

# Selective Trials: A Principal-Agent Approach to Randomized Controlled Experiments

By SYLVAIN CHASSANG, GERARD PADRÓ I MIQUEL, AND ERIK SNOWBERG\*

*We study the design of randomized controlled experiments when outcomes are significantly affected by experimental subjects' unobserved effort expenditure. While standard randomized controlled trials (RCTs) are internally consistent, the unobservability of effort compromises external validity. We approach trial design as a principal-agent problem and show that natural extensions of RCTs—which we call selective trials—can help improve external validity. In particular, selective trials can disentangle the effects of treatment, effort, and the interaction of treatment and effort. Moreover, they can help identify when treatment effects are affected by erroneous beliefs and inappropriate effort expenditure.*

*Keywords: randomized controlled trials, selective trials, blind trials, incentivized trials, marginal treatment effects, mechanism design, selection, heterogeneous beliefs, compliance.*

*JEL: C81, C93, D82, O12.*

This paper studies the design of experimental trials when outcomes depend significantly on unobserved effort decisions taken by subjects (agents).<sup>1</sup> Even in an ideal setting where the experimenter (principal) can randomly and independently assign an arbitrarily large number of agents to the treatment and control groups, unobserved effort limits the informativeness of randomized controlled trials (RCTs). For example, if a technology's measured returns are low, it is difficult to distinguish whether this is because true returns *are low* or because most agents *believe they are low* and therefore expend no effort using the technology. Moreover, to the extent that effort responds to beliefs, and beliefs respond to

\* Chassang: Princeton University, [chassang@princeton.edu](mailto:chassang@princeton.edu). Padró i Miquel: London School of Economics, [g.padro@lse.ac.uk](mailto:g.padro@lse.ac.uk). Snowberg: California Institute of Technology, [snowberg@caltech.edu](mailto:snowberg@caltech.edu). We are particularly indebted to Abhijit Banerjee, Roland Benabou, and Jeff Ely for advice and encouragement. The paper benefited greatly from conversations with Attila Ambrus, Nava Ashraf, Oriana Bandiera, Angus Deaton, Esther Dufo, Pascaline Dupas, Greg Fischer, Kripa Freitas, Drew Fudenberg, Paul Gertler, Justin Grimmer, Rema Hanna, Jim Heckman, Johannes Hörner, Dean Karlan, Michael Kremer, Guido Imbens, John Ledyard, Maggie McConnell, Stephen Morris, Muriel Niederle, Marcin Peski, Nancy Qian, Antonio Rangel, Imran Rasul, Dan Scharfstein, Sam Schulhofer-Wohl, Jesse Shapiro, Monica Singhal, Andy Skrzypacz, Francesco Sobbrío, Lars Stole, Steven Tadelis, Chris Woodruff and Eric Zitzewitz, as well as seminar participants at Berkeley Haas, Boston University, Brown, Caltech, Chicago Booth, Cornell, Harvard/MIT, HEC Lausanne, Johns Hopkins, LSE, MPSA, NYU Stern, Princeton, The Radcliffe Institute, Stanford, Stockholm School of Economics, SWET, UT Austin, Washington University in St. Louis, the World Bank and Yale. Part of this work was done while Chassang visited the Department of Economics at Harvard, and he gratefully acknowledges their hospitality. Paul Scott provided excellent research assistance.

<sup>1</sup>Throughout the paper we call experimental subjects agents, and call the experimenter the principal. Following usual conventions, we assume the principal is female, and agents are male.

information, this makes it difficult to predict the returns to the technology in the same population as it becomes better informed. In other words, unobserved effort is a source of heterogeneity in treatment effects, and is a significant challenge to the external validity of experimental trials.<sup>2</sup>

We propose simple extensions of RCTs—which we call selective trials—that improve the external validity of trial results without sacrificing robustness or internal validity. These experimental designs can be used to determine the extent to which erroneous beliefs or inappropriate effort affect measured treatment effects. We provide a systematic analysis of trial design using a principal-agent framework with both adverse selection—an agent’s type is unobserved—and moral hazard—an agent’s effort is unobserved. However, unlike the standard principal-agent framework, our principal’s goal is to maximize information about a technology’s returns—in the sense of Blackwell—rather than profits. The principal seeks to achieve this objective through single-agent mechanisms that assign agents to treatments of varying sophistication based on the message they send.

These mechanisms improve on RCTs for two reasons. First, they let agents express preferences over their treatment by probabilistically selecting themselves in or out of the treatment group at a cost—hence the name *selective trials*.<sup>3</sup> This makes implicit, unobserved selection an explicit part of the experimental design. Second, these mechanisms allow for treatments of varying richness: in open trials, treatment corresponds to access to the new technology; in blind trials, treatment corresponds to an undisclosed allotment of the technology, and information about the probability of having been allotted the technology; and in incentivized trials, treatment corresponds to access to the technology as well as an incentive, or insurance, contract based on outcomes.

Our results fall into two broad categories. Given a type of treatment (open, blind or incentivized) our first set of results characterize maximally informative mechanisms and examine the sampling patterns such mechanisms induce. We show that a mechanism is maximally informative if and only if it identifies an agent’s preferences over all possible treatment assignments and, given preferences, still assigns each agent to the treatment or control group with positive probability. Thus, our designs encapsulate the data generated by a standard randomized controlled trial. These designs can be implemented in a number of intuitive ways, such as a menu of lotteries, or utilizing the design of Becker et al. (1964), referred to as the BDM mechanism.

While our main focus is on identification, and thus infinite samples, selective trials may impose sampling costs on experimenters. In particular, sampling patterns do not matter when arbitrarily large samples are available, but affect statistical

<sup>2</sup>Unobserved effort is an issue whether a trial is open—agents know their treatment status—or blind—agents’ treatment status is obscured by giving the control group a placebo. See Duflo et al. (2008b) for a more detailed description of RCTs and the external validity issues frequently associated with them.

<sup>3</sup>For simplicity, we focus on monetary costs, but selection could also be based on non-monetary costs. For example, agents could choose between lines with different wait times to place themselves into the treatment group with different probabilities.

power in finite samples. In any mechanism that identifies agents' preferences in a strictly incentive-compatible way, agents with a higher value for the technology must be assigned to the treatment group with higher probability, which can reduce statistical power. However, these sampling costs can be reduced by diminishing incentives for the truthful reporting of preferences. This allows the experimenter to strike a balance between sampling costs and the precision of the preference data that is obtained. As detailed later, these results contribute to recent discussions about the usefulness of charging subjects for access to treatment in RCTs (see, for instance, Cohen and Dupas, 2010; Dupas, 2009; or Ashraf et al., 2010).

Our second class of results characterizes what can be inferred from selective trials, and highlights how they contribute to the ongoing discussion about the external validity of field experiments (Deaton, 2010; Imbens, 2010). By eliciting agents' value for the technology, open selective trials recover the distribution of returns as a function of willingness to pay. As a result, open trials provide a simple and robust way to recover the marginal treatment effects (MTEs) introduced by Heckman and Vytlacil (2005). Identifying MTEs is valuable as they can be used to forecast the effect of policies that change accessibility of the technology, such as subsidies. However, MTEs are typically not sufficient to make projections about interventions that alter beliefs and effort expenditure, such as informational campaigns.<sup>4</sup>

Selective trials go beyond MTEs and identify deep parameters by letting agents express preferences over richer treatments. Specifically, we consider blind trials where treatment status is hidden from agents by providing the control group with a placebo. This allows the principal to vary the information an agent has about his treatment status. This variation can be used to identify the pure effect of treatment and effort, the effect of their interaction, and agents' perceived returns to effort.<sup>5</sup> As blind trials are rarely used in economics—often for want of a convincing, ethical placebo—we extend the analysis to incentivized selective trials in which agents know their treatment status, but receive different transfers conditional on observable outcomes. Under mild assumptions, this produces

<sup>4</sup>In addition, selective trials may alleviate subversions of experimental protocol discussed in Deaton (2010). That is, explicitly allowing the agents to select themselves in and out of treatment may reduce the number of agents in the control group who obtain the treatment by other means, as well as the number of agents in the treatment group that refuse to be treated. Furthermore, the principal may use the information revealed by agents' preferences to increase monitoring of agents who expressed a high value for treatment but were assigned to the control group. Malani (2008) proposes a related solution: a trial design in which agents may select the nature of their control treatment, thus reducing incentives to subvert the experimental protocol.

<sup>5</sup>Although uncommon in economics, blind trials are quite common in medicine. For a brief review of RCTs in medicine see Stolberg et al. (2004). Jadad and Enkin (2007) provides a more comprehensive review. Selective trials nest a class of trial designs referred to as preference trials, in which at least one group of agents is allowed to choose their treatment. These designs have primarily been used in medicine to understand the ethics of randomized controlled trials and facilitate informed consent. Our work shows that eliciting preferences is not incompatible with randomization, and that preferences carry information that facilitates inference from treatment effects. For more on preference trials, see Zelen (1979); Flood et al. (1996); Silverman and Altman (1996); King et al. (2005); Jadad and Enkin (2007); Tilbrook (2008).

information similar to that produced by selective blind trials.

This paper contributes primarily to the literature on treatment effects. Most of this literature, based on a statistical framework quite different from our principal-agent approach, has focused on much simpler effort decisions and the ex post analysis of data. Agents are usually viewed as either taking treatment or not (with some exceptions: see Jin and Rubin, 2008, for a recent example), and more importantly, this decision is assumed to be observable, or sufficiently correlated with exogenous observable variables (Imbens and Angrist, 1994; Angrist et al., 1996; Heckman and Vytlacil, 2005). In contrast, we consider effort decisions which are unobservable and high dimensional. Additionally, most previous approaches, even those that rely—as we do—on decision theory, focus on modeling data from an RCT after it has been run (Philipson and Desimone, 1997; Philipson and Hedges, 1998).<sup>6</sup> We take an ex ante perspective and propose designs for experimental trials that can help understand how beliefs and effort affect treatment effects.

Successful implementation of the trial designs suggested by our principal-agent approach requires addressing a number of challenges. A practical limitation of our approach is that large samples may be needed to estimate all identifiable parameters. This limitation is inherent in any non-parametric estimation of treatment effects conditional on a large set of explanatory variables (see, for example, Pagan and Ullah, 1999). Another challenge is how to extract reliable preference data from agents. Mechanisms that are equivalent in theory, due to the assumption of rationality, may have very different properties in practice. Thus, experimenters may prefer to elicit coarser preference information in order to use simpler designs. We believe that these practical concerns are best resolved through a mix of laboratory and field experiments in well-understood environments. Therefore, it is encouraging that many elements of selective trials are already being evaluated in field settings (see, for example, Karlan and Zinman, 2009; Ashraf et al., 2010; Cohen and Dupas, 2010; Berry et al., 2011). A final set of challenges are more theoretical, and deal with extending our mechanisms to elicit richer information, such as the variation of preferences over time, or beliefs about other participants.

The paper is organized as follows. Section I uses a simple example to illustrate the main points of the paper. Section II defines the general framework. Section III investigates open selective trials. Section IV turns to blind selective trials, and shows how they can be used to identify true and perceived returns to effort. Section V analyzes incentivized trials, which eschew placebos, and shows that under reasonable assumptions they can be as informative as blind selective trials. Section VI concludes with a discussion of the limitations of, and future directions for, our approach to designing randomized controlled experiments.

<sup>6</sup>These studies use information correlated with agents' decisions to comply or not comply with their assigned treatments to refine understanding of treatment effects. This approach, as well as ours, is closely related to the classic Roy (1951) selection model in which selection into treatment reveals information about an agent's type (Heckman and Honoré, 1990; Heckman et al., 1997).

## I. An Example

To illustrate the basic insights underlying selective trials, and the potential usefulness of the data they generate, this section adopts a particularly simple model of the relationships between agents' beliefs, effort decisions, and outcomes. We emphasize that this structure makes inference particularly stark. Subsequent sections study inference in a much more general model that incorporates many important aspects of actual experiments.

To fix ideas, we discuss the example in terms of an experiment evaluating the health effects of a water treatment product.<sup>7</sup>

### A. A Simple Model

There are infinitely many agents indexed by  $i \in \mathbb{N}$ . Each agent has a treatment status  $\tau_i \in \{0, 1\}$ . If  $\tau_i = 1$  agent  $i$  is in the treatment group, and is given the water treatment product. Otherwise,  $\tau_i = 0$ , and the agent is in the control group.

Agent  $i$  obtains a final outcome  $y_i \in \{0, 1\}$  that can be measured by the principal. In our example  $y_i = 1$  indicates that the agent has remained healthy. The probability that an agent remains healthy depends on both treatment and effort

$$\text{Prob}(y_i = 1 | e_i, \tau_i) = q_0 + R e_i \tau_i,$$

where  $e_i \in [0, 1]$  is agent  $i$ 's decision of whether or not to expend effort using the product,  $R \in [R_L, R_H]$  is the technology's return, which is common to all agents, and  $q_0$  is the unknown baseline likelihood of staying healthy over the study period, which will be controlled for using randomization. Agents have different types  $t$  that characterize their beliefs about returns  $R$ . We denote by  $R_t = \mathbb{E}_t R$  the returns expected by an agent of type  $t$ . The distribution  $F_{R_t}$  of expectations  $R_t$  in the population need not be known to the principal or the agents.<sup>8</sup>

We assume throughout that effort is private and cannot be monitored by the principal. In other words, we assume that all observable dimensions of effort are already controlled for, and focus on those dimensions that are not observable. For example, with a water treatment product, an experimenter may be able to

<sup>7</sup>It should be noted that while our main focus is on the use of RCTs in medical, public health, and development contexts, our analysis applies to most environments involving decentralized experimentation. For instance, if a firm wants to try a new way to organize production, specific plant managers will have to decide how much effort to expend implementing it. The firm's CEO is in the same position as the principal in our framework, and must guess the effort expended by his managers when evaluating returns to the new production scheme. Similarly, if a school board wants to experiment with a new program, individual teachers and administrators will have to decide how much effort to expend on implementing the program.

<sup>8</sup>For illustrative purposes, this example focuses on heterogenous beliefs as a source of heterogenous behavior and returns. In this setting, convincingly identifying true returns to treatment has a large effect on behavior, and would be particularly valuable. Moreover, the example allows effort to affect outcomes only in the treatment group. The general framework, described in Section II, allows for: general, idiosyncratic, returns; effort in both the treatment and control group; and effort along an arbitrary number of dimensions, which can accommodate dynamic effort expenditure.

determine whether or not the agent has treated water in his home, but it may be much more difficult to determine if the agent drinks treated water when away from home.<sup>9</sup>

Given effort  $e_i$ , agent  $i$ 's expected utility is given by

$$\mathbb{E}_t[y_i|e_i] - ce_i,$$

where  $c \in (R_L, R_H)$  is the agents' cost of effort. In our example, this may be the cost of remembering to use the product, the social cost of refusing untreated water, or disliking the taste of treated water. In addition, we assume each agent has quasilinear preferences with respect to money. An agent's willingness to pay for treatment is  $V_t = \max\{R_t - c, 0\}$ , which we assume is less than some value  $V_{\max}$  for all agents.

We focus initially on open trials, in which agents know their treatment status, and contrast two ways of running trials: a standard RCT, where agents are randomly assigned to the treatment group with probability  $\pi$ , and a selective open trial that allows agents to express preferences for treatment by selecting their probability of treatment.

The implementation of selective trials we explore here uses the BDM mechanism, and proceeds as follows:

- 1) Each agent sends a message  $m_i \in [0, V_{\max}]$  indicating his willingness to pay for treatment;
- 2) A price  $p_i$  to obtain treatment is independently drawn for each agent from a distribution with convex support and c.d.f.  $F_p$  that satisfies  $0 < F_p(0) < F_p(V_{\max}) < 1$ ; and
- 3) If  $m_i \geq p_i$ , the agent obtains the treatment at price  $p$ , otherwise, the agent is in the control group and no transfers are made.

Note that a higher message  $m$  increases an agent's probability of treatment,  $F_p(m)$ , as well as his expected payment:  $\int_{p \leq m} p dF_p$ . As  $F_p$  has convex support, it is dominant for an agent of type  $t$  to send message  $m = V_t$ .

### B. The Limits of RCTs and the Value of Self-Selection

#### INFERENCE FROM RANDOMIZED CONTROLLED TRIALS

We begin by considering the information produced by an RCT. If agent  $i$  is in the treatment group, he chooses to expend effort  $e = 1$  if and only if  $R_t \geq c$ .

<sup>9</sup>Still, as Duflo et al. (2010) shows, innovative monitoring technologies may be quite effective. To the extent that monitoring is feasible, it should be done.

Hence, the average treatment effect identified by an RCT is

$$\begin{aligned}\Delta^{RCT} &= \mathbb{E}[y|\tau = 1] - \mathbb{E}[y|\tau = 0] \\ &= \mathbb{E}[q_0 + R \times \mathbf{1}_{R_t \geq c} | \tau = 1] - \mathbb{E}[q_0 | \tau = 0] \\ &= R \times \text{Prob}(R_t > c) = R \times (1 - F_{R_t}(c)).\end{aligned}$$

When the distribution of agents' expectations  $F_{R_t}$  is known, then an RCT will identify  $R$ . However, in most cases  $F_{R_t}$  is not known, and the average treatment effect  $\Delta^{RCT}$  provides a garbled signal of the underlying returns  $R$ . If the outcomes in the treatment group are only weakly better than those in the control group, the principal does not know if this is because the water treatment product is not particularly useful, or because the agents did not expend sufficient effort using it.

#### INFERENCE FROM OPEN SELECTIVE TRIALS

We now turn to selective trials and show they are more informative than RCTs.

The selective trial described above elicits agents' willingness to pay and, conditional on a given willingness to pay  $V$ , generates non-empty treatment and control groups. As it is dominant for agents to truthfully reveal their value, an agent with value  $V_t$  has probability  $F_p(V_t)$  of being in the treatment group and probability  $1 - F_p(V_t)$  of being in the control group. Both of these quantities are strictly positive as  $0 < F_p(0) < F_p(V_{\max}) < 1$ .<sup>10</sup>

This trial provides us with the set of local instruments needed by Heckman and Vytlacil (2005) to estimate marginal treatment effects (MTEs). That is, for any willingness to pay  $V$ , we are able to estimate

$$\begin{aligned}\Delta^{MTE}(V) &\equiv \mathbb{E}[y|\tau = 1, V_t = V] - \mathbb{E}[y|\tau = 0, V_t = V] \\ &= \mathbb{E}[y|\tau = 1, m_t = V] - \mathbb{E}[y|\tau = 0, m_t = V],\end{aligned}$$

which can be used to perform policy simulations in which the distribution of types is constant but access to the technology is changed, for example, by subsidies. Moreover, MTEs can be integrated to recover the average treatment effect identified by an RCT.

In the current environment, willingness to pay is a good signal of future use, and thus MTEs can be used to identify the true returns  $R$ . Specifically, all agents with value  $V_t > 0$  have expectations  $R_t$  such that  $R_t - c > 0$ , and expend effort

<sup>10</sup>Note also that agents with higher value are treated with higher probability. This matters for the precision of estimates in actual experiments, where sample size is not infinite. We return to this point in Section III.

$e = 1$  using the technology.<sup>11</sup> Hence, it follows that

$$\begin{aligned}\Delta^{MTE}(V > 0) &= \mathbb{E}[q_0 + R \times e_t | \tau = 1, V_t > 0] - \mathbb{E}[q_0 | \tau = 0, V_t > 0] \\ &= R.\end{aligned}$$

A selective trial identifies the average treatment effect, MTEs, and true returns  $R$ . Hence, it is more informative than an RCT, which only identifies the average treatment effect.

The true returns  $R$ , and the distribution of valuations  $V_t$ , have several policy uses. First, knowing  $R$  allows us to simulate the treatment effect for a population in which all agents expend the appropriate amount of effort. Second, these variables allow us to estimate the returns to increasing usage within a given population. Third, and finally, the data provided by a selective trial can be used to inform agents and disrupt learning traps more effectively than data from an RCT. For example, imagine that true returns to the technology are high, but most agents believe they are low. In that case, an RCT will measure low returns to the treatment and will not convince agents that they should be expending more effort. In contrast, the data generated by a selective trial would identify that true returns are high, and lead agents to efficiently adopt the water treatment product.<sup>12</sup>

### C. Richer Treatments

In the previous subsection, a selective trial identified true returns because willingness to pay was a good predictor of usage. However, as our continuing example shows, this will not always be the case. Thus, MTEs are generally not sufficient to infer true returns, nor whether beliefs are affecting measured treatment effects. However, more sophisticated selective trials, such as blind selective trials or incentivized selective trials, can be used to recover true returns.

We modify the example so that the returns  $R$  to the technology include both baseline returns and returns to effort:  $R = (R_b, R_e) \in \mathbb{R}^2$ . In the context of a water treatment product,  $R_b$  could be the baseline returns to using the product only when it is convenient to do so, and  $R_e$  the additional returns to using it more thoroughly (for example, bringing treated water when away from home). Success

<sup>11</sup>In this environment, the same result can be obtained by charging a price  $p$  for a probability of treatment  $\pi$  such that  $F_{R_t}(\frac{p}{\pi} - c) > 0$ , and evaluating treatment effects only for those willing to pay. The idea that higher prices will select individuals who value the technology more, and use it more intensely, can be traced back to the seminal selection model of Roy (1951). See Oster (1995) for a discussion of related ideas in the context of non-profit organizations.

<sup>12</sup>For empirical work in development economics on the effect of information on behavior, see Thornton (2008), Nguyen (2009) or Dupas (2011). For theoretical work on failures of social learning, see the classic models of Banerjee (1992) or Bikhchandani et al. (1992).



rates given effort and treatment status are:

$$\begin{aligned}\text{Prob}(y = 1|\tau = 0, e) &= q_0 \\ \text{Prob}(y = 1|\tau = 1, e) &= q_0 + R_b + eR_e.\end{aligned}$$

An agent of type  $t$  has expectation  $(R_{b,t}, R_{e,t})$  about returns  $R = (R_b, R_e)$ , and expends effort if and only if  $R_{e,t} \geq c$ . Therefore, an agent's willingness to pay for treatment is given by  $V_t = R_{b,t} + \max\{R_{e,t} - c, 0\}$ .

#### INFERENCE FROM OPEN SELECTIVE TRIALS

We have already shown that open selective trials can identify treatment effects conditional on willingness to pay. However, in the current environment, willingness to pay is no longer a good signal of effort. Indeed, there are now two reasons why an agent might value treatment: he believes that a thorough use of the product has high returns ( $R_{e,t}$  is high)—the channel emphasized in Section I.B—or he believes that a casual use of the water treatment product is sufficient to obtain high returns and that thorough use brings little additional return ( $R_{b,t}$  is high, but  $R_{e,t}$  is low). That is, agents who are willing to pay because they think baseline returns are high need not be the agents who will actually expend effort. Formally, a selective trial still identifies MTEs,

$$\Delta^{MTE}(V) = R_b + R_e \text{Prob}(R_{e,t} \geq c | R_{b,t} + \max\{R_{e,t} - c, 0\} = V),$$

but these are generally not sufficient to recover  $R_b$  and  $R_e$ .<sup>13</sup> As a result, MTEs are insufficient to simulate the returns in a population of agents that all expended appropriate effort, or more generally, the returns from increasing the effort of agents. Nor do MTEs provide the information needed to infer true returns.

#### BLIND SELECTIVE TRIALS

In a blind trial, an agent does not know his treatment status  $\tau \in \{0, 1\}$  at the time of effort, but rather, knows his probability  $\phi \in [0, 1]$  of having been assigned to the treatment group. Open trials are blind trials where  $\phi$  is either 0 or 1.

Given a probability  $\phi$  of being treated, an agent expends effort if and only if  $\phi R_{e,t} - c > 0$ . An agent's expected value for being treated with probability  $\phi$  is

$$V_t(\phi) = \phi R_{b,t} + \max\{\phi R_{e,t} - c, 0\}.$$

We depart from standard blind trials in a simple but fundamental way: while they keep  $\phi$  fixed and do not infer anything from the specific value of  $\phi$  used, we

<sup>13</sup>For instance, it is not possible to distinguish a situation in which returns to effort are equal to  $R_e$  and a proportion  $\eta V$  of agents with value  $V$  expends effort, from a situation in which returns to effort are  $2R_e$  and a proportion  $\frac{\eta}{2} V$  of agents with value  $V$  expends effort.

allow  $\phi$  to vary and use both willingness to pay for, and outcomes at, different values of  $\phi$  for inference.<sup>14</sup>

As with open trials, willingness to pay can be elicited using a BDM-type mechanism. Since willingness to pay  $V_t(\phi)$  now depends on  $\phi$ , the mechanism in Section I.A is implemented after the agent is asked to send a message  $m(\phi)$  for each possible value of  $\phi$ . A value of  $\phi_i$  is then drawn independently for each agent from a c.d.f.  $F_\phi$ , with full support on  $[0, 1]$  and mass points at 0 and 1. Transfer  $p_i$  is independently drawn from a c.d.f.  $F_p$ , as before. If  $m(\phi_i) \geq p_i$ , the agent pays  $p_i$  and is allotted the treatment with probability  $\phi_i$ ; otherwise, the agent is in the control group and no transfers are made.

A first advantage of blind trials is that, unlike open trials, an agent's actual treatment status  $\tau$  and his belief  $\phi$  about his treatment status can be different. This allows for a robust identification of baseline returns  $R_b$ . If an agent is assigned a probability of treatment  $0 < \phi < 1$  low enough that  $\phi R_H < c$ , he will not expend any effort. Still, a proportion  $\phi > 0$  of such agents receive treatment while a proportion  $1 - \phi > 0$  do not. Hence, we can identify  $R_b$  by measuring the effect of treatment for agents known not to expend effort:

$$R_b = \mathbb{E} \left[ y \mid \phi < \frac{c}{R_H}, \tau = 1 \right] - \mathbb{E} \left[ y \mid \phi < \frac{c}{R_H}, \tau = 0 \right].$$

A second advantage of blind trials is that the agents' value mapping  $V_t(\phi)$  allows identification of which agents would expend effort when treated for sure. The amount that an agent with belief  $\phi = 1/2$  is willing to pay to learn his treatment status is

$$\theta_t \equiv \frac{1}{2} [V_t(\phi=1) + V_t(\phi=0)] - V_t(\phi=1/2).$$

If the agent does not intend to expend effort when treated for sure, he will not value information, and  $\theta_t$  will be equal to 0. Inversely, if the agent does intend to expend effort, information is valuable since it allows him to tailor his behavior to his treatment status, and thus  $\theta_t > 0$ .<sup>15</sup> In the current example, provided that a positive measure of agents satisfies  $\theta_t > 0$ , we can identify  $R_e$  using either of the

<sup>14</sup>A similar insight comes from Malani (2006), which identifies placebo effects by examining variation in outcomes associated with variations in the probability of treatment across blinded experiments.

<sup>15</sup>This result holds very generally—see Proposition 5. To verify this relationship in the current example, note that if the agent expends effort conditional on being treated for sure (that is,  $R_{e,t} > c$ ), then

$$\theta_t = \frac{1}{2} [R_{b,t} + R_{e,t} - c] - \frac{1}{2} R_{b,t} - \max \left\{ \frac{1}{2} R_{e,t} - c, 0 \right\} \geq \min \left\{ \frac{R_{e,t} - c}{2}, \frac{c}{2} \right\} > 0.$$

following expressions:

$$\begin{aligned} R_e &= \mathbb{E}[y|\phi=1, \theta_t > 0, \tau=1] - \mathbb{E}[y|\phi=1, \theta_t=0, \tau=1] \\ &= \mathbb{E}[y|\phi=1, \theta_t > 0, \tau=1] - \mathbb{E}\left[y \mid \phi < \frac{c}{R_H}, \theta_t > 0, \tau=1\right]. \end{aligned}$$

#### INCENTIVIZED SELECTIVE TRIALS

We now show that incentivized selective trials can provide the principal with information similar to that produced by blind selective trials. This is useful as in many areas of economic interest, blind trials are not practical due to the lack of suitable, or ethical, placebos.

In an incentivized selective trial, an agent obtains a treatment status  $\tau \in \{0, 1\}$ , makes a fixed transfer  $p$  (which can be positive or negative), and is given a bonus (or penalty)  $w$  in the event that  $y = 1$ . Note that if  $p > 0$  and  $w > 0$ , then an agent is assigned an incentive contract. If, instead,  $p < 0$  and  $w < 0$ , an agent is assigned an insurance contract.

Given a bonus level  $w$ , an agent expends effort if and only if  $(1+w)R_{e,t} - c > 0$ . In turn, an agent's willingness to pay for treatment, given bonus  $w$ , is

$$V_t(w) = (1+w)R_{b,t} + \max\{(1+w)R_{e,t} - c, 0\}.$$

As before, the mapping  $w \mapsto V_t(w)$  can be elicited using a variant of the BDM mechanism (described in Appendix B). Incentivized trials allow us to evaluate baseline returns in a straightforward manner. When offered a full insurance contract  $w_0 = -1$ , an agent will expend effort  $e = 0$  so that

$$R_b = \mathbb{E}[y|w=w_0, \tau=1] - \mathbb{E}[y|w=w_0, \tau=0].$$

In turn, notice that for any type  $t$  with  $R_{e,t} > 0$ , there exists a value  $w_t$  such that whenever  $w > w_t$ , the agent expends effort  $e = 1$ . The value  $w_t$  is identified from the mapping  $w \mapsto V_t(w)$  as

$$\left. \frac{\partial V_t}{\partial w} \right|_{w > w_t} = R_{e,t} + R_{b,t} > R_{b,t} = \left. \frac{\partial V_t}{\partial w} \right|_{w < w_t}.$$

Additionally, this last expression allows us to identify the agent's subjective beliefs about baseline returns and returns to effort ( $R_{b,t}, R_{e,t}$ ). For a value  $\bar{w}$  sufficiently high that it induces some agents to expend effort, returns to effort can be identified by either of the following expressions

$$\begin{aligned} R_e &= \mathbb{E}[y|w=\bar{w}, \bar{w} - w_t > 0, \tau=1] - \mathbb{E}[y|w=\bar{w}, \bar{w} - w_t < 0, \tau=1] \\ &= \mathbb{E}[y|w=\bar{w}, \bar{w} - w_t > 0, \tau=1] - \mathbb{E}[y|w=w_0, \bar{w} - w_t > 0, \tau=1]. \end{aligned}$$

Thus, just like blind selective trials, incentivized selective trials identify true re-

turns  $R = (R_b, R_e)$ .

Altogether, this section suggests that while unobserved effort is an issue for the external validity of standard randomized controlled trials, appropriate ex ante trial design—rather than ex post data treatment—may help in alleviating these concerns.

The rest of the paper explores how these results extend in a general framework that allows for many realistic elements. In particular, this general framework allows for arbitrary heterogeneity among agents, including heterogeneous preferences, beliefs, and returns. Moreover, the general framework allows for multidimensional effort in both the treatment and control group. This allows the model to accommodate complex technologies, dynamic effort expenditure, and attempts by agents in the control group to obtain substitute treatments.

The following sections provide systematic results about which mechanisms are the most informative, what sampling patterns they produce, and what can be inferred from the data they generate.

## II. A General Framework

We now generalize the framework used in our example. Once again, there are infinitely many agents, indexed by  $i \in \mathbb{N}$ . However, returns to the technology are now described by a multidimensional parameter  $R \in \mathcal{R} \subset \mathbb{R}^\kappa$ .

### TYPES

Each agent  $i$  has a type  $t \in \mathcal{T}$ , which includes a belief about returns  $R$ , as well as other factors that might affect behavior and outcomes, such as idiosyncratic costs of effort, idiosyncratic returns, and beliefs about such factors. We assume that agents are exchangeable, so that their types are i.i.d. draws from some distribution  $\chi \in \Delta(\mathcal{T})$ , which is itself a random variable. A profile of types is given by  $\mathbf{t} \in \mathcal{T}^\mathbb{N}$ . For conciseness we omit publicly observable traits, but it is straightforward to allow for them.

### OUTCOMES AND SUCCESS RATES

Agent  $i$  obtains an outcome  $y_i \in \{0, 1\}$ .<sup>16</sup> An agent's true and perceived likelihoods of success (that is,  $\text{Prob}(y = 1)$ ) depend on his type, the aggregate returns to the technology and the agent's effort choice  $e \in E$ , where  $E$  is a compact subset of  $\mathbb{R}^{\kappa'}$ .

Success rates are denoted by

$$\begin{aligned} q(R, t, \tau_i, e_i) &= \text{Prob}(y = 1 | R, t, \tau_i, e_i) \\ q_t(\tau_i, e_i) &= \int_{\mathcal{R}} q(R, t, \tau_i, e_i) dt(R) \end{aligned}$$

<sup>16</sup>As Appendix A shows, binary outcomes simplify notation, but are not essential to our results.

where  $q(R, t, \tau, e)$  is the true success rate of an agent of type  $t$ , which allows for idiosyncratic, heterogeneous returns, while  $q_t(\tau, e)$  is the probability of success perceived by an agent of type  $t$ , which allows for idiosyncratic, heterogeneous beliefs about those returns. We assume that  $q$  and  $q_t$  are continuous with respect to effort  $e$ . Note that as  $e$  can be multidimensional, the model is consistent with dynamic effort expenditure, and agents learning about returns to treatment, or their treatment status (as in Philipson and Desimone, 1997; or Chan and Hamilton, 2006).<sup>17</sup>

#### PREFERENCES

Given effort  $e_i$ , treatment status  $\tau_i$ , monetary transfer  $p_i$ , and final outcome  $y_i$ , agent  $i$ 's utility is  $u(y_i, t_i) - c(e_i, t_i) - p_i$ .

Note that  $p_i$  can be negative and all transfers can be shifted by a fixed amount, for example, when there is compensation for participating in the experiment. Such compensation may be used to increase participation, or relax agents' cash constraints.<sup>18</sup>

#### ASSIGNMENT MECHANISMS

We distinguish three ways of assigning treatment:

- 1) *Open selective trials* are mechanisms  $G_o = (M_o, \mu_o)$  where  $M_o$  is a set of messages and  $\mu_o : M_o \rightarrow \Delta(\{0, 1\} \times \mathbb{R})$  maps individual messages to a probability distribution over treatment status  $\tau_i \in \{0, 1\}$  and transfers  $p_i \in \mathbb{R}$ ;
- 2) *Blind selective trials* are mechanisms  $G_b = (M_b, \mu_b)$  where  $M_b$  is a set of messages and  $\mu_b : M_b \rightarrow \Delta([0, 1] \times \mathbb{R})$  maps messages to a probability distribution over uncertain treatment status  $\phi_i = \text{Prob}(\tau_i = 1)$  and transfers  $p_i$ ; and
- 3) *Incentivized selective trials* are mechanisms  $G_w = (M_w, \mu_w)$  where  $M_w$  is a set of messages and  $\mu_w : M_w \rightarrow \Delta(\{0, 1\} \times \mathbb{R} \times \mathbb{R})$  maps messages to a probability distribution over treatment status  $\tau_i$ , a fixed transfer  $p_i$  from the agent to the principal, and a bonus  $w_i$  transferred from the principal to the agent conditional on  $y_i = 1$ .

<sup>17</sup> For example, it is not enough for agents to just expend effort spreading fertilizer. As Duflo et al. (2008a) highlights, effort is needed to choose the appropriate seeds to go with the fertilizer, to learn how much and when to water the crops, and to learn how much fertilizer gives the highest returns at the lowest cost. In this case, it is natural to think of effort as a vector, where the first component corresponds to choosing the amount of fertilizer, the second to picking the right seeds, the third to properly applying it, and so on.

To accommodate dynamic effort expenditure, different dimensions of the effort vector may indicate contingent effort expenditure depending on realized observables, such as the state of crops, or how they seem to respond to previous fertilizer use.

<sup>18</sup> Appendix A allows for agents with non-quasilinear preferences and thus tradeoffs between treatment and non-monetary costs.

Note that these are single agent mechanisms. Agent  $i$ 's final assignment depends only on his message, and not on messages sent by others. We denote the likelihood of being given the treatment when sending message  $m$  by  $\pi(m) \equiv \text{Prob}(\tau = 1|m)$ . We focus largely on mechanisms  $G$  such that  $\chi$ -almost surely every agent  $i$  has a dominant message  $m_G(t_i)$ . In all these designs, agents can probabilistically select their assignment using messages, hence the name *selective trials*.

#### INFORMATIVENESS OF MECHANISMS

We evaluate mechanisms according to their informativeness in the sense of Blackwell. We say a mechanism  $G$  is at least as informative as a mechanism  $G'$ , denoted by  $G' \preceq G$ , if the data generated by  $G'$  can be simulated using only data generated by  $G$ .

Specifically, denote by  $a_i$  the assignment given to agent  $i$  by whichever mechanism is chosen. The principal observes data  $\mathbf{d}_G = (m_i, a_i, y_i)_{i \in \mathbb{N}}$ . Denote by  $\mathcal{D}_G$  the set of possible data sequences generated by mechanism  $G$ . Mechanism  $G' \preceq G$  if and only if there exists a fixed data manipulation procedure  $h : \mathcal{D}_G \rightarrow \Delta(\mathcal{D}_{G'})$  such that for all  $\mathbf{t} \in \mathcal{T}^{\mathbb{N}}, R \in \mathcal{R}$ ,  $h(\mathbf{d}_G(\mathbf{t}, R)) \sim \mathbf{d}_{G'}(\mathbf{t}, R)$ .

This notion of informativeness is easier to work with in environments with infinite samples, as this focuses on issues of identification rather than issues of statistical power. However, this definition also applies in the case of finitely many agents.<sup>19</sup>

Although our framework is quite general, we intentionally limit our approach in three ways. First, we assume agents are rational, that is, they play undominated strategies, regardless of the complexity of the assignment mechanism. Second, we examine only single-agent mechanisms. Third, despite the fact that effort expenditure may be dynamic, we restrict attention to mechanisms that elicit preferences only once. Note, however, that the timing of this elicitation may be freely chosen by the principal. Specifically, messages could be elicited before agents have any exposure to the technology, or after they have assessed it. Section VI discusses the limitations of assuming rationality and examining only single-agent mechanisms, and the difficulties of eliciting preferences more than once.

### III. Open Selective Trials

In open selective trials an agent is assigned a treatment status  $\tau$  and a transfer  $p$  based on message  $m$ . Given this assignment  $(\tau, p)$ , the indirect utility of an agent with type  $t$  is  $V_t(\tau) - p$  where,

$$V_t(\tau) = \max_{e \in E} q_t(\tau, e)u(y=1, t) + [1 - q_t(\tau, e)]u(y=0, t) - c(e, t).$$

<sup>19</sup>With infinite samples, sampling patterns do not matter. Thus, there is a large equivalence class of most informative mechanisms. When samples are finite, these mechanisms remain undominated in the sense of Blackwell, but need no longer be equivalent.

We normalize the value of being in the control group  $V_t(\tau=0)$  to zero for every type. Hence  $V_t \equiv V_t(\tau=1)$  denotes the agent's willingness to pay for treatment. For simplicity, we assume that there exists a known value  $V_{\max} \in \mathbb{R} > 0$  such that for all  $t \in \mathcal{T}$ ,  $V_t \in (-V_{\max}, V_{\max})$ , and that the distribution over values induced by the distribution of types  $\chi$  admits a density. The optimal effort for type  $t$  given treatment status  $\tau$  is denoted by  $e^*(\tau, t)$ .<sup>20</sup>

#### A. Information Production in Open Selective Trials

Our first result highlights the fact that selective trials are natural extensions of RCTs. An RCT is a mechanism  $G_0 = (\emptyset, \pi_0)$ . As  $M = \emptyset$ , no messages are sent, all agents are assigned to the treatment group with the same probability  $\pi_0 \in (0, 1)$ , and there are no transfers.

**FACT 1 (full support sampling):** *Consider a mechanism  $G = (M, \mu)$ . If there exists  $\xi > 0$  such that for all  $m \in M$ ,  $\pi(m) \in (\xi, 1 - \xi)$ , then, with infinite samples,  $G_0 \preceq G$ .*

**PROOF:**

All proofs can be found in the online appendix.

Recalling that  $\pi(m) \equiv \text{Prob}(\tau = 1|m)$ , Fact 1 shows that if every type has a positive probability of being in the treatment or control group, then mechanism  $G$  is as informative as an RCT. This holds for any  $\xi > 0$  because the sample size is infinite. The assumption of infinite samples—which is maintained throughout the paper—is important for all of our identification results. We discuss sampling issues that arise with finite samples in Section III.B.

As Plott and Zeiler (2005) and others show, information elicited in non-incentive-compatible ways can be unreliable. Moreover, as Kremer and Miguel (2007) and others note, reported beliefs about a technology's return are often uncorrelated with use. Therefore, we focus on *strictly incentive-compatible* assignment mechanisms—assignment mechanisms such that  $\chi$ -almost every agent has a strictly preferred message.<sup>21</sup>

Our next result shows that an open selective trial is a most informative trial if and only if it identifies each agent's value  $V_t$ , and, conditional on any expressed valuation, assigns a positive mass of agents to both the treatment and control group.

**PROPOSITION 1 (most informative mechanisms):** *Any strictly incentive-compatible mechanism  $G$  identifies at most value  $V_t$  (that is,  $V_t = V_{t'} \Rightarrow m_G(t) = m_G(t')$ ).*

<sup>20</sup>At this stage, whether optimal effort is unique or not does not matter. We explicitly assume a unique optimal effort level in Sections IV and V to apply a convenient version of the Envelope Theorem.

<sup>21</sup>Note that the mechanisms we consider can accommodate surveys. Consider the mechanism  $G = (\mathcal{T}, \pi_0)$  where the message space  $M = \mathcal{T}$ , the likelihood of treatment is constant and equal to  $\pi_0$ , and no transfers are made. This is essentially an RCT supplemented with a rich survey. As assignment does not depend on the message, truthful revelation of one's type is weakly dominant. Unfortunately, any other message is also weakly dominant. Hence, data generated by such a mechanism is likely to be unreliable, especially if figuring out one's preferences is costly.

Whenever  $G$  identifies values  $V_t$  (that is,  $m_G(t) = m_G(t') \Rightarrow V_t = V_{t'}$ ) and satisfies full support ( $0 < \inf_m \pi(m)$  and  $\sup_m \pi(m) < 1$ ), then for any strictly incentive-compatible mechanism  $G'$ ,  $G' \preceq G$ .

It follows that open selective trials can identify at most the distribution of returns conditional on agents' valuations, which can be used to construct marginal treatment effects (MTEs). It is important to note that these mechanisms identify MTEs independently of the experimenter's beliefs. Hence, to the extent that elicited values are reliable, these mechanisms identify MTEs with a degree of robustness comparable to that with which RCTs identify average treatment effects.<sup>22</sup>

#### IMPLEMENTING MOST INFORMATIVE TRIALS

Here we exhibit two straightforward implementations of most informative selective trials. The first is the BDM mechanism described in Section I.A, with the expanded message space  $M = [-V_{\max}, V_{\max}]$ . Once again, the principal draws a price  $p_i \in [-V_{\max}, V_{\max}]$  independently for each agent from a common c.d.f.  $F_p$  with support  $[-V_{\max}, V_{\max}]$ . If  $m_i \geq p_i$ , then the agent is assigned  $(\tau = 1, p_i)$ ; otherwise, he is assigned  $(\tau = 0, 0)$ .

**FACT 2 (BDM Implementation):** *Whenever  $F_p$  has full support over  $[-V_{\max}, V_{\max}]$ , an agent with value  $V_t$  sends optimal message  $m_{BDM} = V_t$  and the BDM mechanism is a most informative mechanism.*

A second implementation is a menu of lotteries. Consider mechanism  $G^*$ , where  $M = (-\frac{1}{2}, \frac{1}{2})$ , any agent sending message  $m$  is assigned to the treatment group with probability  $\pi(m) = \frac{1}{2} + m$ , and must make a transfer  $p(m) = V_{\max}m^2$ . One can think of agents as having a baseline probability of being in the treatment group equal to  $\frac{1}{2}$  and deciding by how much they want to deviate from this baseline. An agent with value  $V_t$  chooses message  $m$  to maximize

$$\pi(m)V_t - p(m) = V_t \left( \frac{1}{2} + m \right) - V_{\max}m^2.$$

This problem is concave in  $m$ , and first order conditions yield an optimal message  $V_t/2V_{\max}$ , which identifies  $V_t$ . In addition, every agent is assigned to the treatment and control group with positive probability. Thus  $G^*$  is a most informative mechanism.

Note that  $G^*$  gives agents higher expected utility than an RCT that assigns agents to the treatment and control group with probability  $\frac{1}{2}$ . More generally, for any RCT, a selective trial that assigns price  $p = 0$  for a probability of treatment  $\pi$  equal to that of the RCT must increase the agents' expected utility. Thus,

<sup>22</sup>Note that selective trials also identify higher order moments of the outcome distribution conditional on treatment status and willingness to pay, which may be useful to researchers.



selective trials may help decrease the number of agents who refuse randomization. This is potentially useful as refusals reduce the external validity of treatment effects (Malani, 2008).<sup>23</sup>

### B. The Cost of Running Selective Trials

In equilibrium, the menu of lotteries  $G^*$  yields sampling profile  $\pi(V) = \frac{1}{2} \left(1 + \frac{V}{V_{\max}}\right)$ , which is strictly increasing in value  $V$ . In the BDM mechanism, the sampling profile  $\pi_{BDM}(V) = F_p(V)$  is also increasing in  $V$ . This property holds for any mechanism.

**PROPOSITION 2 (monotonicity):** *Consider a strictly incentive compatible mechanism  $G$ . If agents  $t$  and  $t'$  with values  $V_t > V_{t'}$  send messages  $m_G(t) \neq m_G(t')$ , then it must be that  $\pi(m_G(t)) > \pi(m_G(t'))$ .*

Thus, in any selective trial, agents with high values are over-sampled—they have a higher likelihood of being in the treatment group—and those with low values are under-sampled. In contrast, RCTs have a flat sampling profile. While sampling patterns do not matter when there is an arbitrarily large number of agents, they can significantly affect statistical power when samples are finite.

This issue is related to the recent debate in development economics about charging for treatment in RCTs.<sup>24</sup> If, as in Ashraf et al. (2010), willingness to pay is correlated with product usage, then eliciting willingness to pay might be quite useful in understanding true returns. If, instead, as in Cohen and Dupas (2010), most agents have low values, and willingness to pay is a poor predictor of actual use, then undersampling agents with low values may significantly reduce statistical power. Furthermore, in such a setting, willingness to pay provides little information about intended use.<sup>25</sup>

We make two contributions to this debate. First, we note that when trade-offs between money and treatment are uninformative, selective trials can and should be based on more informative trade-offs. For instance, if most of the heterogeneity in willingness to pay is driven by wealth and credit constraints, then eliciting willingness to wait, or willingness to perform a tedious task—like sitting through multiple information sessions—may be a better indicator of future usage than willingness to pay. As we discuss in Section VI, this requires some knowledge of the agents and their environment.

Second, we show that carefully designed selective trials can reduce the costs of oversampling agents with high values by reducing the slope of the sampling profile.

<sup>23</sup>Jadad and Enkin (2007) reports refusal rates approaching 50 percent in some medical trials.

<sup>24</sup>This literature is motivated by questions of efficiency, and is mostly interested in whether charging for usage improves how well treatment is matched with those who need and use it. This paper takes a slightly different perspective, and is interested in how controlling for willingness to pay improves inference from experimental trials.

<sup>25</sup>As Dupas (2010) shows, this can also hinder social learning.

PROPOSITION 3 (sampling rates and incentives): *For any mechanism  $G = (M, \mu)$  and  $\underline{\rho} < \bar{\rho}$  in  $(0, 1)$ , there exists a mechanism  $G' = (M, \mu')$  such that  $G \preceq G'$ , and for all  $m \in M$ ,  $\pi'(m) \in [\underline{\rho}, \bar{\rho}]$ .*

*The following must also hold. Denoting the expected utility of type  $t$  sending message  $m$  in mechanism  $G'$  (including transfers) by  $U(t|m, G')$ , then*

$$\max_{m_1, m_2 \in M} |U(t|m_1, G') - U(t|m_2, G')| \leq 2(\bar{\rho} - \underline{\rho})V_{\max}.$$

Proposition 3 implies that it is always possible to reduce the slope of a mechanism's sampling profile without affecting identification. Unfortunately, reducing the slope of the sampling profile also reduces incentives for truth-telling. We illustrate this with the family of mechanisms  $(G_\lambda^*)_{\lambda \in (0, 1)}$  which generalize  $G^*$  as follows:  $M = (-\frac{1}{2}, \frac{1}{2})$ ,  $\pi(m) = \frac{1}{2} + \lambda m$  and  $p(m) = \lambda V_{\max} m^2$ . As the slope of the sampling profile  $\lambda$  goes to zero, each agent will be sampled with probability approaching  $\frac{1}{2}$  and will pay an amount approaching zero, irrespective of the message he sends. For any  $\lambda > 0$ ,  $m = V_t/2V_{\max}$  is still a dominant strategy for an agent of type  $t$ . However, if an agent with value  $V_t$  instead sends message  $V/2V_{\max}$  with  $V \neq V_t$ , his expected loss is

$$U(t|m = V_t/2V_{\max}) - U(t|m = V/2V_{\max}) = \frac{\lambda}{4V_{\max}}(V_t - V)^2,$$

which vanishes as the slope of the sampling profile  $\lambda$  goes to 0.

Importantly, although there is a trade-off between oversampling agents with high values and the noisiness of the preference information that may be elicited, the slope of the sampling profile is a free parameter over which the principal can, and should, optimize.

Altogether, this section has shown that open selective trials provide a simple way to identify MTEs and, more generally, the distribution of returns conditional on willingness to pay. In addition, while selective trials systematically oversample high value agents, this issue is negligible when sample size is large or agents are very responsive to incentives. However, as Section I highlights, willingness to pay need not be a good predictor of actual effort, and MTEs may not allow identification of deep parameters of interest. The following sections explore richer treatments that can better identify the role of effort.

## IV. Blind Selective Trials

### A. Framework and Basic Results

In blind trials, an agent is assigned a probability of being in the treatment group,  $\phi \in [0, 1]$ , which is disclosed to the agent, and an actual treatment status,  $\tau \in \{0, 1\}$ , which is known only to the principal. Thus, the pair  $(\tau, \phi)$  can be thought of as a full description of an agent's overall treatment. Blind selective

trials nest both open selective trials, where  $\phi \in \{0, 1\}$ , and standard blind trials, where  $\phi$  is fixed.

#### ASSIGNMENT MECHANISMS

As noted in Section II, selective blind trials are mechanisms  $G = (M, \mu)$  where  $\mu : M \rightarrow \Delta([0, 1] \times \mathbb{R})$ . Given a message  $m$ ,  $\mu$  assigns the agent a likelihood of being treated  $\phi \in [0, 1]$ , and a transfer  $p \in \mathbb{R}$ . An actual, and unrevealed, treatment status  $\tau \in \{0, 1\}$  is drawn according to  $\phi$ .

#### UTILITY AND EFFORT

An agent of type  $t$ 's value for uncertain treatment status  $\phi$  is:

$$(1) \quad V_t(\phi) = \max_{e \in E} \left( \phi q_t(\tau=1, e) + (1-\phi) q_t(\tau=0, e) \right) \left( u(y=1, t) - u(y=0, t) \right) + u(y=0, t) - c(e, t).$$

The corresponding effort decision is  $e^*(\phi, t)$ , which we assume is unique.<sup>26</sup> Consistent with earlier notation, we maintain  $V_t(\phi=0) = 0$ . Note that  $V_t(\phi=1) = V_t$  is the agent's value for treatment in an open trial. Throughout the section, we keep  $\phi$  as an argument of  $V_t(\phi)$  and denote the value of  $V_t(\phi)$  at  $\varphi$  by  $V_t(\phi=\varphi)$ . Thus,  $V_t(\phi)$  denotes the entire mapping:  $\varphi \mapsto V_t(\phi=\varphi)$ . Denoting by  $\mu(\phi|m)$  the distribution of assignments  $\phi$  given message  $m$ , we have:

**PROPOSITION 4** (most informative mechanisms): *Any strictly incentive-compatible blind mechanism  $G$  identifies at most the mapping  $V_t(\phi)$  (that is,  $V_t(\phi) = V_{t'}(\phi) \Rightarrow m_G(t) = m_G(t')$ ).*

*If  $G$  identifies  $V_t(\phi)$  (that is,  $m_G(t) = m_G(t') \Rightarrow V_t(\phi) = V_{t'}(\phi)$ ) and satisfies  $\inf_{\phi, m} \mu(\phi|m) > 0$  then  $G' \preceq G$  for any strictly incentive-compatible mechanism  $G'$ .*

A simple generalization of the BDM mechanism is a most informative blind trial. The blind BDM Mechanism (bBDM) is composed of distributions  $F_\phi$  over  $[0, 1]$ , and  $F_{p|\phi}$  over  $[-V_{\max}, V_{\max}]$  with densities bounded away from 0, and the message space  $M = [-V_{\max}, V_{\max}]^{[0,1]}$ , so that a message  $m$  corresponds to a value function  $V_t(\phi)$ . An agent sends message  $m_i$ , and the principal draws values  $\phi_i = \varphi$  and  $p_i$  according to distributions  $F_\phi$  and  $F_{p|\varphi}$ . If  $m_i(\varphi) \geq p_i$ , the agent is assigned  $(\varphi, p_i)$ . Otherwise, the agent is assigned  $(0, 0)$ . It is straightforward to show that  $m_{bBDM}(t) = V_t(\phi)$ . Additionally, bBDM satisfies the full sampling constraint  $\inf_{\phi, m} \mu(\phi|m) > 0$ .

Blind selective trials have two distinct advantages over open selective trials. First, blind selective trials distinguish an agent's belief  $\phi$  and treatment status  $\tau$ . As detailed in the next subsection, this allows the principal to identify whether

<sup>26</sup>Using the results of Milgrom and Segal (2002) this allows us to apply the usual Envelope Theorem to  $V_t(\phi)$  in Proposition 6. Note that this also implies that  $e^*(\phi, t)$  is continuous in  $\phi$ .

empirical success rates are being driven by the agent's behavior or by the treatment itself. Second, by identifying the value function  $V_t(\phi)$ , blind selective trials provide useful information about an agent's intended behavior and his perceived success rate under different conditions.

### B. The Value of Distinguishing Beliefs and Treatment Status

Changes in success rates due to treatment come from two sources: the effect of the treatment itself, and the effect of behavioral changes induced by treatment. In an open trial, changes in behavior are perfectly correlated with changes in treatment status. As a result, the effect of treatment, and the effect of behavioral changes are difficult to distinguish. In contrast, blind trials allow us to disentangle these two effects by distinguishing an agent's actual treatment status  $\tau$  and his (correct) belief  $\phi$  that he is being treated.

We can disentangle these effects by considering  $\mathbb{E}[y|V_t(\phi), \phi = \varphi, \tau]$ , the measured success rate conditional on the value function  $V_t(\phi)$ , belief  $\phi = \varphi$ , and treatment status  $\tau$ , which is identified by selective blind trials. This allows identification of MTEs conditioned on the entire value function,  $\Delta^{MTE}(V_t(\phi))$ , as well as the pure treatment and behavioral effects  $\Delta^T(V_t(\phi))$  and  $\Delta^B(V_t(\phi))$ :

$$\begin{aligned}\Delta^T(V_t(\phi)) &= \lim_{\substack{\varphi \rightarrow 0 \\ \varphi > 0}} \mathbb{E}[y|V_t(\phi), \phi = \varphi, \tau = 1] - \mathbb{E}[y|V_t(\phi), \phi = \varphi, \tau = 0] \\ \Delta^B(V_t(\phi)) &= \lim_{\substack{\varphi \rightarrow 1 \\ \varphi < 1}} \mathbb{E}[y|V_t(\phi), \phi = \varphi, \tau = 0] - \mathbb{E}[y|V_t(\phi), \phi = 0, \tau = 0].\end{aligned}$$

As  $\varphi$  approaches zero, an agent's effort converges to  $e^*(\tau = 0, t)$ , the effort he would expend if he knew he was not treated.<sup>27</sup> Hence,  $\Delta^T$  identifies the returns to treatment keeping the agent's behavior at its default level  $e^*(\tau = 0, t)$ . Similarly, as  $\varphi$  approaches one, the agent's effort converges to  $e^*(\tau = 1, t)$ , the effort associated with sure treatment. Thus,  $\Delta^B$  is the effect of behavior change alone. Finally,

$$\Delta^I \equiv \Delta^{MTE} - \Delta^T - \Delta^B$$

measures the aggregate treatment effect (conditional on value  $V_t(\phi)$ ), net of the effect of treatment and behavior alone. That is,  $\Delta^I$  measures the interaction effect between behavior and treatment. If  $\Delta^I$  is positive, then treatment and effort changes are complementary in producing successful outcomes. If, instead,  $\Delta^I$  is negative, this suggests that there is a negative interaction between treatment and the perceived optimal effort of agents.<sup>28</sup>

<sup>27</sup>We use a continuity argument because  $\phi = 0$  implies  $\tau = 0$ , hence, there is no treatment group. This is essentially an identification at infinity argument, as in Heckman (1990) or Heckman and Honoré (1990), which entails well-known practical difficulties.

<sup>28</sup>These quantities can also be measured unconditionally across the entire agent population, or conditioned only on the value for sure treatment,  $V_t$ . Moreover,  $\Delta^T$  can be estimated using a standard blind RCT with a sufficiently low value of  $\phi$ .

Being able to identify  $\Delta^T$  and  $\Delta^B$  has important practical implications. Consider, for example, a cholesterol-reducing drug. If agents react to anticipated treatment by eating more fatty foods, then the aggregate effect of treatment could be quite small even if the effect of the drug alone is significant. In this environment,  $\Delta^T$  is the treatment effect purified of changes in behavior, that is, the effect of the drug on people who do not change their diet. Moreover, the sum of the interaction effect and the pure effect of treatment,  $\Delta^I + \Delta^T$ , is the part of treatment effect that could not be obtained without treatment.

When interpreting  $\Delta^B$  and  $\Delta^I$  it is important to keep in mind that these are the direct and interaction effects at the agents' *perceived* optimal effort level  $e^*(\tau = 1, t)$ . Consequently, if  $\Delta^I$  and  $\Delta^B$  are small, this may be because effort does not improve the success rate of treatment, or because the agent is expending inappropriate effort. In order to distinguish these two possibilities, we need additional information on the effort of agents. As the following subsection shows, this is what  $V_t(\phi)$  provides.

### C. The Value of Eliciting Preferences $V_t(\phi)$

As highlighted in Section I.C, the mapping  $V_t(\phi)$  can tell us whether, and by how much, treatment changes an agent's effort. Recalling that  $V_t(\phi = 0) = 0$ , knowledge of the mapping  $V_t(\phi)$  provides the following simple test.

PROPOSITION 5 (a test of "intention to change behavior"):

If  $e^*(\phi = 0, t) = e^*(\phi = 1, t)$ , then for all  $\varphi$ ,  $V_t(\phi = \varphi) = \varphi V_t(\phi = 1)$ .

If  $e^*(\phi = 0, t) \neq e^*(\phi = 1, t)$ , then for all  $\varphi \in (0, 1)$ ,  $V_t(\phi = \varphi) < \varphi V_t(\phi = 1)$ .

When effort changes with  $\tau$ , the agent gets additional surplus from tailoring his behavior to his treatment status. The difference  $\varphi V_t(\phi = 1) - V_t(\varphi)$  is thus the agent's willingness to pay to learn his actual treatment status, which will be zero if effort is independent of treatment.<sup>29</sup> Recalling that  $q_t(\tau, e)$  is the perceived success rate of an agent with type  $t$ , the value function  $V_t(\phi)$  also allows us to estimate an agent's perceived returns to effort.

PROPOSITION 6 (identifying perceived returns to effort): For any value  $\varphi$ ,

$$\left. \frac{\partial V_t(\phi)}{\partial \phi} \right|_{\varphi} = [q_t(\tau = 1, e^*(\varphi, t)) - q_t(\tau = 0, e^*(\varphi, t))] \times [u(y = 1, t) - u(y = 0, t)].$$

Note that selective blind trials can allow for double-blind designs in which the principal has varying beliefs about the likelihood that an agent is being treated. Varying the beliefs of the principal may help identify the treatment effect due to variations in the principal's behavior. A proper analysis of this approach requires a better understanding of the principal's incentive problem, which we abstract away from in this paper.

<sup>29</sup>When  $\varphi = 1/2$  this coincides with test statistic  $\theta_t$  defined in Section I.C.

Note that in a richer decision theoretic framework, agents may have preferences for early revelation of uncertainty, even though their actions do not depend on information (Kreps and Porteus, 1978). In such a framework, an agent's value for information would be a noisy (but still informative) signal of intent to change behavior.

In particular, we can compute the agent's perceived increase in treatment effects when moving from default effort (induced by  $\varphi = 0$ ) to perceived optimal effort given treatment (induced by  $\varphi = 1$ ):

$$\frac{\partial V_t(\phi) \Big|_1}{\partial \phi} \Big/ \frac{\partial V_t(\phi) \Big|_0}{\partial \phi} = \frac{q_t(\tau=1, e^*(\varphi=1, t)) - q_t(\tau=0, e^*(\varphi=1, t))}{q_t(\tau=1, e^*(\varphi=0, t)) - q_t(\tau=0, e^*(\varphi=0, t))}.$$

This data helps evaluate whether under-provision of effort is to blame for low returns to treatment.<sup>30</sup> Returning to the example in Section I, imagine a trial of a water treatment product known to the principal to be effective only if agents use it whenever they drink water. If measured returns to the treatment are low, there are two competing explanations: 1) the treatment is not effective in the agents' disease environment, or 2) agents are not expending appropriate effort using the product. Agents' perceived returns can help distinguish these explanations. If perceived returns to effort are high, then the agent is likely to be expending significant effort, and it is more likely that the treatment is not effective in a particular disease environment. If, instead, perceived returns are low, it becomes more likely that the treatment has an effect that is unmeasured due to agents' lack of effort.

Preference data  $V_t(\phi)$  may also provide some insight into the nature of placebo effects. Under a sufficiently broad definition of behavior—including unconscious or involuntary behavior—behavioral treatment effects  $\Delta^B$  are largely undistinguishable from placebo effects (Malani, 2006). However, because indirect preferences identify whether or not agents intend to change their behavior (Proposition 5), this data provides some indication of whether behavioral effects  $\Delta^B$  are driven by changes in behavior of which the agent is aware. For instance, if agents do not value information ( $V_t(\phi = \varphi) = \varphi V_t(\phi = 1)$ ), and yet exhibit positive behavioral effects ( $\Delta^B > 0$ ), this suggests that changes in behavior the agent is unaware of are driving behavioral effects.

## V. Incentivized Selective Trials

We now show how quantities similar to those identified by blind selective trials can be identified without a placebo. This can be accomplished using an incentivized selective trial, which allows agents to express preferences over contracts.<sup>31</sup>

<sup>30</sup>Identifying these derivatives requires the precise elicitation of an agent's preferences. This relies heavily on the rationality of agents, but not sample size.

Note that the logic underlying Proposition 5 implies  $V_t(\phi)$  must be convex. This follows from the fact that any mean preserving spread in belief  $\phi$  is equivalent to the arrival of a signal about treatment status. As more information is necessarily useful in this setting, this implies that  $V_t(\phi)$  is convex. Thus, in practice, it may be preferable to use simpler mechanisms that elicit  $V_t(\phi)$  for very few values of  $\phi$ , and construct discrete approximations of the desired derivatives. As  $V_t(\phi)$  is convex in  $\phi$ , a few points are sufficient to obtain correct bounds on these derivatives.

<sup>31</sup>For field experiments using explicit incentives, see, for instance, Gertler (2004); Schultz (2004); Volpp et al. (2006, 2008); Thornton (2008); and Kremer et al. (2009). A fully worked-out numerical example illustrating inference from incentivized trials is given in the online appendix.

## A. Framework and Basic Results

## ASSIGNMENT MECHANISMS

As noted in Section II, an incentivized trial is a mechanism  $G = (M, \mu)$ , where  $\mu : M \rightarrow \Delta(\{0, 1\} \times \mathbb{R} \times \mathbb{R})$ . Given a message  $m$ ,  $\mu$  is used to draw a treatment status  $\tau$ , a fixed transfer  $p$  from the agent, as well as a bonus  $w$  transferred to the agent in the event of success. Note that both  $p$  and  $w$  may be negative in the case of insurance. The pair  $(\tau, w)$  can be thought of as an aggregate treatment.

## UTILITY AND EFFORT

The agents' indirect preferences over contracts  $(\tau, w)$ , denoted by  $V_t(\tau, w)$ , are given by

$$(2) \quad V_t(\tau, w) = \max_{e \in E} q_t(\tau, e)[u(y=1, t) + w] + [1 - q_t(\tau, e)]u(y=0, t) - c(e, t).$$

We denote by  $e^*(\tau, w, t)$  the induced effort level, and maintain the normalization  $V_t(\tau=0, w=0) = 0$ .

## INSURANCE

A specific value of  $w$  that will be useful is  $w_{0,t} \equiv -[u(y=1, t) - u(y=0, t)]$ , the utility difference between success ( $y = 1$ ) and failure ( $y = 0$ ) for an agent of type  $t$ . The transfer  $w_{0,t}$  essentially provides an agent with perfect insurance over the outcome  $y$ . When fully insured, an agent will choose  $e$  to minimize the cost of his effort, regardless of his treatment status. We refer to this effort choice as *no effort*. Note that no effort differs from the default behavior of untreated agents in an open trial, as they may still be expending effort to improve their success rate.

We proceed by assuming that  $w_{0,t}$  is known to the principal. At the end of the section we show that under mild assumptions,  $w_{0,t}$  can be inferred from elicited preferences  $V_t(\tau, w)$ .

## B. What can be Inferred from Incentivized Trials?

It is straightforward to extend Propositions 1 and 4, which characterize most informative mechanisms. That is,  $G$  is a most informative incentivized trial if it identifies the mapping  $V_t(\tau, w)$  and, given any message, puts positive density on all possible treatments  $(\tau, w)$ . As before, the BDM mechanism can be adapted to identify  $V_t(\tau, w)$ —Appendix B provides a detailed description. Note that the information produced by incentivized trials nests that produced by open trials. In particular,  $V_t(\tau=1, w=0) = V_t$ .

As in the case of blind selective trials, incentivized selective trials allow us to disentangle the effects of treatment and effort, as well as infer an agent's perception of how effort affects outcomes. Incentivized selective trials recover the

empirical success rate  $\mathbb{E}[y|V_t(\tau, w), \tau, w]$  as a function of preferences, treatment, and incentives. This will be independent of reward  $w$  if effort does not matter for outcomes, or if incentives do not affect effort expenditure.

#### ISOLATING RETURNS TO TREATMENT AND RETURNS TO EFFORT

A contract with transfer  $w_{0,t} \equiv -[u(y=1, t) - u(y=0, t)]$  provides an agent of type  $t$  with perfect insurance. Thus, the agent expends no effort, regardless of his treatment status. Given  $w_{0,t}$ , we can identify two quantities similar to those discussed in Section IV.B:

$$\begin{aligned} \text{Returns to Treatment | No Effort} &= \mathbb{E}[y|V_t(\tau, w), \tau=1, w=w_{0,t}] \\ &\quad - \mathbb{E}[y|V_t(\tau, w), \tau=0, w=w_{0,t}] \\ \text{Returns to Effort | Treatment} &= \mathbb{E}[y|V_t(\tau, w), \tau=1, w=0] \\ &\quad - \mathbb{E}[y|V_t(\tau, w), \tau=1, w=w_{0,t}] \end{aligned}$$

Note that here returns are measured using no effort as a baseline, rather than the default effort level  $e^*(\tau=0, w=0, t)$  expended by agents in the control group of an open trial.<sup>32</sup>

#### IDENTIFYING PERCEIVED RETURNS TO EFFORT

Indirect preferences over contracts  $V_t(\tau, w)$  also provide information on perceived returns to effort. Recall that  $q_t(\tau, e)$  denotes the agent's perceived likelihood of success given treatment status  $\tau$  and effort  $e$ .

PROPOSITION 7 (identifying perceived success rates):

$$\forall \tau, w, \quad \frac{\partial V_t(\tau, w)}{\partial w} = q_t(\tau, e^*(\tau, w, t)).$$

Given knowledge of  $w_{0,t}$ , this allows us to compute subjective returns to treatment and perceived optimal effort:

$$\begin{aligned} \text{Perceived Returns to Treatment} &= q_t(\tau=1, w=w_{0,t}|V_t(\tau, w)) \\ &\quad - q_t(\tau=0, w=w_{0,t}|V_t(\tau, w)) \\ \text{Perceived Returns to Effort} &= q_t(\tau=1, w=0|V_t(\tau, w)) \\ &\quad - q_t(\tau=1, w=w_{0,t}|V_t(\tau, w)). \end{aligned}$$

Note that if perceived returns to effort are low, this can indicate that an agent plans to expend little or no effort using the technology. The principal can use this information in deciding which agents' usage to monitor more closely.

<sup>32</sup>Note that, unlike blind selective trials, identification here does not rely on identification at infinity (see Footnote 27).



The monetary equivalent of the cost of an agent’s optimal effort can be obtained by rearranging (2):

$$c(e^*(\tau, w=0, t)) - c(e^*(\tau, w=w_{0,t}, t)) = -w_{0,t} \times q_t(\tau, e^*(\tau, w=0, t)) - [V_t(\tau, w=0) - V_t(\tau, w=w_{0,t})].$$

Note that all parameters on the right-hand side are identified from data, except perhaps  $w_{0,t}$ .

Identifying the costs incurred by agents can improve inference by allowing the principal to distinguish—among agents who believe that appropriate effort has high returns—those who believe that only a small amount of effort is sufficient to obtain high returns from those who believe that a significant amount of effort is necessary to obtain high returns.

#### IDENTIFYING THE FULL INSURANCE CONTRACT

One drawback of incentivized trials is that they rely on identifying the full insurance contract  $w_{0,t}$ , which may depend on the agent’s type. However,  $w_{0,t}$  can be identified from preference information under mild assumptions.

**FACT 3:** *Assume that outcome  $y = 1$  yields strictly greater utility than  $y = 0$ , that is,  $u(y=1, t) > u(y=0, t)$ , and an agent perceives treatment to be beneficial:*

$$\forall e_0 \in E, \exists e_1 \in E \text{ s.t. } c(e_1, t) \leq c(e_0, t) \quad \text{and} \quad q_t(\tau=0, e_0) < q_t(\tau=1, e_1).$$

$$\text{Then, } w_{0,t} = \max\{w \mid V_t(\tau=1, w) = V_t(\tau=0, w)\}.$$

In words, when treatment facilitates success, the full insurance transfer  $w_{0,t}$  is the highest transfer such that an agent places no value on obtaining treatment. Note that our assumptions rule out cases where an agent believes treatment reduces the likelihood of success, as well as environments where an agent values treatment only for reasons other than its impact on the principal’s outcome of interest. Whenever the assumptions of Fact 3 do not hold,  $w_{0,t}$  must be calibrated from alternative data, for example, the expected amount of wages lost when sick. This is a delicate task, and estimates of  $w_{0,t}$  are likely to be noisy. The corresponding insurance contract would not induce no effort, but rather a small, and slightly uncertain, level of effort. Hence, whenever the full insurance contract  $w_{0,t}$  is estimated with noise, this leads to noisy estimates of treatment effects.

## VI. Discussion

This paper studies inference and external validity when experimental subjects take unobserved decisions that can affect outcomes. As effort expenditure is driven by beliefs, and beliefs can respond to information, the returns measured by an RCT may not be representative of the returns a better informed population

would obtain. We take a principal-agent approach to trial design, where the principal maximizes the informativeness of data. This leads us to study selective trials, which improve on RCTs by allowing agents to express preferences over treatments of varying richness. We show that selective trials can identify whether agents' beliefs are reducing measured treatment effects, as well as separate the returns from treatment, effort, and their interaction.

More generally, this paper advocates a mechanism design approach to randomized controlled experiments, an approach we believe can help build bridges between reduced form methods—largely concerned with robustness and internal validity—and structural methods—which use models to identify deep parameters needed for external validity. While we believe this research agenda can yield many useful applications, successfully implementing its insights requires overcoming a number of practical difficulties. In the remainder of this section we discuss some of these implementation challenges and directions for future work.

#### A. Implementation Issues

In theory, the selective trials described in this paper are robust and require no specific knowledge on the part of the principal. However, our results are obtained under three important sets of assumptions that may not hold in practice.

#### BEHAVIORAL ASSUMPTIONS

The correct elicitation of preferences, which is key to our analysis, relies strongly on the assumption that agents are rational. However, as people often fail to play dominant strategies, BDM-like mechanisms only provide a noisy signal of the agents' underlying valuations (Keller et al., 1993; Bohm et al., 1997). This suggests that running even relatively simple open selective trials, let alone full-fledged blind or incentivized selective trials, is likely to be challenging.

Agents may also be subject to other behavioral biases that are not taken into account by our framework.<sup>33</sup> A specific concern is that the act of making choices may change agents' preferences. For example, it is possible that an agent who expresses a strong desire for, but does not get, treatment, may attempt to obtain treatment by other means, but would not do so if his valuation was never elicited.<sup>34</sup> Another concern is that agents may try to infer the value of treatment from the principal's choice of experimental design. For example, similar to Milgrom and Roberts (1986), if treatment is only available at a high cost, agents may infer that the technology is more valuable. In these environments, a principal should take into account how experimental design influences behavior before

<sup>33</sup>For instance, loss aversion, ambiguity aversion, or even social preferences may play a significant role. A different bias might come from the psychological cost of parting from *any* amount money (Cohen and Dupas, 2010; Ashraf et al., 2010).

<sup>34</sup>A simple way to test for this is to construct a second control group that is never asked to express preferences.

drawing inferences.<sup>35</sup>

Ultimately, we believe the best way to address these concerns is through careful and extensive experimentation, blending both laboratory and field work. As laboratory experiments allow the observation of underlying fundamentals, they are essential to understand which implementations of selective trials produce more reliable data, and what the relevant biases may be. In turn, field experiments—in simple environments where actual behavior is observable, and trustworthy surveys may be conducted—are needed to check that the insights gathered from the laboratory apply in more realistic settings. We anticipate that appropriate implementations should give agents multiple opportunities to learn how the relevant mechanism works before they actually express preferences over treatment (Plott and Zeiler, 2005). Additionally, it may be preferable to use mechanisms that elicit coarse information about preferences, but impose a smaller cognitive burden on agents.<sup>36</sup>

Finally, even if our behavioral assumptions are wrong, the data generated still enriches that obtained through an RCT. Although this invalidates the interpretation of the data put forth in this paper, it does not preclude a more standard analysis focusing on average treatment effects, or a more sophisticated analysis taking into account relevant biases.

#### SAMPLE SIZE

Large samples are likely to be necessary in order to realize the full value of the additional data our mechanisms elicit. Note that the difficulty is not with the data collection process, as the correct elicitation of preferences relies only on rationality. Rather, sample size restricts the ability to compute meaningful estimates of treatment effects conditional on preferences. This issue is inherent to any non-parametric estimation of treatment effects conditional on a rich set of explanatory variables, and existing methodologies apply (see, for instance, Pagan and Ullah (1999)). Given sufficiently large samples, a kernel regression may be practical. In small samples, it may be necessary to bin agents with similar preferences. Alternatively, it may be informative to estimate parametric relationships between treatment effects and preference data.<sup>37</sup>

#### CASH CONSTRAINTS

Eliciting preferences using monetary trade-offs is impractical in the presence of severe cash constraints. When only agents are cash constrained, a possible, but

<sup>35</sup>For example, the design itself could be considered as part of the experimental treatment. This implies the principal should compare agents whose preferences are the same, but have been elicited using different mechanisms.

<sup>36</sup>In the case of open selective trials, one may elicit the agent's preferences over only a few lotteries—see Appendix B for a discussion. In the case of blind selective trials, one may elicit  $V_i(\phi)$  at a few values of  $\phi$  and exploit the fact that  $V_i(\phi)$  is convex to fit simple functional forms.

<sup>37</sup>Note that controlling for preferences may reduce the heterogeneity of treatment effects within each bin. This may alleviate statistical power concerns.

expensive, solution is to give agents a show-up fee that they can use to express preferences.

More fundamentally, monetary trade-offs may be uninformative of intended behavior in environments where there is sizable heterogeneity in the marginal value of income. For example, Cohen and Dupas (2010) finds that willingness to pay for bednets in Kenya is a poor predictor of actual use.<sup>38</sup> In that setting, other trade-offs—such as willingness to wait, willingness to perform tedious tasks, or willingness to return at a later time—may be more informative of agents' intended behavior. The choice of the relevant trade-off is an important degree of freedom that can and should be guided by local knowledge.

In general, it is clear that implementing the ideas advocated in this paper entails complex experimental designs, and the details of an individual experiment may need to be fine tuned with careful, context-dependent, pilot projects. However, we are encouraged by recent field experiments showing that complex designs can be successfully implemented (see Ashraf et al., 2010; Karlan and Zinman, 2009; and particularly Berry et al., 2011, which implements a BDM mechanism in the field). Thus, despite the significant caveats detailed in this section, we are hopeful that our approach will prove useful in guiding future field work.

### B. Theoretical Extensions

Our approach also suggests directions for further theoretical work. We believe these extensions are sufficiently interesting in their own right to deserve independent analyses. We outline two of these extensions, specifying both the challenges they pose and their potential value added.

#### EXTENSION TO DYNAMIC MECHANISMS

While our framework can accommodate learning and dynamic effort expenditure by agents, we focus on mechanisms that elicit agents' preferences only once. This is a significant restriction, as identifying whether, and how, agents change their behavior over time is an important input in the analysis of treatment effects (Philipson and Desimone, 1997; Philipson and Hedges, 1998; Scharfstein et al., 1999; Chan and Hamilton, 2006). However, the timing of elicitation is a free design variable. In particular, it may occur before or after an agent has been exposed to the technology.

For concreteness, consider a technology that requires sustained effort to yield returns, for example, anti-depressants with delayed effects, technologies exhibiting significant learning-by-doing, and so on. Eliciting how preferences change over time would improve inference by helping to distinguish agents exhibiting consistent motivation throughout the trial from agents whose motivation drops

<sup>38</sup>Note that this is not always the case. Ashraf et al. (2010) documents the opposite finding for water treatment products in Zambia.

in the middle. The difficulty is that eliciting preferences in the future necessarily changes an agent's beliefs about future treatment status, and, in turn, changes current effort expenditure. In particular, if an agent is promised treatment in future periods to induce a particular effort level today, then it becomes impossible to elicit preferences in the future without breaking this promise.<sup>39</sup>

#### EXTENSION TO MULTI-AGENT MECHANISMS

The mechanisms considered in this paper are all single-agent mechanisms—an agent's assignment depends only on the message he sends and not on the messages sent by other agents. This allows us to identify an agent's preferences, and thus his beliefs about his own returns to treatment and to effort. Considering multi-agent mechanisms, in which assignment depends on the messages sent by others, can allow us to identify an agent's beliefs about others agents' values, others agents' success rates, and so on.

The information elicited by multi-agent mechanisms may be useful if there are externalities between agents, as in Miguel and Kremer (2004), or to investigate social learning. For example, if we observe that most agents have low value for the technology, but believe that others have high value for the technology, this suggests a specific failure of social learning, and provides us with the means to correct it. Indeed, if most agents do not expend effort using the technology, but believe others do, then they will interpret each others' poor outcomes as a signal that even with high effort the technology does not yield returns. Providing the agents with actual data on others' willingness to pay corrects these inference mistakes and may increase experimentation.

#### REFERENCES

- Abbring, Jaap H. and Gerard J. Van den Berg**, “The Nonparametric Identification of Treatment Effects in Duration Models,” *Econometrica*, September 2003, 71 (5), 1491–1517.
- and —, “Social Experiments and Instrumental Variables with Duration Outcomes,” 2005. Tinbergen Institute Discussion Paper 2005-047/3.
- Angrist, Joshua D., Guido W. Imbens, and Donald B. Rubin**, “Identification of Causal Effects using Instrumental Variables,” *Journal of the American Statistical Association*, June 1996, 91 (434), 444–455.
- Ashraf, Nava, James Berry, and Jesse M. Shapiro**, “Can Higher Prices Stimulate Product Use? Evidence from a Field Experiment in Zambia,” *American Economic Review*, December 2010, 100 (6), 2383–2413.

<sup>39</sup>In the context of labor market experiments, Abbring and Van den Berg (2003, 2005) makes a similar point: if expectations of potential access to treatment change ex ante behavior (for example, investment in human capital), then treatment effects are not identified.

- Banerjee, Abhijit**, “A Simple Model of Herd Behavior,” *The Quarterly Journal of Economics*, August 1992, 107 (3), 797–817.
- Becker, Gordon M., Morris H. DeGroot, and Jacob Marschak**, “Measuring Utility by a Single-Response Sequential Method,” *Behavioral Science*, 1964, 9 (3), 226–232.
- Berry, James, Greg Fischer, and Raymond Guiteras**, “Eliciting and Utilizing Willingness to Pay: Evidence from Field Trials in Northern Ghana,” 2011. London School of Economics, *mimeo*.
- Bikhchandani, Sushil, David Hirshleifer, and Ivo Welch**, “A Theory of Fads, Fashion, Custom, and Cultural Change as Informational Cascades,” *Journal of Political Economy*, 1992, 100 (5), 992–1026.
- Bohm, Peter, Johan Lindén, and Joakin Sonnegård**, “Eliciting Reservation Prices: Becker-DeGroot-Marschak Mechanisms vs. Markets,” *The Economic Journal*, July 1997, 107 (443), 1079–1089.
- Chan, Tat Y. and Barton H. Hamilton**, “Learning, Private Information, and the Economic Evaluation of Randomized Experiments,” *Journal of Political Economy*, 2006, 114 (6), 997–1040.
- Cohen, Jessica and Pascaline Dupas**, “Free Distribution or Cost-Sharing? Evidence from a Randomized Malaria Prevention Experiment,” *Quarterly Journal of Economics*, 2010, 125 (1), 1–45.
- Deaton, Angus**, “Instruments, Randomization, and Learning about Development,” *Journal of Economic Literature*, 2010, 48 (2), 424–455.
- Duflo, Esther, Michael Kremer, and Jonathan Robinson**, “How High Are Rates of Return to Fertilizer? Evidence from Field Experiments in Kenya,” *American Economic Review*, 2008, 98 (2), 482–488.
- , **Rachel Glennerster, and Michael Kremer**, “Using Randomization in Development Economics Research: A Tool Kit,” in T. Paul Schultz and John Strauss, eds., *Handbook of Development Economics, Vol. 4*, Amsterdam: Elsevier, 2008, pp. 3895–3962.
- , **Rema Hanna, and Stephen Ryan**, “Monitoring Works: Getting Teachers to Come to School,” 2010. MIT, *mimeo*.
- Dupas, Pascaline**, “What Matters (and What Does Not) in a Households’ Decision to Invest in Malaria Prevention,” *American Economic Review*, 2009, 99 (2), 224–230.
- , “Short-Run Subsidies and Long-Term Adoption of New Health Products: Experimental Evidence from Kenya,” 2010. University of California, Los Angeles *mimeo*.

- , “Do Teenagers Respond to HIV Risk Information? Evidence from a Field Experiment in Kenya,” *American Economic Journal: Applied Economics*, January 2011, *3* (1), 1–36.
- Flood, A.B., J.E. Wennberg, R.F. Nease, F.J. Fowler, J. Ding, and L.M. Hynes**, “The Importance of Patient Preference in the Decision to Screen for Prostate Cancer,” *Journal of General Internal Medicine*, 1996, *11* (6), 342–349.
- Gertler, Paul**, “Do Conditional Cash Transfers Improve Child Health? Evidence from PROGRESA’s Control Randomized Experiment,” *American Economic Review*, 2004, *94* (2), 336–341.
- Heckman, James J.**, “Varieties of Selection Bias,” *The American Economic Review*, 1990, *80* (2), 313–318.
- **and Bo E. Honoré**, “The Empirical Content of the Roy Model,” *Econometrica*, 1990, *58* (5), 1121–1149.
- **and Edward Vytlacil**, “Structural Equations, Treatment Effects, and Econometric Policy Evaluation,” *Econometrica*, May 2005, *73* (3), 669–738.
- **, Jeffrey Smith, and Nancy Clements**, “Making the Most out of Programme Evaluations and Social Experiments: Accounting for Heterogeneity in Programme Impacts,” *The Review of Economic Studies*, 1997, *64* (4), 487–535.
- Imbens, Guido W.**, “Better LATE Than Nothing: Some Comments on Deaton (2009) and Heckman and Urzua (2009),” 2010. Harvard University, *mimeo*.
- **and Joshua D. Angrist**, “Identification and Estimation of Local Average Treatment Effects,” *Econometrica*, March 1994, *62* (2), 467–475.
- Jadad, Alejandro R. and Murray Enkin**, *Randomized Controlled Trials: Questions, Answers, and Musings*, BMJ Books, 2007.
- Jin, Hui and Donald B. Rubin**, “Principal Stratification for Causal Inference with Extended Partial Compliance,” *Journal of the American Statistical Association*, 2008, *103* (481), 101–111.
- Karlan, Dean S. and Jonathan Zinman**, “Observing Unobservables: Identifying Information Asymmetries with a Consumer Credit Field Experiment,” *Econometrica*, 2009, *77* (6), 1993–2008.
- Keller, L. Robin, Uzi Segal, and Tan Wang**, “The Becker-DeGroot-Marschak Mechanism and Generalized Utility Theories: Theoretical Predictions and Empirical Observations,” *Theory and Decision*, 1993, *34* (2), 83–97.

- King, Michael, Irwin Nazareth, Fiona Lampe, Peter Bower, Martin Chandler, Maria Morou, Bonnie Sibbald, and Rosalind Lai**, “Impact of Participant and Physician Intervention Preferences on Randomized Trials: A Systematic Review,” *Journal of the American Medical Association*, 2005, 293 (9), 1089–1099.
- Kremer, Michael and Edward Miguel**, “The Illusion of Sustainability,” *The Quarterly Journal of Economics*, 2007, 122 (3), 1007–1065.
- , – , and **Rebecca Thornton**, “Incentives to Learn,” *The Review of Economics and Statistics*, 2009, 91 (3), 437–456.
- Kreps, David M. and Evan L. Porteus**, “Temporal Resolution of Uncertainty and Dynamic Choice Theory,” *Econometrica*, 1978, 46 (1), 185–200.
- Malani, A.**, “Identifying Placebo Effects with Data from Clinical Trials,” *Journal of Political Economy*, 2006, 114 (2), 236–256.
- Malani, Anup**, “Patient enrollment in medical trials: Selection bias in a Randomized Experiment,” *Journal of Econometrics*, 2008, 144 (2), 341–351.
- Miguel, Edward and Michael Kremer**, “Worms: Identifying Impacts on Education and Health in the Presence of Treatment Externalities,” *Econometrica*, January 2004, 72 (1), 159–217.
- Milgrom, Paul and Ilya Segal**, “Envelope Theorems for Arbitrary Choice Sets,” *Econometrica*, 2002, 70 (2), 583–601.
- and **John Roberts**, “Price and Advertising Signals of Product Quality,” *The Journal of Political Economy*, 1986, 94 (4), 796–821.
- Nguyen, Trang**, “Information, Role Models and Perceived Returns to Education: Experimental Information, Role Models and Perceived Returns to Education: Experimental Evidence from Madagascar,” 2009. MIT, *mimeo*.
- Oster, Sharon M.**, *Strategic Management for Nonprofit Organizations: Theory and Cases*, Oxford, UK: Oxford University Press, 1995.
- Pagan, Adrian and Aman Ullah**, *Nonparametric Econometrics*, Cambridge University Press, 1999.
- Philipson, Tomas and Jeffrey Desimone**, “Experiments and Subject Sampling,” *Biometrika*, 1997, 84 (3), 619–631.
- and **Larry V. Hedges**, “Subject Evaluation in Social Experiments,” *Econometrica*, 1998, 66 (2), 381–408.



- Plott, Charles R. and K. Zeiler**, “The Willingness to Pay-Willingness to Accept Gap, The ‘Endowment Effect,’ Subject Misconceptions, and Experimental Procedures for Eliciting Valuations,” *American Economic Review*, 2005, 95 (3), 530–545.
- Rothschild, Michael**, “A Two-Armed Bandit Theory of Market Pricing,” *Journal of Economic Theory*, 1974, 9 (2), 185–202.
- Roy, A.D.**, “Some Thoughts on the Distribution of Earnings,” *Oxford Economic Papers*, 1951, 3 (2), 135–146.
- Scharfstein, Daniel O., Andrea Rotnitzky, and James M. Robins**, “Adjusting for Nonignorable Drop-Out Using Semiparametric Nonresponse Models,” *Journal of the American Statistical Association*, 1999, 94 (448), 1096–1120.
- Schultz, T. Paul**, “School Subsidies for the Poor: Evaluating the Mexican Progresa Poverty Program,” *Journal of Development Economics*, 2004, 74 (1), 199–250.
- Silverman, W.A. and D.G. Altman**, “Patients’ Preferences and Randomised Trials,” *The Lancet*, 1996, 347 (8995), 171–174.
- Stolberg, Harald O., Geoffrey Norman, and Isabelle Trop**, “Randomized Controlled Trials,” *American Journal of Roentgenology*, 2004, 183 (6), 1539–1544.
- Thornton, Rebecca**, “The Demand for and Impact of Learning HIV Status: Evidence from a Field Experiment,” *American Economic Review*, 2008, 98 (5), 1829–1863.
- Tilbrook, Helen**, “Patients’ Preferences within Randomised Trials: Systematic Review and Patient Level Meta-analysis,” *British Medical Journal*, 2008, 337, 1864–1871.
- Volpp, Kevin G., Andrea Gurmankin Levy, David A. Asch, Jesse A. Berlin, John J. Murphy, Angela Gomez, Harold Sox, Jingsan Zhu, and Caryn Lerman**, “A Randomized Controlled Trial of Financial Incentives for Smoking Cessation,” *Cancer Epidemiology Biomarkers & Prevention*, 2006, 15 (1), 12.
- , **Leslie K. John, Andrea B. Troxel, Laurie Norton, Jennifer Fassbender, and George Loewenstein**, “Financial Incentive-based Approaches for Weight Loss: A Randomized Trial,” *Journal of the American Medical Association*, 2008, 300 (22), 2631–2637.
- Zelen, Marvin**, “A New Design for Randomized Clinical Trials,” *New England Journal of Medicine*, 1979, 300 (22), 1242–1245.

## EXTENSIONS

## A1. General Outcome Space

Most of the results extend directly to the case where  $y$  takes values in a general outcome space  $Y$ , and is distributed according to some density function  $f_y(R, \tau, e, t)$ . We denote by  $f_{y,t}(\tau, e) \equiv \int_R f_y(R, \tau, e, t) dt(R)$  the subjective distribution of returns from the perspective of an agent of type  $t$ . Values go from being sums of two terms to being integrals, and incentive contracts are now functions  $w : Y \rightarrow \mathbb{R}$ . We have that

$$\begin{aligned} V_t &= \max_{e \in E} \int_y u(y, t) f_{y,t}(\tau = 1, e) dy - c(e, t) \\ V_t(\phi) &= \max_{e \in E} \phi \int_y u(y, t) f_{y,t}(\tau = 1, e) dy + (1 - \phi) \int_y u(y, t) f_{y,t}(\tau = 0, e) dy - c(e, t) \\ V_t(\tau, w) &= \max_{e \in E} \int_y [u(y, t) + w(y)] f_{y,t}(\tau, e) dy - c(e, t). \end{aligned}$$

Propositions 1, 2, 3, 4 and 5 extend directly with these generalized value functions. Propositions 6 and 7, which identify subjective returns to treatment and effort differ as follows. As we have that

$$\forall y_0, \quad \frac{\partial V_t(\tau, w)}{\partial w(y_0)} = f_{y,t}(\tau, e^*(\tau, w, t))(y_0),$$

Proposition 7 extends directly.

Proposition 6, which deals with blind trials, is more difficult to extend, as now we have only a one-dimensional instrument,  $\phi \in [0, 1]$  to identify an entire function  $f_{y,t}$  rather than the single parameter  $q_t$ . We now identify

$$(A1) \quad \frac{\partial V_t(\phi)}{\partial \phi} = \int_y u(y, t) [f_{y,t}(\tau = 1, e^*(\phi, t))(y) - f_{y,t}(\tau = 0, e^*(\phi, t))(y)] dy,$$

which corresponds to a utility weighted subjective treatment effect given subjectively appropriate effort under belief  $\phi$ .

## A2. Eliciting Preferences under Non-Quasilinear Utility

The approach developed in this paper largely extends to the case where preferences are not quasilinear, although we must consider slightly different mechanisms. We now consider utility taking the form  $u(y, e, p, t)$  where  $y \in Y$ ,  $e \in E$ ,  $p \in P$  is now a prize (that is, a bundle of goods which may or may not include monetary transfers), and  $t$  is the agent's type. We focus on the case where there exists an unambiguously most desirable prize  $\bar{p} \in P$ , and an unambiguously least desirable prize,  $\underline{p} \in P$ .

In the case of open trials, indirect preferences take the following form:

$$V_t(\tau, p) = \max_e \int_y u(y, e, p, t) f_{y,t}(\tau, e) dy.$$

Say we want to elicit preference over  $(\tau, p) \in \{0, 1\} \times P$ . We assume for simplicity that for all such  $(\tau, p)$ ,  $V_t(\tau = 0, \underline{p}) \leq V_t(\tau, p) \leq V_t(\tau = 1, \bar{p})$ . We normalize  $V_t(\tau = 0, p = \underline{p}) = 0$  and  $V_t(\tau = 1, p = \bar{p}) = 1$ . Consider the following generalization of the BDM mechanism: an agent sends a message  $m \in \mathbb{R}^{\{0,1\} \times P}$ , which corresponds to a value function; the principal randomly picks  $(\tau, p, \lambda)$  from some continuous distribution over  $\{0, 1\} \times P \times [0, 1]$ ; an agent is assigned  $(\tau, p)$  if  $m(\tau, p) > \lambda$  and the lottery  $\lambda \times (\tau = 1, p = \bar{p}) + (1 - \lambda) \times (\tau = 0, p = \underline{p})$  otherwise. In this setting it is dominant for an agent to send message  $m = V_t$ . Similar mechanisms allow us to identify indirect preferences in the case of blind and incentivized trials.

Propositions 1, 3, 4 and 5 extend directly with these generalized value functions. Again, extending Propositions 6 and 7 requires some more work. Proposition 6—which identifies subjective returns to effort using blind trials—extends as is when  $y \in \{0, 1\}$ , and extends according to (A1) when  $y$  takes values in a general outcome set  $Y$ . Proposition 7 extends as is when preferences are separable in prize  $p$ , that is, when  $u(y, e, p, t) = u_0(y, e, t) - u_1(p, t)$ . When preferences are not separable in prize  $p$ , incentivized trials allow us to identify  $f_{y,t}(y) \frac{\partial u}{\partial w(y)} \Big|_{y,p}$  for all values of  $y$  and  $p$ . Note that when preferences are separable, the multiplicative constant can be identified from the fact that probabilities sum to 1.

## IMPLEMENTATION

### *B1. Implementing Open Selective Trials as a Finite Menu of Lotteries*

The mechanisms described in the paper all use a continuum of messages and elicit the agent's exact willingness to pay. Of course, it is possible to use simpler mechanisms to elicit coarser information. This example shows how to identify which of  $N$  intervals an agent's willingness to pay belongs to.

The principal chooses value thresholds  $-V_{\max} = V_0 < V_1 < \dots < V_N = V_{\max}$ . She can elicit the interval where an agent's value lies by offering a menu of lotteries. This menu is constructed with messages  $M = \{1, \dots, N\}$  and any increasing sequence  $\pi(1) < \pi(2) < \dots < \pi(N)$  of sampling rates. Thus, message  $m \in M$  corresponds to buying the lottery that delivers treatment with probability  $\pi(m)$ . In order to match these messages with the appropriate value interval, the principal simply sets  $p(m)$ , the price of lottery  $m$ , according to:

$$(B1) \quad \forall k > 1, \quad p(k) = p(k-1) + (\pi(k) - \pi(k-1))V_{k-1}.$$

Note that the sequence of prices is entirely determined by  $p(1)$ . Denote by  $G^{\pi,p}$  the mechanism corresponding to this menu of lotteries, then:

FACT 4: Under mechanism  $G^{\tau,p}$  an agent of type  $t$  sends message  $k$  if and only if  $V_t \in [V_{k-1}, V_k]$ .

This emphasizes the many degrees of freedom the principal has when implementing selective trials as menus of lotteries. The value intervals according to which agents are classified, and the rates according to which they obtain treatment are, to a large extent, free parameters. The only restriction is that sampling rates must be increasing in an agent's value (Proposition 2).

### B2. Implementing Incentivized Selective Trials

This section complements Section V by describing how to implement incentivized selective trials as an extension of the BDM mechanism. Let the message space  $M$  be the set of (normalized) possible utility functions  $V_t(\tau, w)$ :

$$M = \left\{ m \in \mathbb{R}^{\{0,1\} \times \mathbb{R}} \text{ s.t. } m(0,0) = 0 \right\}.$$

Let  $F_{\tau,w}$  be a full support probability distribution over  $\{0,1\} \times \mathbb{R}$  and let  $(F_{p|\tau,w})_{(\tau,w) \in \{0,1\} \times \mathbb{R}}$  denote a set of full-support conditional probability distributions over  $p \in \mathbb{R}$ . The mechanism is run as follows: the agent submits a utility function  $m_i$ . A pair  $(\tau_i, w_i)$  and a price  $p_i$  are drawn according to  $F_{\tau,w}$  and  $F_{p|\tau_i,w_i}$ . If  $p_i \leq m_i(\tau_i, w_i)$ , then the agent is given allocation  $(\tau_i, w_i)$  and pays  $p_i$ . If  $p_i > m_i(\tau_i, w_i)$ , the agent is assigned  $(0,0)$  and makes no transfers. Because  $F_{\tau,w}$  as well as  $F_{p|\tau,w}$  have full-support, it is optimal for an agent to send message  $m_i(t) = V_t(\tau, w)$ . In turn, a mechanism is a *most informative incentivized trial* if and only if: (i) it elicits value function  $V_t(\tau, w)$ , and (ii), for any message  $m$ , the induced distribution over  $(\tau, w) \in \{0,1\} \times \mathbb{R}$  has full support.

Note that instead of eliciting preferences over a continuous domain  $\{0,1\} \times \mathbb{R}$ , the same methodology can be used to elicit preferences over a finite grid. The distribution  $F_{\tau,w}$  then needs to have full-support with respect to the grid of interest.

# Online Appendix for: Selective Trials: A Principal-Agent Approach to Randomized Controlled Experiments

By SYLVAIN CHASSANG, GERARD PADRÓ I MIQUEL, AND ERIK SNOWBERG

## PROOFS

**FACT 1** (full support sampling): *Consider a mechanism  $G = (M, \mu)$ . If there exists  $\xi > 0$  such that for all  $m \in M$ ,  $\pi(m) \in (\xi, 1 - \xi)$ , then, with infinite samples,  $G_0 \preceq G$ .*

**PROOF:**

The data  $\mathbf{d}_G$  can be broken in two subsamples,  $(d_G^{\sigma_0(i)})_{i \in \mathbb{N}}$  and  $(d_G^{\sigma_1(i)})_{i \in \mathbb{N}}$ , such that  $\sigma_0, \sigma_1$  are non-decreasing mappings from  $\mathbb{N}$  to  $\mathbb{N}$ , and for all  $i \in \mathbb{N}$ ,  $\tau_{\sigma_0(i)} = 0$  and  $\tau_{\sigma_1(i)} = 1$ . Since  $\forall m, \pi(m) \in [\xi, 1 - \xi]$ , we have that each such subsample is infinite and we can pick  $\sigma_1$  and  $\sigma_0$  to be strictly increasing from  $\mathbb{N}$  to  $\mathbb{N}$ . We define mapping  $h$  (such that  $h(\mathbf{d}_G) \sim \mathbf{d}_{G_0}$ ) as follows.

We use the notation  $h(\mathbf{d}_G) = (d_i^h)_{i \in \mathbb{N}}$ , where  $d_i^h = (m_i^h, p_i^h, \tau_i^h, y_i^h)$ . For every  $i \in \mathbb{N}$ , set  $m_i^h = \emptyset$ ,  $p_i^h = 0$ , and draw  $\tau_i^h$  as a Bernoulli variable of parameter  $\pi_0$ . Finally, set  $y_i^h = y_{\sigma_{\tau_i^h}(i)}$ . It is easy to check that indeed,  $h(\mathbf{d}_G) \sim \mathbf{d}_{G_0}$ .

**PROPOSITION 1** (most informative mechanisms): *Any strictly incentive-compatible mechanism  $G$  identifies at most value  $V_t$  (that is,  $V_t = V_{t'} \Rightarrow m_G(t) = m_G(t')$ ).*

*Whenever  $G$  identifies values  $V_t$  (that is,  $m_G(t) = m_G(t') \Rightarrow V_t = V_{t'}$ ) and satisfies full support ( $0 < \inf_m \pi(m)$  and  $\sup_m \pi(m) < 1$ ), then for any strictly incentive-compatible mechanism  $G'$ ,  $G' \preceq G$ .*

**PROOF:**

The proof of the first claim is very similar to that of Fact 1. Consider a mechanism  $G = (M, \mu_G)$  such that every player has a strictly dominant strategy. An agent with value  $V(t_i)$  chooses a message  $m_i$  to solve

$$\max_{m \in M} \pi(m)V(t_i) - \mathbb{E}_\mu[p_i | m_i = m].$$

This problem is entirely defined by player  $i$ 's value  $V(t_i)$ . Since a.e. player has a strictly optimal message, this problem has a unique solution for a.e. value.

We now construct a mapping  $h : \mathcal{D} \rightarrow \Delta(\mathcal{D})$  such that the data generated by  $G'$  can be simulated from data generated by  $G$  using mapping  $h$ . For simplicity

we describe the mapping  $h$  in the case where  $M$  is finite. Given  $\mathbf{d}_G$ ,  $h(\mathbf{d}_G)$  is generated as follows.

First, we break down the basic data  $\mathbf{d}_G$  in  $2 \times \text{card } M$  subsets, according to treatment  $\tau$  and the message  $m_G(V)$  corresponding to the value declared by the agent. Formally, for all  $m \in M$  and  $\tau \in \{0, 1\}$ , we define  $(d_G^{\sigma_{m,\tau}(i)})_{i \in \mathbb{N}}$  the ordered subsequence such that for all  $i$ ,  $m_G(V_{\sigma_{m,\tau}(i)}) = m$  and  $\tau_{\sigma_{m,\tau}(i)} = \tau$ . Since  $0 < \inf_m \pi(m) < \sup_m \pi(m) < 1$ , all these subsamples are infinite. Hence,  $\sigma_{m,\tau}$  can be chosen to be strictly increasing from  $\mathbb{N} \rightarrow \mathbb{N}$ . We use these subsamples to simulate data  $\mathbf{d}_{G'}$ .

Let us denote  $h(\mathbf{d}_G) = (d_i^h)_{i \in \mathbb{N}}$ . For all  $i \in \mathbb{N}$ ,  $d_i^h = (m_i^h, p_i^h, \tau_i^h, y_i^h)$ . We first set  $m_i^h = m_{G'}(V_i)$ . Then using  $\mu_{G'}(m_i^h)$ , we draw values  $\tau_i^h$  and  $p_i^h$ . Finally we set  $y_i^h = y_{\sigma_{m_i^h, \tau_i^h}(i)}$ . This defines  $h : \mathcal{D} \rightarrow \Delta(\mathcal{D})$ . It is easy to check that  $h(\mathbf{d}_G) \sim \mathbf{d}_{G'}$ .<sup>1</sup> This concludes the proof.

**FACT 2 (BDM Implementation):** *Whenever  $F_p$  has full support over  $[-V_{\max}, V_{\max}]$ , an agent with value  $V_t$  sends optimal message  $m_{BDM} = V_t$  and the BDM mechanism is a most informative mechanism.*

**PROOF:**

The fact that the BDM mechanism elicits values is well-known. Since  $F_p$  has full support over  $[-V_{\max}, V_{\max}]$ , assignment to treatment also satisfies full support and the second part of Proposition 1 implies that  $G_{BDM}$  is a most informative mechanism.

**PROPOSITION 2 (monotonicity):** *Consider a strictly incentive compatible mechanism  $G$ . If agents  $t$  and  $t'$  with values  $V_t > V_{t'}$  send messages  $m_G(t) \neq m_G(t')$ , then it must be that  $\pi(m_G(t)) > \pi(m_G(t'))$ .*

**PROOF:**

Agents of type  $t$  and  $t'$  are such that  $V_t > V_{t'}$  and  $m_G(t) \neq m_G(t')$ . Denote  $\pi(m) = \text{Prob}(\tau = 1|m)$  and  $p_{m_G} = \mathbb{E}_{\mu_G(\cdot|m)}[p]$ . By optimality of the message, it must be that

$$\begin{aligned} \pi(m_G(t))V_t - p_{m_G(t)} &> \pi(m_G(t'))V_t - p_{m_G(t')} \\ \pi(m_G(t'))V_{t'} - p_{m_G(t')} &> \pi(m_G(t))V_{t'} - p_{m_G(t)}. \end{aligned}$$

Adding the two inequalities yields that  $[\pi(m_G(t)) - \pi(m_G(t'))](V_t - V_{t'}) > 0$ , which implies that  $\pi(m_G(t)) > \pi(m_G(t'))$ .

**PROPOSITION 3 (sampling rates and incentives):** *For any mechanism  $G = (M, \mu)$  and  $\underline{\rho} < \bar{\rho}$  in  $(0, 1)$ , there exists a mechanism  $G' = (M, \mu')$  such that  $G \preceq G'$ , and for all  $m \in M$ ,  $\pi'(m) \in [\underline{\rho}, \bar{\rho}]$ .*

<sup>1</sup>Note that for the sake of notational simplicity, this construction ends up wasting data points by not taking consecutive elements from the subsamples. This is inconsequential here since we have infinitely many data points.

The following must also hold. Denoting the expected utility of type  $t$  sending message  $m$  in mechanism  $G'$  (including transfers) by  $U(t|m, G')$ , then

$$\max_{m_1, m_2 \in M} |U(t|m_1, G') - U(t|m_2, G')| \leq 2(\bar{\rho} - \underline{\rho})V_{\max}.$$

PROOF:

We begin with the first assertion. Given mechanism  $G = (M, \mu)$ , we define mechanism  $G' = (M, \mu')$  as follows:

$$\forall m \in M, \quad \mu'(m) = \begin{cases} \tau = 0, p = 0 & \text{with probability } \frac{\rho}{\bar{\rho} - \underline{\rho}} \\ \mu(m) & \text{with probability } \frac{\bar{\rho} - \rho}{\bar{\rho} - \underline{\rho}} \\ \tau = 1, p = 0 & \text{with probability } \frac{\rho}{\bar{\rho}} \end{cases}$$

Clearly, mechanism  $G'$  is strategically equivalent to mechanism  $G$ . The proof that  $G \preceq G'$  is omitted since it is essentially identical to that of Fact 1.

We now turn to the second assertion. Consider two messages  $m_1$  (optimally) sent by a type with value  $V_1$ , and  $m_2$  (optimally) sent by a type with value  $V_2$ . Let  $p_{G'}(m) = \mathbb{E}_{\mu_{G'}(\cdot|m)}[p]$ . We must have that

$$\begin{aligned} \pi_{G'}(m_1)V_1 - p_{G'}(m_1) &\geq \pi_{G'}(m_2)V_1 - p_{G'}(m_2) \\ \pi_{G'}(m_2)V_2 - p_{G'}(m_2) &\geq \pi_{G'}(m_1)V_2 - p_{G'}(m_1) \end{aligned}$$

within mechanism  $G'$ . These two inequalities yield that  $(\pi_{G'}(m_2) - \pi_{G'}(m_1))V_1 \leq p_{G'}(m_2) - p_{G'}(m_1) \leq (\pi_{G'}(m_2) - \pi_{G'}(m_1))V_2$ , which implies that  $|p_{G'}(m_2) - p_{G'}(m_1)| < (\bar{\rho} - \underline{\rho})V_{\max}$ . Hence the difference in utilities between sending two messages  $m_1$  and  $m_2$  for an agent with value  $V \in [-V_{\max}, V_{\max}]$  is  $|(\pi_{G'}(m_1) - \pi_{G'}(m_2))V - p_{G'}(m_1) + p_{G'}(m_2)| \leq 2(\bar{\rho} - \underline{\rho})V_{\max}$ .

PROPOSITION 4 (most informative mechanisms): *Any strictly incentive-compatible blind mechanism  $G$  identifies at most the mapping  $V_t(\phi)$  (that is,  $V_t(\phi) = V_{t'}(\phi) \Rightarrow m_G(t) = m_G(t')$ ).*

*If  $G$  identifies  $V_t(\phi)$  (that is,  $m_G(t) = m_G(t') \Rightarrow V_t(\phi) = V_{t'}(\phi)$ ) and satisfies  $\inf_{\phi, m} \mu(\phi|m) > 0$  then  $G' \preceq G$  for any strictly incentive-compatible mechanism  $G'$ .*

PROOF:

The proof of Proposition 4 is essentially identical to that of Proposition 1 and hence omitted.

PROPOSITION 5 (a test of ‘‘intention to change behavior’’):

*If  $e^*(\phi=0, t) = e^*(\phi=1, t)$ , then for all  $\varphi$ ,  $V_t(\phi=\varphi) = \varphi V_t(\phi=1)$ .*

*If  $e^*(\phi=0, t) \neq e^*(\phi=1, t)$ , then for all  $\varphi \in (0, 1)$ ,  $V_t(\phi=\varphi) < \varphi V_t(\phi=1)$ .*

PROOF:

The proof is given for the general case where there might be multiple optimal effort choices. Let  $V_t(\tau, e)$  denote the expected value of type  $t$  under treatment status  $\tau$  and when expending effort  $e$ . We have that

$$\begin{aligned} V_t(\phi) &= \max_{e \in E} \phi V_t(\tau=1, e) + (1 - \phi) V_t(\tau=0, e) \\ &\leq \phi \max_{e \in E} V_t(\tau=1, e) + (1 - \phi) \max_{e \in E} V_t(\tau=0, e). \end{aligned}$$

If  $\arg \max_{e \in E} V_t(\tau=1, e) \cap \arg \max_{e \in E} V_t(\tau=0, e) \neq \emptyset$ , the inequality is an equality and, since we normalized  $V_t(\phi=0) = 0$  we obtain that  $V_t(\phi) = \phi V_t(\phi=1)$ . Inversely, if  $\arg \max_{e \in E} V_t(\tau=1, e) \cap \arg \max_{e \in E} V_t(\tau=0, e) = \emptyset$ , the inequality is strict and  $V_t(\phi) < \phi V_t(\phi=1)$ .

PROPOSITION 6 (identifying perceived returns to effort): *For any value  $\phi$ ,*

$$\left. \frac{\partial V_t(\phi)}{\partial \phi} \right|_{\phi} = [q_t(\tau=1, e^*(\phi, t)) - q_t(\tau=0, e^*(\phi, t))] \times [u(y=1, t) - u(y=0, t)].$$

PROOF:

The result follows directly from applying the Envelope Theorem to (1).

PROPOSITION 7 (identifying perceived success rates):

$$\forall \tau, w, \quad \frac{\partial V_t(\tau, w)}{\partial w} = q_t(\tau, e^*(\tau, w, t)).$$

PROOF:

The result follows directly from applying the Envelope Theorem to (2).

FACT 3: *Assume that outcome  $y = 1$  yields strictly greater utility than  $y = 0$ , that is,  $u(y=1, t) > u(y=0, t)$ , and an agent perceives treatment to be beneficial:*

$$\forall e_0 \in E, \exists e_1 \in E \text{ s.t. } c(e_1, t) \leq c(e_0, t) \quad \text{and} \quad q_t(\tau=0, e_0) < q_t(\tau=1, e_1).$$

$$\text{Then, } w_{0,t} = \max\{w \mid V_t(\tau=1, w) = V_t(\tau=0, w)\}.$$

PROOF:

Whenever  $w = w_{0,t}$ , the agent is perfectly insured and  $V_t(\tau=1, w) = V_t(\tau=0, w)$  since access to the technology is valuable only in so far as it affects outcomes. We now show that whenever  $w > w_{0,t}$ ,  $V_t(\tau=1, w) > V_t(\tau=0, w)$ . The agent's value is

$$V_t(\tau, w) = \max_{e \in E} q_t(\tau, e)[u(y=1, t) - u(y=0, t) + w] + u(y=0, t) - c(e, t).$$

Let  $e_0^*$  be the agent's optimal effort level if  $\tau = 0$ . By assumption, there exists  $e_1$  such that  $c(e_1, t) \leq c(e_0^*, t)$  and  $q_t(\tau=1, e_1) > q_t(\tau=0, e_0^*)$ . Since  $w >$



$w_{0,t} = u(0, t) - u(1, t)$ , it follows that the agent gets strictly higher value under configuration  $(\tau = 1, e_1)$  than under configuration  $(\tau = 0, e_0^*)$ . This concludes the proof.

FACT 4: *Under mechanism  $G^{\pi,p}$  an agent of type  $t$  sends message  $k$  if and only if  $V_t \in [V_{k-1}, V_k]$ .*

PROOF:

Indeed,  $m_{G^{\pi,p}}(V) = k$  if and only if for all  $k' \neq k$ ,

$$(C1) \quad V\pi_k - p_k > V\pi_{k'} - p_{k'}.$$

For  $k' < k$ , this last condition is equivalent to  $V \geq \max_{k' < k} \{(p_k - p_{k'}) / (\pi_k - \pi_{k'})\}$ , which in turn is equivalent to  $V > V_{k-1}$ . Similarly, for  $k' > k$ , (C1) is equivalent to  $V_k > V$ . This concludes the proof.

#### A NUMERICAL EXAMPLE

This section illustrates the step-by-step process of inference from trial data, starting with a standard RCT, adding data from open selective trials, and concluding by adding both objective and subjective data from an incentivized trial.

We return to a setting where returns are two dimensional:  $R = (R_b, R_e)$ . As before, in the context of a water treatment product,  $R_b$  could be the baseline returns of using the water treatment product only when it is convenient to do so and  $R_e$  the returns to using it more thoroughly (for instance, bringing treated water when away from home). Success rates are given by:

$$q(\tau=0, e) = 0 \quad \text{and} \quad q(\tau=1, e) = R_b + eR_e,$$

where  $e \in \mathbb{R}_+$  is the agent's effort expenditure. An agent with type  $t$  has beliefs  $R_t = (R_{b,t}, R_{e,t})$  and maximizes  $\mathbb{E}_t[y] - c(e)$  where  $c(e) = \frac{e^2}{2}$ . The effort expended in an incentivized trial is thus  $e^*(w, t) = R_{e,t}(1+w)$ , which nests the effort decision of an open trial,  $e^*(w=0, t) = R_{e,t}$ .

Throughout, we illustrate the inference process by considering the case where each parameter has a low and high value:  $R_e, R_{e,t} \in \{1/4, 1/2\}$ ,  $R_b \in \{0, 1/8\}$  and  $R_{b,t} \in \{0, 3/32\}$ . Each element of a selective trial adds data which will narrow down the set of possible values.<sup>2</sup>

#### INFERENCE FROM AN RCT

An RCT identifies the average treatment effect,  $\widehat{\Delta} = R_b + R_e \times R_{e,t}$ . For the numerical values specified above, the possible outcomes are described in the following matrix

<sup>2</sup>For simplicity, we consider priors that put point masses on a few possible states. Unfortunately, such strong priors often result in degenerate inference problems. We computed the states to keep the inference problem well-defined and better reflect the mechanics of inference from a continuous state space. This accounts for our somewhat unusual parameter values.

	$R_e = 1/2$		$R_e = 1/4$	
	$R_{e,t} = 1/2$	$R_{e,t} = 1/4$	$R_{e,t} = 1/2$	$R_{e,t} = 1/4$
$R_b = 1/8$	$\hat{\Delta} = 3/8$	$\hat{\Delta} = 1/4$	$\hat{\Delta} = 1/4$	$\hat{\Delta} = 3/16$
$R_b = 0$	$\hat{\Delta} = 1/4$	$\hat{\Delta} = 1/8$	$\hat{\Delta} = 1/8$	$\hat{\Delta} = 1/16$

As illustrated by the matrix, if  $\hat{\Delta} \in \{1/16, 3/16, 3/8\}$  this identifies the returns of the technology  $(R_b, R_e)$ . However, treatment effects  $\hat{\Delta} \in \{1/8, 1/4\}$  are consistent with multiple true returns.<sup>3</sup> In particular, when  $\hat{\Delta} = 1/4$ , it may be that casual use of the water treatment product is not particularly effective ( $R_b = 0$ ), more thorough use is not particularly effective ( $R_e = 1/4$ ), or more thorough use is effective, but agents don't believe it is, and so do not expend much effort into using the water treatment product more thoroughly ( $R_e = 1/2, R_{e,t} = 1/4$ ).

#### INFERENCE FROM A SELECTIVE OPEN TRIAL

By Fact 1, open selective trials identify treatment effects  $\hat{\Delta}$ . Additionally, by Proposition 1, an open selective trial identifies the agent's willingness to pay for treatment  $V_t = R_{b,t} + R_{e,t}^2/2$ . To illustrate the value of this data, focus on the case where  $\hat{\Delta} = 1/4$ . As shown above, this is consistent with three different vectors of  $(R_b, R_e, R_{e,t})$ . Based on this, we illustrate the six possible values of  $V_t$  in the following matrix:

	$R_b = 0, R_e = 1/2, R_{e,t} = 1/2$	$R_b = 1/8, R_e = 1/2, R_{e,t} = 1/4$	$R_b = 1/8, R_e = 1/4, R_{e,t} = 1/2$
$R_{b,t} = 3/32$	$V_t = 7/32$	$V_t = 1/8$	$V_t = 7/32$
$R_{b,t} = 0$	$V_t = 1/8$	$V_t = 1/32$	$V_t = 1/8$

If  $V_t = 1/32$  the data from selective trials indicates  $R_{e,t} = 1/4 = e^*$ . As the treatment effect is  $\hat{\Delta} = 1/4$  the only consistent returns are  $R_b = 1/8$  and  $R_e = 1/2$ . If  $V_t = 7/32$ , there remains uncertainty, as the data is consistent with both  $(R_b = 0, R_e = 1/2)$  and  $(R_b = 0, R_e = 1/4)$ . Finally if  $V_t = 1/8$ , the data is consistent with any of the states  $(R_b, R_e, R_{e,t})$  that produce  $\hat{\Delta} = 1/4$ . That is to say that even in this limited example, data from a selective open trial (and, hence, MTEs) may not help in identifying underlying returns. We now turn to how incentivized trials allow us to infer whether effort, or returns to effort, are low.

<sup>3</sup>For example,  $(R_b = 0, R_e = 1/2, R_{e,t} = 1/2)$ ,  $(R_b = 1/8, R_e = 1/2, R_{e,t} = 1/4)$  and  $(R_b = 1/8, R_e = 1/4, R_{e,t} = 1/2)$  are all consistent with  $\hat{\Delta} = 1/4$ .

Note that agents' beliefs may be self-confirming. For instance, an agent who believes that effort has high returns,  $R_{e,t} = 1/2$ , who observes  $\hat{\Delta} = 1/4$  will continue to believe returns are high, even though this data could be generated by  $R_e = 1/4$ . Such self-confirming beliefs are frequent in the experimentation and social learning literatures (???)

## INFERENCE FROM AN INCENTIVIZED TRIAL

Incentivized trials yield:

$$\widehat{\Delta}(w) = R_b + R_e \times R_{e,t}(1+w) \quad \text{and} \quad V_t(\tau=1, w) = R_{b,t}(1+w) + \frac{[R_{e,t}(1+w)]^2}{2}.$$

As an open selective trial already identifies  $V_t = V_t(w=0) = R_{b,t} + R_{e,t}^2/2$  and  $\widehat{\Delta} = \widehat{\Delta}(w=0) = R_b + R_e \times R_{e,t}$ , by eliciting valuations and treatment effects for a small  $w$ , the principal can also identify  $\left. \frac{\partial V_t(\tau, w)}{\partial w} \right|_{w=0} = R_{b,t} + R_{e,t}^2$  and  $\left. \frac{\partial \widehat{\Delta}(w)}{\partial w} \right|_{w=0} = R_e \times R_{e,t}$ . With this data the principal can identify:

$$R_{e,t} = \left[ 2 \left( \left. \frac{\partial V_t}{\partial w} \right|_{w=0} - V_t(w=0) \right) \right]^{1/2},$$

and thus, the rest of the unknown parameters:  $R_e = \left. \frac{\partial \widehat{\Delta}(w)}{\partial w} \right|_{w=0} / R_{e,t}$ ,  $R_{b,t} = \left. \frac{\partial V_t(\tau, w)}{\partial w} \right|_{w=0} - R_{e,t}^2$ ,  $R_b = \widehat{\Delta} - R_e \times R_{e,t}$ . The same information can be identified in a mathematically simpler, but more data intensive, way by identifying  $w_{0,t}$  and the empirical quantities associated with that value.

Altogether, incentivized selective trials allow us to identify both the true returns  $(R_b, R_e)$  and the agents' beliefs  $(R_{b,t}, R_{e,t})$ . Thus, in this example, data from a selective incentivized trial allows a principal to determine how effective casual and thorough use of the water treatment product is, without having to observe individual agents' usage. This is possible, as eliciting each agent's indirect preferences over the water treatment product, and bonuses associated with staying healthy, allows the principal to infer the agents' beliefs about the effects of casual and more thorough usage. This, in turn, allows the principal to infer behavior and identify the deep structural parameters determining the product's effectiveness, as well as how beliefs about effectiveness lead to different outcomes.