

# Narrow Inference and Incentive Design<sup>\*</sup>

Alexander Clyde<sup>†</sup>

October 22, 2024

[\[CLICK HERE FOR LATEST VERSION\]](#)

## Abstract

There is evidence that people struggle to do causal inference in complex multidimensional environments. This paper explores the consequences of this in a principal-agent setting. A principal chooses a mechanism to screen an agent. The agent makes choices on multiple dimensions, and infers the effect of each action separately without properly controlling for the other actions. I characterize the principal’s optimal mechanism when facing an agent who does such ‘narrow’ inference, and contrast it with their optimal mechanism when the agent is fully rational. I show that the principal can benefit from narrow inference, and identify cases where this holds.

**Keywords:** Behavioural Mechanism Design, Screening, Narrow Bracketing, Misspecified Models.

**JEL Classification:** D90, D02, D82

---

<sup>\*</sup>I am grateful to my advisors Ran Spiegler and Philippe Jehiel, for their support and guidance at all stages of this project. I also thank Martin Cripps, Nathan Hancart, Deniz Kattwinkel, Vasiliki Skreta, Duarte Gonçalves, Michael Thaler, Mark Armstrong, Nikita Roketskiy, Elena Ashtari Tafti and Dajana Xhani for helpful comments. I gratefully acknowledge financial support from an Economic and Social Research Council studentship (1923039).

<sup>†</sup>Department of Economics, University College London; Address: Drayton House, 30 Gordon Street, London, WC1H 0AX; Email: [alexander.clyde.econ@gmail.com](mailto:alexander.clyde.econ@gmail.com) or [alex.clyde@ucl.ac.uk](mailto:alex.clyde@ucl.ac.uk)

# 1 Introduction

Understanding the incentives we face requires understanding how the many choices we make affect outcomes we care about. For example, in the labour market workers have to form beliefs about how their choices of effort, occupation and education affect their final wage. This can require making sophisticated causal inference from any available data. Economic models often assume that people form these beliefs and makes these choices jointly in one over-arching decision problem. However, work in experimental economics suggests that people both fail to consider their choices jointly and struggle to correctly infer the effects of their actions<sup>1</sup>. In line with this, I formulate a model of bounded rationality in which people form beliefs separately in a piecewise way for each decision they face. I call this model of belief formation ‘narrow inference’.

Taking into account narrow inference matters for how we should design incentives. Consider a policymaker who is determining how much to subsidize university education. A worker has to decide whether to make human capital investments in education and work experience. Assume that the workers beliefs about how these actions affect future net earnings have to be consistent with actual data. A worker who makes narrow inference forms beliefs about the effect of education and work experience on earnings separately. In forming beliefs about the effect of any individual action in this way, they fail to control for other dimensions of action. This leads to a confounding bias that distorts the worker’s perception of the size of earnings benefits of obtaining education and work experience. The extent of this misperception affects how the policymaker wants to design both education subsidies and earnings taxation.

In this paper I analyze such incentive design problems. I consider a principal-agent screening model with a principal who has full understanding of their problem, but an agent who only performs narrow inference. I explore how a principal would design an incentive mechanism if they knew the agent had this form of bounded rationality. I obtain a result characterizing the principal’s optimal mechanism with an agent who makes narrow inference as the solution to a zero-sum game. I use this result to demonstrate in what cases the principal benefits from the agent making narrow rather than fully rational inference,

---

<sup>1</sup>In the literature review I discuss work on ‘narrow bracketing’ from behavioural economics, and work in experimental economics on correlation neglect and causal misperceptions.

and in what cases the principal does not. I then explore what happens when the number of dimensions of action grows large. In doing so, I obtain a result that demonstrates the effect of narrow inference on agent’s perception of incentives in an optimal mechanism can be quantitatively large.

In the screening model, the agent faces a binary decision problem on whether to take an action or not on multiple dimensions. The principal chooses a function mapping the agent’s actions to an outcome, and the agent needs to infer how their actions affect the outcome. There is a large population of agents who differ according to a single dimensional type variable that affects the immediate costs and benefits of the actions. The principal’s choice of outcome function screens different types of agents into choosing different actions. In the absence of bounded rationality this problem is standard, with the outcome providing a zero-sum transfer of utility between the agent and principal. The principal and the agent derive the opposite utility from the outcome, and also derive a potentially different immediate utility from the agent’s actions that is not mediated by the outcome. The immediate utility of the actions for both the agent and the principal is additively separable across dimensions, but the principal can choose an outcome function where there is an interactive effect.

An agent who makes narrow inference calculates the effect of each dimension’s action on the outcome separately. Their beliefs about the effect of an action on a given dimension must be consistent with data on the population level average outcome conditional on that action. The difference between the average population level outcomes between any two actions in a given dimension is then used to estimate the relative effect of each action on the outcome. This is a naive way to estimate the ‘treatment effect’ of any action. It can lead to incorrect expectations if the distribution over actions are correlated across dimensions, something that is possible due to joint dependence on the type variable. The estimated effect of the action on the outcome is then biased from confounding. This holds even when the true outcome function is additively separable, as the inferential failure of the agent involves neglecting the correlation in the data and not just in misspecifying the functional form of the outcome function.

In what follows, I develop the policymaker-worker example to illustrate narrow inference and to demonstrate how the principal might benefit from agents using narrow

inference.

**Example 1.** The agent has two dimensions of action, whether to obtain work experience  $a_1 = 1$  or not  $a_1 = 0$  and whether to study at university  $a_2 = 1$  or not  $a_2 = 0$ . The principal, a policymaker, uses general taxation and a relative subsidy for education to determine a future net earnings schedule that depends jointly on these two actions  $t : A_1 \times A_2 \rightarrow \mathbb{R}$ . There are three types of agent  $s \in \{0, 1, 2\}$ . The probabilities of the types are denoted  $p_0, p_1, p_2 \in (0, 1)$  respectively. The agent's utility depends on their type  $s$ , their actions  $a_1, a_2$  and the outcome  $t$ .

$$t(a_1, a_2) - (3 - s)(a_1 + a_2)$$

Suppose that the principal wants to implement that the type  $s = 0$  chooses neither action, the type  $s = 1$  obtains work experience  $a_1 = 1$  but not a university education  $a_2 = 0$ , while the highest type  $s = 2$  obtains both  $a_1 = a_2 = 1$ . An outcome function that implements this must satisfy the incentive constraints. The first ensures that type  $s = 1$  chooses  $(1, 0)$  over  $(0, 0)$  and the second ensures type  $s = 2$  chooses action  $(1, 1)$  over  $(1, 0)$ .

$$t(1, 0) - 2 \geq t(0, 0)$$

$$t(1, 1) - 2 \geq t(1, 0) - 1$$

Choosing  $t$  such that these incentive constraints bind allows the principal to minimize the earnings paid to types  $s = 1$  and  $s = 2$ . This means  $t(1, 1) > t(1, 0) > t(0, 0)$ . These local incentive constraints binding suffices for all incentive constraints to hold.

Now consider if the agent used narrow inference. They expect the earnings from any action to be the average population level earnings of those who have taken that action. For an agent of type  $s = 1$  making narrow inference to choose work experience thus requires that

$$\frac{p_1}{p_1 + p_2}t(1, 0) + \frac{p_2}{p_1 + p_2}t(1, 1) - 2 \geq t(0, 0)$$

Since  $t(1, 1) > t(1, 0)$ , this narrow perception of the expected earnings benefit of work experience is biased upward from the true effect. It fails to adjust for the fact that a

type  $s = 1$  agent who is on the margin between obtaining experience does not obtain a university education and thus has lower future earnings than the average.

Similarly for the type  $s = 2$  agent to want to obtain an education under narrow inference requires that

$$t(1, 1) - 1 \geq \frac{p_1}{p_1 + p_0} t(1, 0) - \frac{p_0}{p_1 + p_0} t(0, 0)$$

As  $t(1, 0) > t(0, 0)$ , there perception of the earnings from not obtaining education is biased downwards. The expected earnings for those who do not obtain education mixes the earnings of those who get work experience and those who do not. It is therefore less than the earnings obtained by the types on the margin of obtaining an education, type  $s = 2$ , all of whom obtain work experience.

This upward bias in the incentives the agent perceives allows the principal to implement the same action choices for each type while obtaining higher tax revenue by reducing net earnings across the type distribution.

My analysis of the design problem proceeds as follows. First, in order to contrast the principal's optimal mechanism when agents make narrow inference to that with fully rational agents, I first state results describing the principal's optimal mechanism in the rational benchmark. This involves applying results in single-dimension monopolistic screening problems adapted to the multidimensional action setting. Under a standard regularity assumption, the principal's optimal mechanism is fully separable across dimensions. Facing such a mechanism, the agent chooses a strategy that on each dimension selects the action if and only if their type is above a dimension-specific threshold.

I then obtain a result characterizing the principal's optimal mechanism under narrow inference. In this characterization, the principal plays a zero-sum game against an adversarial player who can shrink the immediate utility from actions on any dimension by some constant factor that lies between one and zero, with the shrinkage factors summing to one across dimensions. The principal's optimal mechanism under narrow inference then solves the same problem as in the rational benchmark except with the shrunk immediate utilities. The result allows us to both solve for the principal's optimal mechanism for specific parameterizations, and also enables us to easily compare how the mechanism

differs from that with rational agents.

I show using the characterization result what effect narrow inference has on this threshold strategy and on the principal's welfare relative to the rational benchmark. This involves considering two different cases; one where the agent's actions have immediate utility costs to the agent but benefits to principal on all dimensions, and one where all actions have immediate benefits to the agents but costs to the principal. The first case fits the policymaker-worker example, where human capital investments have costs to the worker but social benefits for the principal. For an example that fits the second case, consider an environmental agency that is designing regulations for a firm. The firm has two divisions, and each division makes a production decision. The production also causes a negative externality via pollution that the regulator wants to abate. The regulator can tax production from each division separately. Although each division has the same objective to maximize the overall profits of the firm, they make separate narrow inference about how their own production decision affects the total production tax that the firm has to pay.

In the case where actions immediately cost the agent but benefit the principal, when facing the principal's optimal mechanism the agent's thresholds are lower under narrow inference than when they are rational and the principal is always at least as well-off under narrow inference. Lower thresholds mean a greater proportion of types take the action on any dimension. On the other hand, when actions immediately cost the principal but benefit the agent then thresholds are higher and the principal is always weakly worse-off under narrow inference. These effects are the consequence of the agent overestimating the causal effects of their actions on outcomes due to confounding bias. Taking the action on one dimension is associated with taking the action on others, as both are chosen by higher type agents, and when agents neglect this their estimates are biased upwards. I show that the principal's gains and losses from narrow inference are purely due to this confounding, as under narrow inference the principal's optimal mechanism has a separable outcome function just as in the rational case.

The effect of narrow inference can be quantitatively large. Assume that the immediate utility of the actions is identical across dimensions for both the principal and agent. In this case, under the principal's optimal mechanism there is a single threshold above which

all types of agent take the action on all dimensions and below which they take no action on any dimension. I show that if total immediate utilities are split over  $n$  dimensions, there is a finite size of the dimensionality above which either all types of agents take the actions on all dimensions, or no agents take the actions on any dimensions. In the first case, the principal is able to implement this strategy at arbitrarily small cost in terms of outcomes. This result has relevance to the organization design question of whether top management wants to split or merge different divisions of their firm.

Finally, I explore extensions that vary the form of the participation constraints and consider the consequences of dropping the regularity assumption. In the main model, I assume the agent also makes narrow dimension-by dimension decisions on whether to participate in the mechanism. This is motivated by an interpretation that non-participation involves resorting to an outside option where the agent takes no action and obtains a default outcome. The agent treats the non-participation decision as they treat decisions within the mechanism. They believe that they can use the outside option narrowly, taking it on some dimensions but continuing to participate in the mechanism on others. I consider a version of the model where the agent instead perceives the participation decision as a joint one, taking the sum of perceived narrow utility across dimensions as the value of participation. I demonstrate how we can modify the characterization result under this alternative participation constraint. The results on the principal’s welfare and on large dimensionality still hold.

## Literature Review

My paper builds on several distinct but related strands of literature. Experimental work in psychology and economics documenting that people make inferential errors similar to narrow inference. [Enke and Zimmermann \(2019\)](#) find subjects fail to adjust for correlation between multiple information sources. Similar logic extends to predictive tasks, in [He and Kučinskas \(2024\)](#) subjects’ forecasting performance deteriorates when information from a single variable is split into two. [Fernbach et al. \(2010\)](#) present evidence suggesting people focus narrowly on a few variables when trying to make causal predictions. In line with this, [Graeber \(2023\)](#) finds subjects ignore the effect of variables that are not directly involved in a predictive task despite these variables containing valuable information.

Narrow inference involves causal misperceptions, but also thinking about decision problems narrowly. The literature on narrow bracketing considers decision makers who break decision problems into smaller sub-problems without accounting for how these decisions interact in the larger joint problem. Work in this area; [Tversky and Kahneman \(1981\)](#), [Thaler \(1985\)](#), [Thaler \(1999\)](#), [Read et al. \(1999\)](#), [Rabin and Weizsäcker \(2009\)](#), has both documented evidence for and explored the theoretical implications of narrow decision making. Recent work exploring theoretical foundations for narrow behaviour includes [Kőszegi and Matějka \(2020\)](#), who use a model of costly information acquisition to explain both mental accounting and naive diversification. [Lian \(2021\)](#) builds a theory of ‘narrow thinking’, which models decision makers as playing an incomplete information game between multiple-selves. In both these papers, narrow behaviour arises from coordination frictions or costs for agents who have otherwise rational beliefs. Under the type of narrow inference considered in this paper, the agent is able to perfectly coordinate their actions, but has distorted expectations about the effect of actions due to having a misspecified narrow causal model.

In modelling agents with narrow causal perceptions, this paper builds on work studying decision making by agents using misspecified models of how action choices map into consequences. There is a growing literature on the Berk-Nash Equilibrium of [Esponda and Pouzo \(2016\)](#), a solution concept founded as the limit of a process of misspecified learning; [Heidhues et al. \(2018\)](#), [Frick et al. \(2020\)](#), [Bohren and Hauser \(2021\)](#), [Fudenberg et al. \(2021\)](#). Another connected literature is that on modelling causal misperceptions using Bayesian Networks; [Spiegler \(2016\)](#), [Eliaz and Spiegler \(2020\)](#). [Schumacher and Thyssen \(2022\)](#) use this Bayesian Network approach in a principal-agent moral hazard problem where the agent has causal misperceptions of how their actions map into output. In [Eliaz and Spiegler \(2024\)](#) a Bayesian Network formalism is used to model the design of narratives for misspecified news consumers by media organizations.

Earlier work on design when agents misperceive incentives by [Rubinstein \(1993\)](#) and [Piccione and Rubinstein \(2003\)](#) explores monopolistic pricing when customers have a coarse misperception of any pricing strategy. For more general settings, [Jehiel \(2005\)](#) develops an equilibrium concept for extensive form games —Analogy Based Expectation Equilibrium (ABEE)— in which players have coarse misperceptions of other players’



strategies. The behaviour of the principal and the agent under narrow inference can be formulated as an ABEE of an extensive form game, and I discuss this in more detail in Section 5.3. The first papers to explicitly apply the ABEE concept to design problems are Jehiel (2011) and Jehiel and Mierendorff (2024). In Jehiel (2011), an auction designer manipulates bidders who do not perceive how the distribution of bids varies with different auction formats and the identities of different bidders. Similarly, in Jehiel and Mierendorff (2024) a proportion of bidders form beliefs about how signals of their own valuation of an item vary with opponents bids in a way that neglects correlation between their own signal and the signals of the other bidders.

In contributing to the small literature on mechanism design where agents use mis-specified models, this paper also contributes to a larger literature on mechanism design that takes into account agents' limited rationality in a variety of other dimensions. A detailed review can be found in Kőszegi (2014). This includes work in contract theory (Eliaz and Spiegler, 2006), (Heidhues and Kőszegi, 2010), (Herweg et al., 2010) and optimal taxation (O'Donoghue and Rabin, 2006), (Spinnewijn, 2015), (Farhi and Gabaix, 2020), (Lockwood, 2020).

## 2 Model

An agent faces a multidimensional decision. Let  $A = \{0, 1\}^n$  be the agent's set of feasible action profiles. I refer to  $i \in \{1, \dots, n\} \equiv N$  as a dimension, such that  $a_i$  is the agent's action in dimension  $i$ . The agent has a type that lies in a bounded interval  $s \in S \equiv [0, 1]$ . This type is drawn from an atomless distribution that admits a density  $p(s)$  such that  $p(s) > 0$  for all  $s \in S$ . Denote the cdf of the distribution by  $P(s) = \int_0^s p(\tilde{s})d\tilde{s}$ .

The dimension  $i$  action  $a_i$  generates an immediate utility  $v_i(s)a_i$ , where  $v_i(s)$  is strictly increasing, continuously differentiable in  $s$  and can be positive or negative. In addition to the immediate utility, the agent receives utility from an outcome  $t$  that needs to be inferred. The utility of an agent of type  $s$ , choosing action  $a$  with outcome  $t \in \mathbb{R}$  is

$$u(s, a, t) = \sum_{i \in N} v_i(s)a_i + t \tag{1}$$

The principal derives benefit/costs  $w_i a_i$  from an action  $a_i$ , where  $w_i \in \mathbb{R}$ . The outcome  $t \in \mathbb{R}$  represents a zero-sum transfer of surplus between the agent and the principal. The principal's payoff given actions  $a$  and outcome  $t \in \mathbb{R}$  is

$$W(a, t) = -t + \sum_{i \in N} w_i a_i \quad (2)$$

Throughout the paper, I make the following standard regularity assumption on the type distribution. In Section 5.2, I explore the implications of dropping this assumption.

**Assumption 1.** For every dimension  $i \in N$

$$\phi_i(s) = v_i(s) - \frac{1 - P(s)}{p(s)} v'_i(s)$$

is strictly increasing in  $s \in S$ . We refer to this property as increasing virtual values (IVV).

## 2.1 Mechanisms

I focus on a natural class of indirect mechanisms. Before the agent takes any actions the principal commits to a mechanism, which consists of a function mapping actions to outcomes  $t : A \rightarrow \mathbb{R}$ . After learning their type, the agent chooses a distribution over actions according to a strategy  $g : S \rightarrow \Delta(A)$ . The marginal distribution over actions in dimension  $i$  is denoted by  $g_i(a_i|s) = \sum_{a_{-i} \in A_{-i}} g(a_i, a_{-i}|s)$ .

For an outcome function  $t \in \mathbb{R}^A$ , given a strategy  $g$  the expected payoff for the principal is

$$W(t, g) = \int_0^1 \sum_{a \in A} [-t(a) + \sum_{i \in N} w_i a_i] g(a|s) p(s) ds \quad (3)$$

The restriction to this class of mechanisms is simple to reconcile with narrow inference. Under narrow inference, the agent perceives the outcome as measurable only with respect to their own actions. Suppose the principal could choose a more general mechanism in which the outcome function varied with an arbitrary message space as well as the actions. The principal could present information on how the outcome varies with more finely grained messages, drawing the agent's attention to the joint multidimensional nature of

their problem and undoing the narrow inference.

In Section 3, I show the restriction makes no difference to the analysis of the principal's optimal mechanism in the rational case. Under the optimal mechanism the agent chooses a strategy that is deterministic, and as such can be implemented with an outcome function that only depends on the chosen action. I consider a more general class of mechanisms where the principal can force the agent to randomize over actions in Section 5.1.

## 2.2 Model Interpretations

The model allows actions to have both a positive and negative effect on payoffs. The sign of  $v_i(s)$  determines whether a type  $s$  agent has immediate positive utility from action  $a_i = 1$  or immediate disutility. Likewise, the direct effect of actions on the principal's payoff can be positive ( $w_i \geq 0$ ) or negative ( $w_i \leq 0$ ).

The outcome function  $t$  can be interpreted as the division of the social surplus or costs from actions. The total social surplus or costs of the actions  $a$  is given by  $\sum_{i \in N} w_i a_i$ . The outcome function is then the share of that social surplus the agents receive, or the share of the social cost they must bear.

This framework can capture the stories given in the introduction. Suppose the agent is a multi-divisional firm and the principal a regulator. In this story,  $v_i(s) > 0$  represents the immediate payoff benefit to the type  $s$  firm of producing output  $a_i = 1$ , while  $w_i < 0$  represents the social cost of pollution caused by production. The regulator then chooses a tax on production  $t$ .

In another story the agent is a worker and the principal a policymaker. Here  $v_i(s) < 0$  is the cost of studying or gaining work experience and  $w_i > 0$  is the ultimate social benefits of the human capital the worker acquires. The policymaker then chooses a net earnings schedule  $t$ , which they can determine using both general taxation and the relative tax or subsidy on education.

## 2.3 Rational Inference

Given a strategy  $g$ , write the expected utility of an agent of type  $s$  as

$$U(s) = \sum_{a \in A} g(a|s) u(s, a, t(a)) = \sum_{a \in A} g(a|s) \left[ \sum_{i \in N} v_i(s) a_i + t(a) \right] \quad (4)$$

Incentive Compatibility (IC) of strategy  $g$  under outcome function  $t \in \mathbb{R}^A$  requires that  $g$  is a best response to  $t$ . This means for any  $s \in S$ ,  $a \in \text{supp}\{g(\cdot|s)\}$  and any  $a' \in A$

$$\sum_{i \in N} v_i(s) a_i + t(a) \geq \sum_{i \in N} v_i(s) a'_i + t(a') \quad (5)$$

The agent always has the option of not participating in the mechanism, taking the actions  $a = 0$  and obtaining a baseline outcome. Normalize the utility of this baseline outcome to zero. We then have the following participation constraint; for all  $s \in S$

$$U(s) \geq 0 \quad (6)$$

## 2.4 Narrow Inference

Given a strategy  $g$ , an unconditional distribution over actions in  $A$  is induced as follows.

$$g(a) = \int_0^1 g(a|s) p(s) ds \quad (7)$$

Let the marginal over an action in dimension  $i$  be denoted  $g_i(a_i) = \sum_{a_{-i} \in A_{-i}} g(a_i, a_{-i})$ . I use the terms action distribution and strategy interchangeably throughout the paper.

The agent forms narrow perceptions of the mechanism's outcome function. In particular an agent believes when taking a decision in dimension  $i$  that in expectation they will receive  $\bar{t}_i(a_i)$  if they take action  $a_i$ . When  $g_i(a_i) > 0$  we require that this expectation is consistent with the actual conditional expectation of outcomes given  $a_i$ .

$$\bar{t}_i(a_i) = \sum_{a_{-i} \in A_{-i}} \frac{g(a_i, a_{-i})}{g_i(a_i)} t(a_i, a_{-i}) \quad (8)$$

Denote  $\bar{t}(a) = (\bar{t}_i(a_i))_{i=1}^n \in \mathbb{R}^n$  and  $\bar{t} = (\bar{t}(a))_{a \in A} \in \mathbb{R}^{2n}$ . When  $g_i(a_i) = 0$ ,

we allow  $\bar{t}_i(a_i)$  to take on arbitrary values. Analogously to the rational benchmark, these ‘off-path’ actions do not affect the principal’s objective and  $\bar{t}$  can be set to ensure incentive compatibility. Henceforth, I refer to agents performing narrow inference as ‘narrow agents’.

A narrow agent imposes an additively separable form on their estimate of the outcome function using  $\bar{t}$ . This gives them the following perceived expected utility from strategy  $g$  when they are of type  $s \in S$ .

$$\bar{U}(s) = \sum_{i \in N} \sum_{a_i \in A_i} g_i(a_i|s)[v_i(s)a_i + \bar{t}_i(a_i)] = \sum_{i \in N} \bar{U}_i(s) \quad (9)$$

Where

$$\bar{U}_i(s) = \sum_{a_i \in A_i} g_i(a_i|s)[v_i(s)a_i + \bar{t}_i(a_i)] \quad (10)$$

denotes the narrow perceived expected utility of type  $s$  in dimension  $i$ . A strategy  $g$  is *narrow incentive compatible* (NIC) if for any dimension  $i \in N$ , type  $s \in S$  and actions  $a_i \in \text{supp}\{g_i(\cdot|s)\}$ ,  $a'_i \in A_i$

$$v_i(s)a_i + \bar{t}_i(a_i) \geq v_i(s)a'_i + \bar{t}_i(a'_i) \quad (11)$$

On each dimension, there is a baseline outcome that the agent can always obtain even when both not participating in the mechanism and taking the zero action. Normalize the utility of this baseline outcome to zero. We then have dimension by dimension *narrow participation constraints*; for all  $s \in S$ ,  $i \in N$

$$\bar{U}_i(s) \geq 0 \quad (12)$$

This fits the following interpretation; the agent believes they can reject any additional effect on the outcome resulting from participation in the mechanism separately on each dimension, whilst still obtaining the effect from participation on other dimensions. This is in line with the agent believing the true outcome function is additive. Although this appears to add participation constraints relative to the rational benchmark, in practice

it does not. In Section 3, I show that the outcome function in the principal's optimal mechanism with a rational agent is additively separable and thus also satisfies these dimension-by-dimension constraints. I consider alternative participation constraints in Section 5.1.

### 3 Rational Benchmark

With rational agents, we have a screening problem with a single dimension of type but multiple dimensions of action. I restate existing results adapted to our setting<sup>2</sup>.

It will be shown that the agent's strategy under the principal's optimal mechanism with rational agents takes a *threshold form* where there is a potentially different threshold  $\hat{s}_i \in S$  on each dimension such that  $g_i(1|s) = \mathbb{1}\{s \geq \hat{s}_i\}$ . Let the vector of thresholds across dimensions be denoted  $\hat{s} = (\hat{s}_i)_{i \in N} \in S^n$ . We can characterize the principal's problem in terms of choosing these thresholds. Denote the value of the principal's objective under threshold strategy  $\hat{s}$  by  $W(\hat{s})$ .

**Proposition 1.** *Assume the IVV assumption holds. The principal maximizes their objective over all IC mechanisms that satisfy the participation constraint if and only if they choose an outcome function implementing a threshold strategy that solves the following problem.*

$$\max_{\hat{s} \in S^n} W(\hat{s}) = \sum_{i \in N} \int_{\hat{s}_i}^1 (\phi_i(s) + w_i) p(s) ds \quad (13)$$

*The principal's value under an objective maximizing mechanism can be achieved by an additively separable outcome function*

$$t(a_1, \dots, a_n) = \sum_{i \in N} t^i(a_i) \quad (14)$$

$$t^i(0) = 0, t^i(1) = -v_i(\hat{s}_i) \text{ for all } i \in N \quad (15)$$

*Proof.* In Appendix □

Thus, the principal's optimal outcome function can be treated as the sum of separate outcome functions, one for each dimension. This does not result directly from IC, but

---

<sup>2</sup>In particular Proposition 3.1 of [Carroll \(2017\)](#).

rather from the optimality for the principal of implementing a threshold strategy when IVV holds. We will see in Section 5.2 that without IVV it can be optimal for the principal to choose a non-separable outcome function, and that under such an outcome function a non-threshold strategy is implemented and thus IC.

## 4 Narrow Agents

The solution to the principal's design problem with narrow agents can be characterized as a zero-sum game between the principal and an adversarial player. In this game, the principal faces the same design problem as in the rational benchmark except the immediate utilities are scaled by some factor  $\beta_i \in [0, 1]$  in each dimension. The immediate utility in dimension  $i$  is then  $\beta_i v_i(s)$ , and the scaling factors sum to one across dimensions  $\sum_{i \in N} \beta_i = 1$ . The principal chooses a mechanism to maximize their objective while the adversarial player simultaneously chooses the scaling factors to minimize the value of the objective. The result shows that the agent's strategy and value of the principal's objective under the principal's optimal mechanism with narrow agents coincide with those that arise as the solution to this zero sum game with rational agents.

### 4.1 Main Characterization Result

The characterization result is stated as follows.

**Theorem 1.** *Assume the IVV assumption holds. The principal maximizes their objective over all NIC mechanisms that satisfy the narrow participation constraints if and only if they choose an outcome function that implements a threshold strategy that solves*

$$\min_{\beta \in [0,1]^n: \sum_{i \in N} \beta_i = 1} \max_{\hat{s} \in S^n} \overline{W}(\hat{s}; \beta) = \max_{\hat{s} \in S^n} \min_{\beta \in [0,1]^n: \sum_{i \in N} \beta_i = 1} \overline{W}(\hat{s}; \beta) \quad (16)$$

with the value of the principal's objective given by

$$\overline{W}(\hat{s}; \beta) = \sum_{i \in N} \int_{\hat{s}_i}^1 (\beta_i \phi_i(s) + w_i) p(s) ds \quad (17)$$

*Proof.* [In Appendix](#)

□

To see some intuition for this result, consider the case with two dimensions and symmetric immediate utility  $v_i(s) = v(s)$  and principal's direct utility  $w_i = w$  from actions across all dimensions  $i \in N$ . With a rational agent, from Proposition 1 the principal's optimum sets a threshold  $\hat{s}$  so that the action is taken on all dimensions for all types above and the zero action is taken on all dimensions for all types below. From Proposition A.2, this optimum is induced with an outcome function that is additive and identical across dimensions  $t(a_1, a_2) = t^1(a_1) + t^2(a_2)$  and  $t^1(1) = t^2(1) = \tilde{t}(1)$ ,  $t^1(0) = t^2(0) = 0$ . With a narrow agent, under the same strategy and outcome function the agents double count the effect of each action on the outcome. In each dimension, they believe that the outcome resulting from taking the action  $a_i = 1$  is  $2 \cdot \tilde{t}(1)$  and the outcome resulting from  $a_i = 0$  is 0. This double-counting is the result of confounding neglect; the agent fails to adjust for the fact that every type who takes action  $a_1 = 1$  also takes action  $a_2 = 1$ . The principal then has to half the size of the the difference in outcomes in order to maintain the same thresholds  $\frac{1}{2}\tilde{t}(1)$ . This has the same effect as scaling the immediate utilities down by  $\frac{1}{2}$  in each dimension when the agent is rational.

The result extends this logic to asymmetric cases. It allows us to both solve for the principal's optimal mechanism with narrow agents and also demonstrates the connection between any given problem with narrow agents to the rational benchmark. I use the characterization to obtain additional results. I give conditions under which the principal does and does not benefit from facing narrow over rational agents, and how the implemented strategy changes between the two cases. I then explore the effect of symmetric immediate utility across dimensions and what happens when the number of dimensions grows large. First, I present some preliminaries that are used in obtaining the characterization.

## 4.2 Preliminaries for Characterization Result

I first obtain a result characterizing how narrow expected utilities relate to the implemented strategy. In particular we obtain a version of the envelope theorem for each dimension separately. I then consider which beliefs can be induced by an outcome function. I show that beliefs must satisfy a statistical correctness constraint, and that for any deterministic threshold strategy there is a valid additively separable outcome function that induces beliefs such that the strategy is NIC.



Since beliefs  $\bar{t}$  can only depend on  $a$ , given a fixed strategy  $g$  we cannot necessarily obtain any value of the narrow expected utility in dimension  $i$ ;  $\bar{U}_i(s)$  from (11). We say that narrow expected utilities  $(\bar{U}_i(s))_{s \in S, i \in N}$  can be *achieved* given strategy  $g$  if there exists an outcome function that induces beliefs  $\bar{t}$  such that for each dimension  $i \in N$  and every  $s \in S$

$$\sum_{a_i \in A_i} g_i(a_i|s) \bar{t}_i(a_i) = -v_i(s) \sum_{a_i \in A_i} g_i(a_i|s) a_i + \bar{U}_i(s) \quad (18)$$

**Lemma 1.** *A strategy  $g$  and narrow expected utilities  $(\bar{U}(s))_{s \in S}$  that can be achieved given  $g$  are NIC if and only if*

1. *The strategy is monotonic on each dimension; that is for all  $i \in N$  we have that*

$$\sum_{a_i \in A_i} a_i g_i(a_i|s) \quad (19)$$

*is increasing in  $s \in S$ .*

2. *On each dimension  $i \in N$ ,  $\bar{U}_i(s)$  is increasing in  $s \in S$ .*

3. *On each dimension  $i \in N$ , the following envelope condition holds for any two types  $s, s' \in S$*

$$\bar{U}_i(s) = \bar{U}_i(s') + \int_{s'}^s v'_i(z) \sum_{a_i \in A_i} a_i g_i(a_i|z) dz \quad (20)$$

*Proof.* [In Appendix](#) □

The requirement that an NIC strategy has increasing immediate utility in type separately on each dimension means that any deterministic strategy in an NIC mechanism must have a threshold form. This differs from the rational case where a threshold strategy is optimal for the principal under the IVV assumption, but is not an implication of IC.

It will be useful in characterizing the principal's optimal mechanism under NIC to show how beliefs and the outcome function relate for any fixed distribution over actions. The following result shows when we can write the outcome distribution in terms of the beliefs over the expected outcome in either dimension. It gives a standard statistical correctness result that applies to beliefs formed using Bayesian Networks under perfect Direct Acyclic Graphs (DAGs) ([Spiegler, 2020](#)).

**Lemma 2** (Statistical Correctness). *Given any distribution over actions  $g$ , for any two dimensions  $i, j \in N$  we have that beliefs  $\bar{t}_i, \bar{t}_j$  satisfy*

$$\sum_{a_i \in A_i} g_i(a_i) \bar{t}_i(a_i) = \sum_{a_j \in A_j} g_j(a_j) \bar{t}_j(a_j) \quad (21)$$

*Proof.* Rearranging the expected outcome gives us the first part. For any  $i \in N$ :

$$\sum_y g(a) t(a) = \sum_{a_i} g_i(a_i) \sum_{a_{-i}} \frac{g(a_i, a_{-i})}{g_i(a_i)} t(a_i, a_{-i}) = \sum_{a_i} g_i(a_i) \bar{t}_i(a_i)$$

□

This statistical correctness is necessary but not sufficient for an outcome function to exist that induces given beliefs for a fixed strategy and action distribution. For an example of beliefs that satisfy the statistical correctness constraint but cannot be induced, consider the case with  $N = \{1, 2\}$  and  $g(1, 1) = g(0, 0) = \frac{1}{2}$ . If beliefs do not also satisfy  $\bar{t}_1(1) = \bar{t}_2(1)$  and  $\bar{t}_1(0) = \bar{t}_2(0)$ , then there is no outcome function implementing these beliefs under this action distribution.

The principal's objective can be written in terms of beliefs. This means it is useful to work directly with beliefs rather than the underlying outcome function when we characterize the principal's optimal mechanism. Although the statistical correctness constraint is not sufficient, I now show that any action distribution taking a deterministic threshold form;  $g_i(1|s) = \mathbb{1}\{s \geq \hat{s}_i\}$  for some  $\hat{s}_i \in [0, 1]$  for all  $i \in N$ , can be made NIC by an outcome function that is additive across dimensions.

**Proposition 2.** *For any deterministic threshold strategy  $g$ , we can construct an outcome function  $t$  that implements beliefs  $\bar{t}$  so that  $g$  is NIC. The constructed outcome function is additive; for any  $a_{-i}, \tilde{a}_{-i} \in A_{-i}$  we have that*

$$t(1, \tilde{a}_{-i}) - t(0, \tilde{a}_{-i}) = t(1, a_{-i}) - t(0, a_{-i}) \quad (22)$$

Moreover, any outcome function  $\tilde{t}$  such that  $g$  is NIC can only differ from this additive  $t$  at action combinations that do not occur under  $g$ ;  $t(a) \neq \tilde{t}(a)$  only if  $g(a|s) = 0$  for all  $s \in S$ .

*Proof.* [In Appendix](#)

□

Proposition 2 means that if the principal implements a deterministic threshold strategy, they have no payoff gain from implementing an outcome function that is not additive. The principal can exploit two features of the narrow agents misperception, that they can only perceive of the outcome function as additive and that their beliefs do not account for confounding bias. The principal only exploits the second misperception: in Section 5.1.1 I show that only deterministic threshold strategies are NIC and in the proof of Theorem 1 the principal does not want to implement any other form of strategy even if they can force the agent to randomize over actions.

The proof of Theorem 1 works as follows. First it uses Lemma 1 to write both the principal's objective and the statistical correctness constraint from Lemma 2 only in terms of the strategy and the narrow perceived utility of the type taking action zero;  $\overline{U}_i(0)$ . It then shows that a minimax upper bound to this constrained problem is solved by implementing deterministic threshold strategies on each dimension, with the  $\beta$  weights in the proof coming from a rewriting of the Lagrange multipliers from the statistical correctness constraint. Under IVV, we can apply a standard minimax theorem argument to obtain a saddle point for this upper bound problem. Finally, we can show that we can achieve this minimax upper bound with an outcome function that solves the full problem using Proposition 2.

### 4.3 Effect of Narrow Agents on Principal's Welfare

Using the characterization of the principal's optimal mechanism, I obtain a result on how the principal's optimal thresholds differ when we move to narrow agents from the rational benchmark. We see that when facing narrow inference, if actions have an immediate cost to the agent but benefit to the principal, the principal implements a strategy with a lower type threshold for taking the action on any dimension than in the rational benchmark. The opposite holds in the immediate utility benefit, principal's loss case where the thresholds are higher when the agent is narrow.

**Proposition 3.** *Assume the IVV assumption holds.*

1. When for every  $i \in N$  we have  $w_i > 0$  and  $v_i(s) \leq 0$  for all  $s \in S$ , then on each dimension the objective-maximizing thresholds are weakly lower with narrow agents than under the rational benchmark, so  $a_i = 1$  is taken by a larger proportion of types for all  $i \in N$ .
2. When for every  $i \in N$  we have  $w_i < 0$  and  $v_i(s) \geq 0$  for all  $s \in S$ , then on each dimension the objective-maximizing thresholds are weakly greater with narrow agents than under the rational benchmark, so  $a_i = 1$  is taken by a smaller proportion of types for all  $i \in N$ .

*Proof.* [In Appendix](#) □

The intuition for this result is as follows. When the agent's action is costly in terms of immediate utility, for any threshold strategy the principal transfers positive utility from outcomes to the agent. The single dimension of type results in actions in one dimension being positively correlated with actions on any other dimension. Since actions result in higher transfers to the agent, this leads a narrow agent to overestimate the transfer they will get from the principal. Rational agents adjust for the fact their type is on the margin between taking an action or not, so will receive a lower overall transfer than the average obtained by agents taking that action. The overestimation of the transfer by narrow agents means less transfer has to be given to higher type agents in order to implement any given strategy. This reduces the marginal cost to the principal of implementing that any given proportion of agents take the action on any dimension. Given the fixed benefits of the actions to the principal, this lower marginal cost means they want a higher proportion of agents to take the action.

When actions have an immediate utility benefit to the agent, the principal is a net receiver of utility from outcome transfers from the agent. In this case the converse logic holds as the principal has a lower marginal benefit from a higher proportion of agents taking the action, but a fixed cost. This is clearest if we consider the principal is a profit-maximizing seller of multiple goods and the agent is a customer. Narrow inference results in the agent overestimating the price of the good, which hurts the seller as it means the sell to fewer types at any price.

Both these intuitions also apply to the effect of narrow inference on the welfare of

the principal. With immediate utility costs of the actions to the agent, under narrow inference any fixed strategy requires less transfer of outcome utility to the agent, while with immediate utility benefits it requires more transfer to the agent. This gives the following result.

**Proposition 4.** *Assume the IVV assumption holds.*

1. *When for every  $i \in N$  we have  $v_i(s) \leq 0$  for all  $s \in S$ , the principal can obtain at least as high an objective value when the agent is narrow compared to the rational benchmark.*
2. *When for every  $i \in N$  we have  $v_i(s) \geq 0$  for all  $s \in S$ , the principal obtains at least as high an objective value in the rational benchmark compared to when the agent is narrow.*

*Proof.* [In Appendix](#)

□

## 4.4 Effect of Greater Dimensionality

I now consider what happens when the number of dimensions grows large, with symmetry across dimensions. In this setting when actions have a direct benefit for the principal, the principal is able to incentivize types to take the actions at vanishing cost. This is because they only have to pay the transfer of utility from outcomes to the agent on one dimension in this symmetric narrow agent setting. When actions have a direct cost to the principal, the opposite is true and it becomes too costly for the principal to extract transfers from the agent. In this case, narrow agents overestimate the cost to themselves of taking the action for any given outcome function.

Define a *symmetric dimension space of size  $n$*  as follows. For any  $n$ ,  $v_i^{(n)}(s) = \frac{1}{n}v(s)$ , and  $w_i = \frac{1}{n}w$ . Both immediate utility and the principal's direct utility from actions are decreasing as the dimensionality of the action space  $n$  grows, but such that the total effect of actions  $\sum_{i \in N} v_i^{(n)}(s) = v(s)$ ,  $\sum_{i \in N} w_i^{(n)} = w$  is constant. Let  $\hat{s}_i^{(n)}$  be the solution to the principal's problem with narrow agents when there is a symmetric dimension space of size  $n$ .

**Proposition 5.** *Assume the IVV assumption holds. Consider a sequence as  $n \rightarrow \infty$  of symmetric dimension spaces of size  $n$ .*

1. *When  $w > 0$ , there exists an  $\bar{n}$  such that for any  $n \geq \bar{n}$ , we have that the principal's optimal mechanism with narrow agents implements a strategy such that all types take the action on all dimensions; for all  $i \in N$ ,  $\hat{s}_i^n = 0$ .*
2. *When  $w < 0$ , there exists an  $\bar{n}$  such that for any  $n \geq \bar{n}$ , we have that the principal's optimal mechanism with narrow agents implements a strategy such that all types take no action on all dimensions; for all  $i \in N$ ,  $\hat{s}_i^n = 1$ .*

*Proof.* [In Appendix](#) □

The result follows from the logic discussed in the intuition for the characterization result in Section 4.1. To implement symmetric thresholds, the principal must scale down the outcome utility from taking the action on any dimension by  $\frac{1}{n}$  relative to the rational benchmark to maintain the same thresholds as narrow incentive compatible. As  $n$  grows large, the contribution of the outcome to the principal's utility then shrinks and the outcome transfers to or from the agent become smaller. Eventually for some  $n$ , the direct immediate utility to the principal dominates their objective. When  $w > 0$  this means they want all types of agent to take the beneficial action while when  $w < 0$  they want no types to take the action.

If we consider a variant of the multi-dimensional firms story then Proposition 5 provides a theory of organization design. Let the principal be the CEO or top management of the firm. If divisions take actions that are costly to themselves but benefit the firms overall profitability, then under narrow inference the CEO benefits from splitting divisions. Conversely, if the divisions take actions that benefit themselves but reduce profitability, then the CEO would want to merge divisions under narrow inference.

## 4.5 Symmetric Immediate Utility

If the immediate utilities of each action to the agent are the same across dimensions, we have that the threshold strategy implemented by the principal's optimal mechanism has an identical threshold in every dimension. This means the agent perfectly correlates their

actions across dimension; above some threshold type the agent chooses action  $a_i = 1$  for all  $i \in N$  and below this threshold the agent chooses  $a_i = 0$ . This contrasts with the principal's optimal mechanism in the rational benchmark, where the principal generally implements thresholds that differ across dimensions even when immediate utilities are symmetric.

**Proposition 6.** *Assume the IVV assumption holds. If immediate utility of the actions is the same across all dimensions;  $v_i(s) = v_j(s) = v(s)$  for any  $i, j \in N$ , and either*

1. *For every  $i \in N$ , we have  $w_i > 0$  and  $v(s) \leq 0$  for all  $s \in S$ .*
2. *For every  $i \in N$ , we have  $w_i < 0$  and  $v(s) \geq 0$  for all  $s \in S$ .*

*Then the principal's optimal mechanism with narrow agents implements a threshold strategy such that for any two dimensions  $i, j \in N$ ,  $\hat{s}_i = \hat{s}_j$ .*

*Proof.* [In Appendix](#) □

The result holds because with symmetric immediate utilities, then in both cases the narrow participation constraints have to bind on all dimensions. The statistical correctness constraint then pins down that the thresholds have to be equal across dimensions.

Intuitively, if two thresholds differ then under symmetric immediate utilities the principal can either equalise the higher threshold to the lower one when they benefit from the actions ( $w_i > 0$ ), or raise the the lower threshold to the higher level if the actions are costly to the principal ( $w_i < 0$ ). Under the assumptions on the agent's immediate utilities for both cases, this can be done in a way that improves the principal's direct utility without changing the expected outcome utility costs or benefits to the principal.

## 4.6 Illustrative Example

The following example illustrates the rational benchmark and the narrow agent results.

**Example 2.** A company has two divisions  $N = \{1, 2\}$ . Each division has a binary action decision  $A_i = \{0, 1\}$ , corresponding to whether to produce output or not. The profitability of production for companies varies with a uniformly distributed type variable;  $s \sim U[0, 1]$ ,  $P(s) = s$ . The divisions obtain the same profits from each unit of production

across dimensions;  $v_1(s) = v_2(s) = r \cdot s$  with  $r > 0$ . Production causes potentially different levels of pollution across divisions;  $w_1 < 0, w_2 < 0$ . However, the benefits of production exceed the cost of pollution for the highest types;  $-w_1 < r, -w_2 < r$ .

For these parameters, we have that for each  $i \in N$

$$\phi_i(s) = v_i(s) - \frac{1 - P(s)}{p(s)} v'_i(s) = r(2s - 1) \quad (23)$$

We have increasing virtual values and the strategy implemented by the principal's optimal mechanism has a threshold form  $g_i(1|s) = \mathbb{1}\{s \geq \hat{s}_i^*\}$  in each dimension, with thresholds

$$\hat{s}_i^* = \frac{1}{2} - \frac{w_i}{2r} \quad (24)$$

This gives the principal objective value  $W^{rat} = \sum_{i \in N} r(\frac{1}{2} + \frac{w_i}{2r})^2$

With an agent who does narrow inference, given  $\beta$ , the thresholds that solve the problem  $\hat{s}^*(\beta) = \arg \max_{\hat{s} \in S^n} \overline{W}(\hat{s}, \beta)$  are

$$\hat{s}_i(\beta_i) = \begin{cases} \frac{1}{2} - \frac{w_i}{2r\beta_i} & \text{if } -w_i \leq \beta_i r \\ 1 & \text{if } -w_i > \beta_i r \end{cases} \quad (25)$$

for  $i \in N$ .

When  $-(w_1 + w_2) \leq r$ , we have that the value of the principal's objective is

$$\begin{aligned} \min_{\beta \in [0,1]^2, \beta_1 + \beta_2 = 1} \overline{W}(\hat{s}_i(\beta_i), \beta) &= \min_{\beta \in [0,1]^2, \beta_1 + \beta_2 = 1} \sum_{i \in N} r\beta_i \left(\frac{1}{2} + \frac{w_i}{2\beta_i r}\right)^2 \\ &= \sum_{i \in N} r \frac{w_i}{w_1 + w_2} \left(\frac{1}{2} + \frac{1}{2} \frac{w_1 + w_2}{r}\right)^2 \end{aligned} \quad (26)$$

so  $\beta_i = \frac{w_i}{w_1 + w_2}$  (which verifies that  $-w_i \leq \beta_i r$  for both dimensions) and the thresholds in the objective maximizing strategy are

$$\hat{s}_1 = \hat{s}_2 = \frac{1}{2} - \frac{w_1 + w_2}{2r} \quad (27)$$

When  $-(w_1 + w_2) > r$ , then  $\hat{s}_i = 1$  for some dimension  $i$ . Without loss of generality,



let  $-w_1 \leq \beta_1 r$ ,  $-w_2 > \beta_2 r$ . Then we have

$$\overline{W}(\hat{s}_i(\beta_i), \beta) = \frac{1}{4}(r\beta_1 + 2w_1 + \frac{w_1^2}{r\beta_1}) \quad (28)$$

Which is minimized by  $\beta_1 = \min\{1, \frac{-w_1}{r}\}$ , giving symmetric thresholds  $\hat{s}_1 = \hat{s}_2 = 1$  and setting the value of the principal's objective to zero.

This demonstrates how when we have symmetric immediate utilities across dimensions, Proposition 6 applies and we have symmetric thresholds across dimensions. We see also that the thresholds are higher with narrow agents, as Proposition 3 implies.

The welfare of the principal under narrow inference is then

$$W^{nar} = r(\frac{1}{2} + \frac{w_1 + w_2}{2r})^2 \mathbb{1}\{-(w_1 + w_2) \leq r\} \quad (29)$$

Comparing with the principal's objective value under the rational benchmark  $W^{rat}$  in Example 2, we see that it is lower with narrow agents than in the rational benchmark when  $W^{nar} \leq W^{rat} \Leftrightarrow w_1 w_2 \leq \frac{1}{2} r^2$ . This always holds when  $-(w_1 + w_2) \leq r$ , as the maximum value of  $w_1 w_2$  subject to this constraint is  $\frac{r^2}{4}$ . When  $-(w_1 + w_2) > r$  the objective value with narrow agents is zero which is again lower than the positive objective value than is achieved with rational agents. This demonstrates Proposition 4.

## 5 Extensions

### 5.1 Alternative Participation Constraints

In this section I consider two alternatives to narrow dimension-by-dimension participation constraints. I first analyze the consequences of an 'ex-post' participation constraint where the agent can opt out of the mechanism after learning the true realization of their utility. This can be interpreted as either the agent truly having an ex-post right to withdraw from the mechanism, or that the principal has concerns for the welfare of the agent or their own reputation that mean they do not want to reduce the true welfare of the agent below a particular level.

Under this true welfare participation constraint (henceforth TWPC), the principal can benefit from both committing to a non-separable outcome function and forcing the agent to randomize over actions. This is because in combination they allow the principal to ‘redistribute’ the utility from outcomes between different types of agents without affecting the narrow participation constraints, which can relax the true welfare participation constraint. This contrasts with the principal’s optimal mechanism under the narrow participation constraints, where the principal optimally chooses an additive mechanism and where we can show there is no benefit to the principal from forcing randomization.

I then consider a participation constraint that requires that the sum of narrow utilities  $\sum_{i \in N} \bar{U}_i(s)$  is greater than the value of the outside option for all types  $s \in S$ . This is a relaxation of the narrow participation constraint considered earlier. I show how we can modify the analysis of the principal’s optimal mechanism for this setting.

### 5.1.1 True Welfare Participation Constraint

Under TWPC, given a strategy  $g$  and an outcome function  $t$  we require that for all types  $s \in S$ , for any  $a \in \sup p\{g(\cdot|s)\}$

$$u(s, a, t(a)) = \sum_{i \in N} v_i(s) a_i + t(a) \geq 0 \quad (30)$$

Now consider an extension of the basic model where the Principal can force randomization by agents. They do this by forcing agents to choose from a restricted set of lotteries on each dimension,  $G_i \subseteq \Delta(A_i)$ . An outcome function  $t$  and strategy  $g$  are *narrow incentive compatible given restrictions* (NICR) if there exist  $\{G_i\}_{i \in N}$  such that for all  $i \in N$  and  $s \in S$  we have that for any  $\tilde{g}_i \in G_i$

$$\sum_{a_i \in A_i} g_i(a_i|s) [a_i v_i(s) + \bar{t}_i(a_i)] \geq \sum_{a_i \in A_i} \tilde{g}_i(a_i) [a_i v_i(s) + \bar{t}_i(a_i)] \quad (31)$$

with  $g_i \in G_i$  on all  $i \in N$ .

Clearly every NIC strategy is NICR, but there are NICR strategies that are not NIC. I now show that any NICR strategy must take a *random threshold form*. A strategy has a random threshold form if for all  $i \in N$  there exists a threshold  $\hat{s}_i \in [0, 1]$  such that for

almost every type  $s \in S$ ,  $g_i(1|s) = q_i^h \cdot \mathbb{1}\{s \geq \hat{s}_i\} + q_i^l \cdot \mathbb{1}\{s < \hat{s}_i\}$ , where  $1 \geq q_i^h \geq q_i^l \geq 0$ . This nests the concept of a deterministic threshold strategy as a special case with  $q_i^h = 1$  and  $q_i^l = 0$ .

**Lemma 3.** *Every NICR strategy  $g$  takes a **random threshold form**. Every NIC action strategy takes a **deterministic threshold form** with  $q_i^h = 1$  and  $q_i^l = 0$ .*

*Any strategy  $g$  that takes a **random threshold form** is NICR if there exists an outcome function  $t$  and a set of restrictions  $\{G_i\}_{i \in N} = \{(q_i^h, 1 - q_i^h), (q_i^l, 1 - q_i^l)\}_{i \in N}$  such that  $t$  together with  $g$  induces beliefs that for all  $i \in N$  satisfy*

$$\bar{t}_i(1) = \bar{t}_i(0) - v_i(\hat{s}_i) \quad (32)$$

*Proof.* [In Appendix](#) □

In the example that follows, the principal can benefit from selecting a mechanism with both a non-separable outcome function and that forces randomization by the agent. This is something that is not the case for the narrow participation constraints considered earlier. Under those constraints, the principal's optimal mechanism can be implemented with an additive outcome function and allowing NICR strategies does not benefit the principal<sup>3</sup>.

For intuition, consider the case where  $N = \{1, 2\}$ . Focus on action distributions that have full support;  $g(a) > 0$  for all  $a \in A$ . Given  $g$ , for any action  $a \in A$ , denote the lowest type choosing that action by  $\underline{s}(a, g) = \inf\{s \in S : g(a|s) > 0\}$ . This is well defined if  $g$  has full support. Under a full support action distribution, there are multiple outcome functions that can implement the same beliefs. We can write any outcome function that

---

<sup>3</sup>The proof of Theorem 1 does not use the fact that  $g$  is restricted to being a deterministic threshold distribution under NIC. Implementing a deterministic threshold distribution is optimal for the principal in this case even when random threshold strategies are allowed.

supports beliefs  $\bar{t}$  as

$$\begin{aligned} t(1, 1) &= \frac{g_1(1)}{g(1, 1)} \bar{t}_1(1) - \frac{g_2(0)}{g(1, 1)} \bar{t}_2(0) + \frac{g(0, 0)}{g(1, 1)} t(0, 0) \\ t(1, 0) &= \frac{g_2(0)}{g(1, 0)} \bar{t}_2(0) - \frac{g(0, 0)}{g(1, 0)} t(0, 0) \\ t(0, 1) &= \frac{g_1(0)}{g(0, 1)} \bar{t}_1(0) - \frac{g(0, 0)}{g(0, 1)} t(0, 0) \end{aligned} \tag{33}$$

The principal can use the degree of freedom by selecting  $t(0, 0)$  to redistribute utility between types taking action combinations that share an action in some dimension. For example, the principal can equalize the welfare of the lowest types choosing  $(1, 0)$  and  $(0, 0)$  by setting.

$$t(0, 0) = \bar{t}_2(0) + \frac{g_2(0)}{g(0, 0)} v_1(\underline{s}((1, 0), g))$$

The outcome function that achieves this can be non-separable. I now demonstrate with an example how under TWPC, the principal can benefit from this ‘hidden redistribution’ channel by selecting a mechanism that is both non-separable and induces a small amount of randomization in the strategy.

I take the highest value of the principal’s objective that can be obtained from a deterministic strategy under TWPC and show that under some parameters there is a random strategy with a non-separable outcome function that betters it.

**Example 3.** Let the agent have actions  $N = \{1, 2\}$  that result in benefits to the principal  $w_1 > w_2 > 0$  but have a disutility cost. For this example, we assume the immediate disutility cost are symmetric across dimensions  $v(s) = -r(1 - s)$ ,  $r > 0$ , and depend on the uniformly distributed type variable  $s \sim U[0, 1]$ ,  $P(s) = s$ .

Consider parameters  $w_1 = 1, w_2 = 0.3, h = 0.5$ . In Appendix A.3 we calculate if the principal is restricted to implement a deterministic threshold strategy, they choose thresholds  $\hat{s}_1^* = 0.2, \hat{s}_2^* = 0.6$ . Let  $W_y(\hat{s}_1, \hat{s}_2)$  be the principal’s welfare if they implemented a strategy with thresholds  $(\hat{s}_1, \hat{s}_2)$  and only had to consider the TWPC for the worst-off agent choosing actions  $a$ . The principal objective value when all TWPC hold

at optimal thresholds  $(\hat{s}_1^*, \hat{s}_2^*)$  is

$$\begin{aligned} \min\{W_{0,0}(\hat{s}_1^*, \hat{s}_2^*), W_{1,0}(\hat{s}_1^*, \hat{s}_2^*), W_{1,1}(\hat{s}_1^*, \hat{s}_2^*)\} &= W_{1,0}(\hat{s}_1^*, \hat{s}_2^*) \\ &= \min\{0.6, 0.56, 0.64\} = 0.56 \end{aligned} \quad (34)$$

Now consider what happens if we switch to a random threshold strategy  $\tilde{g}_i(1|s) = (1 - \epsilon)\mathbb{1}\{s \geq \hat{s}_i^*\}$  some  $\epsilon \in (0, 1)$ . Choose  $\epsilon = 0.001$ ,  $t(0, 0) = 0$ , and  $t(1, 0)$ ,  $t(0, 1)$ ,  $t(1, 1)$  to satisfy the equations in (33). This is NICR as it induces beliefs such that  $\bar{t}_1(1) - \bar{t}_1(0) = 0.4 = r(1 - \hat{s}_1^*)$ ,  $\bar{t}_2(1) - \bar{t}_2(0) = 0.2 = r(1 - \hat{s}_2^*)$ . It also satisfies the true welfare participation constraint;  $t(1, 1) - 2r(1 - \hat{s}_2^*) \approx 0.053$ ,  $t(1, 0) - r(1 - \hat{s}_1^*) = t(0, 0) = 0$ ,  $t(0, 1) - r(1 - \hat{s}_2^*) \approx 13.11$ ,

The value of the principal's objective under this random threshold strategy is

$$- \sum_{a \in A} g(a)t(a) + (1 - \epsilon)((1 - P(s_1))w_1 + (1 - P(s_2))w_2) \approx 0.573$$

This is greater than the highest objective value that can be obtained from a deterministic threshold strategy. The outcome function is given by  $t(0, 0) = 0$ ,  $t(1, 0) = 0.4$ ,  $t(0, 1) \approx 13.31$ ,  $t(1, 1) \approx 0.453$ . Thus  $t(1, 1) - t(1, 0) \neq t(0, 1) - t(0, 0)$  and unlike under the narrow participation constraints, the principal can benefit from choosing a non-separable outcome function.

### 5.1.2 Sum of Narrow Utility Participation Constraint

Now I consider a participation constraint such that the agent's sum of narrow utilities across dimensions must exceed an outside option worth zero. That is, for all  $s \in S$

$$\sum_{i \in N} \bar{U}_i(s) = \sum_{i \in N} \sum_{a_i \in A_i} g_i(a_i|s)[v_i(s)a_i + \bar{t}_i(a_i)] \geq 0 \quad (35)$$

We call this a *sum-narrow participation constraint*. This constraint reflects that the agent makes a joint decision across dimensions on whether to participate or not. The agent has the option to take the zero action on all dimensions and reject any outcome given by the mechanism. Under the sum-narrow participation constraint, the agent understands that

the outcome function may be interactive across dimensions and they cannot just reject the outcome on each dimension individually. This contrasts with their beliefs formed from narrow inference for different actions within the mechanism, which could be correct only if the outcome function was additive.

I modify the proof of Theorem 1 to obtain the following result.

**Theorem 2.** *Assume IVV holds. The principal maximizes their objective over all NIC mechanisms that satisfy the sum-narrow participation constraint if and only if they choose an outcome function that implements a threshold strategy that solves*

$$\max_{\hat{s} \in S^n} \overline{W}(\hat{s}; \beta) = \max_{\hat{s} \in S^n} \left\{ \sum_{i \in N} \int_{\hat{s}_i}^1 \left( \frac{1}{n} \phi_i(s) + w_i \right) p(s) ds \right\} \quad (36)$$

*Proof.* [In Appendix](#) □

The only subsequent result that doesn't hold under sum-narrow participation constraints is Proposition 6. Propositions 3, 4 and 5 continue to hold, and their proofs require no modification.

## 5.2 Dropping the IVV assumption

I now explore what happens to the principal's optimal mechanism without the IVV assumption. I first present an example where IVV does not hold, and show the principal's optimal mechanism under the rational benchmark is no longer separable and no longer implements a deterministic threshold strategy. Thus, since only deterministic threshold strategies are NIC, we cannot make the same neat comparisons between the principal's optimal mechanism under narrow and rational inference that we made in Propositions 3 and 4. Without IVV, the fact that there are strategies that are IC but not NIC becomes relevant.

However, despite the reduction in the set of strategies that are implementable we can prove that an analogous result to Proposition 4 still holds. The principal is still always better off when agents make narrow compared to rational inference if  $v_i(1) \leq 0$ ,  $w_i > 0$  for all  $i \in N$  and the principal is restricted to implement a *deterministic interval strategy*. In a deterministic interval strategy, for some  $k \in \mathbb{N}$  there is a partition of the type space

$z_0 = 0 < z_1 < \dots < z_{k-1} < z_k = 1$  where for any  $s \in [z_{l-1}, z_l)$  for  $l \in \{1, \dots, k\}$  we have  $g(a|s) = 1$  for some  $a \in A$ .

Finally, I show how the characterization result Theorem 1 can be modified to deal with the absence of IVV.

### 5.2.1 Example without IVV

There are two dimensions  $N = \{1, 2\}$  and the type distribution is uniform  $P(s) = s$ . Let the immediate utility take the following form on either dimension, for some  $d_i, b_i$  and  $r_i \in (0, 1)$

$$v_i(s) = \begin{cases} (d_i - b_i)(1 - r_i) - \frac{d_i}{2}(1 - s) - \frac{d_i - b_i}{2} \frac{(1 - r_i)^2}{1 - s} & \text{if } s \in [0, r_i) \\ -\frac{b_i}{2}(1 - s) & \text{if } s \in [r_i, 1] \end{cases} \quad (37)$$

this has a continuous derivative equal to

$$v'_i(s) = \begin{cases} \frac{d_i}{2} - \frac{d_i - b_i}{2} \left( \frac{1 - r_i}{1 - s} \right)^2 & \text{if } s \in [0, r_i) \\ \frac{b_i}{2} & \text{if } s \in [r_i, 1] \end{cases} \quad (38)$$

Assume that  $d_i(1 - (1 - r_i)^2) + b_i(1 - r_i)^2 > 0$  and  $b_i > 0$  so that  $v_i(\cdot)$  is strictly increasing.

We can write the virtual values associated with these immediate utilities.

$$v_i(s) - \frac{1 - P(s)}{p(s)} v'_i(s) = \begin{cases} (d_i - b_i)(1 - r_i) - d_i(1 - s) & \text{if } s \in [0, r_i) \\ -b_i(1 - s) & \text{if } s \in [r_i, 1] \end{cases} \quad (39)$$

These virtual values can be decreasing for some interval of types depending on the parameters. Let  $d_1 = b_1 = 0.54$ ,  $r_1 = 0.6$ ,  $d_2 = -0.12$ ,  $b_2 = 1.1$ ,  $r_2 = 0.66$ ,  $w_1 = 0.4266$  and  $w_2 = 0.32$ . I show that under these parameters the principal's optimal separable mechanism is dominated by a non-separable mechanism implementing a non-threshold strategy.

First calculate the best case mechanism for the principal facing a rational agent if they were restricted to separable mechanisms. For these parameters, the optimal thresholds

are interior and solve  $\phi_i(\hat{s}_i) + w_i = 0$  for  $i \in N$ .

$$\begin{aligned}\hat{s}_1^* &= 1 - \frac{w_1}{b_1} = 0.21 \\ \hat{s}_2^* &= 1 - \frac{w_2}{b_2} = \frac{39}{55}\end{aligned}$$

The value of the principal's objective under this mechanism is then  $W(\hat{s}_1^*, \hat{s}_2^*) \approx 0.215$ .

The following mechanism is better for the principal. It implements the deterministic interval strategy  $g^{int}$

$$\begin{aligned}g^{int}(a_1, a_2|s) \\ = \mathbb{1}\{s \in [0, \hat{s}_1^*)\}(1 - a_1) \cdot a_2 + \mathbb{1}\{s \in [\hat{s}_1^*, \hat{s}_2^*)\}a_1 \cdot (1 - a_2) + \mathbb{1}\{s \in [\hat{s}_2^*, 1]\}a_1 \cdot a_2\end{aligned}$$

This gives the principal payoff

$$\begin{aligned}W(g^{int}) &= \int_0^{\hat{s}_1^*} (\phi_2(s) + w_2) p(s) ds + \int_{\hat{s}_1^*}^{\hat{s}_2^*} (\phi_1(s) + w_1) p(s) ds \\ &\quad + \int_{\hat{s}_2^*}^1 ((\phi_1(s) + w_1) + (\phi_2(s) + w_2)) p(s) ds \\ &\approx 0.217\end{aligned}$$

which is greater than the payoff from the best-case for a separable mechanism. The strategy  $g^{int}$  can be made both IC and to satisfy the participation constraint by the following non-separable outcome function.

$$t(0, 0) = 0, t(0, 1) = t(0, 0) - v_1(0), t(1, 0) = t(0, 1) - v_1(\hat{s}_1^*), t(1, 1) = t(1, 0) - v_2(\hat{s}_2^*)$$

### 5.2.2 Welfare of the Principal without IVV

I now show that Proposition 4—which distinguishes cases where the principal benefits from facing a narrow agent—also holds without the IVV assumption. In the case where the principal is worse-off under narrow inference, the result is trivial as the principal may now benefit from implementing a larger set of strategies in the rational case compared to the narrow case. The interesting case is the one where the principal gains from narrow



inference, where actions have an immediate cost to the agent but a direct benefit to the principal.

The proof for the second case works as follows. For a fixed interval strategy, for each dimension  $i \in N$  take the lowest type  $\underline{s}_i$  that both takes the action  $a_i = 1$  and would produce positive value to the principal if implemented as a threshold;  $(v_i(\underline{s}_i) + w_i)(1 - P(\underline{s}_i)) \geq 0$ . We then take the dimension  $i^*$  at which such a threshold has the greatest cost to the principal in terms of the outcome function required for IC. We can obtain an upper bound on the loss of payoff to the principal if they moved to a threshold strategy where only the action on dimension  $i^*$  is taken, and only by types above threshold  $\underline{s}_{i^*}$ . Under narrow inference, for exactly the same cost the principal has to pay to implement this solo action strategy under the rational benchmark, the principal can implement a threshold strategy where substantial payoff from other dimensions is obtained for free. The gain from this move to narrow inference exceeds the upper bound on the loss from moving to the solo action strategy. This shows the principal's payoff from any interval strategy in the rational benchmark is lower than their payoff from some threshold strategy under narrow inference.

**Proposition 7.** *Assume the principal is restricted to implementing a deterministic interval strategy*

1. *When for every  $i \in N$  we have  $v_i(s) \leq 0$  for all  $s \in S$ , the principal can obtain at least as high an objective value when the agent is narrow compared to the rational benchmark.*
2. *When for every  $i \in N$  we have  $v_i(s) \geq 0$  for all  $s \in S$ , the principal obtains at least as high an objective value in the rational benchmark compared to when the agent is narrow.*

*Proof.* [In Appendix](#)

□

### 5.2.3 Minmax characterization without IVV

In this section I show how we can adapt Theorem 1 when IVV does not necessarily hold. This works by applying the ironing procedure of [Myerson \(1981\)](#). Let  $P^{-1}(s)$  be the

quantile function for the type distribution<sup>4</sup>. For any dimension  $i \in N$  and any  $x \in [0, 1]$ , define

$$\Phi_i(x) = \int_x^1 \phi_i(P^{-1}(u)) du \quad (40)$$

Let  $to(\Phi_i)$  be convex hull of the graph of  $\Phi_i$ , the upper concave envelope of  $\Phi_i(\cdot)$  is then defined as

$$\widehat{\Phi}_i(x) = \sup\{z : (z, x) \in co(\Phi_i)\} \quad (41)$$

Then let  $\widehat{\Phi}_i(s) = -\widehat{\Phi}_i'(P(s))$ , this is well defined and increasing for all  $s \in S$  by definition of  $\widehat{\Phi}_i(\cdot)$ <sup>5</sup>. For any  $\beta \in [0, 1]^n$  and dimension  $i \in N$ , define the threshold<sup>6</sup>

$$\hat{s}_i^*(\beta) = \min\{\tilde{s} \in \arg \max_{\hat{s}_i \in S} \int_{\hat{s}_i}^1 (\beta_i \widehat{\Phi}_i(s) + w_i) p(s) ds\}$$

and let  $\hat{s}^*(\beta) = (\hat{s}_i^*(\beta))_{i \in N}$ . We can now state the result.

**Theorem 3.** *The principal maximizes their objective over all NIC mechanisms that satisfy the narrow participation constraints if and only if they choose an outcome function that implements a threshold strategy  $\hat{s}^*(\beta^*)$ , where*

$$\beta^* \in \arg \min_{\beta \in [0, 1]^n : \sum_{i \in N} \beta_i = 1} \sum_{i \in N} \int_{\hat{s}_i^*(\beta)}^1 (\beta_i \widehat{\Phi}_i(s) + w_i) p(s) ds \quad (42)$$

and the value of the principal's objective is given by

$$\widehat{W}(\hat{s}^*(\beta^*); \beta^*) = \sum_{i \in N} \int_{\hat{s}_i^*(\beta^*)}^1 (\beta_i^* \widehat{\Phi}_i(s) + w_i) p(s) ds \quad (43)$$

*Proof.* [In Appendix](#) □

---

<sup>4</sup>This is the inverse of cdf  $P$  as  $P$  is strictly increasing, and thus  $P^{-1}$  is also strictly increasing.

<sup>5</sup>We have that by concavity  $\widehat{\Phi}_i'(P(\cdot))$  is defined for all but a countable set of points in  $S$ , and by right continuity we can extend it to all  $S$ .

<sup>6</sup>This is well defined due to continuity of  $\int_{\hat{s}_i}^1 (\beta_i \widehat{\Phi}_i(s) + w_i) p(s) ds$ .

### 5.3 Connection to ABEE

It is possible to express behaviour under narrow inference as an Analogy Based Expectation Equilibrium (ABEE) (Jehiel, 2005). Under ABEE, each player in a game has an ‘analogy partition’ of the set of histories where other players move. For any cell in the partition, a player believes that the strategy of the other players is the average of the true strategies for histories in that cell.

Take a game with  $n + 2$  players; consisting of the principal,  $n$  different ‘selves’ of the agent and a player of nature. Each of the  $n$  selves corresponds to one of the  $n$  actions available to the agent, so that self  $i \in \{1, \dots, n\}$  controls action  $a_i$ . All selves share identical preferences over the actions and outcome. The timing of the game is as follows; first the principal chooses an outcome function  $t$ . Then the common type of the agent’s selves is drawn. After learning this common type then, moving in any order, each of the  $n$  selves choose either an action from the set they control or to not participate in the mechanism. Finally, the player of nature implements the outcome function chosen by the principal.

Although each of the agent’s selves have common preferences, they differ according to their analogy partitions. Each self partitions the history at which the player of nature moves, with each cell in the partition corresponding to a different action chosen by the self. Thus their beliefs about the expected outcome from each action is the average outcome obtained among all types of agents choosing that action. This coincides with the beliefs under narrow inference. Behaviour under narrow inference then coincides with an ABEE of this multi-selves game.

## 6 Conclusion

This paper takes a step towards understanding how errors in causal inference might affect how we should design economic incentives. I consider a model boundedly rational belief formation I call narrow inference. I then explore the consequences of this model for economic design. In Theorem 1, I show how we can solve for a principal’s optimal mechanism when facing an agent who makes narrow inference, and how this mechanism contrasts with the principal’s optimal mechanism when they face an agent who is fully

rational. I demonstrate how differences in the underlying environment affect both whether the principal benefits or not from causal inferential errors and the shape of the principal's favoured mechanism. I also demonstrate the robustness of these conclusions to some variations of the underlying model.

One can imagine many additional variations and extensions of the model of bounded rationality explored in this paper. In particular, it seems interesting to consider how the principal could shape the extent of agent's departure from rational beliefs. For example, the principal could provide data on how additional variables correlate with the outcomes or otherwise frame the mechanism in a way that influences inference by the agent. The analysis of this paper suggests that in some cases this could be as important a margin of design as the size of material incentives.

It could also be interesting to explore in other economic contexts the underlying idea that people can have different mental models of the world for different decisions that they face. People may not use the same underlying mental model for forecasting inflation as they use for forming impressions about the economic competence of policymakers. Our underlying ethical perspectives may differ in an inconsistent way when making choices in the workplace vs when voting in elections.

# A Appendices

## A.1 Rational Agents

We analyze the principal's optimal mechanism in the rational benchmark. The next result provides a standard characterization of all IC strategies and expected utilities. We say that expected utilities  $\{U(s)\}_{s \in S}$  can be *achieved* given strategy  $g$  if there exists an outcome function  $t$  such that for every type  $s \in S$ ,  $U(s)$  is the expected utility.

**Lemma A.1.** *A strategy  $g$  and expected utilities  $\{U(s)\}_{s \in S}$  that can be achieved given  $g$  are IC if and only if*

1. *Weak monotonicity condition: For any  $s, s' \in S$*

$$\sum_{i \in N} (v_i(s) - v_i(s')) \sum_{a_i \in A_i} g_i(a_i|s) a_i \geq \sum_{i \in N} (v_i(s) - v_i(s')) \sum_{a_i \in A_i} g_i(a_i|s') a_i \quad (44)$$

2. *The expected utility  $U(s)$  is increasing in  $s \in S$ .*

3. *The following envelope condition holds for any two types  $s, s' \in S$*

$$U(s) = U(s') + \sum_{i \in N} \int_{s'}^s v'_i(z) \sum_{a_i \in Y} a_i g(a_i|z) dz \quad (45)$$

*Proof.* For any types  $s, s' \in S$ , the rational incentive constraints (5) require that

$$U(s) \geq U(s') + \sum_{i \in N} (v_i(s) - v_i(s')) \sum_{a_i \in A_i} g_i(a_i|s') a_i$$

Together with the IC from interchanging  $s, s'$  in the above, we get that the weak monotonicity condition must hold. This then implies the second condition. Using the rewritten ICs, the envelope condition holds from the Lipschitz continuity arguments in Theorems 1 and 2 of [Milgrom and Segal \(2002\)](#).

For the converse, the envelope formula implies reduces the rewritten IC for any  $s, s'$

to the following

$$\sum_{i \in N} \int_{s'}^s v'_i(z) \sum_{a_i \in Y} a_i g(a_i|z) dz \geq \sum_{i \in N} (v_i(s) - v_i(s')) \sum_{a_i \in A_i} g_i(a_i|s') a_i$$

which holds under the weak monotonicity condition.  $\square$

As noted, a strategy where the expected utility of actions in individual dimensions is non-monotonic in type can be IC as long as the weak monotonicity condition is satisfied. This is in contrast to the narrow agent model where the action has to be monotonic in type for each dimension. We shall see that under IVV, this does not matter as the principal's optimal mechanism implements a threshold strategy that is monotone on each dimension anyway.

The example in Section 3.2 of [Carroll \(2017\)](#) shows that in we can have non-separability in an optimal selling mechanism with a co-monotonic type distribution. This is due to different monotonicity condition when we have multiple goods relative to when we have a single good. With a single good, we have that under IC higher types must get the good with higher probability, while with multiple goods we can trade-off probabilities across goods without violating IC. In the example of Section 5.2.1, we show that this also applies in our model when IVV does not hold.

We use Lemma A.1 to prove the next lemma, showing that any deterministic threshold strategy can be implemented with an additively separable outcome function.

**Lemma A.2.** *Let  $g^*$  be a deterministic threshold strategy and  $U(0)$  be the expected utility of the type  $s = 0$ . The strategy is IC and achieves the expected utility  $U(0)$  for type  $s = 0$  under outcome function*

$$t(a_1, \dots, a_n) = \sum_{i \in N} t^i(a_i) \tag{46}$$

$$t^i(0) = \frac{1}{n} U(0), t^i(1) = -v_i(\hat{s}_i) + \frac{1}{n} U(0) \text{ for all } i \in N \tag{47}$$

*Proof.* Combining the envelope formula and the expression for expected utility gives us

that the requirement for  $t$  to implement the correct expected utilities is

$$\begin{aligned}
\sum_{a \in A} g^*(a|s)t(a) &= - \sum_{i \in N} \sum_{a_i \in A_i} a_i g_i^*(a_i|s) v_i(s) + U(s) \\
&= - \sum_{i \in N} v_i(s) \sum_{a_i \in A_i} g_i^*(a_i|s) a_i + \sum_{i \in N} \int_0^s v_i'(z) \sum_{a_i \in A_i} a_i g_i^*(a_i|z) dz + U(0) \\
&= \sum_{i \in N} -v_i(s) \mathbb{1}\{s \geq \hat{s}_i\} + \mathbb{1}\{s \geq \hat{s}_i\} \int_{\hat{s}_i}^s v_i'(t) dt + U(0) \\
&= \sum_{i \in N} -\mathbb{1}\{s \geq \hat{s}_i\} v_i(\hat{s}_i) + U(0)
\end{aligned}$$

Since  $g_i^*(1|s) = \mathbb{1}\{s \geq \hat{s}_i\}$  we can write this in terms of actions as  $t(a) = \sum_{i \in N} -\mathbb{1}\{a_i = 1\} v_i(\hat{s}_i) + \frac{1}{n} U(0)$ . The outcome function in the proposition function satisfies this and achieves expected utility  $U(0)$  for type  $s = 0$ , so by Lemma A.1 we have the result.  $\square$

#### A.1.1 Proof of Proposition 1

*Proof.* We can rewrite the principal's objective (2) using the envelope formula (45) from Lemma A.1 and the expression for the outcome function in terms of utilities in the direct

mechanism.

$$\begin{aligned}
W(t, g) &= \int_0^1 \sum_{a \in A} [-t(a) + \sum_{i \in N} w_i a_i] g(a|s) p(s) ds \\
&= \sum_{i \in N} \int_0^1 \sum_{a_i \in A_i} [a_i v_i(s) + w_i a_i] g_i(a_i|s) p(s) ds - \int_0^1 U(s) p(s) ds \\
&= \int_0^1 \sum_{a \in A} [a_i v_i(s) + \sum_{i \in N} w_i a_i] g_i(a_i|s) p(s) ds \\
&\quad - \sum_{i \in N} \int_0^1 \left[ \int_0^s v_i'(z) \sum_{a_i \in Y} a_i g_i(a_i|z) dz \right] p(s) ds - U(0) \\
&= \int_0^1 \sum_{a \in A} \left[ s \sum_{i \in N} v_i(a_i) + \sum_{i \in N} w_i a_i \right] g(a|s) p(s) ds \\
&\quad - \sum_{i \in N} \int_0^1 \left[ \int_t^1 p(s) ds \right] v_i'(z) \sum_{a_i \in A_i} g_i(a_i|z) dz - U(0) \\
&= \sum_{i \in N} \int_0^1 \sum_{a_i \in A_i} [s v_i(a_i) w_i a_i] g_i(a_i|s) p(s) ds \\
&\quad - \sum_{i \in N} \int_0^1 [1 - P(t)] v_i'(z) \sum_{a_i \in A_i} g_i(a_i|z) dz - U(0) \\
&= \sum_{i \in N} \int_0^1 \sum_{a_i \in A_i} [\phi_i(s) a_i + w_i a_i] g_i(a_i|s) p(s) ds - U(0)
\end{aligned}$$

Where the first line follows from expressing the outcome function in terms of expected utility and the last two lines follow from a standard switching of the order of integration and rewriting in terms of marginal strategies.

Clearly it is optimal to set the expected utility of the lowest type to zero. We now consider a relaxed version of the Principal's problem where we ignore the weak monotonicity constraints from Lemma A.1 and that expected utilities might not be achieved given  $g$ . We show that under IVV, the mechanism that solves this relaxed problem implements a deterministic threshold strategy. Under a threshold strategy we have that for all  $i \in N$ ,  $\sum_{a_i \in A_i} a_i g_i(a_i|s)$  is increasing in type  $s$ . Thus the weak monotonicity condition of Lemma A.1 is satisfied.

By Lemma A.2 we can find an outcome function that implements the threshold strategy and achieves any given expected utility for the lowest type. Thus the solution to the



relaxed problem coincides with the solution to the full problem.

$$\max_{g_i \in \Delta(A_i)^S} \int_0^1 \left[ \sum_{i \in N, a_i \in A_i} (\phi_i(s) + w_i) a_i g_i(a_i|s) \right] p(s) ds$$

This problem can be solved pointwise by strategy  $g_i(a_i|s) = 1$  if and only if  $a_i \in \arg \max_{\tilde{a}_i \in A_i} (\phi_i(s) + w_i) \tilde{a}_i$ . BY IVV,  $\phi_i(s) + w_i$  is strictly increasing in  $s \in S$ . Thus either we have  $\phi_i(\hat{s}_i) + w_i = 0$  for some  $\hat{s}_i \in [0, 1]$ , or either  $\phi_i(s) + w_i < 0$  or  $\phi_i(s) + w_i > 0$  for all  $s \in S$ . In the first case the maximizing strategy on dimension  $i \in N$  is  $g_i(a_i|s) = \mathbb{1}\{s \geq \hat{s}_i\}$ , where without loss of generality we set  $g_i(a_i|\hat{s}_i) = 1$  since  $\hat{s}_i \in S$  has measure zero. In the other cases we can write the maximizing strategy as having a threshold form with thresholds  $\hat{s}_i = 0$  and  $\hat{s}_i = 1$  respectively.  $\square$

## A.2 Proofs

I first present proofs for results that are not stated in the main body but are used in the proofs for some of the results.

### Lemma A.3

To obtain the principal's optimal deterministic mechanism in Example 3, we use the following lemma. Given a strategy  $g$ , let the set of all outcome functions that are NICR be denoted  $\mathcal{M}^{NICR}(g)$ . For any vector of reservation values  $\underline{t} = (\underline{t}_i)_{i \in N} \in \mathbb{R}^n$ , denote by  $\mathcal{M}^{NICR}(g, \underline{t})$  the set of all outcome functions  $t \in \mathcal{M}^{NICR}(g)$  such that the induced beliefs satisfy  $\bar{t}_i(0) \geq \underline{t}_i$  for all  $i \in N$ .

**Lemma A.3.** *For a fixed random threshold strategy  $g$ , if  $\mathcal{M}^{NICR}(g, \underline{t}) \neq \emptyset$  the principal's optimal outcome function in  $\mathcal{M}^{NICR}(g, \underline{t})$  results in objective value*

$$\max_{c \in \mathcal{M}^{NICR}(g, \underline{t})} W(t, g, S) = \min_{i \in N, a_i \in A_i} \{ -\underline{t}_i + g_i(1) v_i(\hat{s}_i) + \sum_{j \in N} g_j(1) w_j \} \quad (48)$$

*Proof.* The principal's objective can be written in terms of the beliefs in any dimension

$i \in N$ .

$$\begin{aligned}
W(g, t) &= \int_0^1 \sum_{a \in A} [-t(a) + \sum_{i \in N} w_i a_i] g(a|s) p(s) ds \\
&= - \sum_{a_i \in A_i} g_i(a_i) \bar{t}_i(a_i) + \sum_{j \in N} \sum_{a_j \in A_j} g_j(a_j) w_j a_j \\
&= -\bar{t}_i(0) + g_i(1) v_i(\hat{s}_i) + \sum_{j \in N} g_j(1) w_j
\end{aligned}$$

Therefore we want to set  $\bar{t}_i(0)$  as low as possible, which for a fixed dimension would mean  $\bar{t}_i(0) \geq \underline{t}_i$ . By Lemma 2, the statistical correctness constraint requires that for any  $i, j \in N$

$$-\bar{t}_i(0) + g_i(1) v_i(\hat{s}_i) = -\bar{t}_j(0) + g_j(1) v_j(\hat{s}_j)$$

This can be violated if we set  $\bar{t}_k(0) = \underline{t}_k$  for every  $k \in N$ . The principal cannot set  $\bar{t}_k(0) = \underline{t}_k$  for any  $k \notin \arg \min_{i \in N} \{-\underline{t}_i + g_i(1) v_i(\hat{s}_i)\}$ , as then the statistical correctness constraint requires that for any  $j \in \arg \min_{i \in N} \{-\underline{t}_i + g_i(1) v_i(\hat{s}_i)\}$

$$\bar{t}_j(0) = \underline{t}_k - g_k(1) v_k(\hat{s}_k) + g_j(1) v_j(\hat{s}_j) < \underline{t}_j$$

violating the dimension  $j$  relaxed participation constraint. The statistical correctness constraint can be satisfied by setting  $\bar{t}_j(0) = \underline{t}_j$  for any  $j \in \arg \min_{i \in N} \{-\underline{t}_i + g_i(1) v_i(\hat{s}_i)\}$ , in which case we have  $t_i(0) \geq \underline{t}_i$  for all  $i \in N \setminus \{j\}$ .  $\square$

I now present proofs from the main body.

### Proof of Lemma 1

For any dimension  $i \in N$  and any two types  $s, s' \in S$  NIC requires

$$\bar{U}_i(s) \geq \bar{U}_i(s') + (v_i(s) - v_i(s')) \sum_{a_i \in A_i} a_i g_i(a_i|s')$$

From this we can obtain that the first two conditions are necessary. The envelope formula on each dimension then holds via the usual Lipschitz continuity arguments as in [Milgrom and Segal \(2002\)](#).

Conversely, combining the envelope formula with the rewritten NIC gives

$$\int_{s'}^s v'_i(z) \sum_{a_i \in A_i} a_i g_i(a_i|z) dz \geq (v_i(s) - v_i(s')) \sum_{a_i \in A_i} a_i g_i(a_i|s')$$

which holds under the monotonicity condition.

## Proof of Proposition 2

By the envelope formula (20) of Lemma 1, the beliefs inducing any deterministic threshold strategy must satisfy

$$\bar{t}_i(1) = \bar{t}_i(0) - v_i(\hat{s}_i)$$

for each dimension  $i \in N$ . Thus NIC and the thresholds pin down beliefs, and any threshold strategy can be rendered NIC by some beliefs.

We now show that for any two action combinations that occur under a deterministic threshold strategy, one of the action vectors is weakly larger on all dimensions than the other, with strict inequality for one of the dimensions.

**Lemma A.4.** *Let  $\tilde{g}$  be a deterministic threshold strategy. Then for any  $s > s'$  and  $a'', a'$  such that  $a'' \neq a'$ ,  $\tilde{g}(a''|s'') > 0$  and  $\tilde{g}(a'|s') > 0$  only if  $a''_j \geq a'_j$  for all  $j \in N$ .*

*Proof.* Since  $\tilde{g}$  is a deterministic threshold distribution, for any two dimensions  $i, j \in N$  there is an  $\hat{s}_i$  such that  $\tilde{g}_i(1|s) = \mathbb{1}\{s \geq \hat{s}_i\}$  and an  $\hat{s}_j$  such that  $\tilde{g}_j(1|s) = \mathbb{1}\{s \geq \hat{s}_j\}$ . Suppose for that we can find  $a'' \neq a'$ ,  $\tilde{g}(a''|s'') > 0$  and  $\tilde{g}(a'|s') > 0$  for some  $s'', s' \in S$  such that  $a''_i = 1 > a'_i = 0$  for some dimension  $i \in N$  but  $a''_j = 0 < a'_j = 1$  for another dimension  $j \in N \setminus \{i\}$ . But then  $s \geq \hat{s}_i > s''$  and  $s'' \geq \hat{s}_j > s$ , a contradiction.  $\square$

This implies that the any two action combinations occurring with positive probability under a threshold strategy can be ranked. Denote the set of all action combinations that have positive probability under  $g$  by  $A(g) = \{a \in A : \exists s \in S \text{ such that } g(a|s) > 0\}$ . Denote the projection of  $A(g)$  on dimension  $i \in N$  by  $A_i(g)$ . Define the order  $\succ$  so that  $a'' \succ_i a'$  if and only if  $a''_j \geq a'_j$  for all  $j \in N$  with strict inequality for at least one such  $j$ . By Lemma A.4 this is a strict total order.

Given our NIC deterministic threshold strategy  $g$ , we enumerate the set  $A(g) = \{1, \dots, |A(g)|\}$  so that  $k > l$  means that for  $a^k, a^l \in A(g)$ ,  $a^k \succ a^l$ . Now we can form a partition of  $A(g)$  for each dimension  $i \in N$ . For each action  $a_i \in A_i$ , define the set  $\mathcal{A}_i(a_i) = \{(a_i, \tilde{a}_{-i}) \in A(g)\}$ . This is a partition as  $\emptyset \notin \mathcal{A}_i(a_i)$  for any  $a_i \in A_i(g)$ ,  $\cup_{\tilde{a}_{-i} \in A_i} \mathcal{A}_i(\tilde{a}_{-i}) = A(g)$  and  $\mathcal{A}_i(a_i'') \cap \mathcal{A}_i(a_i') = \emptyset$  for any  $a_i'' \neq a_i' \in A_i$ .

We can then show that any action combination that occurs with positive probability under an NIC deterministic threshold strategy must be maximal in the partition cell according to the order  $\succ$  for at least one dimension.

**Lemma A.5.** *Given an NIC deterministic threshold strategy  $g$ , any  $a \in A(g)$  is such that for at least one dimension  $j \in N$ ,  $a = (a_j, a_{-j}) \succ \tilde{a} = (a_j, \tilde{a}_{-j})$  for all  $\tilde{a} \in \mathcal{A}_j(a_j)$ .*

*Proof.* Suppose this does not hold, then for all  $i \in N$ , we can find a  $\tilde{a}(i) \in \mathcal{A}(a_i)$  such that  $\tilde{a}(i) \succ a$ . The finite set  $\{a(1), \dots, a(n)\}$  must contain a member that is minimal in the strict total order  $\succ$ . Denote this element  $a(k)$  for some  $k \in N$ . Then on all dimensions  $j \in N$ ,  $a(k)_j \leq a_j$ , as if  $a(k)_j > a_j$  then  $a(j) \not\succ a(k)$  which would contradict the minimality of  $a(k)$  in  $\{a(1), \dots, a(n)\}$ . However,  $a(k)_j \leq a_j$  for all  $j \in N$  contradicts that  $a(k) \succ a$ .  $\square$

Now using this, for any action combination  $a^l \in A(g) = \{1, \dots, |A(g)|\}$  assign a dimension  $\pi(l) \in N$  so that  $a^l \succ \tilde{a}$  for any  $\tilde{a} \in \mathcal{A}_{\pi(l)}(g)$ . Then we can recursively define an outcome function  $t$  from the beliefs  $\bar{t}$  that render  $g$  NIC. For any  $k \in \{2, \dots, |A(g)|\}$

$$t(a^1) = \bar{t}_{\pi(1)}(a_{\pi(1)}^1)$$

$$t(a^k) = \frac{\sum_{a_{-\pi(k)} \in A_{-\pi(k)}} g(a_{\pi(k)}^k, a_{-\pi(k)})}{g_{\pi(k)}(a_{\pi(k)}^k)} \bar{t}_{\pi(k)}(a_{\pi(k)}^k) - \sum_{y \in A_{\pi(k)}(a_{\pi(k)}^k)} \frac{g(a)}{g_{\pi(k)}(a_{\pi(k)}^k)} t(a)$$

This outcome function is well defined as for any  $a^k \in A(g)$  all  $a \in A_{\pi(k)}(a_{\pi(k)}^k)$ ,  $t(a)$  has been defined at an earlier stage as  $a^k$  is maximal in  $\succ$  on dimension  $\pi(k)$ . Since  $a^1$  is minimal in  $A(g)$ , we have that  $a_1 = A_{\pi(1)}(a_{\pi(1)}^1)$ , so the first equation is in fact a special case of the second. We now show that  $t$  is additive for the action combinations  $a \in A(g)$  at which it is well defined.

**Lemma A.6.** *For any  $a_{-i} \in A_{-i}$  with  $(1, a_{-i}), (0, a_{-i}) \in A(g)$ , there exists no  $\tilde{a}_{-i} \neq a_{-i}$  such that  $(1, \tilde{a}_{-i}) \in A(g)$  and  $(0, \tilde{a}_{-i}) \in A(g)$ .*

*Proof.* Suppose for contradiction that there is a  $\tilde{a}_{-i} \neq a_{-i}$  such that  $(1, \tilde{a}_{-i}) \in A(g)$  and  $(0, \tilde{a}_{-i}) \in A(g)$ . As we have strict total order  $\succ$  on  $A(g)$ , we have two cases. In the first case  $(0, a_{-i}) \succ (0, \tilde{a}_{-i})$ . This means that  $a_j \geq \tilde{a}_j$  for all  $j \in N \setminus \{i\}$  with strict inequality for some such  $j$ . Then neither  $(0, a_{-i}) \succ (1, \tilde{a}_{-i})$  nor  $(1, \tilde{a}_{-i}) \succ (0, a_{-i})$ , a contradiction as Lemma A.4 implies since  $(1, \tilde{a}_{-i}) \neq (0, a_{-i})$  they must be ranked.

Similarly if  $(0, \tilde{a}_{-i}) \succ (0, a_{-i})$ , then  $\tilde{a}_j \geq a_j$  for all  $j \in N \setminus \{i\}$  with strict inequality for some such  $j$ . Then neither  $(1, a_{-i}) \succ (0, \tilde{a}_{-i})$  or  $(0, \tilde{a}_{-i}) \succ (1, a_{-i})$ .  $\square$

Therefore we cannot have  $t(1, a_{-i}) - t(0, a_{-i}) = t(1, \tilde{a}_{-i}) - t(0, \tilde{a}_{-i})$  and  $(1, a_{-i}), (0, a_{-i}), (1, \tilde{a}_{-i}), (0, \tilde{a}_{-i}) \in A(g)$ . This means the outcome function is additive for all  $a \in A(g)$ .

We can additively extend the outcome function  $t$  defined above to all  $a \in A$ . For any  $\tilde{a} \in Y$ ,  $\tilde{a} \neq a^1$  such that  $a_j^1 \geq \tilde{a}_j$ , define  $t(a) = t(a^1)$ . For each dimension  $i \in N$ , we will define  $t^i(a_i)$  for each  $a_i \in A_i$  so that  $t(a) = \sum_{i \in N} t^i(a_i)$ . First set  $t^i(a_i^1) = \frac{1}{n} t(a^1)$  for all  $i \in N$ , and  $t^i(1) = t^i(0)$  if  $a_i^1 = 1$ . Now move through the elements  $a^l \in \{2, \dots, |A(g)|\} \subset A(g)$  in  $\succ$  order. If  $a^l$  differs in one dimension  $j$  from  $a^{l-1}$ , then by Lemma A.6 there is a unique  $a_{-j} \in A_{-j}$  such that  $(1, a_{-j}), (0, a_{-j}) \in A(g)$ , and we can write  $t^j(1) - t^j(0) = t(1, a_{-j}) - t(0, a_{-j})$ . If  $a^l$  differs from  $a^{l-1}$  on multiple dimensions  $N^l$ , choose one arbitrarily  $j \in N^l$  and set  $t^j(1) - t^j(0) = t(a^l) - t(a^{l-1})$  and set  $t^k(1) = t^k(0)$  for all other  $k \in N^l \setminus \{j\}$ . For the remaining dimensions, which are such that  $(1, a_{-i}) \notin A(g)$  for every  $a_{-i} \in A_{-i}$ , set  $t^i(1) = t^i(0)$ .

For the final part, any  $a \in A \setminus A(g)$  is such that  $g(a) = 0$  and thus  $t(a)$  does not affect on-path beliefs  $\bar{t}(a_i) = \sum_{a_{-i} \in A_{-i}} \frac{g(a_i, a_{-i})}{g_i(a_i)} t(a)$ . For all  $a \in A(g)$  and  $t'$  that implements the beliefs  $\bar{t}$  must match our recursive construction  $t$ . To see this take  $a^1 \in A(g)$ , at this action combination  $t'(a^1) = t(a^1)$  is pinned down by the beliefs only. This is because  $a^1$  is the minimal action in  $A(g)$  according to  $\succ$  but is also maximal in  $\succ$  on one of the dimension partitions, so on this dimension  $t'(a^1)$  and  $t(a^1)$  must both be equal to the belief. Then any expression that implements  $\bar{t}$  for  $t'(a^2)$  is also pinned down in terms of beliefs according to the recursive formula given for  $t(a_2)$ . Continuing up the order  $\succ$ , we

have  $t'(a^l) = t(a^l)$  for every  $a^l \in A(g)$ .

### Proof of Theorem 1

We break the proof into 4 steps.

**Step 1:** We write the principal's problem in a virtual value form. From the statistical correctness constraint for any dimension  $i \in N$  we can write

$$\int_0^1 g(a|s)t(a) = \sum_{a_i} g_i(a_i)\bar{t}_i(a_i) = \int_0^1 \sum_{a_i \in A_i} \bar{t}_i(a_i)g_i(a_i|s)p(s)ds$$

We then have for any  $i, j \in N$  that

$$\begin{aligned} \int_0^1 \sum_{a_i \in A_i} \bar{t}_i(a_i)g_i(a_i|s)p(s)ds &= \int_0^1 \sum_{a_j \in A_j} \bar{t}_j(a_j)g_j(a_j|s)p(s)ds \\ \Leftrightarrow \int_0^1 \bar{U}_i(s)p(s)ds - \int_0^1 v_i(s) \sum_{a_i \in A_i} a_i g_i(a_i|s)p(s)ds \\ &= \int_0^1 \bar{U}_j(s)p(s)ds - \int_0^1 v_j(s) \sum_{a_j \in A_j} a_j g_j(a_j|s)p(s)ds \end{aligned}$$

We can then write the principals objective in terms of the perceived utility in one of the dimensions, and then use the envelope formula to write in terms of the utility of the

lowest type and the marginal strategy.

$$\begin{aligned}
W(t, g) &= \int_0^1 \sum_{a \in A} [-t(a) + \sum_{i \in N} w_i a_i] g(a|s) p(s) ds \\
&= - \int_0^1 \overline{U}_i(s) p(s) ds + \int_0^1 v_i(s) \sum_{a_i \in A_i} a_i g_i(a_i|s) p(s) ds \\
&\quad + \sum_{j \in N} \sum_{a_j \in A_j} w_j \int_0^1 a_j g_j(a_j|s) p(s) ds \\
&= -\overline{U}_i(0) - \int_0^1 \left( \int_0^s v_i(z) \sum_{a_i \in A_i} a_i g_i(a_i|z) dz \right) p(s) ds \\
&\quad + \sum_{a_i \in A_i} \int_0^1 v_i(s) \sum_{a_i \in A_i} a_i g_i(a_i|s) p(s) ds + \sum_{j \in N} \sum_{a_j \in A_j} w_j \int_0^1 a_j g_j(a_j|s) p(s) ds \\
&= -\overline{U}_i(0) + \sum_{a_i \in A_i} \int_0^1 \left( v_i(s) - \frac{1-P(s)}{p(s)} v'_i(s) \right) g_i(a_i|s) p(s) ds \\
&\quad + \sum_{j \in N} \sum_{a_j \in A_j} w_j \int_0^1 a_j g_j(a_j|s) p(s) ds
\end{aligned}$$

Where we use a switch of the order of integration in the final equation. We can also use the envelope formula to write the statistical correctness constraints for any  $i, j \in N$  as

$$\begin{aligned}
&- \overline{U}_i(0) + \sum_{a_i \in A_i} \int_0^1 \left( v_i(s) - \frac{1-P(s)}{p(s)} v'_i(s) \right) g_i(a_i|s) p(s) ds \\
&= -\overline{U}_j(0) + \sum_{a_j \in A_j} \int_0^1 \left( v_j(s) - \frac{1-P(s)}{p(s)} v'_j(s) \right) g_j(a_j|s) p(s) ds
\end{aligned}$$

**Step 2:** We solve problem that is an upper bound to the principal's full problem and show that the solution to this upper bound implements a deterministic threshold strategy.

The principal wants to maximize  $W(t, g)$  and must implement beliefs that satisfy the statistical correctness constraints. The Lagrangian of the problem of maximizing this objective given this constraint can be written as follows, remembering that  $\phi_i(s) = v_i(s) - \frac{1-P(s)}{p(s)} v'_i(s)$  and denoting Lagrange multipliers by  $\lambda_j \in \mathbb{R}$  for the  $j$ th of the  $n-1$  statistical correctness constraints

$$\begin{aligned}
& \overline{W}(g, \overline{U}(0), \lambda) \\
&= -\overline{U}_i(0) + \sum_{a_i \in A_i} \int_0^1 \Phi_i(s) a_i g_i(a_i|s) p(s) ds + \sum_{i \in N} \sum_{a_i \in A_i} w_i \cdot a_i \int_0^1 g_i(a_i|s) p(s) ds \\
&+ \sum_{j \in N \setminus \{i\}} \lambda_j [-\overline{U}_j(0) + \sum_{a_j \in A_j} \int_0^1 \Phi_j(s) a_j g_j(a_j|s) p(s) ds \\
&\quad + \overline{U}_i(0) - \sum_{a_i \in A_i} \int_0^1 \Phi_i(s) a_i g_i(a_i|s) p(s) ds] \\
&= -(1 - \sum_{j \in N \setminus \{i\}} \lambda_j) [\overline{U}_i(0) + \sum_{a_i \in A_i} \int_0^1 \Phi_i(s) a_i g_i(a_i|s) p(s) ds] \\
&+ \sum_{j \in N \setminus \{i\}} \lambda_j [-\overline{U}_j(0) + \sum_{a_j \in A_j} \int_0^1 \Phi_j(s) a_j g_j(a_j|s) p(s) ds \\
&+ \sum_{i \in N} \sum_{a_i \in A_i} w_i \cdot a_i \int_0^1 g_i(a_i|s) p(s) ds]
\end{aligned}$$

Note that we can write the  $n - 1$  lagrange multipliers as  $\beta = (\beta_1, \dots, \beta_n) \in \mathbb{R}^n$  with  $\sum_{j \in N} \beta_j = 1$  by setting  $\beta_j = \lambda_j$  for  $j \in N \setminus \{i\}$  and  $\beta_i = 1 - \sum_{j \in N \setminus \{i\}} \lambda_j$ . We can use this to write a relaxed version of the principal's problem

$$\sup_{g \in \Delta(A)^S, \overline{U}(0) \in \mathbb{R}^n} \min_{\beta \in \mathbb{R}^n, \sum_{j \in N} \beta_j = 1} \overline{W}(g, \overline{U}(0), \beta) \quad \text{subject to } \overline{U}_j(0) \geq 0 \text{ for } j \in N$$

The problem is relaxed as it does not satisfy the following constraints that must hold in the full problem. These are that the strategy must satisfy the monotonicity requirement implied by NIC in Lemma 1. In addition since  $t$  only depend on actions, for a given strategy  $g$  there might not exist a  $t$  inducing beliefs  $\bar{t}$  such that any arbitrary expected utilities  $\overline{U}_i(s) = \sum_{a_i \in A_i} g_i(a_i|s) [a_i v_i(s) + \bar{t}_i(a_i)]$  can be achieved.

We will show that the solution to the relaxed problem implements a strategy that has a deterministic threshold form. Such a strategy satisfies the monotonicity requirement and will be able to induce the expected utilities and beliefs in the solution. Thus both type of additional constraints hold in the solution of the relaxed problem that ignores them.

Restricting the domain of  $\beta$  so that  $\beta_j \in [0, 1]$  for all  $j \in N$  in the minimization



problem gives us the following upper bound.

$$\begin{aligned}
& \min_{\tilde{\beta} \in [0,1]^n: \sum_{i \in N} \tilde{\beta}_i = 1} \sup_{g \in \Delta(A)^S, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \overline{W}(g, \overline{U}(0), \beta) \\
& \geq \sup_{g \in \Delta(A)^S, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \min_{\tilde{\beta} \in [0,1]^n: \sum_{i \in N} \tilde{\beta}_i = 1} \overline{W}(g, \overline{U}(0), \beta) \\
& \geq \sup_{g \in \Delta(A)^S, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \min_{\tilde{\beta} \in \mathbb{R}^n: \sum_{i \in N} \tilde{\beta}_i = 1} \overline{W}(g, \overline{U}(0), \beta)
\end{aligned}$$

For fixed  $\beta \in [0, 1]^n$  with  $\sum_{j \in N} \beta_j = 1$ , we have

$$\overline{W}(g, \overline{U}(0), \beta) = \sum_{j \in N} [-\beta_j \overline{U}_j(0) + \sum_{a_j \in A_j} \int_0^1 (\phi_j(s) + w_j) a_j g_j(a_j | s) p(s) ds]$$

The strategy  $g \in \Delta(A)^S$  that maximizes this expression must satisfy  $g_i(a_i | s) = 1$  if and only if  $a_i \in \arg \max_{\tilde{a}_i \in A_i} (\beta_i \phi_i(s) + w_i) \tilde{a}_i$  for every  $i \in N$ . As in Proposition A.2, this strategy takes a threshold form WLOG and can be induced by the following beliefs that give the correct expected utilities. Given a dimension  $i \in N$ , and threshold  $\hat{s}_i \in [0, 1]$  the following beliefs implement the truthful reporting for the threshold strategy  $g$ .

$$\begin{aligned}
\bar{t}_i(0) &= -\overline{U}_i(0) \\
\bar{t}_i(1) &= \bar{t}_i(0) - v_i(\hat{s}_i)
\end{aligned}$$

By Proposition 2, we can construct an outcome function  $t$  that implements these beliefs given the threshold strategy.

**Step 3:** We show that the objective function in our upper bound problem satisfies the conditions of the minimax theorem. This allows us to interchange the min and max operator and means we have a saddle point solution.

Denote the vector of thresholds by  $\hat{s} = (\hat{s}_i)_{i \in N} \in S^n$ . We can now write the objective in terms of the thresholds.

$$\overline{W}(\hat{s}, \overline{U}(0), \beta) = \sum_{j \in N} [-\beta_j \overline{U}_j(0) + \int_{\hat{s}_j}^1 (\phi_j(s) + w_j) p(s) ds]$$

Define the quantile function  $P^{-1}(s)$ . Since  $P(s)$  is strictly increasing, this is just

the inverse and is also strictly increasing. For any vector  $x \in [0, 1]^n$ , we can write  $P^{-1}(x_i) = \hat{s}_i$ . Letting  $P^{-1}(x) = (P^{-1}(x_i))_{i \in N}$ , we use this to rewrite the objective.

$$\overline{W}(P^{-1}(x), \overline{U}(0), \beta) = \sum_{j \in N} [-\beta_j \overline{U}_j(0) + \int_{x_j}^1 (\phi_j(P^{-1}(u)) + w_j) du]$$

Taking derivatives of  $\int_{x_j}^1 (\phi_j(P^{-1}(u)) + w_j) du$  with respect to the threshold  $x_j$  gives

$$-(\phi_j(P^{-1}(u)) + w_j)$$

By the IVV assumption, this is decreasing and thus  $\int_{x_j}^1 (\phi_j(P^{-1}(u)) + w_j) du$  is concave. We have that  $-\beta_j \overline{U}_j(0)$  is also concave. The sum of concave functions defined on different domains is also concave. Thus for fixed  $\beta$ ,  $\overline{W}(P^{-1}(x), \overline{U}(0), \beta)$  is concave in  $\overline{U}(0)$  and the quantiles  $x$ . Since  $\overline{W}(P^{-1}(x), \overline{U}(0), \beta)$  is convex in  $\beta$  for fixed  $\overline{U}(0), x$  we can apply the minimax theorem (Sion, 1958) and obtain that

$$\begin{aligned} \min_{\beta \in [0, 1]^n, \sum_{i \in N} \beta_i = 1} \sup_{x \in [0, 1]^n, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \overline{W}(P^{-1}(x), \overline{U}(0), \beta) &= \sup_{x \in [0, 1]^n, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \overline{W}(P^{-1}(x), \overline{U}(0), \beta^*) \\ &= \sup_{x \in [0, 1]^n, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \min_{\beta \in [0, 1]^n, \sum_{i \in N} \beta_i = 1} \overline{W}(P^{-1}(x), \overline{U}(0), \beta) \end{aligned}$$

where  $\beta^*$  is the minimizer. Using  $\hat{s}_i = P^{-1}(x_i)$ , we then have

$$\begin{aligned} \min_{\beta \in [0, 1]^n, \sum_{i \in N} \beta_i = 1} \sup_{\hat{s} \in S^n, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \overline{W}(\hat{s}, \overline{U}(0), \beta) &= \sup_{\hat{s} \in S^n, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \overline{W}(\hat{s}, \overline{U}(0), \beta^*) \\ &= \sup_{\hat{s} \in S^n, \overline{U}(0) \in \mathbb{R}_{\geq 0}^n} \min_{\beta \in [0, 1]^n, \sum_{i \in N} \beta_i = 1} \overline{W}(\hat{s}, \overline{U}(0), \beta) \end{aligned}$$

**Step 4:** We can show we can attain the value of the upper bound in the full problem by the following argument. Take the minimizing  $\beta^*$ . If  $\beta_i^* > 0$  we must have that when comparing the dimension  $i$  and any other dimension  $j$  that

$$-\overline{U}_i(0) + \int_{\hat{s}_i}^1 (v_i(s) - \frac{1 - P(s)}{p(s)} v'_i(s)) p(s) ds \leq -\overline{U}_j(0) + \int_{\hat{s}_j}^1 (v_j(s) - \frac{1 - P(s)}{p(s)} v'_j(s)) p(s) ds$$

So that the dimension  $i$  term is minimal. At the solution the dimension  $i$  participation

constraint binds  $\overline{U}_i(0) = 0$  as otherwise the principal could increase the value of the objective by reducing  $\overline{U}_i(0)$ . If  $\beta_j^* > 0$ , then we have that this inequality must hold with equality with  $\overline{U}_j(0) = 0$ , which means the statistical correctness and participation constraint hold. If  $\beta_j^* = 0$ , then we can increase the value of  $\overline{U}_j(0)$  without affecting the principal's objective value. From the above inequality, this can be done so that the statistical correctness constraint holds without violating the participation constraint. Thus the solution for the upper bound can be attained by a NIC strategy and  $t$  that satisfy the statistical correctness and participation constraints, and therefore solves the full problem. As  $\beta_i^* > 0$ , we have that  $\overline{U}_i(0) = 0$  so the equality can be achieved with bounded  $\overline{U}(0)$ . Therefore we can replace the supremum with a maximum in the saddle point problem.

### Proof of Proposition 3

For every  $i \in N$ , let  $\hat{s}_i^{rational}$  be the threshold such that  $g_i(1|s) = \mathbb{1}\{s_i \geq \hat{s}_i^{rational}\}$  is the strategy that solves the principal's problem in the rational benchmark. Let  $\hat{s}_i^{narrow}$  be the solution to the principal's problem a narrow agent, as given by the solution to the minimax problem in Theorem 1. Let  $\beta^*$  be the saddle point distribution over dimensions from that problem.

When  $v_i(1) \leq 0$  we have  $\phi_i(s) = v_i(s) - \frac{1-P(s)}{p(s)}v_i'(s) \leq 0$  for all  $s \in [0, 1]$ , then  $\phi_i(s) + w_i \leq \beta_i^* \phi_i(s) + w_i$ . Suppose that for some dimension  $j \in N$  we have  $\hat{s}_j^{narrow} > \hat{s}_j^{rational}$ . Optimality of  $\hat{s}_j^{rational}$  as a threshold implies  $\phi_j(s) + w_j > 0$  for all  $s \in (\hat{s}_j^{rational}, 1]$ . However, then  $\hat{s}_j^{narrow} > \hat{s}_j^{rational}$  cannot be optimal for the principal as  $0 < \phi_j(s) + w_j \leq \beta_j^* \phi_j(s) + w_j$  for  $s \in (\hat{s}_j^{rational}, \hat{s}_j^{narrow}]$ , a contradiction.

For the second case we again assume that  $\beta_i^* > 0$ , as otherwise  $\hat{s}_i^{narrow} = 1$  in which case the result holds. Then  $\phi_i(s) = v_i(s) - \frac{1-P(s)}{p(s)}v_i'(s) \geq 0$  for any  $s \in [\hat{s}_i^{narrow}, 1]$ . Otherwise there is a  $\tilde{s}_i \in (\hat{s}_i^{narrow}, 1]$  such that by IVV  $\beta_i^* \phi_i(s) + w_i < 0$  for all  $s \in [\hat{s}_i^{narrow}, \tilde{s}_i]$ , and the principal could then switch to implementing  $g_i(0|s) = 1$  for all  $s \in [\hat{s}_i^{narrow}, \tilde{s}_i]$  and obtain a higher payoff.

Now for contradiction assume  $\hat{s}_i^{rational} > \hat{s}_i^{narrow}$ . We have that  $\phi_i(s) \geq 0$  and thus  $\phi_i(s) + w_i \geq \beta_i^* \phi_i(s) + w_i$  for any  $s \in [\hat{s}_i^{narrow}, \hat{s}_i^{rational}]$ . Then optimality of  $\hat{s}_i^{rational}$  is contradicted by the fact that optimality of  $\hat{s}_i^{narrow}$  requires  $0 < \beta_i^* \phi_i(s) + w_i \leq \phi_i(s) + w_i$

for all  $s \in (\hat{s}_i^{narrow}, \hat{s}_i^{rational}]$ , where the first strict inequality follows from IVV.

### Proof of Proposition 4

For any fixed threshold strategy and fixed  $\beta$  the difference in the principal's objective can be written as

$$\begin{aligned}
W(\hat{s}) - \overline{W}(\hat{s}; \beta) &= \sum_{i \in N} \int_{\hat{s}_i}^1 [(\Phi_i(s) + w_i) - \beta_i(\Phi_i(s) + w_i)] p(s) ds \\
&= \sum_{i \in N} (1 - \beta_i) \int_{\hat{s}_i}^1 (v_i(s) - \frac{1 - P(s)}{p(s)} v_i'(s)) p(s) ds \\
&= \sum_{i \in N} (1 - \beta_i) v_i(\hat{s}_i) (1 - P(\hat{s}_i))
\end{aligned}$$

where the last line follows from the fact that  $\int_{\hat{s}_i}^1 (v_i(s) - \frac{1 - P(s)}{p(s)} v_i'(s)) p(s) ds = v_i(\hat{s}_i) (1 - P(\hat{s}_i))$ . Then as  $\hat{s}_i \in [0, 1]$  for the first case where  $v_i(s) \leq 0$  clearly we have  $W(\hat{s}) \leq \overline{W}(\hat{s}; \beta)$  and for the second case where  $0 \leq v_i(s)$  we have  $W(\hat{s}) \geq \overline{W}(\hat{s}; \beta)$ .

### Proof of Proposition 6

Let  $(\hat{s}^*, \beta^*)$  be the saddle point solution to the characterization problem (16). If for any pair of dimensions  $i, j \in N$ ,  $\beta_i^*, \beta_j^* \in (0, 1)$  then

$$\int_{\hat{s}_i}^1 (v(s) - \frac{1 - P(s)}{p(s)} v'(s)) p(s) ds = \int_{\hat{s}_j}^1 (v(s) - \frac{1 - P(s)}{p(s)} v'(s)) p(s) ds$$

as otherwise we would have  $\beta_k^* = 0$  for one of the dimensions  $k = \{i, j\}$  as putting any weight on that dimension would not be minimizing.

The remaining case is when  $\beta_i^* = 0$  for some  $i \in N$ . Then it is optimal for the principal to implement threshold  $\hat{s}_i^* = 0$  when  $w_i > 0$  and  $\hat{s}_i^* = 1$  when  $w_i < 0$ . But then in both cases either all dimensions either choose the same marginal strategy as on  $i$ , in which case the result holds, or the solution to the minimization problem would be to have  $\beta_i^* > 0$ , a contradiction. To see this, in the first case for any  $j \in N$  unless  $\hat{s}_j^* = 0$

we have.

$$\int_0^1 \left( v(s) - \frac{1-P(s)}{p(s)} v'(s) \right) p(s) ds < \int_{\hat{s}_j^*}^1 \left( v(s) - \frac{1-P(s)}{p(s)} v'(s) \right) p(s) ds = h(\hat{s}_j^*) (1 - P(\hat{s}_j^*))$$

as  $v(s) < h(1) \leq 0$  for all  $s \in S \setminus \{1\}$ . For the second case, for all  $j \in N$  unless  $\hat{s}_j^* = 1$

$$0 = \int_1^1 \left( v(s) - \frac{1-P(s)}{p(s)} v'(s) \right) p(s) ds < \int_{\hat{s}_j^*}^1 \left( v(s) - \frac{1-P(s)}{p(s)} v'(s) \right) p(s) ds = h(\hat{s}_j^*) (1 - P(\hat{s}_j^*))$$

as  $0 \leq h(0) < v(s)$  for all  $s \in S \setminus \{0\}$ .

We see that we must have that  $\beta_i^* \in (0, 1)$  for all  $i \in N$ . Rearranging the implied equality, we have that for any  $i, j \in N$

$$\begin{aligned} \int_{\hat{s}_i^*}^1 \left( v(s) - \frac{1-P(s)}{p(s)} v'(s) \right) p(s) ds &= \int_{\hat{s}_j^*}^1 \left( v(s) - \frac{1-P(s)}{p(s)} v'(s) \right) p(s) ds \\ \Leftrightarrow \int_{\hat{s}_i^*}^{\hat{s}_j^*} \left( v(s) - \frac{1-P(s)}{p(s)} v'(s) \right) p(s) ds &= 0 \end{aligned}$$

which implies  $\hat{s}_i^* = \hat{s}_j^*$  by IVV.

## Proof of Proposition 5

For any fixed  $n$ , by Proposition 6 we have that  $\hat{s}_i^{(n)} = \hat{s}^{(n)}$  for all  $i \in N$  under the symmetric dimension space assumption. Given the principal's optimal thresholds  $\hat{s}^{(n)}$ , the saddle point problem in Theorem 1 is

$$\begin{aligned} &\min_{\beta \in [0,1]^n, \sum_{i \in N} \beta_i = 1} \overline{W}(\hat{s}^{(n)}, \beta) \\ &= \sum_{i \in N} [\beta_i \int_{\hat{s}^{(n)}}^1 \left( \frac{1}{n} v(s) - \frac{1-P(s)}{p(s)} \frac{1}{n} v'(s) \right) p(s) ds + \frac{1}{n} w(1 - P(\hat{s}^{(n)}))] \end{aligned}$$

This has solution  $\beta_i = \frac{1}{n}$  for all  $i \in N$ . The optimal threshold  $\hat{s}^n$  can then be characterized by the following variational inequality. For all  $s \in S$  we must have

$$\begin{aligned} (s - \hat{s}^{(n)}) \left( \frac{1}{n} \left( \frac{1}{n} v(s) - \frac{1 - P(s)}{p(s)} \frac{1}{n} v'(s) \right) + \frac{1}{n} w \right) &\geq 0 \\ \Leftrightarrow \\ (s - \hat{s}^{(n)}) \left( \left( \frac{1}{n} v(s) - \frac{1 - P(s)}{p(s)} \frac{1}{n} v'(s) \right) + w \right) &\geq 0 \end{aligned}$$

If  $w > 0$ , there exists an  $\underline{n}$  such that for all  $\tilde{n} \geq \underline{n}$  since  $h(s)$  is bounded.

$$\frac{1}{\tilde{n}} \left( h(0) - \frac{1 - P(\hat{0})}{p(\hat{0})} v'(0) \right) + w > 0$$

By IVV, we then have that  $\left( \frac{1}{n} v(s) - \frac{1 - P(s)}{p(s)} \frac{1}{n} v'(s) \right) + w > 0$  for all  $s \in S$ . Thus the only solution to the variational inequality is  $\hat{s}^{(\tilde{n})} = 0$ , as if  $\hat{s}^{(n)} > 0$  then the inequality is violated for all  $s \in [0, \hat{s}^{(n)}]$ .

We can make an analogous argument when  $w < 0$  to show the second part.

### Proof of Lemma 3

Let  $\hat{U}_i(s, a_i) = v_i(s)a_i + \bar{t}_i(a_i)$ . Clearly  $\hat{U}_i(s, 1) - \hat{U}_i(s, 0) = v_i(s) + \bar{t}_i(1) - \bar{t}_i(0)$  is increasing in  $s$ . Let  $q_i^h$  be the lottery in  $G_i$  that puts maximal probability on  $a_i = 1$ , and  $q_i^l$  be the analogous lottery for  $a_i = 0$ . If  $G_i = \Delta(A_i)$  then we must have  $q_i^h = 1$ ,  $q_i^l = 0$ . There is a threshold  $\hat{s}_i$  such that  $\hat{U}_i(s, 1) > \hat{U}_i(s, 0)$  for all  $s > \hat{s}_i$ , in which case type  $s$  will choose the lottery  $q_i^h$  from  $G_i$ . The same logic applies for types  $s < \hat{s}_i$  who will choose lottery  $q_i^l$ . The types who are indifferent have measure zero, so it is without loss for them to choose  $q_i^h$  also.

Given a random threshold strategy, the beliefs in the proposition statement ensure that the agent is indifferent between taking either action at the threshold.

$$v_i(\hat{s}_i) + \bar{t}_i(1) = \bar{t}_i(0)$$

They then prefer to take the lottery with the higher probability of  $a_i = 1$  if and only if their type is above the threshold. The restrictions contain all lotteries that are chosen by

some type under the random threshold strategy.

### Proof of Theorem 2

The same proof as in Theorem 1 applies up to Step 3. Step 3 still applies but the upper bound problem modified so that we no longer have  $\overline{U}(0) \in \mathbb{R}_{\geq 0}^n$ , but instead  $\overline{U}(0) \in \{U \in \mathbb{R}^n : \sum_{i \in N} U_i = 0\} \equiv \mathcal{U}^{SN}$ . The upper bound problem is then

$$\begin{aligned} \min_{\beta \in [0,1]^n, \sum_{i \in N} \beta_i = 1} \sup_{\hat{s} \in S^n, \overline{U}(0) \in \mathcal{U}^{SN}} \overline{W}(\hat{s}, \overline{U}(0), \beta) &= \sup_{\hat{s} \in S^n, \overline{U}(0) \in \mathcal{U}^{SN}} \overline{W}(\hat{s}, \overline{U}(0), \beta^*) \\ &= \sup_{\hat{s} \in S^n, \overline{U}(0) \in \mathcal{U}^{SN}} \min_{\beta \in [0,1]^n, \sum_{i \in N} \beta_i = 1} \overline{W}(\hat{s}, \overline{U}(0), \beta) \end{aligned}$$

with

$$\overline{W}(\hat{s}, \overline{U}(0), \beta) = \sum_{j \in N} [-\beta_j \overline{U}_j(0) + \int_{\hat{s}_i}^1 (\phi_i(s) + w_i) p(s) ds]$$

The saddle point values  $\beta^*$  must be such that  $\beta_i^* = \frac{1}{n}$  for all  $i \in N$ . Otherwise if  $\beta_j^* > \frac{1}{n}$  for some  $j \in N$  then for any  $U < 0$ , we can choose  $\overline{U}_j(0) = U$ , and for all  $l \in N \setminus \{j\}$   $\overline{U}_l(0) = -\frac{1}{n-1} U$ . This satisfies the sum-narrow participation constraint and allows us to obtain an arbitrarily large payoff by choosing  $U$ .

With  $\beta^* = (\frac{1}{n}, \dots, \frac{1}{n})$ , we can attain the value of the upper bound in the full problem by setting  $\overline{U}(0)$  such that  $\sum_{i \in N} \overline{U}_i(0) = 0$  and for any  $i, j \in N$  the statistical correctness constraint  $\overline{U}_i(0) - v_i(s)(1 - P(\hat{s}_i)) = \overline{U}_j(0) - v_j(s)(1 - P(\hat{s}_j))$  holds. This can be achieved by

$$\overline{U}_i(0) = v_i(s)(1 - P(\hat{s}_i)) - \frac{1}{n} \sum_{j \in N} v_j(s)(1 - P(\hat{s}_j)) \text{ for all } i \in N$$

### Proof of Proposition 7

For the second case, Proposition 4 the result holds if the principal is restricted to implement a threshold strategy. Since the principal can also implement a non-threshold strategy in the rational benchmark, the result holds without IVV.

For the first case, let  $g^{int}$  be an interval strategy with  $k$  intervals. Let  $N_l \subseteq N$  be the

subset of dimensions such that the agent takes the action  $a_i = 1$  for some type in interval  $l$ ; if  $i \in N_l$  then  $g_i^{int}(a_i|s) = 1$  for all  $s \in [z_{l-1}, z_l)$ .

For each dimension  $i \in N$ , we define the smallest type  $\underline{s}_i$  that both takes the action  $a_i = 1$  under  $g^{int}$  and for which the principal would get a positive payoff if all types above where to take the action  $a_i = 1$ ;  $\underline{s}_i = \min\{s \in S : g_i^{int}(1|s) = 1 \text{ and } w_i + v_i(s) \geq 0\}$ . Each  $\underline{s}_i$  is in one of the  $k$  intervals of  $g^{int}$ , for each  $i \in N$  denote this interval by  $[z_{l_i-1}, z_{l_i})$ . We have that  $\underline{s}_i$  is well defined because the lower bound of any  $[z_{l_i-1}, z_{l_i})$  is closed. The set of dimensions at which action 1 is taken in this interval is denoted  $N_{l_i}$  and includes  $i$ . For each dimension  $i$ , denote the interval at which action  $a_i = 1$  is first taken by  $m_i$ . The lowest type that takes action  $a_i = 1$  is then  $z_{m_i-1}$ .

Let  $i^* = \arg \min_{i \in N} v_i(\underline{s}_i)(1 - P(\underline{s}_i))$ , we define a threshold strategy  $\tilde{g}$  such that the action 1 is only ever taken on dimension  $i^*$ , and it is only taken by types above  $\underline{s}_{i^*}$ ;  $\tilde{g}_{i^*}(1|s) = \mathbb{1}\{s \geq \underline{s}_{i^*}\}$  and  $\tilde{g}_k(1|s) = 0$  for all  $s \in S$  and  $k \in N \setminus \{i^*\}$ .

We now calculate maximum loss to the principal from switching from  $g^{int}$  to  $\tilde{g}$ . We can write the principal's payoff from action  $a_i = 1$  in interval  $[z_{l-1}, z_l)$  as

$$\begin{aligned} & (w_i + v_i(z_{l-1}))(1 - P(z_{l-1})) - (w_i + v_i(z_l))(1 - P(z_l)) \\ &= (w_i + v_i(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_i(z_l) - v_i(z_{l-1}))(1 - P(z_l)) \end{aligned}$$

The loss in dimension  $i^*$  from switching to  $\tilde{g}$  is

$$\begin{aligned} & \sum_{l=m_{i^*}}^k \mathbb{1}\{i^* \in N_l\} [(w_{i^*} + v_{i^*}(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_{i^*}(z_l) - v_{i^*}(z_{l-1}))(1 - P(z_l))] \\ & - (w_{i^*} + v_{i^*}(\underline{s}_{i^*}))(1 - P(\underline{s}_{i^*})) \end{aligned}$$

By splitting the summation in the first line and rewriting the second line in terms of a



sum over intervals we have that this is equal to

$$\begin{aligned}
& \sum_{l=m_{i^*}}^{l_{i^*}-1} \mathbb{1}\{i^* \in N_l\} [(w_{i^*} + v_{i^*}(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_{i^*}(z_l) - v_{i^*}(z_{l-1}))(1 - P(z_l))] \\
& + \sum_{l=l_{i^*}}^k \mathbb{1}\{i^* \in N_l\} [(w_{i^*} + v_{i^*}(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_{i^*}(z_l) - v_{i^*}(z_{l-1}))(1 - P(z_l))] \\
& - \sum_{l=l_{i^*}}^k [(w_{i^*} + v_{i^*}(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_{i^*}(z_l) - v_{i^*}(z_{l-1}))(1 - P(z_l))] \\
& + (w_{i^*} + v_{i^*}(z_{l_{i^*}-1}))(P(\underline{s}_{i^*}) - P(z_{l_{i^*}-1})) - (v_{i^*}(\underline{s}_{i^*}) - v_{i^*}(z_{l_{i^*}-1}))(1 - P(\underline{s}_{i^*}))
\end{aligned}$$

This has the following upper bound, as if  $\underline{s}_{i^*} > t_{l_{i^*}-1}$  then  $w_{i^*} + v_{i^*}(z_{l_{i^*}-1}) < 0$ , and for all  $m_{i^*} \leq l < l_{i^*}$ ,  $w_{i^*} + v_{i^*}(z_{l-1}) < 0$  meaning the first and last lines are negative.

$$\begin{aligned}
& \sum_{l=l_{i^*}}^k \mathbb{1}\{i^* \in N_l\} [(w_{i^*} + v_{i^*}(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_{i^*}(z_l) - v_{i^*}(z_{l-1}))(1 - P(z_l))] \\
& - \sum_{l=l_{i^*}}^k [(w_{i^*} + v_{i^*}(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_{i^*}(z_l) - v_{i^*}(z_{l-1}))(1 - P(z_l))]
\end{aligned}$$

Adding and subtracting terms then gives us

$$\begin{aligned}
& \sum_{l=l_{i^*}}^k (P(z_l) - P(z_{l-1})) \left[ \sum_{j \in N_l} (w_j + v_j(z_{l-1})) - \sum_{i \in N^*} (w_i + v_i(z_{l-1})) \right] \\
& - \sum_{l=l_{i^*}}^k (1 - P(z_l)) \left[ \sum_{j \in N_l} (v_j(z_l) - v_j(z_{l-1})) - \sum_{i \in N^*} (v_i(z_l) - v_i(z_{l-1})) \right] \\
& - \sum_{l=l_{i^*}}^k \sum_{j \in N_l \setminus \{i^*\}} [(w_j + v_j(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_j(z_l) - v_j(z_{l-1}))(1 - P(z_l))] \\
& + \sum_{l=l_{i^*}}^k \sum_{i \in N^* \setminus \{i^*\}} [(w_i + v_i(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_i(z_l) - v_i(z_{l-1}))(1 - P(z_l))]
\end{aligned}$$

Subtracting out terms in the first line and last lines (the new last line is terms that do not subtract out), removing the last negative terms in the last line and rearranging the

second line then gives us a new upper bound

$$\begin{aligned}
& \sum_{l=l_{i^*}}^k (P(z_l) - P(z_{l-1})) \sum_{j \in N_l \setminus \{i^*\}} (w_j + v_j(z_{l-1})) \\
& - \sum_{l=l_{i^*}}^k (P(z_l) - P(z_{l-1})) \left[ \sum_{m=l_{i^*}}^l \sum_{j \in N_m} (v_j(z_m) - v_j(z_{m-1})) - \sum_{m=l_{i^*}}^l \sum_{i \in N^*} (v_i(z_m) - v_i(z_{m-1})) \right] \\
& - \sum_{l=l_{i^*}}^k \sum_{j \in N_l \setminus \{i^*\}} [(w_j + v_j(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_j(z_l) - v_j(z_{l-1}))(1 - P(z_l))] \\
& - \sum_{l=l_{i^*}}^k (P(z_l) - P(z_{l-1})) \mathbb{1}_{\{i^* \notin N_l\}} (w_{i^*} + v_{i^*}(z_{l-1})) \\
& \leq \sum_{l=l_{i^*}}^k (P(z_l) - P(z_{l-1})) \left[ \sum_{j \in N_l \setminus \{i^*\}} (w_j + v_j(z_{l-1})) \right] \\
& - \sum_{l=l_{i^*}}^k (P(z_l) - P(z_{l-1})) \left[ \sum_{m=l_{i^*}}^l \sum_{j \in N_m} (v_j(z_m) - v_j(z_{m-1})) - \sum_{m=l_{i^*}}^l \sum_{i \in N^*} (v_i(z_m) - v_i(z_{m-1})) \right] \\
& - \sum_{l=l_{i^*}}^k \sum_{j \in N_l \setminus \{i^*\}} [(w_j + v_j(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_j(z_l) - v_j(z_{l-1}))(1 - P(z_l))]
\end{aligned}$$

Where the inequality follows as the only interval  $l \geq l_{i^*}$  for which we can have  $w_{i^*} + v_{i^*}(z_{l-1}) < 0$  is  $l = l_{i^*}$ , and  $i^* \in N_{l_{i^*}}$  by definition.

For  $g^{int}$  to be IC, Lemma A.1 requires that for any intervals  $l, m \in \{1, \dots, k\}$  and any types  $s \in [z_{l-1}, z_l)$ ,  $s' \in [z_{m-1}, z_m)$

$$\sum_{i \in N_l} (v_i(s) - v_i(s')) \geq \sum_{i \in N_m} (v_i(s) - v_i(s'))$$

From this and continuity of  $v_i(\cdot)$ , we have that for any  $l \in \{1, \dots, k\}$

$$\sum_{m=l}^k \sum_{i \in N_m} (v_i(z_m) - v_i(z_{m-1})) \geq \sum_{m=l}^k \sum_{j \in N_l} (v_j(z_m) - v_j(z_{m-1})) = \sum_{j \in N_l} (v_j(z_m) - v_j(z_{l-1}))$$

Combining this with our previous inequality gives that the principal's worst case loss from switching the actions on dimension  $i^*$  is

$$\begin{aligned}
& \sum_{l=l_{i^*}}^k (P(z_l) - P(z_{l-1})) \sum_{j \in N_l \setminus \{i^*\}} (w_j + v_j(z_{l-1})) \\
& - \sum_{l=l_{i^*}}^k \sum_{j \in N_l \setminus \{i^*\}} [(w_j + v_j(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_j(z_l) - v_j(z_{l-1}))(1 - P(z_l))]
\end{aligned}$$

The loss from switching to  $\tilde{g}$  on dimensions  $j \in N \setminus \{i^*\}$  is

$$\begin{aligned}
& \sum_{l=1}^k \sum_{j \in N_l \setminus \{i^*\}} [(w_j + v_j(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_j(z_l) - v_j(z_{l-1}))(1 - P(z_l))] \\
& \leq \sum_{l=l_{i^*}}^k \sum_{j \in N_l \setminus \{i^*\}} [(w_j + v_j(z_{l-1}))(P(z_l) - P(z_{l-1})) - (v_j(z_l) - v_j(z_{l-1}))(1 - P(z_l))] \\
& + \sum_{l=1}^{l_{i^*}-1} (P(z_l) - P(z_{l-1})) \sum_{j \in N_l \setminus \{i^*\}} (w_j + v_j(z_{l-1}))
\end{aligned}$$

So on all dimensions the loss from switching is at most

$$\sum_{l=1}^k (P(z_l) - P(z_{l-1})) \sum_{j \in N_l \setminus \{i^*\}} (w_j + v_j(z_{l-1}))$$

Since for any  $j \in N$ , on any interval  $l$  such that  $j \in N_l$  and  $z_{l-1} < \underline{s}_j$  we have  $w_j + v_j(z_{l-1}) < 0$ , the worst case loss to the principal is bounded below  $\sum_{j \in N \setminus \{i^*\}} w_j (1 - P(\underline{s}_j))$ . We now show that under narrow inference we can implement a deterministic threshold strategy  $g^*$  such that  $g_i(1|s) = \mathbb{1}\{s \geq \underline{s}_i\}$  for all  $i \in N$ , with the same cost to the principal as the outcome function that implements the deterministic strategy  $\tilde{g}$  under rational inference. This gives our result as we have lost at most  $\sum_{j \in N \setminus \{i^*\}} w_j (1 - P(\underline{s}_j))$  from switching from  $g^{int}$  to  $\tilde{g}$  in the rational benchmark, and we can gain this for free under narrow inference.

The cost of the outcome function that implements  $\tilde{g}$  under rational inference is  $v_{i^*}(\underline{s}_{i^*})(1 - P(\underline{s}_{i^*}))$ . If we take beliefs such that  $\bar{t}_{i^*}(0) = 0$ ,  $\bar{t}_{i^*}(1) = -v_{i^*}(\underline{s}_{i^*})$ , and for all  $j \in N \setminus \{i^*\}$

$$\begin{aligned}
\bar{t}_j(0) &= v_j(\underline{s}_j)(1 - P(\underline{s}_j)) - v_{i^*}(\underline{s}_{i^*})(1 - P(\underline{s}_{i^*})) \geq 0 \\
\bar{t}_j(1) &= \bar{t}_j(0) - v_j(\underline{s}_j)
\end{aligned}$$

then from Proposition 2 there is an outcome function that implements these beliefs since  $g^*$  is a deterministic threshold strategy.

### Proof of Theorem 3

Step 1 of Theorem 1 works as before. From Step 2 onwards, we replace the objective by

$$\widehat{W}(x, \overline{U}(0), \beta) = \sum_{j \in N} [-\beta_j \overline{U}_j(0) + \int_{x_j}^1 (\widehat{\Phi}_j(P^{-1}(u)) + w_j) du]$$

where  $P^{-1}(\cdot)$  is the strictly increasing quantile function for the type distribution as before. This is an upper bound on the original objective as by definition of the upper concave envelope,  $\int_x^1 (\widehat{\Phi}_i(P^{-1}(u)) du = \widehat{\Phi}_i(x) \geq \Phi_i(x) = \int_x^1 (\phi_i(P^{-1}(u)) du$  for all  $i \in N$ ,  $x \in [0, 1]$ . Since by Lemma 3 only deterministic threshold strategies are NIC and all deterministic threshold strategies can be made NIC by some outcome function, the new objective remains an upper bound of the full problem even without increasing  $\phi_i(\cdot)$ .

Since  $\widehat{\Phi}_i(\cdot)$  is increasing in  $s$  the new objective is concave for fixed  $\beta$ , and we can make the same argument as we made for Theorem 1 and obtain a saddle point  $(\hat{s}^*, \overline{U}(0)^*, \beta^*)$ . This means  $(\hat{s}^*(\beta^*), \overline{U}(0)^*, \beta^*)$  is also a saddle point. In  $\widehat{W}(\hat{s}, \overline{U}(0), \beta)$  the term for  $\overline{U}(0)$  is separable, therefore  $\hat{s}^*(\beta) \in \arg \max_{\hat{s} \in [0, 1]^n} \widehat{W}(\hat{s}, \overline{U}(0), \beta)$  for all  $\overline{U}(0)$ . Thus, these thresholds maximize the new objective for fixed  $\beta$ . Since  $\hat{s}^* \in \arg \max_{s \in S} \widehat{W}(\hat{s}^*, \overline{U}(0)^*, \beta^*)$ , we must have  $\widehat{W}(\hat{s}^*(\beta^*), \overline{U}(0)^*, \beta^*) = \widehat{W}(\hat{s}^*, \overline{U}(0)^*, \beta^*)$ .

We now apply the same argument as Myerson (1981). For any  $i \in N$

$$\begin{aligned} \int_{\hat{s}_i}^1 (\phi_i(s) - \widehat{\Phi}_i(s)) p(s) ds &= \int_{P(\hat{s}_i)}^1 (\phi_i(P^{-1}(u)) - \widehat{\Phi}_i(P^{-1}(u))) du \\ &= \int_0^1 (\phi_i(P^{-1}(u)) - \widehat{\Phi}_i(P^{-1}(u))) du - \int_0^{P(\hat{s}_i)} (\phi_i(P^{-1}(u)) - \widehat{\Phi}_i(P^{-1}(u))) du \\ &= -(\Phi_i(1) - \widehat{\Phi}_i(1)) + (\Phi_i(0) - \widehat{\Phi}_i(0)) + \Phi_i(P(\hat{s}_i)) - \widehat{\Phi}_i(P(\hat{s}_i)) \\ &= \Phi_i(P(\hat{s}_i)) - \widehat{\Phi}_i(P(\hat{s}_i)) \end{aligned}$$

where we use the fact that  $\Phi_i(1) = \widehat{\Phi}_i(1)$  and  $\Phi_i(0) = \widehat{\Phi}_i(0)$ <sup>7</sup>.

We can show that  $\Phi_i(P(\hat{s}_i^*(\beta))) = \widehat{\Phi}_i(P(\hat{s}_i^*(\beta)))$  for any  $\beta$ . When  $\hat{s}_i^*(\beta) \in \{0, 1\}$  this is clear. If  $\hat{s}_i^*(\beta) \in (0, 1)$  then  $\beta_i \widehat{\Phi}_i(\hat{s}_i^*(\beta)) + w_i = 0$ , and by definition  $\hat{s}_i^*(\beta)$  is the smallest type satisfying this. Since  $\widehat{\Phi}_i(P(s)) > \Phi_i(P(s))$  only in intervals  $s \in [\underline{s}, \overline{s})$  where  $\widehat{\Phi}_i(s)$  is constant, at  $\hat{s}_i^*(\beta)$  we must have  $\Phi_i(P(\hat{s}_i^*(\beta))) = \widehat{\Phi}_i(P(\hat{s}_i^*(\beta)))$  otherwise

<sup>7</sup>This follows from continuity of  $\Phi_i(\cdot)$ .

we can find a smaller threshold in the maximizing set.

For any  $\beta$  we can write the old objective as

$$\begin{aligned}
& \overline{W}(\hat{s}, \overline{U}(0), \beta) \\
&= \sum_{j \in N} [-\beta_j \overline{U}_j(0) + \int_{\hat{s}_j}^1 (\beta_j \Phi_j(s) + w_j) p(s) ds] \\
&= \sum_{j \in N} [-\beta_j \overline{U}_j(0) + \int_{\hat{s}_j}^1 (\beta_j \widehat{\Phi}_j(s) + w_j) p(s) ds + \beta_j \int_{\hat{s}_j}^1 (\Phi_j(s) - \widehat{\Phi}_j(s)) p(s) ds] \\
&= \sum_{j \in N} [-\beta_j \overline{U}_j(0) + \int_{\hat{s}_j}^1 (\beta_j \widehat{\Phi}_j(s) + w_j) p(s) ds + \beta_j (\Phi_j(P(\hat{s}_j)) - \widehat{\Phi}_j(P(\hat{s}_j)))]
\end{aligned}$$

At  $\hat{s}^*(\beta^*)$ , since  $\Phi_i(P(\hat{s}_i^*(\beta^*))) = \widehat{\Phi}_i(P(\hat{s}_i^*(\beta^*)))$  the value the old and new objectives are identical for any  $\beta$ ;  $\overline{W}(\hat{s}^*(\beta^*), \overline{U}(0), \beta) = \widehat{\overline{W}}(\hat{s}^*(\beta^*), \overline{U}(0), \beta)$ . Thus,  $(\hat{s}^*(\beta^*), \overline{U}(0)^*, \beta^*)$  is a saddle point for the old objective also as for any  $(\hat{s}, \beta)$

$$\begin{aligned}
\overline{W}(\hat{s}^*(\beta^*), \overline{U}(0)^*, \beta) &\geq \min_{\beta \in [0,1]^n: \sum_{i \in N} \beta_i = 1} \widehat{\overline{W}}(\hat{s}^*(\beta), \overline{U}(0)^*, \beta) = \overline{W}(\hat{s}^*(\beta^*), \overline{U}(0)^*, \beta^*) \\
&= \widehat{\overline{W}}(\hat{s}^*(\beta), \overline{U}(0)^*, \beta^*) \geq \widehat{\overline{W}}(\hat{s}, \overline{U}(0)^*, \beta^*) \geq \overline{W}(\hat{s}, \overline{U}(0)^*, \beta^*)
\end{aligned}$$

We can then take apply Step 4 from Theorem 1 to show that this objective value can be achieved in the full problem, which completes the proof.

### A.3 Details for Examples

#### Example 3

We calculate the principal's optimal mechanism under the true welfare participation constraint when the principal is restricted to implementing a deterministic threshold strategy. Using the construction in Proposition 2, we can write the outcome function

that implements the beliefs for a threshold strategy where  $\hat{s}_1 < \hat{s}_2$  as

$$t(0, 0) = \bar{t}_1(0) \quad (49)$$

$$t(1, 0) = \frac{P(\hat{s}_2)}{P(\hat{s}_2) - P(\hat{s}_1)} \bar{t}_2(0) - \frac{P(\hat{s}_1)}{P(\hat{s}_2) - P(\hat{s}_1)} \bar{t}_1(0) \quad (50)$$

$$t(1, 1) = \bar{t}_2(1) \quad (51)$$

This means that the true welfare participation constraints in terms of beliefs are

$$\bar{t}_1(0) \geq 0 \quad (52)$$

$$\frac{P(\hat{s}_2)}{P(\hat{s}_2) - P(\hat{s}_1)} \bar{t}_2(0) - \frac{P(\hat{s}_1)}{P(\hat{s}_2) - P(\hat{s}_1)} \bar{t}_1(0) + v_1(\hat{s}_1) \geq 0 \quad (53)$$

$$\bar{t}_2(1) + v_1(\hat{s}_2) + v_2(\hat{s}_2) \geq 0 \quad (54)$$

The beliefs must also satisfy the statistical correctness constraint by Lemma 2. This requires that

$$\begin{aligned} P(\hat{s}_1) \bar{t}_1(0) + (1 - P(\hat{s}_1)) \bar{t}_1(1) &= P(\hat{s}_2) \bar{t}_2(0) + (1 - P(\hat{s}_2)) \bar{t}_2(1) \\ \Leftrightarrow \bar{t}_1(0) - (1 - P(\hat{s}_1)) v_1(\hat{s}_1) &= \bar{t}_2(0) - (1 - P(\hat{s}_2)) v_2(\hat{s}_2) \end{aligned} \quad (55)$$

The requirement that both these sets of constraints holds can be reduced to

$$\bar{t}_1(0) \geq 0 \quad (56)$$

$$\bar{t}_2(0) \geq \min\left\{-v_1(\hat{s}_1), -\frac{P(\hat{s}_2)}{P(\hat{s}_1)} v_1(\hat{s}_1) + \frac{P(\hat{s}_1)}{P(\hat{s}_2) - P(\hat{s}_1)} [(1 - P(\hat{s}_1)) v_1(\hat{s}_1) - (1 - P(\hat{s}_2)) v_2(\hat{s}_2)]\right\} \quad (57)$$

Under the parametric assumptions

$$\begin{aligned} &\min\left\{-v_1(\hat{s}_1), -\frac{P(\hat{s}_2)}{P(\hat{s}_1)} v_1(\hat{s}_1) + \frac{P(\hat{s}_1)}{P(\hat{s}_2) - P(\hat{s}_1)} [(1 - P(\hat{s}_1)) v_1(\hat{s}_1) - (1 - P(\hat{s}_2)) v_2(\hat{s}_2)]\right\} \\ &= \min\{r(1 - \hat{s}_2), r(1 - \hat{s}_1) + 2r\hat{s}_1 - r\hat{s}_1(\hat{s}_1 + \hat{s}_2)\} \end{aligned}$$

Therefore by Lemma A.3 and the parametric assumptions, the objective value of the

principal is

$$\begin{aligned}
\max_{c \in \mathcal{M}^{NICR}(g)} W(t, g, S) &= \min\{W_{0,0}(\hat{s}_1, \hat{s}_2), W_{1,0}(\hat{s}_1, \hat{s}_2), W_{1,1}(\hat{s}_1, \hat{s}_2)\} \\
W_{0,0}(\hat{s}_1, \hat{s}_2) &= -r(1 - \hat{s}_1)^2 + (1 - \hat{s}_1)w_1 + (1 - \hat{s}_2)w_2 \\
W_{1,0}(\hat{s}_1, \hat{s}_2) &= -r(1 - \hat{s}_1) + 2r\hat{s}_1 - r\hat{s}_1(\hat{s}_1 + \hat{s}_2) - r(1 - \hat{s}_2)^2 + (1 - \hat{s}_1)w_1 + (1 - \hat{s}_2)w_2 \\
W_{1,1}(\hat{s}_1, \hat{s}_2) &= -r(1 - \hat{s}_2) - r(1 - \hat{s}_2)^2 + (1 - \hat{s}_1)w_1 + (1 - \hat{s}_2)w_2
\end{aligned} \tag{58}$$

We consider parameters at which there is a saddle point where only the true welfare participation constraint for the worst-off type choosing  $(1, 0)$  binds. Consider the thresholds  $\hat{s}_1 < \hat{s}_2$  that maximize

$$W_{1,0}(\hat{s}_1, \hat{s}_2) = -r(1 - \hat{s}_1) - 2r\hat{s}_1 + r\hat{s}_1(\hat{s}_1 + \hat{s}_2) - r(1 - \hat{s}_2)^2 + (1 - \hat{s}_1)w_1 + (1 - \hat{s}_2)w_2 \tag{59}$$

When the maximizing thresholds are interior, they are equal to

$$\begin{aligned}
\hat{s}_1^* &= \frac{4}{3} - \frac{2}{3} \frac{w_1}{r} + \frac{1}{3} \frac{w_2}{r} \\
\hat{s}_2^* &= \frac{1}{3} - \frac{2}{3} \frac{w_2}{r} + \frac{1}{3} \frac{w_1}{r}
\end{aligned}$$

Under parameters  $w_1 = 1, w_2 = 0.3, r = 0.5$ . We then have that these thresholds and the corresponding value of the principal's objective are

$$\begin{aligned}
\hat{s}_1^* &= 0.2, \hat{s}_2^* = 0.6 \\
\min\{W_{0,0}(\hat{s}_1^*, \hat{s}_2^*), W_{1,0}(\hat{s}_1^*, \hat{s}_2^*), W_{1,1}(\hat{s}_1^*, \hat{s}_2^*)\} &= W_{1,0}(\hat{s}_1^*, \hat{s}_2^*) \\
&= \min\{0.6, 0.56, 0.64\} = 0.56
\end{aligned} \tag{60}$$

The case where  $\hat{s}_1 > \hat{s}_2$  is symmetric except that the objective value is lower since  $w_1 > w_2$ . We now consider the symmetric threshold case  $\hat{s}_1 = \hat{s}_2 = \hat{s}$ . Here the principal

chooses a single threshold to maximize

$$\begin{aligned}
\max_{c \in \mathcal{M}^{NICR}(g)} W(t, g, S) &= \min\{W_{0,0}(\hat{s}), W_{1,1}(\hat{s})\} = W_{1,1}(\hat{s}) \\
W_{0,0}(\hat{s}) &= -r(1 - \hat{s})^2 + (1 - \hat{s})w_1 + (1 - \hat{s})w_2 \\
W_{1,1}(\hat{s}) &= -r(1 - \hat{s}) - r(1 - \hat{s})^2 + (1 - \hat{s})w_1 + (1 - \hat{s})w_2
\end{aligned} \tag{61}$$

The optimal threshold under our parameters is  $\hat{s} = \frac{3}{2} - \frac{w_1 + w_2}{2r} = 0.2$ , which gives an objective value of  $\max_{c \in \mathcal{M}^{NICR}(g)} W(t, g, S) = W_{1,1}(\hat{s}) = 0.32$ . This is less than the objective value obtained by  $\hat{s}_1^* = 0.2, \hat{s}_2^* = 0.6$ , which is therefore optimal for the principal as it is a saddle point deterministic threshold strategy.



## References

- [1] Bohren, J. A. and D. N. Hauser (2021). Learning with heterogeneous misspecified models: Characterization and robustness. *Econometrica* 89(6), 3025–3077. [8](#)
- [2] Carroll, G. (2017). Robustness and separation in multidimensional screening. *Econometrica* 85(2), 453–488. [14](#), [38](#)
- [3] Eliaz, K. and R. Spiegler (2006). Contracting with diversely naive agents. *The Review of Economic Studies* 73(3), 689–714. [9](#)
- [4] Eliaz, K. and R. Spiegler (2020). A model of competing narratives. *American Economic Review* 110(12), 3786–3816. [8](#)
- [5] Eliaz, K. and R. Spiegler (2024). News media as suppliers of narratives (and information). *arXiv preprint arXiv:2403.09155*. [8](#)
- [6] Enke, B. and F. Zimmermann (2019). Correlation neglect in belief formation. *The review of economic studies* 86(1), 313–332. [7](#)
- [7] Esponda, I. and D. Pouzo (2016). Berk–nash equilibrium: A framework for modeling agents with misspecified models. *Econometrica* 84(3), 1093–1130. [8](#)
- [8] Farhi, E. and X. Gabaix (2020). Optimal taxation with behavioral agents. *American Economic Review* 110(1), 298–336. [9](#)
- [9] Fernbach, P. M., A. Darlow, and S. A. Sloman (2010). Neglect of alternative causes in predictive but not diagnostic reasoning. *Psychological Science* 21(3), 329–336. [7](#)
- [10] Frick, M., R. Iijima, and Y. Ishii (2020). Misinterpreting others and the fragility of social learning. *Econometrica* 88(6), 2281–2328. [8](#)
- [11] Fudenberg, D., G. Lanzani, and P. Strack (2021). Limit points of endogenous misspecified learning. *Econometrica* 89(3), 1065–1098. [8](#)
- [12] Graeber, T. (2023). Inattentive inference. *Journal of the European Economic Association* 21(2), 560–592. [7](#)
- [13] He, S. and S. Kučinskas (2024). Expectation formation with correlated variables. *The Economic Journal* 134(660), 1517–1544. [7](#)
- [14] Heidhues, P. and B. Köszegi (2010). Exploiting naivete about self-control in the credit market. *American Economic Review* 100(5), 2279–2303. [9](#)
- [15] Heidhues, P., B. Köszegi, and P. Strack (2018). Unrealistic expectations and misguided learning. *Econometrica* 86(4), 1159–1214. [8](#)
- [16] Herweg, F., D. Müller, and P. Weinschenk (2010). Binary payment schemes: Moral hazard and loss aversion. *American Economic Review* 100(5), 2451–2477. [9](#)
- [17] Jehiel, P. (2005). Analogy-based expectation equilibrium. *Journal of Economic Theory* 123(2), 81–104. [8](#), [35](#)

- [18] Jehiel, P. (2011). Manipulative auction design. *Theoretical economics* 6(2), 185–217. [9](#)
- [19] Jehiel, P. and K. Mierendorff (2024). Auction design with data-driven misspecifications: Inefficiency in private value auctions with correlation. *Theoretical Economics*, (Forthcoming). [9](#)
- [20] Kőszegi, B. (2014). Behavioral contract theory. *Journal of Economic Literature* 52(4), 1075–1118. [9](#)
- [21] Kőszegi, B. and F. Matějka (2020). Choice simplification: A theory of mental budgeting and naive diversification. *The Quarterly Journal of Economics* 135(2), 1153–1207. [8](#)
- [22] Lian, C. (2021). A theory of narrow thinking. *The Review of Economic Studies* 88(5), 2344–2374. [8](#)
- [23] Lockwood, B. B. (2020). Optimal income taxation with present bias. *American Economic Journal: Economic Policy* 12(4), 298–327. [9](#)
- [24] Milgrom, P. and I. Segal (2002). Envelope theorems for arbitrary choice sets. *Econometrica* 70(2), 583–601. [37](#), [42](#)
- [25] Myerson, R. B. (1981). Optimal auction design. *Mathematics of operations research* 6(1), 58–73. [33](#), [60](#)
- [26] O’Donoghue, T. and M. Rabin (2006). Optimal sin taxes. *Journal of Public Economics* 90(10–11), 1825–1849. [9](#)
- [27] Piccione, M. and A. Rubinstein (2003). Modeling the economic interaction of agents with diverse abilities to recognize equilibrium patterns. *Journal of the European economic association* 1(1), 212–223. [8](#)
- [28] Rabin, M. and G. Weizsäcker (2009). Narrow bracketing and dominated choices. *American Economic Review* 99(4), 1508–1543. [8](#)
- [29] Read, D., G. Loewenstein, M. Rabin, G. Keren, and D. Laibson (1999). Choice bracketing. *Journal of Risk and Uncertainty* 19(1/3), 171–197. [8](#)
- [30] Rubinstein, A. (1993). On price recognition and computational complexity in a monopolistic model. *Journal of Political Economy* 101(3), 473–484. [8](#)
- [31] Schumacher, H. and H. C. Thysen (2022). Equilibrium contracts and boundedly rational expectations. *Theoretical Economics* 17(1), 371–414. [8](#)
- [32] Sion, M. (1958). On general minimax theorems. *Pacific Journal of Mathematics* 8(1), 171 – 176. [50](#)
- [33] Spiegel, R. (2016). Bayesian networks and boundedly rational expectations. *The Quarterly Journal of Economics* 131(3), 1243–1290. [8](#)

- [34] Spiegel, R. (2020). Behavioral implications of causal misperceptions. *Annual Review of Economics* 12(1), 81–106. [17](#)
- [35] Spinnewijn, J. (2015). Unemployed but optimistic: Optimal insurance design with biased beliefs. *Journal of the European Economic Association* 13(1), 130–167. [9](#)
- [36] Thaler, R. (1985). Mental accounting and consumer choice. *Marketing science* 4(3), 199–214. [8](#)
- [37] Thaler, R. H. (1999). Mental accounting matters. *Journal of Behavioral decision making* 12(3), 183–206. [8](#)
- [38] Tversky, A. and D. Kahneman (1981). The framing of decisions and the psychology of choice. *science* 211(4481), 453–458. [8](#)