), the marginal tax rate (r), and other factors such as preferences for leisure (e):

$$H=f(w, N, r, e).$$

For simplicity of exposition, we assume a constant marginal tax rate, though this is not crucial (see Hausman 1985). A more restrictive assumption is that the individual labor supply function is linear and additively separable in non-labor income and the net-of-tax hourly wage:

$$H_i = \gamma_i + w_i(1-r_i) \delta_i + N_i \eta.$$

Now consider a simple experiment that assigns some individuals to a control group where r_i and N_i are not manipulated, and others to a treatment group that receives an additional baseline transfer D and faces an increment to the tax rate t . Then, adopting the earlier potential outcomes framework, each individual has two potential outcomes:

$$H_{i0} = \gamma_i + w_i(1-r_i) \delta_i + N_i \eta_i \text{ and}$$

$$H_{i1} = \gamma_i + w_i(1-r_i - t) \delta_i + (N_i + D) \eta_i.$$

With random assignment, the difference in mean labor supply between treatment and control groups equals

$$E[H_i | D_i = 1] - E[H_i | D_i = 0] = -t E[w_i \delta_i] + D E[\eta_i].$$

The first term here represents substitution effects, while the second represents income effects. But the simple experiment identifies only the combination of them.

Fortunately, the NIT studies involved multiple treatment arms, with various combinations of transfers and tax rates. Consider a simple extension of the above structure, with two treatment groups 1 and 2 and associated parameters $\{D_1, t_1\}$ and $\{D_2, t_2\}$. Now each individual has three potential outcomes associated with assignment to the control group and each of the treatment groups, $H_0, H_1,$ and H_2 . Two distinct treatment-control contrasts can be computed:

$$E[H_i | D_i = 1] - E[H_i | D_i = 0] = -t_1 E[w_i \delta_i] + D_1 E[\eta_i] \text{ and}$$

$$E[H_i | D_i = 2] - E[H_i | D_i = 0] = -t_2 E[w_i \delta_i] + D_2 E[\eta_i].$$

This is a system of two linear equations and two unknowns. So long as the system has full rank – here, as long as $(D_1/D_2 \neq t_1 / t_2)$ – it can be solved for the mean income elasticity of labor supply, $E[\eta_i]$, and for $E[w_i \delta_i]$. The latter can be divided by the mean wage rate, $E[w_i]$, to obtain a wage-rate-weighted mean substitution elasticity. (With a large enough sample, the mean substitution elasticity, $E[\delta_i]$, could be identified by stratifying the treatment-control comparison by the wage rate.)

A number of studies used the NIT experiment data to estimate the parameters of the labor supply function in basically this way, accounting for additional complications that we neglect here (e.g., participation decisions, non-linear tax schedules, etc.) and often using more complex labor supply functions. See, e.g., Moffitt (1979). But this was by no means universal: In the late 1970s, the experimental paradigm was not as well developed, and many of the studies that used the experimental data did not rely solely on the randomly assigned components of non-labor income and tax rates for identification (e.g., Keeley et al., 1978).

In the above simple model the mean income and labor supply elasticities are just identified with two treatment arms. With more than two arms – the Seattle/Denver experiment alone had 11 – the model is over-identified. This opens the possibility of performing over-identification tests of the restrictions imposed when specifying the labor supply function. Ashenfelter and Plant (1990) estimate separate treatment effects of each treatment arm, but we are not aware of studies that investigate formally whether the pattern of effects is consistent with a posited labor supply function.

Even absent multiple treatment arms, sometimes statistical or theoretical models and assumptions can enable researchers to learn about mechanisms that generate a program effect. For example, Card and Hyslop (2005) [henceforth CH] analyze the data from the Canadian Self Sufficiency Program (SSP) RCT. SSP, a welfare-to-work program, combined a strong, temporary work incentive for participating workers with a fixed initial time period during which welfare

recipients had to establish eligibility in the program by working full time. As a result of this two-tiered structure, the simple experiment analysis does not distinguish the effects of the various components of the program. This makes it difficult to compare the effects of SSP with other welfare-to-work programs, to assess how SSP worked, and to draw lessons for similar programs. CH use a parametric statistical model to separately identify the effect of the different incentives inherent in the SSP program. In contrast to static evaluations of welfare-to-work programs, CH focus on the dynamic labor supply incentives inherent in the program.

One cannot directly analyze the effect of the subsidy (which in the following we will refer to as the SSP program) for those who became eligible because of selection in the eligibility decision. One can, however, model eligibility as a type of imperfect compliance, permitting the estimation of the LATE of SSP on total employment or on the fraction employed at any given point in time. When one turns to dynamic analyses, potential differential changes in the nature of selection in the treatment and control groups make it impossible to estimate the dynamic responses of hazard rates or wages just based on the RCT.⁴⁴ In addition, as in other welfare evaluations, endogenous employment decisions make an analysis of wage outcomes problematic. Another issue is that in the short run the strong work incentive arising from the option value in the eligibility period is potentially confounded with the effect of the subsidy.

⁴⁴ CH use a standard search theory to model the incentives of SSP, and capture the effect of eligibility and the SSP subsidy on labor supply incentives via their effects on the reservation wage. The search model clarifies that in the presence of heterogeneity, the pool of workers employed at any given point in time may be selected, whether or not there also is sample selection arising from employment decisions (e.g., Ham and Lalonde 1996).

To address these difficulties, CH proceed by developing a logistic model with random effects and heterogeneity to estimate a benchmark for welfare transitions in the absence of SSP (i.e., for the control group). This model is then combined with parametric specifications of the treatment effects over different ranges of the program spell, as implied by incentives inherent in SSP. This step includes modeling the participation decision and welfare transitions as functions of the SSP subsidy and current and lagged welfare status. A key assumption thereby is that the chosen controls for heterogeneity and the functional form restrictions are sufficient to control for the dynamic selection bias introduced by the eligibility window. CH experiment with different specifications of heterogeneity, and provide ample discussion of the goodness of fit of the model. As a result of this exercise, they are able to obtain separate effects of eligibility and SSP. This allows them to simulate the effects of different components of the program and counterfactual policy changes relating to the time path of the subsidy.

The approach and finding in CH suggest that one may not need a structural model to separately identify multiple treatment effects, the dynamic effects of a program, or to simulate the effect of alternative policies. However, an assumption on functional form is required, as well as harder-to-assess assumptions on the form of underlying heterogeneity.

To estimate mechanisms underlying the effect of experimental or policy variation, other papers have used insights from theory to aid identification without estimating a structural model. For example, Schmieder, von Wachter, and Bender (2016) use insights from the standard search model to estimate the effect of

unemployment duration on wages. A recurring question in the analysis and evaluation of welfare and unemployment programs has been the effect of employment and unemployment on productivity and wages. If wages rise with employment duration, welfare-to-work programs can lead to sustained labor force participation. In contrast, if longer nonemployment duration reduces wages, and hence the disincentive to work, more generous benefits can lead to a welfare trap.

Card and Hyslop (2005) find that increased employment in the course of the Canadian Self-Sufficiency Program did little to increase wages. In contrast, Grogger (2005) finds positive wage impacts of employment in the context of a randomized evaluation of Florida's welfare-to-work program.

Few papers have directly analyzed the effect of unemployment duration on wages.⁴⁵ The question is difficult for at least two reasons. First, as in Card and Hyslop (2005), even with exogenous variation in incentives at the group level, the type of worker employed at any given point in the unemployment spell may differ between the treatment and control groups.⁴⁶ In other words, it is difficult to find a valid instrument for the duration of unemployment. A second complication arises because even if such variation was available, a change in wages might arise either because of a change in wage offers or due to a change in reservation wages.

To address these difficulties, Schmieder, von Wachter, and Bender (2016) use the fact that the canonical search model has the strong prediction that forward-

⁴⁵ An exception is Addison and Blackburn (2000), who discuss some of the issues that arise. A larger number of papers has addressed the question of duration dependence in unemployment spells. See Kroft, Lange, and Notowidigdo (2013) and references therein.

⁴⁶ This bias arises even in the absence of differences in participation.

looking individuals valuing future unemployment insurance benefits will respond to a benefit extension by raising their reservation wage well before benefit exhaustion. Unless reservation wages do not bind, this implies that extensions in UI durations should lead to increases in observed reemployment wages throughout the spell. In contrast to this prediction, Schmieder et al. (2016) find in the context of discontinuous increases in unemployment insurance durations in Germany, that reemployment wages at different points of the unemployment spells are unaffected. They deduce that reservation wages likely had little effect on observed wages and hence that the effect of an increase in UI benefit durations on wages arose from an effect of the rise in nonemployment durations on offered wages. In this case, an exogenous increase in UI benefit durations can be used as an instrument to estimate the effect of nonemployment duration on wages.⁴⁷

Another study incorporating theoretical insights from search theory into an empirical study of unemployment insurance is that of Della Vigna, Lindner, Reizer, and Schmieder (2016), who analyze a change in the time path of UI benefits in Hungary that kept benefits in the final tier unchanged. They use this variation to structurally estimate key parameters of a model with reference dependence, and find the model does quite well compared to an alternative model that explains the pattern based on (unspecified) heterogeneity. The incorporation of non-standard

⁴⁷ The authors argue that their test excludes any affect of the worker's outside option on wages, and hence the findings are not specific to the particular model.

behavioral assumptions into the evaluation of labor market program is still in its infancy, but is an important avenue for future research.⁴⁸

A closely related topic to the question of mechanisms is the extrapolation of experimental evidence to consider the impacts of new policies, not included in the original evaluation. The value of such extrapolations has long been one of the primary arguments in favor of structural modeling (and against reliance on purely experimental evidence), but some scholars have found out ways to synthesize the approaches. The main challenge here is to bridge between the relatively few parameters that are cleanly identified by an experiment and the larger set of parameters that are needed to characterize most structural models.

One way to do this is to start with a characterization of structural behavior that is simple enough to be captured within the experimental evidence. For example, if one assumes that the labor supply function is characterized by constant income and (compensated) substitution elasticities, then the estimates of these parameters that are identified by the NIT experiments are sufficient to identify the effects of alternative NIT parameters that were not included in the experimental treatments. A draw back of such an approach is that the range of policies that can be examined is limited. The approach can be extended, of course, to estimate a more complex structural model that either relies on additional statistical and theoretical assumptions, additional non-experimental moments, or both. In any event, this sort

⁴⁸ For some exceptions, see, Lemieux and MacLeod (2000), DellaVigna and Paserman (2005), Oreopoulos (2007); more recently, Chan (2014) examines the role of time-inconsistency in the context of the randomized evaluation of Florida Transition Program. Babcock, Congdon, Katz, and Mullainathan (2012) give an overview of the potential importance of behavioral assumption for the evaluation of public programs

of exercise is on more solid ground when trying to interpolate to values within the range of tax parameters included in the experiment than when these parameters need to be extrapolated outside of that range.

A more recent, closely related approach is known as the “sufficient statistics” approach (Chetty 2009). Here, the goal is to characterize optimal policy. Starting with a fully characterized (but usually not overly complex) structural model, it is often possible to derive expressions for social welfare, or for the optimal policy, that depend only on a small number of reduced-form parameters. For example, the Baily-Chetty (Baily 1978, Chetty 2006) formula for optimal unemployment insurance benefits expresses the optimal benefit level in terms of the elasticity of unemployment duration with respect to UI benefits, and the income and substitution effects on the exit hazard from unemployment. If one had experimental evidence regarding these effects, one could use the formula to derive the optimal policy (e.g., Chetty 2008, Card, Chetty, and Weber 2007).

Of course, any sufficient statistics approach is dependent upon the validity of the underlying structural model – there is no assurance that the true structural model generates the same sufficient statistics as does the one posited by the researcher. In some cases, this may include a relevant class of models and hence provide a degree of robustness. For example, Chetty (2009) gives the example of heterogeneity in treatment effects, where the optimal policy depends only on the mean effect. Yet, it can be hard to know which assumptions in the structural model matter, and generally the assumptions needed to derive the sufficient statistics are fairly strong. At a practical level, conclusions about optimal policies may involve

extrapolating very far from the range of policy variation included in the experiment, which means relying strongly on the validity of the theoretical model. In this context, a potential drawback of sufficient statistics is that in contrast to explicitly structural work the empirical fit of the model against the data cannot be assessed.

An alternative approach to obtain a framework for policy extrapolation based on experimental variation is to estimate, or calibrate, a full structural model, using experimental evidence to aid in identifying (some of) the necessary parameters. One approach is to fix individual parameters at the values indicated by experiments, then calibrate or structurally estimate the remainder. This approach is pursued, for example, by Davidson and Woodbury (1997), who use the Illinois reemployment bonus experiment to estimate the parameters of a search cost function, then combine this function with calibrated values, derived from non-experimental data, for other parameters of their model of optimal UI benefits. Another approach is to use experimental data to fit a full structural model, but keep the model sufficiently simple such that the main parameters of the model are identified by the available variation, as for example in DellaVigna, Lindner, Reizer, and Schmieder (2016). An alternative is to estimate the structural model solely with non-experimental data to estimate a structural model, then use experimental evidence to validate predictions that the model makes for particular reduced-form comparisons (e.g., Todd and Wolpin 2006).⁴⁹

⁴⁹ Another approach to extrapolation that can be viewed as a hybrid between structural and reduced form approaches is use experimental variation in the incentive to take up a program to effectively estimate a structural model of the compliance rate (e.g., Heckman and Vytlačil 2005). As described in

ii. *Addressing the issue ex ante through the design of the experiment*

In some cases, the experimental design can be structured to help uncover the mechanisms underlying the treatment effect of the program. Economic theory may be particularly useful here in connecting fundamental parameters and mechanisms to the types of impacts that can be measured with experiments. One approach is to design an experiment that targets a particular mechanism of interest, rather than identifying the effect of a well-defined program that might be implemented. Kling, Congdon, Ludwig, and Mullainathan (this volume) refer to this as a “mechanism experiment,” distinguishing it from a program evaluation. Standard models in labor economics or other fields may provide useful characterizations of the behavioral mechanisms to be tested. For example, models of human capital investment have implications for the factors determining take up and success of training or schooling programs that may be useful in structuring the experimental design.

A closely related approach is to introduce multiple treatment arms, with program variation among them that can help uncover underlying parameters. The NIT experiments discussed above present a straightforward example of a congenial marriage of classic (static) labor supply theory and the experimental design. As discussed above, as long as both income and substitution effects are linear in the relevant tax measure, multiple treatments manipulating both the base transfer and

Section 5.d, under certain circumstances this allows one to obtain the full distribution of marginal treatment effects and hence to extrapolate.

the marginal tax rate can be used to separately estimate the income and substitution effects.⁵⁰

The evaluation of the SSP program discussed above is a good example of an experiment that would have benefited from a second treatment arm. Such a treatment might have randomly varied the incentive to become eligible for the (randomly assigned) work subsidy in the main phase of the program. More generally, decisions and programs involving inter-temporal tradeoffs may be an area in which more complex experiments can be particularly insightful. For example, typical UI systems involve expiring benefits, or JSA programs involve sanctions; the timing of benefit exhaustion, reemployment bonuses, or sanctions has been shown to have important empirical effects on reemployment rates (e.g., Meyer 1995, Black, Smith, Berger, and Noel 2003, Schmieder, von Wachter, and Bender 2012). Hence, experiments that try and get at the underlying behavioral mechanisms may provide important insights into how these programs affect labor supply choices. Knowledge of such mechanisms is also a crucial input in optimizing the delivery of insurance or assistance in the labor market. For example, this could involve a reemployment bonus that declines over time, or one that is available only to those who survive to a specified point. By randomly varying the amount, slope, or intervals, one may gain insights into the nature of inter-temporal decision making relevant for these programs. Inter-temporal choice is also an area where theory is likely to be helpful to provide identifying structure. For example, if the goal would

⁵⁰ Multiple treatments may not be necessary. For example, with appropriate data and assumptions, one could in principle experimentally vary *compensated* wage changes to identify the compensated substitution effect. This more closely resembles a mechanism experiment.

be to learn about potential behavioral biases, a model of the effect of particular biases can yield insightful predictions for job search behavior (e.g., Della Vigna, Lindner, Reizer and Schmieder 2016).⁵¹

The usefulness of theory in informing experimental designs hinges, of course, on the model being correct. To mitigate the reliance on particular assumptions (e.g., on functional forms) in principle one could use revealed preference arguments to generate robust predictions from theory that are then used in design of an experiment. E.g., one could use results obtained by Pinto (2015) or Kline and Tartari (2016) to devise multiple treatment arms to test the implied restrictions. However, a model may not be necessary to enrich the experimental design to study underlying channels. The SSP example shows that a basic understanding of the incentives and the nature of the program can be sufficient to design an RCT that uncovers the potentially complex mechanisms underlying the simple SSP evaluation.

V. Conclusion

Because they allow researchers to control assignment into treatment, randomized controlled trials are the Gold Standard for program evaluation. But while random assignment solves the selection problem, there are a broad range of additional relevant design issues that arise routinely in the analysis of central economic questions that are not solved by random assignment on its own. In this

⁵¹ As already mentioned in the discussion of heterogeneous treatment effects, another area where theory is likely to be useful is to understand the determination of compliance rates. As discussed above, the main idea is to experimentally manipulate the incentive to participate and use the variation to trace out the marginal treatment effect (MTE) curve. Theoretical considerations can tell us how to realistically vary the cost of compliance and hence be able to estimate the full range of treatment effects.

chapter, we have discussed six such design issues in depth, including (1) spillover effects and interactions between individuals, leading to a failure of SUTVA; (2) impacts on outcomes that are only observed conditional on individual choices and hence are endogenous, such as wages, hours worked, or participation in a follow-up survey; (3) heterogeneity in treatment effects between experimental sites and observed population groups, or (4) imperfect compliance and heterogeneity in unobserved characteristics, both of which can make it hard to interpret treatment effects and extrapolate to other programs; (5) hidden treatment effects arising because controls also receive versions of the treatment; and (6) the understanding of the mechanisms behind the treatment effect, in particular in the presence of multiple treatment.

We discuss these design issues and solutions in the context of social experiments in the United States labor market, which have provided most of what we know about the functioning of the main labor market programs. Of course, the labor economics literature has been well aware about the limitations of experiments in general and some of these design issues in particular. We have reviewed approaches that can be used to address the design issues in the context of randomized experiments. This includes approaches that can be applied once randomization is completed, and ways to modify the experiments itself to address the concerns we identify.

While we discuss design issues in the context of experiments in the labor market, these issues can arise in all areas that have seen active experimental activities, including field experiments discussed elsewhere in this volume. Hence the

solutions we identify can be applied to a broad range of questions and should be useful for a wide range of researchers interested in harnessing the power of randomized controlled trials.

We close with a brief discussion of recent trends in labor market social experiments, several of which highlight the need to pay more attention to the potential design issues in experimental evaluations that we discuss. One overarching trend, cutting across several areas of research, is that academic economists have become more involved with the implementation of experiments. In labor economics, for example, this has meant a shift away from randomized controlled trials implemented by large, specialized policy consulting firms (e.g., Mathematica, MDRC, or Abt Associates). For example, several experiments have evaluated take up of actual government programs within the context of services provided by H&R Block (e.g., Bettinger, Long, Oreopoulos, and Sanbonmatsu 2012). Another example is the increasing number of randomized trials evaluating the role of economic incentives for teachers (e.g., Fryer, Levitt, List, and Sadoff 2012; Fryer 2013; Springer et al. 2010). Similarly, experiments taking place within private businesses have also been quite successful (e.g., Bandiera, Barankay, and Rasul 2009).

The greater involvement of academic economists harbors both upside potential, if researchers implement state-of-the-art techniques to address additional design issues, and challenges, as there is a broad range of issues that must be considered and monitored when implementing an experimental evaluation of an existing program or a new, complex treatment in a real-world setting. We hope the

discussion of the design issues in this chapter, as well as our summary of the practical aspects of implementing social experiments, will provide a useful guide for those interested in implementing such social experiments.

A second, related trend has been a movement toward evaluating topics in personnel economics (e.g., the response of teachers to incentive pay programs) as distinct from government social programs. These are often conducted within particular firms, and implicate a number of the design issues we discuss, most notably issues of site effects and heterogeneity.

A third important trend has been the use of the actual online labor market, for what amount to field experiments in the taxonomy we set out at the outset (e.g., Pallais 2014). The Internet may well provide a useful resource for future social experiments as well. A key advantage may be that researchers may be able to better control the environment, perhaps allowing them to implement more complex study designs that address some of the issues we pose.

References

- Addison, J. T., Blackburn, M. L. 2000. The effects of unemployment insurance on postunemployment earnings. *Labour Economics*, 7(1), 21-53.
- Ahn, H., Powell, J. L. 1993. Semiparametric estimation of censored selection models with a nonparametric selection mechanism. *Journal of Econometrics*, 58(1), 3-29.
- Alcott, H. 2015. Site selection bias in program evaluation. *Quarterly Journal of Economics*, 130 (3), 1117-1165.
- Altonji, J. G., Blank, R. M. 1999. Race and gender in the labor market. *Handbook of Labor Economics*, 3 (3), 3143-3259.
- Anderson, M. 2008. Multiple inference and gender differences in the effects of early intervention: A reevaluation of the Abecedarian, Perry Preschool, and Early

- Training projects. *Journal of the American Statistical Association*, 103(484), 1481-1495.
- Angrist, J.D., Hull, P., Pathak, P. A., Walters, C. 2015. Leveraging lotteries for school value-added: Testing and estimation. (Working Paper 21748). National Bureau of Economic Research.
- Angrist, J. D., Imbens, G. W. 1995. Two-stage least squares estimation of average causal effects in models with variable treatment intensity. *Journal of the American Statistical Association*, 90 (430), 431-442.
- Angrist, J. D., Imbens, G. W., Rubin, D. B. 1996. Identification of causal effects using instrumental variables. *Journal of the American Statistical Association*, 91(434), 444-455.
- Angrist, J. D., Krueger, A.B. 1999. Empirical strategies in labor economics." *Handbook of Labor Economics*, 3, 1277-1366.
- Ashenfelter, O., Ashmore, D., Deschênes O. 2004. Do unemployment insurance recipients actively seek work? Evidence from randomized trials in four US states. *Journal of Econometrics*, 125(1-2), 53-75.
- Ashenfelter, O., Plant, M. W. 1990. Nonparametric estimates of the labor-supply effects of negative income tax programs. *Journal of Labor Economics*, 8 (1), S396-S415.
- Athey, S., Imbens, G. 2016. The econometrics of randomized experiments. *Handbook of Field Experiments* (forthcoming).
- Babcock, L., Congdon, W. J., Katz, L. F., Mullainathan, S. 2012. Notes on behavioral economics and labor market policy. *IZA Journal of Labor Policy*, 1(2), 1-14.
- Baily, M. N. 1978. Some aspects of optimal unemployment insurance. *Journal of Public Economics*, 10(3), 379-402.
- Baird, S., Bohren, A., McIntosh, C., Ozler, B. 2015. Designing experiments to measure spillover effects, second version (Working Paper 15-021). Penn Institute for Economic Research.
- Bandiera, O., Bankaray, I., Rasul, I. 2009. Social connections and incentives in the workplace: Evidence from personnel data. *Econometrica*, 77(4): 1047-1094.
- Barnes, M. S., Benus, J., Cooper J., Dugan, M.K., Kirsch M.P., Johnson, T. 2014. U.S. Department of Labor Jobs Corps Process Study Final Report. U.S. Department of Labor. [Available at: http://wdr.doleta.gov/research/keyword.cfm?fuseaction=dsp_resultDetails&pub_id=2538&mp=y].

- Barnow, B. S. 2000. Exploring the relationship between performance management and program impact: A case study of the Job Training Partnership Act. *Journal of Policy Analysis and Management*, 19(1), 118-141.
- Becerra, R. M., Lew, V., Mitchell, M. N., Ono, H. 1998. Final report: California Work Pays Demonstration Project, report of the first forty-two months. School of Public Policy and Social Research, University of California-Los Angeles, Los Angeles.
- Beecroft, E., Lee, W., Long, D., Holcomb, P. A., Thompson, T. S., Pindus, N., O'Brien, C., Bernstein, J. 2003. The Indiana welfare reform evaluation: Five-year impacts, implementation, costs and benefits. Abt Associates: Cambridge, MA.
- Bell, S. H., Bloom, H. S., Cave, G., Doolittle, F., Lin, W., Orr, L. L. 1994. The National JTPA Study: Overview: Impacts, benefits, and costs of Title II-A. Abt Associates: Cambridge MA
- Bell, S. H., Orr, L. L., Burstein, N. R. 1987. Evaluation of the AFDC Homemaker-Home Health Aide Demonstrations: Overview of evaluation results. Abt Associates: Cambridge MA.
- Benus, J., Yamagata, E. P., Wang, Y., Blass, E. 2008. Reemployment and Eligibility Assessment (REA) study: FY 2005 initiative: Final report. IMPAQ International, 1-173.
- Bertrand, M., Mullainathan, S. 2004. Are Emily and Greg more employable than Lakisha and Jamal? A field experiment on labor market discrimination. *American Economic Review*, 94(4), 991-1013.
- Bettinger, E., Long, B. T., Oreopoulos, P., Sanbonmatsu, L. 2012. The role of application assistance and information in college decisions: Results from the H&R Block FAFSA experiment. *Quarterly Journal of Economics*, 127(3), 1205-1242.
- Bitler, M. P., Gelbach, J. B., Hoynes, H. W. 2006. What mean impacts miss: Distributional effects of welfare reform experiments. *The American Economic Review*, 96 (4), 988-1012.
- Black, D. A., Galdo, J., Smith, J. A. 2007. Evaluating the Worker Profiling and Reemployment Services System using a regression discontinuity approach. *The American Economic Review*, 97 (2), 104-107.
- Black, D. A., Smith, J. A., Berger, M. C., Noel B. J. 2003. Is the threat of reemployment services more effective than the services themselves? Evidence from random assignment in the UI system. *American Economic Review*, 93(4), 1313-1327.

- Bloom, H.S., Hill, C. J., Riccio J. A. 2005. Modeling cross-site experimental differences to find out why program effectiveness varies. In Bloom, H.S., ed., *Learning more from social experiments: Evolving analytic approaches*. Russell Sage Foundation, 37-74.
- Bloom, D., Kemple, J. J., Morris, P., Scrivener, S., Verma, N., Hendra, R. 2000. Final report on Florida's initial time-limited welfare program. Manpower Demonstration Research Corporation: New York, December.
- Bloom, H. S., Orr, L.L., Bell, S.H., Cave, G., Doolittle, F., Lin, W., Bos, J.M. 1997. The benefits and costs of JTPA Title II-A programs: Key findings from the National Job Training Partnership Act Study. *Journal of Human Resources*, 32 (3), 549-576.
- Bloom, D., Scrivener, S., Michalopoulos, C., Morris, P., Hendra, R., Adams-Ciardullo, D., Walter, J. 2002. *Jobs First: Final report on Connecticut's welfare reform initiative*. Manpower Demonstration Research Corporation.
- Blundell, R., Bozio, A., Laroque, G. 2011. Labor supply and the extensive margin. *The American Economic Review*, 101(3), 482-486.
- Blundell, R., Dias, M. C., Meghir, C., Reenen, J. V. 2004. Evaluating the employment impact of a mandatory job search program. *Journal of the European Economic Association*, 2(4), 569-606.
- Brinch, C., Mogstad, M., Wiswall, M. Forthcoming . Beyond LATE with a discrete instrument. *Journal of Political Economy*.
- Buchinsky, M. 1994. Changes in the US wage structure 1963-1987: Application of quantile regression. *Econometrica: Journal of the Econometric Society*, 62 (2), 405-458.
- Burghardt, J., Schochet, P. Z., McConnell, S., Johnson, T., Gritz, R. M., Glazerman, S., Homrighausen, J., Jackson, R. 2001. *Does Job Corps work? Summary of the National Job Corps Study*. Mathematica Policy Research: Princeton, NJ.
- Card, D., Chetty, R., Weber, A. 2007. Cash-on-hand and competing models of intertemporal behavior: New evidence from the labor market. *The Quarterly Journal of Economics*, 122(4), 1511-1560.
- Card, D., Hyslop, D. R. 2005. Estimating the effects of a time-limited earnings subsidy for welfare-leavers. *Econometrica*, 73(6), 1723-1770.
- Card, D., Kluve, J., Weber, A. 2010. Active labor market programs: A meta-analysis. *The Economic Journal*, 120(548), F452-477.

- Cave, G., Bos, H., Doolittle, F., Toussaint, C. 1993. JOBSTART. Final report on a program for school dropouts. Manpower Demonstration Research Corp: New York.
- Cerqua, A., Pellegrini, G. 2014. Do subsidies to private capital boost firms' growth? A multiple regression discontinuity design approach. *Journal of Public Economics*, 109 (C), 114-126.
- Chan, M. K. 2014. Welfare dependence and self-control: An empirical analysis. Working paper, Economics Discipline Group, UTS Business School, University of Technology, Sydney.
- Chetty, R. 2006. A general formula for the optimal level of social insurance. *Journal of Public Economics*, 90(10), 1879-1901.
- Chetty, R. 2008. Moral hazard versus liquidity and optimal unemployment insurance. *Journal of Political Economy*, 116(2), 173-234.
- Chetty, R. 2009. Is the taxable income elasticity sufficient to calculate deadweight loss? The implications of evasion and avoidance. *American Economic Journal: Economic Policy*, 1(2), 31-52.
- Chetty, R., Friedman, J. N., Rockoff, J.E. 2014. Measuring the impacts of teachers I: Evaluating bias in teacher value-added estimates. *American Economic Review*, 104(9), 2593-2632.
- Chodorow-Reich, G. , Karababounis, L. 2016. The limited macroeconomic effects of unemployment benefit extensions (Working Paper 22163). National Bureau of Economic Research.
- Coglianesi, J. J. (Working Paper). 2015. Do unemployment insurance extensions reduce employment? Mimeo, Harvard University.
- Corson, W., Decker, P., Dunstan, S. M., Kerachsky, S. 1991. Pennsylvania reemployment bonus demonstration: Final report (Unemployment Insurance Occasional Paper 92-1). U.S. Department of Labor: Washington, DC.
- Corson, W., Long, D., Nicholson, W. 1984. Evaluation of the Charleston Claimant Placement and Work Test Demonstration. Mathematica Policy Research.
- Crépon, B., Duflo, E., Gurgand, M., Rathelot, R., Zamora, P. 2013. Do labor market policies have displacement effects? Evidence from a clustered randomized experiment. *The Quarterly Journal of Economics*, 128(2), 531-580.
- Davidson, C., Woodbury, S. A. 1997. Optimal unemployment insurance. *Journal of Public Economics*, 64(3), 359-387.

- Deaton, A. 2010. Instruments, randomization, and learning about development. *Journal of Economic Literature*, 48(2), 424-455.
- Dehejia, R. H., Wahba, S. 2002. Propensity score-matching methods for nonexperimental causal studies. *Review of Economics and Statistics*, 84(1), 151-161.
- DellaVigna, S., Lindner, A., Reizer, B., Schmieder, J. F. 2016. Reference-dependent job search: evidence from Hungary (Working Paper 22257). National Bureau of Economic Research.
- DellaVigna, S., Paserman, M. D. 2005. Job search and impatience. *Journal of Labor Economics*, 23(3), 527-588.
- DiNardo, J., Fortin, N. M., Lemieux, T. 1996. Labor market institutions and the distribution of wages, 1973-1992: A Semiparametric Approach. *Econometrica*, 64(5), 1001-1044.
- Dorsett, R., Hendra, R., Robins, P. K., Williams, S. 2013. Can post-employment services combined with financial incentives improve employment retention for welfare recipients? Evidence from the Texas Employment Retention and Advancement Evaluation. NIESR Discussion Paper No. 409.
- Farber H. S., Silverman, D., Wachter, T. 2015. Factors determining callbacks to job applications by the unemployed: An audit study (Working Paper 21689). National Bureau of Economic Research.
- Fein, D. J., Beecroft, E., Blomquist, J. D. 1994. Ohio Transitions to Independence Demonstration. Final impacts for JOBS and work choice. Abt Associates: Cambridge, MA.
- Feller, A., Grindal, T., Miratrix, L. W., Page, L. C. 2014. Compared to what? Variation in the impacts of early childhood education by alternative care-type settings. Working paper.
- Ferracci, M., Jolivet, G., van den Berg, G. J. 2010. Treatment evaluation in the case of interactions within markets (No. 4700). Working paper, Institute for the Study of Labor (IZA).
- Fraker, T., Maynard, R. 1987. The adequacy of comparison group designs for evaluations of employment-related programs. *Journal of Human Resources*, 22(2), 194-227.
- Freedman, S., Friedlander, D., Riccio, J. 1994. GAIN: Benefits, costs, and three-year impacts of a welfare-to-work program. Manpower Demonstration Research Corp.

- Freedman, S., Knab, J. T., Gennetian, L. A., Navarro, D. 2000. The Los Angeles Jobs-First GAIN Evaluation: Final report on a work first program in a major urban center. Manpower Demonstration Research Corporation: New York.
- Fryer, R., 2013. Teacher incentives and student achievement: Evidence from New York City public schools. *Journal of Labor Economics*, 31(2), 373-427.
- Fryer, R., Levitt, S. D., List, J., Sadoff, S. 2012. Enhancing the efficacy of teacher incentives through loss aversion: A field experiment (Working Paper 18237). National Bureau of Economic Research.
- Gautier, P. A., Muller, P., Rosholm, M., Svarer, M., van der Klaauw, B. 2012. Estimating equilibrium effects of job search assistance (No. 9066). CEPR Discussion Papers.
- Gold, S.F., 1971. The failure of the Work Incentive (WIN) program. *University of Pennsylvania Law Review*, 119(3), 485-501.
- Greenberg, D. H., Robins, P. K. 1986. The changing role of social experiments in policy analysis. *Journal of Policy Analysis and Management*, 5(2), 340-362.
- Greenberg, D. H., Shroder, M. 2004. The digest of social experiments. The Urban Institute, 3rd edition.
- Greenberg D. H., Shroder, M., Onstott, M. 1999. The social experiment market. *The Journal of Economic Perspectives*, 13(3), 157-172.
- Grogger, J. 2005. Welfare reform, returns to experience, and wages: Using reservation wages to account for sample selection bias. *The Review of Economics and Statistics*, 91(3), 490-502.
- Gronau, R. 1973. The effect of children on the housewife's value of time. *Journal of Political Economy*, 81 (2), S168-S199.
- Grossman, J.B., Roberts, J., 1989. Welfare savings from employment and training programs for welfare recipients. *The Review of Economics and Statistics*, 71(3), 532-537.
- Gueron, J. Forthcoming. The politics and practice of social experiments: seeds of a revolution. *Handbook of Field Experiments*.
- Hagedorn, M., Karahan, F., Manovskii, I., Mitman, K. 2015. Unemployment benefits and unemployment in the Great Recession: the role of macro effects. Federal Reserve Bank of New York Staff Report 646, revised February 2015.
- Hagedorn, M., Manovskii, I., Mitman, K. 2015. The impact of unemployment benefit extensions on employment: The 2014 employment miracle? (Working Paper 20884). National Bureau of Economic Research.

- Ham, J. C., LaLonde, R. J. 1996. The effect of sample selection and initial conditions in duration models: Evidence from experimental data on training. *Econometrica: Journal of the Econometric Society*, 64(1), 175-205.
- Ham, J. C., Li, X., Reagan, P. B. 2011. Matching and semi-parametric IV estimation, a distance-based measure of migration, and the wages of young men. *Journal of Econometrics*, 161(2), 208-227.
- Hamilton, G., Freedman, S., Gennetian, L., Michalopoulos, C., Walter, J. 2001. National evaluation of welfare-to-work strategies: How effective are different welfare-to-work approaches? Five-year adult and child impacts for eleven programs. US Department of Health and Human Services and US Department of Education: Washington, DC.
- Hamilton, G. and S. Scrivener. 2012. Increasing employment stability and earnings for low-wage workers lessons from the Employment Retention and Advancement (ERA) project. Office of Planning, Research and Evaluation Report 2012-19. Administration for Children and Families, U.S. Department of Health and Human Services.
- Harrison, G. W., List, J. A. 2004. Field experiments. *Journal of Economic Literature*, 42 (4), 1009-1055.
- Hausman, J. A. 1985. The econometrics of nonlinear budget sets. Fisher-Shultz lecture for the Econometric Society, Dublin: 1982. *Econometrica*, 53(6), 1255-1282.
- Hausman, J. A., Wise, D. A. 1979. Attrition bias in experimental and panel data: The Gary Income Maintenance Experiment. *Econometrica*, 47(2), 455-73.
- Heckman, J. J. 1979. Sample selection bias as a specification error. *Econometrica*, 47(1), 153-61.
- Heckman, J. J. 2010. Building bridges between structural and program evaluation approaches to evaluating policy. *Journal of Economic Literature*, 48(2), 356-98.
- Heckman, J., Hohmann, N., Smith, J., Khoo, M. 2000. Substitution and dropout bias in social experiments: A study of an influential social experiment. *The Quarterly Journal of Economics*, 115(2), 651-694.
- Heckman, J. J., Hotz, V. J. 1989. Choosing among alternative nonexperimental methods for estimating the impact of social programs: The case of manpower training. *Journal of the American statistical Association*, 84(408), 862-874.
- Heckman, J. J., LaLonde, R. J., Smith, J. A. 1999. The economics and econometrics of active labor market programs. *Handbook of Labor Economics*, 3, 1865-2097.

- Heckman, J. J., Smith, J. A. 1995. Assessing the case for social experiments. *The Journal of Economic Perspectives*, 9(2), 85-110.
- Heckman, J. J., Smith, J., Clements, N. 1997. Making the most out of programme evaluations and social experiments: accounting for heterogeneity in programme impacts. *Review of Economic Studies*, 64(4), 487-535.
- Heckman, J. J., Vytlacil, E. 2005. Structural equations, treatment effects, and econometric policy evaluation. *Econometrica*, 73(3), 669-738.
- Herrem, J.W., Schmitt, L.C. 1983. Eligibility review pilot project handbook. Wisconsin Department of Industry, Labor, and Human Relations: Madison, WI.
- Holland, P. W. 1986. Statistics and causal inference. *Journal of the American Statistical Association*, 81(396), 945-960.
- Horowitz, J. L., Manski, C. F. 2000. Nonparametric analysis of randomized experiments with missing covariate and outcome data. *Journal of the American Statistical Association*, 95(449), 77-84.
- Hotz, J. 1992. Recent experience in designing evaluations of social programs: The case of the National JTPA study. In Garfinkel, I., Manski, C., eds., *Evaluating welfare and training programs*, Cambridge, MA: Harvard University Press: 76-114.
- Hotz, J., Imbens, G., Klerman, J. 2006. Evaluating the differential effects of alternative welfare-to-work training components: A reanalysis of the california GAIN program. *Journal of Labor Economics*, 24(3), 521-566.
- Jackson, K. C., Rockoff, J. E., Staiger, D. O. 2014. Teacher effects and teacher-related policies. *Annu. Rev. Econ*, 6(1), 801-825.
- Jacobson, L. S. 2009. Strengthening one-stop career centers: Helping more unemployed workers find jobs and build skills. Hamilton Project Discussion Paper 2009-01, April: The Brookings Institution, Washington DC.
- Jaggers, M. 1984. ERP pilot project final report. Wisconsin Department of Industry, Labor, and Human Relations: Madison, WI.
- Johnson, T.R., Pfiester, J.M., West, R.W., Dickinson, K.P. 1984. Design and implementation of the claimant placement and work test demonstration. SRI International: Menlo Park, CA.
- Johnson, W., Kitamura, Y., Neal, D. 2000. Evaluating a simple method for estimating black-white gaps in median wages. *American Economic Review*, 90(2), 339-343.

- Johnston, A. C., Mas, A. 2015. Potential unemployment insurance duration and labor supply: The individual and market-level response to a benefit cut. Unpublished working paper. Princeton University.
- Kane, T. J., Staiger, D. O. 2008. Estimating teacher impacts on student achievement: An experimental evaluation (Working Paper 14607). National Bureau of Economic Research.
- Keane, M. P. 2010. Structural vs. atheoretic approaches to econometrics. *Journal of Econometrics*, 156(1), 3-20.
- Keeley, M. C., Robins, P. K., Spiegelman, R. G., West, R. W. 1978. The estimation of labor supply models using experimental data. *The American Economic Review*, 68 (5), 873-887.
- Kehrer, K. C., Moffitt, R. A., eds. 1976. *The Gary income maintenance experiment: Initial findings report*. Indiana University: Gary, Ind.
- Kemple, J. J., Friedlander, D., Fellerath V. 1995. *Florida's Project Independence. Benefits, costs, and two-year impacts of Florida's JOBS program*. Manpower Demonstration Research Corporation: New York.
- Kershaw, D., Fair, J. 1976. *The New Jersey income maintenance experiment. Volume 1: Operations, Surveys and Administration*. Academic Press: New York.
- Klepinger, D. H., Johnson, T. R., Joesch, J. M., Benus, J. M. 1997. *Evaluation of the Maryland unemployment insurance work search demonstration (Unemployment Insurance Occasional Paper 98-2)*. U.S. Department of Labor, Employment and Training Administration, Unemployment Insurance Service: Washington DC.
- Klepinger, D.H., Johnson, T.R. Joesch, J.M., 2002. *Effects of unemployment insurance work-search requirements: The Maryland experiment*. *Industrial & Labor Relations Review*, 56(1), pp.3-22.
- Klerman, J. A., Minzner, A., Harkness, J., Mills, S., Cook, R., Savidge-Wilkins, G. 2013. *Design report: Impact evaluation of reemployment and eligibility assessment Program*. Abt Associates: May 7.
- Kline, P., Tartari, M. 2016. *Bounding the labor supply responses to a randomized welfare experiment: A revealed preference approach*. *American Economic Review*, 106 (4), 972-1014.
- Kline, P., Walters, C. 2014. *Evaluating public programs with close substitutes: The case of Head Start*. UC Berkeley Institute for Research on Labor and Employment Working Paper #123-14.

- Kling, J. R., Liebman, J. B., Katz, L. F. 2007. Experimental analysis of neighborhood effects. *Econometrica*, 75(1), 83-119.
- Kling, J.R., J. Ludwig, B. Congdon and S. Mullainathan. Social policy: Mechanism experiments and policy evaluations. *Handbook of Field Experiments* (forthcoming).
- Knox, V. W., Miller, C., Gennetian, L. A. 2000. Reforming welfare and rewarding work: A summary of the final report on the Minnesota Family Investment Program (Vol. 8). Manpower Demonstration Research Corporation, New York.
- Kornfeld, R., Bloom, H. S. 1999. Measuring program impacts on earnings and employment: Do unemployment insurance wage reports from employers agree with surveys of individuals? *Journal of Labor Economics*, 17(1): 168-97.
- Kroft, K., Lange, F., Notowidigdo, M. J. 2013. Duration dependence and labor market conditions: Evidence from a field experiment. *The Quarterly Journal of Economics*, 128(3), 1123-1167.
- Krueger, A. B., Mueller, A. I. 2016. A contribution to the empirics of reservation wages. *American Economic Journal: Economic Policy*, 8(1), 142-179.
- LaLonde, R. J. 1986. Evaluating the econometric evaluations of training programs with experimental data. *The American Economic Review*, 76(4), 604-620.
- Landais, C., Michaillat, P., Saez, E. 2015. A macroeconomic theory of optimal unemployment insurance (Working Paper 16526). National Bureau of Economic Research.
- Lee, D. S. 2009. Training, wages, and sample selection: estimating sharp bounds on treatment effects. *The Review of Economic Studies*, 76(3), 1071-1102.
- Lemieux, T., MacLeod, W. B. 2000. Supply side hysteresis: The case of the Canadian unemployment insurance system. *Journal of Public Economics*, 78(1), 139-170.
- List, J. A., Rasul I. 2011. Field experiments in labor economics. *Handbook of Labor Economics*, 4(4), 103-228.
- Maguire, S., Freely, J., Clymer, C., Conway, M., Schwartz, D. 2010. Tuning in to local labor markets: Findings from the sectoral employment impact study. Public/Private Ventures: New York.
- Manpower Demonstration Research Corporation Board of Directors. 1980. Summary and findings of the national supported work demonstration. Ballinger Publishing Company: Cambridge, MA.

- Meyer, B. D. 1995. Lessons from the US unemployment insurance experiments. *Journal of Economic Literature*, 33 (1), 91-131.
- Mihaly, K., MaCaffrey D. F., Staiger D. O., Lockwood J. R. 2013. A composite estimator of effective teaching. Met Project. [Available at: [http://www.metproject.org/downloads/MET Composite Estimator of Effective Teaching Research Paper.pdf](http://www.metproject.org/downloads/MET_Composite_Estimator_of_Effective_Teaching_Research_Paper.pdf)].
- Miller, C., Van Dok, M., Tessler, B. L., Pennington, A. 2012. Strategies to help low-wage workers advance: Implementation and final impacts of the Work Advancement and Support Center (WASC) demonstration. Manpower Demonstration Research Corp: New York.
- Minnesota Department of Jobs and Training. 1990. Re-employ Minnesota. In Johnson, E.R., eds., Reemployment services to unemployed workers having difficulty becoming reemployed (Unemployment Insurance Occasional Paper 90-2). U.S. Department of Labor, Employment and Training Administration, Unemployment Insurance Service: Washington, DC.
- Moffitt, R. A. 1979. The labor supply response in the Gary experiment. *Journal of Human Resources*, 14(4), 477-487.
- Newey, W., Powell, J. L., Walker, J. R. 1990. Semiparametric estimation of selection models: Some empirical results. *American Economic Review*, 80(2), 324-28.
- O'Leary, C. J. (2006). State UI job search rules and reemployment services. *Monthly Labor Review*, 129(6), 27-37. <http://research.upjohn.org/jrnarticles/3>.
- Oreopoulos, P. 2007. Do dropouts drop out too soon? Wealth, health and happiness from compulsory schooling. *Journal of Public Economics*, 91, 2213-2229.
- Pallais, A. 2014. Inefficient hiring in entry-level labor markets. *American Economic Review*, 104(11), 3565-3599.
- Palmer, J. L., Pechman, J. A. 1978. Welfare in rural areas: the North Carolina-Iowa income maintenance experiment. Brookings Institution: Washington, DC.
- Perez-Johnson, I., Q. Moore, and R. Santillano. 2011. Improving the effectiveness of individual training accounts: Long-term findings from an experimental evaluation of three service delivery models. Final Report. Mathematica, Inc.
- Pinto, R. 2015. Selection bias in a controlled experiment: The case of Moving to Opportunity. Mimeo., University of Chicago.
- Poe-Yamagata, E., J. Benus, N. Bill, H. Carrington, M. Michaelides, and T. Shen. 2011. Impact of the Reemployment and Eligibility Assessment (REA) initiative. Impaq International.

- Powell, J. L. 1984. Least absolute deviations estimation for the censored regression model. *Journal of Econometrics*, 25(3), 303-325.
- Robins, P. K. 1985. A Comparison of the labor supply findings from the four negative income tax experiments. *Journal of Human Resources*, 20(4) 567-582.
- Rothstein, J. 2010. Teacher quality in educational production: Tracking, decay, and student achievement. *Quarterly Journal of Economics*, 125(1), 175-214.
- Rothstein, J. 2016. Revisiting the impacts of teachers. Unpublished working paper. http://eml.berkeley.edu/~jrothst/workingpapers/rothstein_cfr.pdf.
- Schmieder, J. F., von Wachter, T., Bender, S. 2012. The effects of extended unemployment insurance over the business cycle: Evidence from regression discontinuity estimates over 20 years. *Quarterly Journal of Economics*, 127(2), 701-752.
- Schmieder, J. F., von Wachter, T., Bender, S. 2016. The effect of unemployment benefits and nonemployment durations on wages. *American Economic Review*, 106(3), 739-777.
- Schochet, P. Z., Burghardt, J. A. 2008. Do Job Corps performance measures track program impacts? *Journal of Policy Analysis and Management*, 27(3), 556-576.
- Schochet, P., Burghardt, J., McConnell, S. 2008. Does Job Corps work? Impact findings from the national job corps study. *Mathematica Policy Research*.
- Smith, J. A., Todd, P.E. 2005. Does matching overcome LaLonde's critique of nonexperimental estimators? *Journal of Econometrics*, 125(1), 305-353.
- Spiegelman, R. G., O'Leary, C. J., Kline, K. J. 1992. The Washington Reemployment Bonus experiment: Final report (Unemployment Insurance Occasional Paper 92-6). U.S. Department of Labor: Washington, DC.
- Springer, Matthew G., Dale Ballou, Laura S. Hamilton, Vi-Nhuan Le, J.R. Lockwood, Daniel F. McCaffrey, Matthew Pepper, and Brian M. Stecher. 2010. Teacher pay for performance: Experimental evidence from the Project on Incentives in Teaching. Conference paper, National Center on Performance Incentives.
- SRI International. 1983. Final report of the Seattle-Denver income experiment, Volume I: Design and results. U.S. Department of Health and Human Services: Washington, DC.
- Steinman, J. P. 1978. The Nevada claimant placement program. *Employment Security Research*, Nevada Employment Security Department.

- Todd, P. E., Wolpin, K. I. 2006. Assessing the impact of a school subsidy program in Mexico: Using a social experiment to validate a dynamic behavioral model of child schooling and fertility. *American Economic Review*, 96(5), 1384-1417.
- US Department of Health, Education, and Welfare. 1976. Summary report: Rural income maintenance experiment. Government Printing Office: Washington, DC.
- Vytlacil, E. 2002. Independence, monotonicity, and latent index models: An equivalence result. *Econometrica*, 70(1), 331-341.
- Walters, C. 2014. Inputs in the production of early childhood human capital: Evidence from Head Start (Working Paper 20639). National Bureau of Economic Research.
- Watts, H.W., Rees, A., 1977a. The New Jersey Income Maintenance Experiment, Vol. II: Labor supply responses. Academic Press: New York.
- Watts, H.W., Rees, A., 1977b. The New Jersey Income Maintenance Experiment, Vol. III: Expenditures, health, and social behavior, and the quality of the evidence. Academic Press: New York.
- Woodbury, S. A., Spiegelman, R. G. 1987. Bonuses to workers and employers to reduce unemployment: Randomized trials in Illinois. *The American Economic Review*, 77(4), 513-530.

Table 1: Details on Selected Randomized Controlled Trials of Welfare Programs and Other Labor Supply Incentives for Low-Income Workers in the United States

	Target Population	Primary Intervention	Secondary Intervention	Experiment Title	Start Date	Cost (nominal \$)	Sample Size	Treatment	Funding Source	Outcomes of Interest
(1)	Total family income not exceeding 150 percent of the poverty level	Negative income tax		New Jersey Income Maintenance Experiment	1968	\$ 7,800,000	725 - Treatment 632 - Control 1,357 - Total	Eight combinations of income guarantees and tax rates on other income.	OEO	(1) Reduction in work effort and (2) Lifestyle changes
(2)	Rural, low-income families	Negative income tax		Rural Income Maintenance Experiment	1970	\$ 6,100,000	269 - Treatment 318 - Controls 587 - Total	Five negative income tax plans.	The Ford Fdn., OEO Office of Economic Opportunity	(1) Work behavior; (2) Health, school, and other effects on poor children; and (3) Savings and consumption behavior
(3)	Family earning less than \$11,000 in 1971 dollars	Negative income tax	Vocational training	Seattle-Denver Income Maintenance Experiment	1970	\$77,500,000	1,801-Treatment 1 946-Treatment 2 1,012- Treatment 3 1,041- Control	Two types of treatment: a negative income tax plan and a subsidy to vocational training.	HEW, HHS	(1) Effects on labor supply; (2) Marital stability; and (3) Other lifestyle changes.
(4)	Black families with at least one child under the age of 18	Negative income tax		Gary Income Maintenance Experiment	1971	\$20,300,000	1,028 - Treatment 771 - Control 1,799 - Total	Four combinations of guarantee and tax.	HEW	(1) Employment; (2) Schooling; (3) Infant mortality and morbidity; (4) Educational achievement; and (5) Housing consumption
(5)	One- and two-parent families receiving AFDC	Earned income disregard		California Work Pays Demonstration Program (CWPPDP)	1993	\$ 4,500,000	6,278 -Treatment 1 3,471 -Treatment 2 3,276 - Control 1 1,695 - Control 2 14,720 - Total	The treatment involved changing two provisions of the AFDC program. The "\$30 and one-third" provision applied to all AFDC families and allowed welfare recipients to keep the first \$30 and one-third of the remaining wages before welfare grant determinations were made. However, it expired after the recipient had been in the program for four months, and there-after dollar-for-dollar reductions in grant occurred for every dollar of earnings. Under the 100-hour rule, which applied only to two-parent families, the total work hours per month for the primary wage earner could not exceed 100 hours without loss of eligibility. Experimentals received a waiver of the time limit on the \$30 and one-third income disregard, and a waiver of the 100-hour rule. However, the cash grants of experimentals were reduced by 8.5 percent. Controls were subject to the general AFDC rules, with expiring disregards, ineligibility after 100 hours, and higher benefits.	CA Dept of Social Services	(1) Employment; (2) Earnings; and (3) Welfare receipt

(6)	Families on AFDC	Earned income disregard	Individual job search assistance Case management	Florida Family Transition Program (FTP)	1994	\$11,200,000	1,400 - Treatment 1,400 - Control 2,800 - Total	Limited welfare benefits unless "job-ready", enhanced earnings disregard, and intensive case management	FL Dept of Children and Families US Department of Health and Human Services	(1) Earnings; (2) Welfare benefit receipts; and (3) Outcomes or children
(7)	AFDC recipient and recent applicant families	Reemployment bonus Earned income disregard	Job search inventive Child care services	Minnesota Family Investment Program (MFIP)	1994	\$ 5,090,300	5,275 -Treatment 1 1,933 -Treatment 2 5,634 -Treatment 3 1,797 - Control 14,639 - Total	MFIP provided a 20 percent grant increase when recipients became employed, increased the level of income that would be disregarded in grant calculation, an paid the child care subsidy directly to caregiver. Two-parent families were not subject to work history requirements or to the 100-hour rule. Both single-parent and two-parent families assigned to MFIP were subject to mandatory participation in employment services. Rules and procedures were simplified by combining Food Stamps, AFDC, and Minnesota's Family General Assistance (FGA) to form a single cash benefit program. Subjects assigned to the MFIP incentives-only group received identical benefits as MFIP, but were not required to participate in training services. Two other groups.	MN Dept of Human Services; Ford Fdn.; HHS; US Department of Agriculture; Charles Stewart Mott Fdn.; Annie E Casey Fdn.; McKnight Fdn.; Northwest Area Fdn.	(1) Employment; (2) Earnings; (3) Welfare receipt; (4) Total family income; and (5) Other measures of child and family well-being
(8)	AFDC recipients	Earned income disregard Time limit	Job search incentives Vocational training	Connecticut Jobs First	1996	\$ 5,400,000	2138 - Treatment 1821 - Control 3959 - Total	Earnings disregarded below the federal poverty level and required to participate in Job Search Skills Training.	CT Dept of Social Services	(1) Employment; (2) Earnings; (3) Benefit receipt; and (4) Other measures of child well-being
(9)	UI claimants	Reemployment bonus		Illinois Unemployment Insurance Incentive Experiment	1984	\$ 800,000	4,186 - Treatment (claimants) 3,963 - Treatment (employers) 3,963 - Control 12,112 - Total	Unemployed were offered a \$500 bonus if found a job within 11 weeks and held it for 4 months.	IL Dept of Employment Security; WE Upjohn Institute for Employment Research	(1) Reductions in unemployment spells and (2) Net program savings.
(10)	UI claimants	Reemployment bonus	Job search workshop	Pennsylvania Reemployment Bonus Demonstration	1988	\$ 990,000	14,086 - Treatment 3,392 - Control 17,478 - Total	Five combinations of bonus amount and qualification period.	DOL	(1) UI receipt; (2) Employment; and (3) Earnings
(11)	UI claimants	Reemployment bonus		Washington State Reemployment Bonus Experiment	1988	\$ 450,000	12,451 - Treatment 3,083 - Control 15,534 - Total	6 variations of reemployment bonus amount and qualification periods.	Alfred P Sloan Fdn. US DOL, ETA	(1) Weeks of insured unemployment and (2) UI receipt

Abbreviations: DOL = US Department of Labor; ETA = Employment and Training Administration; Fdn. = Foundation; OEO = Office of Economic Opportunity; HEW = US Department of Health, Education, and Welfare; HHS = US Department of Health and Human Services.

Sources: (1) Kershaw and Fair, 1976; Watts and Rees, 1977a and 1977b; (2) US Department of Health, Education, and Welfare 1976; Palmer and Pechman, 1978; (3) SRI International, 1983; (4) Kehrer, McDonald, and Moffit, 1980; (5) Becerra, Lew, Mitchell, and Ono, 1998; (6) Bloom, Kemple, Morris, Scrivener, Verma, and Hendra, 2000; (7) Knox, Miller, and Gennetian, 2000; (8) Bloom, Scrivener, Michalopoulos, Morris, Hendra, Adams-Ciardullo, and Walter, 2002; (9) Woodbury and Spiegelman, 1987; (10) Corson, Decker, Dunstan, and Kerachsky, 1991; (11) Spiegelman, O'Leary, and Kline, 1992.

Table 2: Details on Selected Randomized Controlled Trials of Programs Offering Job Training and Work Experience for Low-Income Individuals in the United States

Target Population	Primary Intervention	Secondary Intervention	Experiment Title	Start Date	Cost (nominal \$)	Sample Size	Treatment	Funding Source	Outcomes of Interest
(1) AFDC recipients, ex-offenders, substance abusers, and high school dropouts	Work experience		National Supported Work Demonstration (NSWD)	1975	\$ 82,400,000	3,214 -Treatment 3,402 - Control 6,616 - Total	Employment in a structured work experience program involving peer group support, a graduated increase in work standards, and close sympathetic supervision, for 12 to 18 months.	DOL, ETA; DOJ; Law Enforcement Assistance Administration; HHS; National Institute on Drug Abuse; HUD; US Department of Commerce; Ford Fdn.	(1) Increases in post-treatment earnings; (2) Reductions in criminal activity; (3) Reductions in transfers payments; and (4) Reductions in drug abuse
(2) AFDC recipients	Work experience		AFDC Homemaker--Home Health Aide Demonstrations	1983	\$ 8,000,000	4,750 -Treatment 4,750 - Control 9,500 - Total	Experimental AFDC subjects (trainees) received a four- to eight-week training course to become a homemaker-home health aide, followed by a year of subsidized employment. Control subjects did not receive this training, nor did they receive subsidized employment.	Health Care Financing Administration	(1) Employment; (2) Earnings; and (3) AFDC and food stamp payments and receipt
(3) Eligible Job Training Partnership Act Title II adults and out-of-school youth	Vocational training General education Work experience On-the-job-training	Individual job search assistance	National Job Training Partnership Act (JTPA) Study	1987	\$ 23,000,000	20,602	Classroom training, on-the-job training, job search assistance, basic education, and work experience.	DOL	(1) Earnings; (2) Employment; (3) Welfare receipt; and (4) Attainment of educational credentials and occupational competencies
(4) AFDC recipients	Vocational training General education Work experience	Individual job search assistance	Greater Avenues for Independence (GAIN)	1988		24,528-Treatment 8,223 - Control 32,751 - Total	basic education, job search activities, assessments, skills training, and work experience.	California Department of Social Services (CDSS)	(1) Participation in employment-related activities; (2) Earnings; (3) Welfare receipt; and (4) Employment
(5) All recipients of ADC (Ohio's AFDC program)	Work experience General education	Individual job search assistance	JOBS	1989	\$ 3,000,000	24,120-Treatment 4,371 - Control	Mandatory employment and training services, which included basic and post-secondary education, community work experience, and job search assistance.	OH Dept of Human Services	(1) Employment; (2) Earnings; and (3) Welfare receipt
(6) low-income, disadvantaged workers and job seekers	Vocational training	Individual job search assistance	Sectoral Employment Impact Study	2003		1,286 -Total	Industry-specific training programs that prepared unemployed and underskilled workers for skilled positions and connect them with employers seeking to fill such vacancies. Sectoral programs employ various approaches depending on the organization leading the effort and local employers' needs.	Charles Stewart Mott Fdn.	(1) Earnings; (2) Employment; and (3) Quality of jobs

(7)	low-wage workers	Vocational training On-the-job training	Case Management	Work Advancement and Support Center (WASC) Demonstration	2005		1,176 - Dayton 971 - San Diego 705 - Bridgeport 2,852 - Total	The program offered participating workers intensive employment retention and advancement services, including career coaching and access to skills training. It also offered them easier access to work supports, in an effort to increase their incomes in the short run and help stabilize their employment. Finally, both services were offered in one location — in existing One-Stop Career Centers created by the Workforce Investment Act (WIA) of 1998 — and by co-located teams of workforce and welfare staff.	State of Ohio; County of San Diego Health and Human Services Agency; DOL ETA; U.S. Department of Agriculture, Food and Nutrition Service; HHS; Administration for Children and Families; Ford Fdn.; Rockefeller Fdn.; Annie E. Casey Fdn.; David and Lucile Packard Fdn.; The William and Flora Hewlett Fdn.; Joyce Fdn.; James Irvine Fdn.; Charles Stewart Mott Fdn.; Robert Wood Johnson Fdn.	(1) Employment and (2) Earnings (along with many other outcome measures)
(8)	school dropouts aged 17-21 years	General education Vocational training	Individual job search assistance	JOBSTART	1985	\$ 6,200,000	1,163 - Treatment 1,149 - Control 2,312 - Total	Education and vocational training, support services, and job placement assistance.	DOL; Rockefeller Fdn.; Ford Fdn.; Charles Stewart Mott Fdn.; William and Flora Hewlett Fdn.; more foundations.	(1) Educational attainment; (2) Employment; (3) Earnings; and (4) Welfare receipt
(9)	16-24 year olds	General education Vocational training	Health care services Housing services	National Job Corps Study	1994	\$ 21,587,202	9,409 - Treatment 5,977 - Control 15,386 - Total	Treatment group allowed to enroll in Job Corps group. Job Corps centers provide vocational training, academic instruction, health care, social skills training, and counseling.	DOL, ETA	(1) Employment; (2) Earnings; (3) Education and job training; (4) Welfare receipt; (5) Criminal behavior; (6) Drug use; (7) Health factors; and (8) Household status

Abbreviations: DOJ = US Department of Justice; HHS = US Department of Health and Human Services; HUD = US Department of Housing and Urban Development; DOL = US Department of Labor; ETA = Employment and Training Administration; Fdn. = Foundation.

Sources: (1) MDRC Board of Directors, 1980; (2) Bell, Burstein, and Orr, 1987; (3) Bell, Bloom, Cave, Doolittle, and Orr, 1994; Bloom, Orr, Bell, Cave, Doolittle, Lin, and Bos, 1997; (4) Freedman, Friedlander, Riccio, 1994; (5) Fein, Beecroft, and Blomquist, 1994; (6) Maguire, Freely, Clymer, Conway, and Schwartz, 2010; (7) Miller, Van Dok, Tessler, and Pennington, 2012; (8) Cave, Bos, Doolittle, and Toussaint, 1993; (9) Burghardt, Schochet, McConnell, Johnson, Gritz, Glazerman, Homrighausen, and Jackson, 2001.

Table 3: Details on Selected Randomized Controlled Trials of Job Search Assistance Programs for Low-Income Individuals and Unemployed Workers in the United States

Target Population	Primary Intervention	Secondary Interventions	Experiment Title	Start Date	Cost (nominal \$)	Sample Size	Treatment	Funding Source	Outcomes of Interest
(1) Single-parent heads of household who were required to participate in the program (recipients of AFDC)	Job Club	General education Vocational training	Project Independence--Florida	1990	\$ 3,600,000	13,513 - Treatment 4,274 - Control 17,787 - Total	The experimental group was eligible to receive Project Independence services and was subject to a participation mandate. Services included independent job search, job club, assessment, basic education, and training. The control group was not eligible for these services and was not subject to a participation mandate.	Florida Department of Health and Rehabilitative Services Ford Fdn. US Department of Health and Human Services	(1) Employment; (2) Earnings; and (3) AFDC receipt
(2) Single-parent welfare recipients	Job Club Case Management	General education Vocational training	National Evaluation of Welfare-to-Work Strategies (NEWWS)	1991	\$ 31,700,000	44,569 - Total	Eleven programs, broadly defined as either employment-focused or education-focused, were tested in seven sites across the US.		(1) Employment; (2) Earnings; (3) Welfare receipt; (4) Cost-effectiveness; and (5) Child well-being
(3) Families on welfare	Individual job search assistance	Earned income disregard Work experience	Indiana Welfare Reform Evaluation	1995	\$ 23,200,000	63,223 - Treatment 1 3,863 - Treatment 2 3,217 - Control 1 1,091 - Control 2 71,394 - Total	Experimentals were subject new welfare reform policies: assisted job search, broader mandatory work participation, earned income disregard, time limits for case assistance, a revised system of child care provision, family benefit cap, and parental responsibility (such as immunizing children). Controls continued under the traditional AFDC policies	Indiana Family and Social Services Administration US Department of Health and Human Services	(1) Employment; (2) Earnings; (3) Welfare receipt; (4) Income; (5) Health insurance; and (6) Parental responsibility
(4) Single-parent (AFDC-FG) and two-parent (AFDC-U) welfare families in Los Angeles County	Job Club Individual job search assistance job search workshop		LA Jobs-First GAIN Evaluation	1995	\$ 29,900,000	11,521 - Treatment 1 4,039 - Treatment 2 4,162 - Control 1 1,009 - Control 2 20,731 - Total	Members of the treatment group were enrolled in Jobs-First GAIN. These subjects were required to participate in at least one of the job search activities, including job clubs and other informational services and job search training sessions. Experimentals were also exposed to Jobs-First GAIN's intensive work-first message. Sanctions were imposed, usually in the form of partial reductions in welfare benefits, for failure to participate. Controls were not exposed to any of Jobs-First GAIN's services, the intensive work-first message, or sanctions. Controls could still receive assistance from other agencies and were subject to existing welfare rules.	Los Angeles Department of Public Social Services US Department of Health and Human Services Ford Fdn.	(1) Employment; (2) Earnings; (3) Welfare benefits; (4) Outcomes for children; and (5) Incremental effects compared with previous LA GAIN program
(5) UI claimants	Individual job search assistance	Vocational training	Nevada Claimant Placement Program (NCP)	1977		3,500	More staff attention and more referrals, weekly interviews and eligibility checks, all services from same ES/UI team which coordinated their efforts		(1) Weeks of benefits; (2) Earnings; (3) Enforcement of work search rules; (4) Job searches; and (5) Referrals and placements

(6)	UI claimants	Job search incentives Individual job search assistance		Claimant Placement and Work Test Demonstration	1983	\$ 225,000	1,485 - Treatment 1 1,493 - Treatment 2 1,666 - Treatment 3 1,277 - Treatment 4	Job search and placement services	US Department of Health and Human Services Ford Fdn.	(1) Employment and (2) UI payments reductions
(7)	UI claimants indefinitely separated from most recent job	Individual job search assistance		Wisconsin Eligibility Review Pilot Project (ERP)	1983		5000	6-hour job search workshop conducted by ES staff; also tried 3-hour job search workshop		(1) Weeks of benefits; (2) Earnings; (3) Enforcement of work search rules; (4) Job searches; and (5) Referrals and placements
(8)	Unemployed	Case management Individual job search assistance Job search workshop		Reemploy Minnesota (REM)	1988	\$ 835,000	4,212 - Treatment unknown - Control (roughly 10 times treatment)	More personalized and intensive unemployment insurance (UI) services, including case management, intensive job search assistance and job matching, claimant targeting for special assistance, and a job-seeking skills seminar. The control group received regular UI services.	Unemployment Insurance Contingent Account of the Minnesota Department of Jobs and Training	(1) Duration of UI benefits and (2) Amount of UI benefits
(9)	UI claimants	Individual job search assistance	Vocational training	Kentucky Worker Profiling and Reemployment Services (WPRS) Experiment	1994	\$ 15,000	1,236 - Treatment 745 - Control 1,981 - Total	Structured job search activities, employment counseling, and retraining	Kentucky Department of Employment Services	(1) Earnings; (2) Length of benefit receipt; and (3) Amount of UI benefits received
(10)	UI claimants	Alternative work search policies		Maryland Unemployment Insurance Work Search Demonstration	1994	\$ 250,000	3,510 - Treatment 1 3,455 - Treatment 2 3,680 - Treatment 3 3,400 - Treatment 4 4,812 - Control 1 4,901 - Control 2 23,758 - Total	4 different rules changes to Maryland UI eligibility rules	US DOL ETA	(1) UI payments in terms of weeks and dollars; (2) Continuing eligibility; (3) Employment; and (4) Earnings

(11)	UI claimants	Individual job search assistance Case management	Vocational training	Reemployment and Eligibility Assessment (REA)	2013	<p>(1) Current REA Program: assistance--defined as the provision of labor market information, developing an individual reemployment plan, a referral to reemployment services, and direct provision of reemployment services + enforcement (see below)</p> <p>(2) Enforcement Only: the requirement that claimants appear for the REA meeting and that REA program staff verify claimants' eligibility and their participation in work search activities, with referral to adjudication and possible suspension of UI benefits for those who do not participate</p>	US DOL ETA	(1) UI benefit receipt; (2) Employment; and (3) Earnings
------	---------------------	--------------------------------------------------------	---------------------	-----------------------------------------------------	------	---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	------------	----------------------------------------------------------

Abbreviations: DOL = US Department of Labor; ETA = Employment and Training Administration; Fdn. = Foundation.

Sources: (1) Kemple, Friedlander, and Fellerath, 1995; (2) Hamilton, Freedman, Gennetian, Michalopoulos, and Walter, 2001; (3) Beecroft, Lee, Long, Holcomb, Thomson, Pindus, O'Brien, and Bernestin, 2003; (4) Freedman, Knab, Gennetian, and Navarro, 2000; (5) Steinman, 1978; (6) Johnson, Pfister, West, and Dickinson, 1984; Corson, Long, and Nicholson, 1984; (7) Herrem and Schmidt, 1983; Jagers, 1984 (8) Minnesota Department of Jobs and Training, 1990; (9) Black, Smith, Berger, and Noel, 2003; (10) Klepinger, Johnson, Joesch, and Benus, 1997; (11) Klerman, Minzner, Harkness, Mills, Cook, and Savidge-Wilkins, 2013.