

**ROBUST IMPLEMENTATION:  
THE ROLE OF LARGE SPACES**

**By**

**Dirk Bergemann and Stephen Morris**

**June 2005**

**COWLES FOUNDATION DISCUSSION PAPER NO. 1519**



**COWLES FOUNDATION FOR RESEARCH IN ECONOMICS  
YALE UNIVERSITY  
Box 208281  
New Haven, Connecticut 06520-8281**

**<http://cowles.econ.yale.edu/>**

# Robust Implementation: The Role of Large Type Spaces\*

Dirk Bergemann<sup>†</sup>

Stephen Morris<sup>‡</sup>

First Version: March 2003

This Version: June 2005

## Abstract

A social choice function is robustly implemented if every equilibrium on every type space achieves outcomes consistent with a social choice function. We identify a *robust monotonicity* condition that is necessary and (with mild extra assumptions) sufficient for robust implementation.

Robust monotonicity is strictly stronger than both Maskin monotonicity (necessary and almost sufficient for complete information implementation) and ex post monotonicity (necessary and almost sufficient for ex post implementation). It is equivalent to Bayesian monotonicity on all type spaces. It requires that there not be too much interdependence of types.

We characterize robust monotonicity for some interesting economic environments. We identify conditions where, if robust implementation is possible, it is possible in a direct mechanism. We identify conditions where, if robust implementation is not possible, virtual robust implementation is not possible either.

KEYWORDS: Mechanism Design, Implementation, Robustness, Common Knowledge, Interim Equilibrium, Iterative Deletion, Dominant Strategies.

JEL CLASSIFICATION: C79, D82

---

\*This research is supported by NSF Grant #SES-0095321. The first author gratefully acknowledges support through a DFG Mercator Research Professorship at the Center of Economic Studies at the University of Munich. Some of the results reported in this paper appeared in early versions of Bergemann and Morris (2004); earlier versions of this paper contained results on ex post implementation which now appear in Bergemann and Morris (2005a); this paper contains an interim report on an ongoing research project.

<sup>†</sup>Department of Economics, Yale University, 28 Hillhouse Avenue, New Haven, CT 06511, dirk.bergemann@yale.edu.

<sup>‡</sup>Department of Economics, Yale University, 30 Hillhouse Avenue, New Haven, CT 06511, stephen.morris@yale.edu.

## 1 Introduction

The mechanism design literature provides a powerful characterization of which social choice functions can be achieved when the designer has incomplete information about agents' types. If we assume a commonly known common prior over the possible types of agents, the revelation principle establishes that if the social choice function can arise as an equilibrium in some mechanism, then it will arise in a truth-telling equilibrium of the direct mechanism (where each agent truthfully reports his type and the designer chooses an outcome assuming they are telling the truth). Thus the Bayesian incentive compatibility constraints characterize whether a social choice function is implementable in this sense.

There are two important limitations of Bayesian incentive compatibility analysis. First, the analysis typically assumes a commonly known common prior over the agents' types. This assumption may be too stringent in practise. In the spirit of the "Wilson doctrine" (Wilson (1987)), we would like implementation results that are *robust* to different assumptions about what players do or do not know about other agents' types. Second, the revelation principle only establishes that the direct mechanism has *an* equilibrium that achieves the social choice function. In general, there may be other equilibria that deliver undesirable outcomes. We would like to achieve *full* implementation, i.e., show the existence of a mechanism all of whose equilibria deliver the social choice function. We studied the first "robustness" problem in an earlier work, Bergemann and Morris (2004). The second "full implementation" problem has been the subject of a large literature. In the incomplete information context, key full implementation references are Postlewaite and Schmeidler (1986), Palfrey and Srivastava (1989) and Jackson (1991). In this paper, we study "robust implementation" where we require robustness and full implementation simultaneously. Requiring both simultaneously adds extra structure to the problem and enables us to derive distinctive new economic results.

Interim implementation on all type spaces is possible if and only if it is possible to implement the social choice function using an iterative deletion procedure. We refer to the resulting notion as *iterative implementation*. We fix a mechanism and iteratively delete messages for each payoff type that are strictly dominated by another message for each payoff type profile and message profile that has survived the procedure. This observation about iterative deletion illustrates a general point well-known from the literature on epistemic foundations of game theory (e.g., Brandenburger and Dekel (1987), Battigalli and Siniscalchi (2003)): equilibrium solution concepts only have bite if we make strong assumptions about type spaces, i.e., we assume small type spaces where the common prior assumption holds.

We exploit this equivalence between robust and iterative implementation to obtain necessary and

sufficient conditions for robust implementation in general environments. Our necessity argument is conceptually novel, exploiting the iterative characterization. The necessary conditions for robust implementation are ex post incentive compatibility of the social choice function and a condition - *robust monotonicity* - that is equivalent to requiring interim monotonicity on every type space. Suppose that we fix a "deception" specifying, for each payoff type  $\theta_i$  of each agent, a set of types that he might misreport himself to be. We require that for some agent  $i$  and a type misreport of agent  $i$  under the deception, for every misreport  $\theta'_{-i}$  that the other agents might make under the deception, there exists an outcome  $y$  which is strictly preferred by agent  $i$  to the outcome he would receive under the social choice function for *every* possible payoff type profile that might misreport  $\theta'_{-i}$ ; where this outcome  $y$  satisfies the extra restriction that no payoff type of agent  $i$  prefers outcome  $y$  to the social choice function if the other agents were really types  $\theta'_{-i}$ . This condition - while a little convoluted - is easier to interpret than the interim (Bayesian) monotonicity conditions. It is very strong and implies both Maskin monotonicity and ex post monotonicity conditions (but is strictly weaker than dominant strategies). We will present examples to illustrate the relationship between these monotonicity conditions.

The sufficiency argument requires only a modest strengthening of the necessary condition by guaranteeing the existence of a "bad" outcome. With the existence of a bad outcome, we show that the necessary conditions are also sufficient for robust implementation. The sufficient conditions guarantee robust implementation in pure, but more generally also in mixed strategies. Our robust analysis thus removes the frequent gap between pure and mixed strategy implementation in the literature.

The iterative characterization comes with the additional benefit that tight implementation results can be proved via a fixed point of a contraction mapping. In particular, we consider a general class of interdependent preferences in which the private types of the agents can be linearly aggregated. In this environment we show that the social choice function can be robustly implemented if and only if the interdependence is not too large. If  $\gamma$  is the weight of the type of agent  $j$  (relative to the type of agent  $i$ ) for the utility of agent  $i$ , then the robust implementation condition can simply be stated as:  $\gamma < 1/(I - 1)$ , where  $I$  is the number of agents. Surprisingly, the converse result for  $\gamma > 1/(I - 1)$  even extends to a robust version of virtual utility. In other words, we also show that if  $\gamma > 1/(I - 1)$ , then not only robust implementation, but even robust virtual implementation fails. We further illustrate the strength of the contraction mapping idea in the implementation context with two examples, one of a public good and one of a private good (single unit auction) allocation problem with quasilinear utility. An important paper of Chung and Ely (2001) analyzed the single (and multi-unit) auction with interdependent valuations with dominance solvability (elimination of weakly rather than strictly dominated actions). In a linear and symmetric setting, they reported

sufficient conditions for direct implementation that coincide with the ones derived here. We show that in the environment with linear aggregation, under strict incentive compatibility, the basic insight extends from the single unit auction model to general allocations models, with elimination of strictly dominated actions only (thus Chung and Ely (2001) require deletion of weakly dominated strategies only because incentive constraints are weak). We also prove a converse result: if there is too much interdependence, then neither the direct nor any augmented mechanism can robustly implement the social choice function (this result will also hold with deletion of weakly dominated strategies).

In the implementation literature, it is a standard practice to obtain the sufficiency results with augmented mechanisms. By augmenting the direct mechanism with additional messages, the designer may elicit additional information about undesirable equilibrium play by the agents. Yet, in many environments common to applied mechanism designs, such a single crossing or supermodular preferences, the structure of the preferences may already permit direct implementation (see a companion paper, Bergemann and Morris (2005), for direct implementation results in ex post equilibrium). We thus provide necessary and sufficient conditions for robust implementation in the direct mechanism. In the direct mechanism, the agents can alert the designer only by a report of their type. In consequence, the incentive compatibility conditions for the rewards are identical to the truth-telling constraints, and the necessary and sufficient conditions for robust implementation coincide.

The results in this paper concern full implementation. An earlier paper of ours, Bergemann and Morris (2004), addresses the analogous questions of robustness to rich type spaces, but looking at the question of partial implementation, i.e., does there exist a mechanism such that *some* equilibrium implements the social choice function. We showed that ex post (partial) implementation of the social choice function is a necessary and sufficient condition for partial implementation on all type spaces.<sup>1</sup> This paper establishes that an analogous result does *not* hold for full implementation.

In a companion paper, Bergemann and Morris (2005), we therefore investigate the notion of ex post implementation. The necessary and sufficient conditions then straddle the Nash and Bayesian implementation conditions as an ex post equilibrium is a Nash equilibrium at every incomplete information (Bayesian) type profile. However in contrast to the iterative argument pursued here, the basic reasoning in Bergemann and Morris (2005) invokes more traditional equilibrium arguments. By comparing the conditions for ex post and robust implementation, it becomes apparent that robust implementation typically imposes additional constraints on the allocation problem. In Bergemann and Morris (2005), we showed that in single crossing environments, the same single

---

<sup>1</sup>This result does not extend to social choice correspondences.

crossing conditions which guarantee incentive compatibility also guarantee full implementation. In contrast, in the linear aggregation environment, we show that robust implementation imposes a strict bound on the interdependence of the preferences, which is not required by the truthtelling conditions. The contraction mapping behind the iterative argument directly pointed at the source of the restriction of the interaction term.

In this paper, we follow the classic implementation literature in allowing for arbitrary mechanisms, including modulo and integer games. By allowing for these mechanisms, we are able to make tight connections with the existing implementation literature. Allowing for these badly behaved mechanisms does complicate our analysis: for example, we must allow for transfinite iterated deletion of best responses in our definition of iterative implementation. Of course, when direct implementation turns out to be possible, we do not need badly behaved mechanisms. We also report some results on what happens with "nice" mechanisms where rationalizable messages always exist and best responses are well defined. We can further strengthen our necessary conditions in this case.

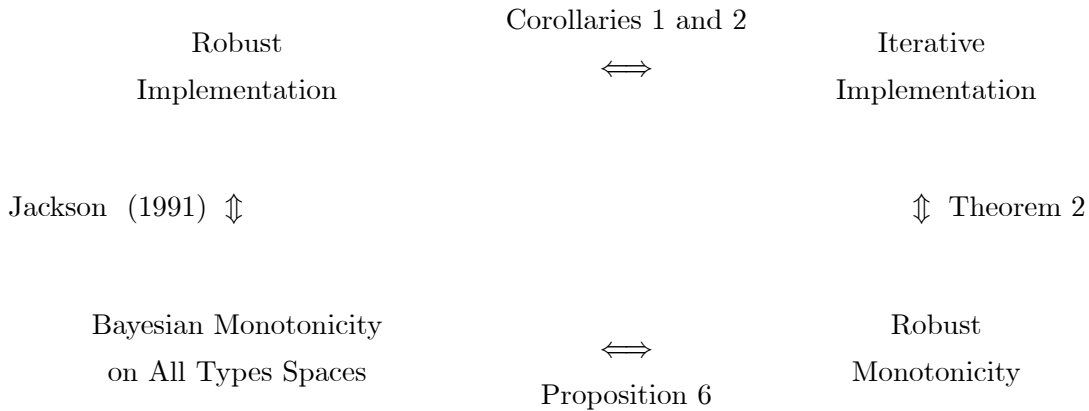


Figure 1: Relationship between Bayesian and Robust Implementation / Monotonicity

Our results extend the classic literature on Bayesian implementation due to Postlewaite and Schmeidler (1986), Palfrey and Srivastava (1989) and Jackson (1991). We focus in this paper on an indirect approach to extending these results. We first note the equivalence between robust implementation and iterative implementation. We then exploit the equivalence to report a direct argument showing that robust monotonicity is a necessary and almost sufficient condition for iterative implementation. But in the light of the classic literature, we know that a necessary and almost sufficient condition for robust implementation must be Bayesian monotonicity on all type spaces. We confirm and clarify our results by directly checking that robust monotonicity is equivalent to

Bayesian (or interim) monotonicity on all type spaces. Figure 1 gives a stylized account of the connection between these alternative approaches.

The remainder of the paper is organized as follows. Section 2 describes the formal environment and solution concepts. Section 3 considers four examples that illustrate the main results of the paper. Section 4 establishes the relation between iterative and robust implementation. Section 5 establishes necessary conditions for robust implementation in the direct mechanism. Section 6 reports our main result on the necessary and sufficient conditions for robust implementation. Section 7 explores the link between robust and virtual implementation. Section 8 considers the preference environment with a linear aggregation of the types and obtains sharp implementation results. Section 9 discusses extensions and variations of our implementation results, examining the role of lotteries, pure strategies and "nice" mechanisms. Section 10 concludes. The Appendix contains some additional examples and proofs.

## 2 Setup

### 2.1 The Payoff Environment

We consider a finite set of agents,  $1, 2, \dots, I$ . Agent  $i$ 's *payoff type* is  $\theta_i \in \Theta_i$ . We write  $\theta \in \Theta = \Theta_1 \times \dots \times \Theta_I$ . There is a set of outcomes  $Z$ . We assume that each  $\Theta_i$  and  $Z$  are countable. Each individual has von Neumann Morgenstern utility function  $u_i : Z \times \Theta \rightarrow \mathbb{R}$ . Thus we are in the world of interdependent types, where an agent's utility depends on other agents' payoff types. We allow for lotteries over deterministic outcomes.<sup>2</sup> Let  $Y = \Delta(Z)$  and extend  $u_i$  to the domain  $Y \times \Theta$  in the usual way:

$$u_i(y, \theta) = \sum_{z \in Z} y(z) u_i(z, \theta).$$

A social choice function is a mapping  $f : \Theta \rightarrow Y$ . If the true payoff type profile is  $\theta$ , the planner would like the outcome to be  $f(\theta)$ . In this paper, we restrict our analysis to the implementation of a social choice function rather than a social choice correspondence or set.<sup>3</sup>

### 2.2 Type Spaces

We are interested in analyzing behavior in a variety of type spaces, many of them with a richer set of types than payoff types. For this purpose, we shall refer to agent  $i$ 's *type* as  $t_i \in T_i$ , where

<sup>2</sup>The role of the lottery assumption and what happens when we drop it are discussed in Section 9.1.

<sup>3</sup>The extension to social choice correspondences might not be straightforward. One reason is that with social choice correspondences, the incentive compatibility conditions that arise from requiring partial implementation would typically be weaker than ex post incentive compatibility, as shown by examples in Bergemann and Morris (2004).

$T_i$  is a countable set.<sup>4</sup> A type of agent  $i$  must include a description of his payoff type. Thus there is a function  $\hat{\theta}_i : T_i \rightarrow \Theta_i$  with  $\hat{\theta}_i(t_i)$  being agent  $i$ 's payoff type when his type is  $t_i$ . A type of agent  $i$  must also include a description of his beliefs about the types of the other agents; thus there is a function  $\hat{\pi}_i : T_i \rightarrow \Delta(T_{-i})$  with  $\hat{\pi}_i(t_i)$  being agent  $i$ 's *belief type* when his type is  $t_i$ . Thus  $\hat{\pi}_i(t_{-i})[t_i]$  is the probability that type  $t_i$  of agent  $i$  assigns to other agents having types  $t_{-i}$ . A *type space* is a collection:

$$\mathcal{T} = \left( T_i, \hat{\theta}_i, \hat{\pi}_i \right)_{i=1}^I.$$

### 2.3 Mechanisms

A planner must choose a *game form* or *mechanism* for the agents to play in order to determine the social outcome. Let  $M_i$  be the countably infinite set of messages available to agent  $i$ .<sup>5</sup> Let  $g(m)$  be the distribution over outcomes if action profile  $m$  is chosen. Thus a mechanism is a collection

$$\mathcal{M} = (M_1, \dots, M_I, g(\cdot)),$$

where  $g : M \rightarrow Y$ .

### 2.4 Solution Concepts

Now holding fixed the payoff environment, we can combine a type space  $\mathcal{T}$  with a mechanism  $\mathcal{M}$  to get an incomplete information game  $(\mathcal{T}, \mathcal{M})$ . The payoff of agent  $i$  if message profile  $m$  is chosen and type profile  $t$  is realized is then given by

$$u_i(g(m), \hat{\theta}(t)).$$

A pure strategy for agent  $i$  in the incomplete information game  $(\mathcal{T}, \mathcal{M})$  is given by

$$s_i : T_i \rightarrow M_i.$$

A (behavioral) strategy is given by

$$\sigma_i : T_i \rightarrow \Delta(M_i).$$

The objective of this paper is to obtain implementation results for interim, or Bayesian Nash,

---

<sup>4</sup>The countable types restriction clarifies the relation to the existing literature. In Section 9.3, we discuss what happens if we allow for uncountable type spaces.

<sup>5</sup>This assumption clarifies the relation with the existing literature. We discuss in Section 9.2 what happens if we restrict attention to finite messages.



equilibria on all possible types spaces.<sup>6</sup> The notion of interim equilibrium for a given type space  $\mathcal{T}$  is defined in the usual way.

**Definition 1 (Interim equilibrium)**

A strategy profile  $\sigma = (\sigma_1, \dots, \sigma_I)$  is an interim equilibrium of the game  $(\mathcal{T}, \mathcal{M})$  if, for all  $i$ ,  $t_i$  and  $m_i$  with  $\sigma_i(m_i|t_i) > 0$ ,

$$\begin{aligned} & \sum_{t_{-i} \in T_{-i}} \sum_{m_{-i} \in M_{-i}} \left( \prod_{j \neq i} \sigma_j(m_j|t_j) \right) u_i \left( g(m_i, m_{-i}), \hat{\theta}(t) \right) \hat{\pi}_i(t_{-i}) [t_i] \\ & \geq \sum_{t_{-i} \in T_{-i}} \sum_{m_{-i} \in M_{-i}} \left( \prod_{j \neq i} \sigma_j(m_j|t_j) \right) u_i \left( g(m'_i, m_{-i}), \hat{\theta}(t) \right) \hat{\pi}_i(t_{-i}) [t_i] \end{aligned}$$

for all  $m'_i$ .

The concern for robustness, expressed by the qualifying condition, for all type spaces, pulls the interim equilibrium in the direction of rationalizability. Consequently we define a message correspondence profile  $S = (S_1, \dots, S_I)$ , where each  $S_i : \Theta_i \rightarrow 2^{M_i}$  and we write  $\mathcal{S}$  for the collection of message correspondence profiles. The collection  $\mathcal{S}$  is a lattice with the natural ordering:  $S \leq S'$  if  $S_i(\theta_i) \subseteq S'_i(\theta_i)$  for all  $i$  and  $\theta_i$ . The largest element is  $\bar{S} = (\bar{S}_1, \dots, \bar{S}_I)$ , where  $\bar{S}_i(\theta_i) = M_i$  for each  $i$  and  $\theta_i$ . The smallest element is  $\underline{S} = (\underline{S}_1, \dots, \underline{S}_I)$ , where  $\underline{S}_i(\theta_i) = \emptyset$  for each  $i$  and  $\theta_i$ .

We define an operator  $b$  to iteratively eliminate never best responses. To this end, we denote the belief of agent  $i$  over message and payoff type profiles of the remaining agents by  $\lambda_i \in \Delta(M_{-i} \times \Theta_{-i})$ . The operator  $b : \mathcal{S} \rightarrow \mathcal{S}$  is now defined as:

$$b_i(S) [\theta_i] = \left\{ m_i \in M_i \left| \begin{array}{l} \exists \lambda_i \text{ s.th.:} \\ \text{(1)} \quad \lambda_i(m_{-i}, \theta_{-i}) > 0 \Rightarrow m_j \in S_j(\theta_j), \forall j \neq i \\ \text{(2)} \quad \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})) \\ \geq \\ \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})), \forall m'_i \in M_i \end{array} \right. \right\}.$$

We observe that  $b$  is increasing by definition: i.e.,  $S \leq S' \Rightarrow b(S) \leq b(S')$ . By Tarski's fixed point theorem, there is a largest fixed point of  $b$ , which we label  $S^{\mathcal{M}}$ . Thus (i)  $b(S^{\mathcal{M}}) = S^{\mathcal{M}}$  and

<sup>6</sup>We label these "interim" equilibria rather than "Bayesian" equilibria in light of the fact that our type space does not have a common prior. Dekel, Fudenberg and Levine (2002) argue that learning justifications that support the equilibrium assumption will - under reasonable assumptions - also imply a common prior on types. We nonetheless maintain the equilibrium solution concept for comparison with the literature. But note that in any case we end up using the fact that with rich type spaces, equilibrium does not have any bite above iterated deletion.

(ii)  $b(S) = S \Rightarrow S \leq S^{\mathcal{M}}$ . We can also construct the fixed point  $S^{\mathcal{M}}$  by starting with  $\bar{S}$  - the largest element of the lattice - and iteratively applying the operator  $b$ . If the message sets are finite, we have

$$S_i^{\mathcal{M}}(\theta_i) \triangleq \bigcap_{n \geq 1} b_i(b^n(\bar{S}))[\theta_i].$$

But because the mechanism  $\mathcal{M}$  may be infinite, transfinite induction may be necessary to reach the fixed point.<sup>7</sup> Thus  $S_i^{\mathcal{M}}(\theta_i)$  are the set of messages surviving (transfinite) iterated deletion of never best responses. If message sets are finite (or compact), a well known duality argument implies that never best responses are equivalent to strictly dominated actions. However, the equivalence does not hold with infinite (non-compact) message sets.<sup>8</sup> In a compact message analysis, Chung and Ely (2001) consider a version of this solution concept in an incomplete information mechanism design context with dominated (not strictly dominated) messages deleted at each round. We observe that the solution concept defined through the iterative application of the operator  $b$  is weaker than the notion of interim rationalizability for a given type space  $\mathcal{T}$ .<sup>9</sup> Under  $b$ , every agent  $i$  is allowed to hold arbitrary beliefs about  $\Theta_{-i}$  and is not restricted to a particular posterior distribution over  $\Theta_{-i}$ . On the other hand, if the type space  $\mathcal{T}$  were the universal type space, then  $S_i^{\mathcal{M}}(\theta_i)$  would be equal to the union of all interim rationalizable actions of agent  $i$  over all types  $t_i \in T_i$  whose payoff type profile coincides with  $\theta_i$ , or  $\hat{\theta}_i(t_i) = \theta_i$ .

For brevity and for lack of a better expression, we refer to the messages  $m_i \in S_i^{\mathcal{M}}(\theta_i)$  as *rationalizable actions*.

## 2.5 Implementation

We now define the notions of interim, robust and iterative implementation.

### Definition 2 (Interim Implementation)

*Social choice function  $f$  is interim implemented on type space  $\mathcal{T}$  by mechanism  $\mathcal{M}$  if the game*

---

<sup>7</sup>Lipman (1994) contains a formal description of the transfinite induction required (for the case of complete information, but nothing important changes with incomplete information). As he notes "we remove strategies which are never a best reply, taking limits where needed".

<sup>8</sup>The following simple example (suggested to us by Andrew Postlewaite) illustrates the non-equivalence. Players 1 and 2 each choose a non-negative integer,  $k_1$  and  $k_2$  respectively. The payoff to player 1 from  $k_1 = 0$  is 1. The payoff to player 1 from action  $k_1 \geq 1$  is 2 if  $k_1 > k_2$ , 0 otherwise. For any belief that player 1 has about 2's actions, there is a (sufficiently high) action from player 1 that gives him a payoff greater than 1. Thus action 0 is never a best response for player 1. However, for any mixed strategy of player 1, there is a (sufficient high) action of player 2 such that action 0 is a better response for player 1 than the mixed strategy. Thus action 0 is not strictly dominated.

<sup>9</sup>For the notion of interim rationalizability, see Battigalli and Siniscalchi (2003) and Dekel, Fudenberg, and Morris (2005).

$(\mathcal{T}, \mathcal{M})$  has an equilibrium and every equilibrium  $\sigma$  of the game  $(\mathcal{T}, \mathcal{M})$  satisfies

$$\sigma(m|t) > 0 \Rightarrow g(m) = f(\widehat{\theta}(t)).$$

We note that a tradition in the implementation literature commonly restricts attention to pure strategy equilibria, but we allow mixed strategy equilibria.<sup>10</sup>

**Definition 3 (Robust Implementation)**

*Social choice function  $f$  is robustly implemented by mechanism  $\mathcal{M}$  if, for every  $\mathcal{T}$ ,  $f$  is interim implemented on type space  $\mathcal{T}$  by mechanism  $\mathcal{M}$ .*

We observe that the notion of robust implementation requires that we can find a mechanism  $\mathcal{M}$  which implements  $f$  for every type space  $\mathcal{T}$ . A weaker requirement would be to ask that for every type space  $\mathcal{T}$  there exists a, possibly different, mechanism  $\mathcal{M}$  such that  $f$  is implemented. This weaker notion would still lead to the same necessary condition as the stronger implementation version we pursue here, and we believe that it would not lead to a substantial change in the sufficiency conditions either.

We shall establish the necessary and sufficient conditions for robust implementation via the iterative deletion process outlined above.

**Definition 4 (Iterative Implementation)**

*Social choice function  $f$  is iteratively implemented by mechanism  $\mathcal{M}$  if, for all  $\theta$ ,  $S^{\mathcal{M}}(\theta) \neq \emptyset$  and if for all  $\theta$  and  $m$ ,  $m \in S^{\mathcal{M}}(\theta) \Rightarrow g(m) = f(\theta)$ .*

### 3 Examples

We precede the formal results with four examples which are meant to illustrate the main insights of the paper. At the same time, they will facilitate a brief review of the key results in the implementation literature.

The first example is a model of majority rule introduced by Palfrey and Srivastava (1989). It highlights the difficulty of Bayesian implementation in a world of interdependent values. Ex post implementation and virtual implementation on some type spaces are possible, but interim implementation is impossible for some type spaces and thus robust implementation is impossible.

The second example builds around a simple coordination game. It shows that robust implementation, even though it is strong requirement, is substantially weaker than dominant strategy implementation. It also highlights the role of augmented mechanisms in achieving implementation.

---

<sup>10</sup>We discuss the pure strategy / mixed strategy modelling choice in Section 9.1.

The third example involves the provision of a public good with quasilinear utility. It demonstrates that robust implementation can frequently be achieved in the direct rather than the augmented mechanism. In addition, the example illustrates the relationship between robust implementation and rationalizability. In fact, we obtain a tight condition, in terms of the interdependence for robust implementation. Conversely, if the conditions fails, we show that even robust virtual implementation is impossible.

The fourth and final example investigates a single unit auction with symmetric bidders. The generalized Vickrey-Groves-Clark mechanism only satisfies weak rather than strict incentive compatibility constraints. We therefore propose an  $\varepsilon$ -efficient allocation rule with strict ex post incentive constraints. The  $\varepsilon$  efficient allocation rule can also be interpreted as virtual implementation of the efficient rule. In either case, we show that the rule itself can be either robustly implemented or robust virtually implemented, respectively, if there is not too much interdependence among the payoff types.

### 3.1 Majority Rule

In the first example, introduced by Palfrey and Srivastava (1989), there are three agents and each has two possible "payoff types",  $\theta_a$  or  $\theta_b$ . There are two possible choices for society,  $a$  or  $b$ . All agents have identical preferences. If a majority of agents (i.e., at least two) are of type  $\theta_y$ , then every agent gets utility 1 from outcome  $y$  and utility 0 from the other outcome. The social choice function agrees with the common preferences of the agents. Thus  $f : \{\theta_a, \theta_b\}^3 \rightarrow \{a, b\}$  satisfies  $f(\theta) = y$  if and only if  $\#\{i : \theta_i = \theta_y\} \geq 2$ . It is useful to use this simple example to review the existing implementation literature and understand the role of interdependent types.

Clearly, incentive compatibility is not a problem in this example. The problem is that in the "direct mechanism" - where all agents simply announce their types - there is the possibility that all agents will choose to always announce  $\theta_a$ . Since no agent expects to be pivotal, he has no incentive to truthfully announce his type when he is in fact  $\theta_b$ . What happens if we allow more complicated mechanisms?

If there were complete information about agents' preferences, then the social choice function is clearly implementable: the social planner could pick an agent, say agent 1, and simply follow that agent's recommendation.

But suppose instead that there is incomplete information about agents' preferences. In particular, suppose it is common knowledge that each agent's type is  $\theta_b$  with independent probability  $q$ , with  $q^2 > \frac{1}{2}$ . This example fails the Bayesian monotonicity condition of Postlewaite and Schmeidler (1986) and Jackson (1991). Palfrey and Srivastava (1989) observe that it is also not possible

to implement in undominated Bayesian Nash equilibrium in this example.

Bergemann and Morris (2005) have analyzed the alternative "more robust" solution concept of ex post equilibrium in this context. It is easy to construct an augmented mechanism whose only ex post equilibrium delivers the social choice function. Let each agent send a message  $m_i \in \{\theta_a, \theta_b\} \times \{\text{truth, lie}\}$ , with the interpretation that an agent is announcing his own type and also sends the message "truth" if he thinks that others are telling the truth and sends the message "lie" if he thinks that someone is lying. Outcome  $y$  is implemented if a majority claim to be type  $\theta_y$  and all agents announce "truth"; or if either 1 or 3 agents claim to be type  $\theta_y$  and at least one agent reports lying.

There is a truthtelling ex post equilibrium where each agent truthfully announces his type and also announces "truth". Now suppose there exists an ex post equilibrium such that at some type profile, the desired outcome is not chosen. Note that whatever the announcements of the other agents, each agent always has the ability to determine make the outcome  $y$ , by sending the message "lie" and - given the announcements of the other agents - choosing his message so that an odd number of players have claimed to be type  $\theta_y$ . So this is not consistent with ex post equilibrium.

Serrano and Vohra (2005) show that virtual Bayesian implementation is possible in this example. The idea is that it is possible to exploit common knowledge of  $q$  to design a bet as a function of  $q$  that will give each player an incentive to truthfully announce his type. This bet may lead to a small but positive probability that the wrong outcome is realized in equilibrium. But virtual implementation is possible.

Robust implementation is impossible in this example. Consider the type space where there is common knowledge that whenever an agent is type  $\theta_y$ , he assigns probability  $\frac{1}{2}$  to both of the other agents being type  $y' \neq y$  and probability  $\frac{1}{2}$  to one being type  $y$  and the other being  $y'$ . Thus every type of every player thinks there is a 50% chance that outcome  $a$  is better and a 50% chance that  $b$  is better. Evidently, there is no way of designing a mechanism that ensures that agents do not fully pool. But if they fully pool, robust implementation is not possible.

### 3.2 Coordination

This example establishes that although robust implementation is a strong requirement, it is weaker than dominant strategies. There are two agents. Each agent  $i$  has two possible types,  $\theta_i$  and  $\theta'_i$ . There are six possible outcomes:  $Z = \{a, b, c, d, z, z'\}$ . The payoffs of the agents are a function of

the allocation and the true payoff type profile, given by:

<b>a</b>	$\theta_2$	$\theta'_2$
$\theta_1$	3, 3	0, 0
$\theta'_1$	0, 0	1, 1

<b>b</b>	$\theta_2$	$\theta'_2$
$\theta_1$	0, 0	3, 3
$\theta'_1$	1, 1	0, 0

<b>c</b>	$\theta_2$	$\theta'_2$
$\theta_1$	0, 0	1, 1
$\theta'_1$	3, 3	0, 0

<b>d</b>	$\theta_2$	$\theta'_2$
$\theta_1$	1, 1	0, 0
$\theta'_1$	0, 0	3, 3

and

<b>z</b>	$\theta_2$	$\theta'_2$
$\theta_1$	2, 2	2, 0
$\theta'_1$	2, 2	2, 0

<b>z'</b>	$\theta_2$	$\theta'_2$
$\theta_1$	2, 0	2, 2
$\theta'_1$	2, 0	2, 2

The social choice function is given by the efficient outcome:

<b>f</b>	$\theta_2$	$\theta'_2$
$\theta_1$	a	b
$\theta'_1$	c	d

Clearly, the social choice function is strictly ex post incentive compatible. But in the "direct mechanism" where each agent simply reports his type, there will always be an equilibrium where each type of each agent misreports his type, and each agent gets a payoff of 1. This is also strictly ex post incentive compatible. The social choice function  $f$  which selects among  $\{a, b, c, d\}$  embeds a coordination game. We further observe that the payoff for agent 1 from allocations  $z$  and  $z'$  are equal and constant for all type profiles. On the other hand, the payoff of agent 2 from  $z$  and  $z'$  depends on his type but not on the type of the other agent.

We now consider the following augmented simple mechanism:

<b>g</b>	$\theta_2$	$\theta'_2$
$\theta_1$	a	b
$\theta'_1$	c	d
$\zeta$	z	z'

The corresponding incomplete information game has the following payoffs:

	type	$\theta_2$		$\theta'_2$	
type	action	$\theta_2$	$\theta'_2$	$\theta_2$	$\theta'_2$
$\theta_1$	$\theta_1$	3, 3	0, 0	0, 0	3, 3
	$\theta'_1$	0, 0	1, 1	1, 1	0, 0
	$\zeta$	2, 2	2, 0	2, 0	2, 2
$\theta'_1$	$\theta_1$	0, 0	1, 1	1, 1	0, 0
	$\theta'_1$	3, 3	0, 0	0, 0	3, 3
	$\zeta$	2, 2	2, 0	2, 0	2, 2

Suppose we iteratively remove actions for each type that could never be a best response given the type action profiles remaining. Thus in the first round, we would observe that type  $\theta_1$  would never send message  $\theta'_1$  and type  $\theta'_1$  would never send message  $\theta_1$ . Knowing this, we could conclude that type  $\theta_2$  would never send message  $\theta'_2$  and type  $\theta'_2$  would never send message  $\theta_2$ . This in turn implies that neither type of agent 1 will ever send message  $\zeta$ . Thus the only remaining message for each type of each agent is truth-telling. But now they must behave this way in any equilibrium on any type space.

### 3.3 Public Good

The third example describes the provision of a public good with quasilinear utility. The utility of each agent is given by:

$$u_i(\theta, x, y) = \left( \theta_i + \gamma \sum_{j \neq i} \theta_j \right) x + y_i,$$

where  $x$  is the level of public good provided and  $y_i$  is the monetary transfer to agent  $i$ . The utility of agent  $i$  depends on his own type  $\theta_i \in [0, 1]$  and the type profile of other agents, with  $\gamma \geq 0$ . The cost of establishing the public good is given by  $c(x) = \frac{1}{2}x^2$ . The planner must choose  $(x, y_1, \dots, y_I) \in \mathbb{R}_+ \times \mathbb{R}^I$  to maximize social welfare, i.e., the sum of gross utilities minus the cost of the public good:

$$\left( (1 + \gamma(I - 1)) \sum_{i=1}^I \theta_i \right) x - \frac{1}{2}x^2.$$

The socially optimal level of the public good is therefore equal to

$$f_0(\theta) = (1 + \gamma(I - 1)) \sum_{i=1}^I \theta_i.$$

We choose essentially unique (up to a constant) transfers that give rise to ex post incentive compatibility:

$$f_i(\theta) = -(1 + \gamma(I - 1)) \left( \gamma \left( \sum_{j \neq i} \theta_j \right) \theta_i + \frac{1}{2} \theta_i^2 \right).$$

We first argue that if  $\gamma < \frac{1}{I-1}$ , the social choice function  $f$  is robustly implementable in the *direct mechanism* where each agent reports his payoff type  $\theta_i$  and the planner chooses outcomes according to  $f$  on the assumption that agents are telling the truth. Consider an iterative deletion procedure. Let  $S^0(\theta_i) = [0, 1]$  and, for each  $k = 1, 2, \dots$ , let  $S^k(\theta_i)$  be the set of reports that agent  $i$  might send, for some conjecture over his opponents' types and reports, with the only restriction on his conjecture being that each type  $\theta_j$  of agent  $j$  sends a message in  $S^{k-1}(\theta_j)$ .

Suppose that agent  $i$  has payoff type  $\theta_i$ , has a point conjecture that other agents have type profile  $\theta_{-i}$  and report their types to be  $\theta'_{-i}$ , and he reports himself to be type  $\theta'_i$ . Then his expected payoff is a constant  $(1 + \gamma(I - 1))$  times

$$\left( \theta_i + \gamma \sum_{j \neq i} \theta_j \right) \left( \theta'_i + \sum_{j \neq i} \theta'_j \right) - \left( \gamma \left( \sum_{j \neq i} \theta'_j \right) \theta'_i + \frac{1}{2} (\theta'_i)^2 \right).$$

The first order condition with respect to  $\theta'_i$  is then

$$\theta_i + \gamma \sum_{j \neq i} \theta_j - \gamma \left( \sum_{j \neq i} \theta'_j \right) - \theta'_i = 0,$$

so he would wish to set

$$\theta'_i = \theta_i + \gamma \sum_{j \neq i} (\theta_j - \theta'_j).$$

Note that this calculation verifies the strict ex post incentive compatibility of  $f$ . The quadratic payoff / linear best response nature of this problem means that we can characterize  $S^k(\theta_i)$  restricting attention to such point conjectures. In particular, we will have

$$S^k(\theta_i) = \left[ \underline{\beta}^k(\theta_i), \bar{\beta}^k(\theta_i) \right],$$

where

$$\begin{aligned} \bar{\beta}^k(\theta_i) &= \min \left\{ 1, \theta_i + \gamma \max_{\{(\theta'_{-i}, \theta_{-i}) : \theta'_j \in S^k(\theta_j) \text{ for all } j \neq i\}} \sum_{j \neq i} (\theta_j - \theta'_j) \right\} \\ &= \min \left\{ 1, \theta_i + \gamma \max_{\theta_{-i}} \sum_{j \neq i} (\theta_j - \underline{\beta}^{k-1}(\theta_j)) \right\}. \end{aligned}$$

Analogously,

$$\underline{\beta}^k(\theta_i) = \max \left\{ 0, \theta_i - \gamma \max_{\theta_{-i}} \sum_{j \neq i} (\bar{\beta}^{k-1}(\theta_j) - \theta_j) \right\}.$$

Thus

$$\bar{\beta}^k(\theta_i) = \min \left\{ 1, \theta_i + (\gamma(I - 1))^k \right\},$$

and

$$\underline{\beta}^k(\theta_i) = \max \left\{ 0, \theta_i - (\gamma(I - 1))^k \right\}.$$

Thus  $\theta'_i \neq \theta_i \Rightarrow \theta'_i \notin S^k(\theta_i)$  for sufficiently large  $k$ , provided that  $\gamma < \frac{1}{I-1}$ .



On the other hand if  $\gamma > \frac{1}{I-1}$ , then we argue that there exist type spaces where the social choice function  $f$  is not virtually implementable. Consider a type space where it is common knowledge whenever agent  $i$  has type  $\theta_i$ , he is convinced that the types of other players  $\theta_{-i}$  are such that

$$\sum_{j \neq i} \theta_j = \frac{1}{\gamma} \left( \frac{1}{2} - \theta_i \right).$$

Observe that  $\gamma > \frac{1}{I-1}$  implies that we can choose the  $\theta_j$  to be in the interval  $[0, 1]$ . Agent  $i$ 's preferences are independent of his type on this type space. Now fix any mechanism and restrict each player to a pooling strategy, i.e., sending the same message independent of his type. Since all types now have identical preferences over outcomes, this pooling strategy is an equilibrium.<sup>11</sup>

### 3.4 Private Good

The final example is a single unit auction with symmetric bidders. There are  $I$  agents and agent  $i$ 's payoff type is  $\theta_i \in [0, 1]$ . If the profile type profile is  $\theta$ , agent  $i$ 's valuation of the object is

$$\theta_i + \gamma \sum_{j \neq i} \theta_j,$$

where  $0 \leq \gamma \leq 1$ .

In Bergemann and Morris (2005), we showed that Maskin monotonicity fails in this example, so complete information implementation is impossible. This in turn implies that robust implementation is impossible. But Bergemann and Morris (2005) also showed that ex post implementation occurs in the direct mechanism if there are at least three agents.

We next argue that virtual robust implementation is possible in this example if  $\gamma < \frac{1}{I-1}$ . An allocation rule is a function  $x : \Theta \rightarrow [0, 1]^I$ , where  $x_i(\theta)$  is the probability that agent  $i$  gets the object and so  $\sum_i x_i(\theta) \leq 1$ . The symmetric efficient allocation rule is the following:

$$x_i^*(\theta) = \begin{cases} \frac{1}{\#\{j: \theta_j \geq \theta_k \text{ for all } k\}}, & \text{if } \theta_i \geq \theta_k \text{ for all } k, \\ 0, & \text{if otherwise.} \end{cases}$$

A symmetric  $\varepsilon$ -efficient allocation rule is the following:

$$x_i^{**}(\theta) = \varepsilon \frac{\theta_i}{I} + (1 - \varepsilon) x_i^*(\theta).$$

In the Appendix, we verify that the resulting generalized VCG transfers satisfy strict ex post incentive compatibility and show that this  $\varepsilon$ -efficient allocation is robustly implementable. Under

---

<sup>11</sup>This argument is complete only if we know that an equilibrium exists in the game where players are restricted to pooling strategies. This will be true for well-behaved mechanisms. Our general negative result holds even if this is not guaranteed.

this allocation rule, the object is not allocated with probability  $\frac{\varepsilon}{2}$ . At the cost of some additional algebra, we could replace this rule with

$$x_i^{**}(\theta) = \varepsilon \frac{\theta_i}{\sum_j \theta_j} + (1 - \varepsilon) x_i^*(\theta)$$

which allocates the object with probability 1.

However, if  $\frac{1}{T-1} < \gamma \leq 1$ , only constant allocations are robust implementable, by the same argument as in the public good case: we can construct beliefs for each type such that types are indistinguishable.

## 4 Interim Equilibrium and Iterative Elimination

The notion of robust implementation requires that a social choice function  $f$  can be interim implemented for all type spaces  $\mathcal{T}$ . As we look for necessary and sufficient conditions for robust implementation, conceptually there are (at least) two approaches to obtain the conditions.

One approach would be to simply look at the interim implementation conditions for every possible type space  $\mathcal{T}$  and then try to characterize the intersection or union of these conditions for all type spaces. This is the approach we initially pursued, and it works in brute force kind of way. In Section 9.1, we review what happens under this approach.

But we focus our analysis on a second, more elegant, approach. We first establish an equivalence between robust and iterative implementation and then derive the necessary conditions for robust implementation as an implication of iterative implementation. The advantage of the second approach is that after establishing the equivalence, we do not need to argue in terms of large type spaces, but rather derive the results from a novel argument using the iterative elimination process.

A complicating element in the implementation context is the fact that the augmented mechanisms often have infinite message spaces and that best responses may not exist. These complications are inherent to the entire implementation literature and we therefore have to carefully address these issues before we establish the implementation results.

### 4.1 Best Response

We start with the fixed point  $S^M$  of the iterative elimination procedure. Recall that by definition,  $S^M$  is a fixed point of  $b$ , and thus for all  $m_i \in S_i^M(\theta_i)$ , there exists  $\lambda_i \in \Delta(M_{-i} \times \Theta_{-i})$  such that

$$\begin{aligned} (1) \quad & \lambda_i(m_{-i}, \theta_{-i}) > 0 \Rightarrow m_j \in S_j^M(\theta_j) \text{ for each } j \neq i \\ (2) \quad & \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) \begin{bmatrix} u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})) \\ -u_i(g(m'_i, m_{-i}), (\theta_i, \theta_{-i})) \end{bmatrix} \geq 0 \text{ for all } m'_i \in M_i \end{aligned} \quad (1)$$

**Proposition 1 (Rationalizable Actions)**

$m_i \in S^{\mathcal{M}}(\theta_i)$  if and only if there exists a type space  $\mathcal{T}$ , an interim equilibrium  $\sigma$  of the game  $(\mathcal{T}, \mathcal{M})$  and a type  $t_i \in T_i$  such that (i)  $\sigma_i(m_i|t_i) > 0$  and (ii)  $\widehat{\theta}_i(t_i) = \theta_i$ .

**Proof.** ( $\Rightarrow$ ) Suppose  $m_i^* \in S^{\mathcal{M}}(\theta_i^*)$ . Now consider the following type space  $\mathcal{T}$  defined through:

$$T_i = \{(m_i, \theta_i) : m_i \in S_i^{\mathcal{M}}(\theta_i)\}.$$

Let

$$\widehat{\theta}_i(m_i, \theta_i) \triangleq \theta_i.$$

By (1), we know that for each  $m_i \in S_i^{\mathcal{M}}(\theta_i)$ , there exists  $\lambda_i^{m_i, \theta_i} \in \Delta(M_{-i} \times \Theta_{-i})$  such that:

$$\lambda_i^{m_i, \theta_i}(m_{-i}, \theta_{-i}) > 0 \Rightarrow m_j \in S_j^{\mathcal{M}}(\theta_j) \text{ for each } j \neq i;$$

and

$$\sum_{m_{-i}, \theta_{-i}} \lambda_i^{m_i, \theta_i}(m_{-i}, \theta_{-i}) [u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})) - u_i(g(m'_i, m_{-i}), (\theta_i, \theta_{-i}))] \geq 0, \quad \forall m'_i \in M_i.$$

Let

$$\widehat{\pi}_i(m_{-i}, \theta_{-i}) [m_i, \theta_i] \triangleq \lambda_i^{m_i, \theta_i}(m_{-i}, \theta_{-i}).$$

Now by construction, there is a pure strategy equilibrium  $s$  with  $s_i(m_i, \theta_i) = m_i$ . But now  $s_i(m_i^*, \theta_i^*) = m_i^*$  and  $\widehat{\theta}_i(m_i^*, \theta_i^*) = \theta_i^*$ .

( $\Leftarrow$ ) Suppose there exists a type space  $\mathcal{T}$ , an equilibrium  $\sigma$  of  $(\mathcal{T}, \mathcal{M})$ , and  $m_i^* \in M_i$  and  $t_i^* \in T_i$  such that  $\sigma_i(m_i^*|t_i^*) > 0$  and  $\widehat{\theta}_i(t_i^*) = \theta_i^*$ . Let

$$S_i(\theta_i) = \left\{ m_i : \sigma_i(m_i|t_i) > 0 \text{ and } \widehat{\theta}_i(t_i) = \theta_i \text{ for some } t_i \in T_i \right\}.$$

Now interim equilibrium conditions ensure that  $b(S) \geq S$ . Thus  $S \leq S^{\mathcal{M}}$ . Thus  $m_i^* \in S_i^{\mathcal{M}}(\widehat{\theta}_i(t_i^*))$ , which concludes the proof. ■

Brandenburger and Dekel (1987) showed an equivalence for finite action complete information games between the set of actions surviving iterated deletion of strictly dominant actions and the set of actions that could be played in a subjective correlated equilibrium. Proposition 1 is a straightforward generalization of Brandenburger and Dekel (1987) to incomplete information and infinite actions. The infinite action extension (for complete information) was shown in Lipman (1994). The finite action incomplete information extension is reported in a recent paper of Battigalli and Siniscalchi (2003) (following an earlier analysis in Battigalli (1999)).

Notice that there is no guarantee that  $S^{\mathcal{M}}(\theta_i)$  is non-empty or that the game  $(\mathcal{T}, \mathcal{M})$  has an equilibrium: the Proposition holds vacuously in this case. But for implementation results, we care

about existence. We have the following two conditions that relate existence of equilibrium on all type spaces to the actions surviving iterated deletion.

**Definition 5 (Ex Post Best Response)**

Message correspondence  $S$  satisfies the ex post best response property if, for all  $i$  and  $\theta_i \in \Theta_i$ , there exists  $m_i^* \in S_i(\theta_i)$  such that

$$m_i^* \in \arg \max_{m_i \in M_i} u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})),$$

for all  $\theta_{-i}$  and  $m_{-i} \in S(\theta_{-i})$ .

We observe that for  $S$  to satisfy the ex post best response property,  $S_i$  must be non-empty for all  $i$  and all  $\theta_i$ .

**Definition 6 (Interim Best Response)**

Message correspondence  $S$  satisfies the interim best response property if, for all  $i$  and  $\psi_i \in \Delta(\Theta_{-i})$ , there exists  $\lambda_i \in \Delta(M_{-i} \times \Theta_{-i})$  such that:

1.  $\lambda_i(m_{-i}, \theta_{-i}) > 0 \Rightarrow m_j \in S_j(\theta_j)$  for each  $j \neq i$ ;
2. for all  $\theta_{-i} \in \Theta_{-i}$ :

$$\sum_{m_{-i} \in M_{-i}} \lambda_i(m_{-i}, \theta_{-i}) = \psi_i(\theta_{-i});$$

3. for all  $\theta_i \in \Theta_i$  there exists  $m_i^* \in S_i(\theta_i)$  such that

$$m_i^* \in \arg \max_{m_i \in M_i} \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})).$$

The interim best response property only requires that for every conjecture over payoff type spaces, there exists some beliefs over messages consistent with the message correspondence  $S$ , such that a best response is in the message correspondence. In particular, it does not require that a best response exists for all possible beliefs over message profiles. Note that the ex post best response property is a stronger requirement than the interim best response property. Also note that the interim best response property implies that  $S_i^M(\theta_i)$  is non-empty for all  $i$  and  $\theta_i$ .

Proposition 1 links every action profile in the set of rationalizable actions to an equilibrium action for some type space  $\mathcal{T}$ . Proposition 2 strengthens the relationship between rationalizable and equilibrium actions, after imposing some structure on the best response property of rationalizable and equilibrium actions, respectively.

**Proposition 2 (Best Response Properties)**

1. If  $S^{\mathcal{M}}$  satisfies the ex post best response property, then  $(\mathcal{T}, \mathcal{M})$  has an equilibrium for each  $\mathcal{T}$ .
2. If  $(\mathcal{T}, \mathcal{M})$  has an equilibrium for each  $\mathcal{T}$ , then  $S^{\mathcal{M}}$  satisfies the interim best response property.

**Proof. (1.)** By the ex post best response property, there exists, for each  $i$ ,  $s_i^* : \Theta_i \rightarrow M_i$  such that

$$s_i^*(\theta_i) \in \arg \max_{m_i \in M_i} u_i(g(m_i, s_{-i}^*(\theta_{-i})), (\theta_i, \theta_{-i}))$$

for all  $\theta_{-i}$ . Now fix any type space. The strategy profile  $s$  with

$$s_i(t_i) = s_i^*(\widehat{\theta}_i(t_i))$$

is an equilibrium of the game  $(\mathcal{T}, \mathcal{M})$ .

**(2.)** Suppose  $(\mathcal{T}, \mathcal{M})$  has an equilibrium for each  $\mathcal{T}$ . Fix any  $i$  and  $\psi_i \in \Delta(\Theta_{-i})$ . Fix any type space  $\mathcal{T}$  with, for each  $\theta_i \in \Theta_i$ , a type  $t_i^*(\theta_i)$  such that (a)  $\widehat{\theta}_i(t_i^*(\theta_i)) = \theta_i$  for each  $\theta_i$ , (b)  $\widehat{\pi}_i(t_i^*(\theta_i)) = \pi_i$  for all  $\theta_i$  and (c)

$$\sum_{\{t_{-i} : \widehat{\theta}_{-i}(t_{-i}) = \theta_{-i}\}} \pi_i(t_{-i}) [t_i^*] = \psi_i(\theta_{-i}) \quad (2)$$

for all  $\theta_i$  and  $\theta_{-i}$ . The game has an equilibrium  $\sigma$ . Let  $m_i$  be any message with  $\sigma_i(m_i | t_i^*(\theta_i)) > 0$ . Let

$$\lambda_i(m_{-i}, \theta_{-i}) = \sum_{\{t_{-i} \in \mathcal{T}_{-i} : \widehat{\theta}_{-i}(t_{-i}) = \theta_{-i}\}} \sigma_{-i}(m_{-i} | t_{-i}) \pi_i(t_{-i}) [t_i^*].$$

Now  $\sigma_i(m_i | t_i^*(\theta_i)) > 0$  implies

$$m_i(\theta_i) \in \arg \max_{m_i \in M_i} \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})).$$

Proposition 1 implies that every message profile  $m_j$  which is played in equilibrium by type  $\theta_j$  is part of the set  $S^{\mathcal{M}}$ , or that:

$$\lambda_i(m_{-i}, \theta_{-i}) > 0 \Rightarrow m_j \in S_j^{\mathcal{M}}(\theta_j) \text{ for each } j \neq i.$$

By construction of the type space  $\mathcal{T}$ , in particular property (c) as expressed by (2), this implies that

$$\sum_{m_{-i} \in M_{-i}} \lambda_i(m_{-i}, \theta_{-i}) = \psi_i(\theta_{-i}) \text{ for all } \theta_{-i} \in \Theta_{-i}.$$

Since these properties hold for arbitrary  $i$  and  $\psi_i \in \Delta(\Theta_{-i})$ ,  $S^{\mathcal{M}}$  satisfies the interim best response property, which concludes the proof. ■

It is unfortunate that there is a gap between the necessary and sufficient conditions in this Proposition. However, an example in the Appendix shows that it is possible to construct (admittedly silly) mechanisms where  $(\mathcal{T}, \mathcal{M})$  has an equilibrium for each  $\mathcal{T}$ , but  $S^{\mathcal{M}}$  fails the ex post best response property. The ex post best response property must be satisfied if the mechanism is nice, i.e., best responses always exist.

## 4.2 Implementation

The first part of the definition of implementation is the requirement that all outcomes that occur in equilibrium are consistent with the social choice function. The first definition checks if this is true for some fixed type space  $\mathcal{T}$  and mechanism  $\mathcal{M}$ .

### Definition 7 (Interim Material Implementation)

*Social choice function  $f$  is interim materially implemented on type space  $\mathcal{T}$  by mechanism  $\mathcal{M}$  if every equilibrium  $\sigma$  of the game  $(\mathcal{T}, \mathcal{M})$  satisfies*

$$\sigma(m|t) > 0 \Rightarrow g(m) = f(\hat{\theta}(t)),$$

*for all  $t$ .*

The next definition checks if the same property holds for every type space:

### Definition 8 (Robust Material Implementation)

*Social choice function  $f$  is robustly materially implemented by mechanism  $\mathcal{M}$  if, for every  $\mathcal{T}$ ,  $f$  is interim materially implemented on type space  $\mathcal{T}$  by mechanism  $\mathcal{M}$ .*

Finally, we ask if iterated deletion delivers acceptable outcomes:

### Definition 9 (Iterative Material Implementation)

*Social choice function  $f$  is iteratively materially implemented by mechanism  $\mathcal{M}$  if, for all  $\theta$ ,  $m \in S^{\mathcal{M}}(\theta) \Rightarrow g(m) = f(\theta)$ .*

Now Proposition 1 immediately implies an equivalence between robust and iterative implementation.

### Corollary 1 (Equivalence)

*Social choice function  $f$  is iteratively materially implemented by  $\mathcal{M}$  if and only if  $f$  is robustly materially implemented by mechanism  $\mathcal{M}$ .*

In all the above definitions, we qualified that implementation as being "material" because the premise of the definition might be vacuous: the mechanism  $\mathcal{M}$  might have the property that on any type space, there is no equilibrium and no messages surviving iterated deletion.<sup>12</sup>

Proposition 2 gave the slightly messier result relating equilibrium existence and properties of messages surviving iterated deletion. The following corollary gives the immediate implications for our implementation definitions:

**Corollary 2 (Necessary Conditions)**

1. *If social choice function  $f$  is iteratively materially implemented by mechanism  $\mathcal{M}$  and  $S^{\mathcal{M}}$  satisfies the ex post best response property, then  $f$  is robustly implemented by  $\mathcal{M}$ .*
2. *If  $f$  is robustly implemented by  $\mathcal{M}$ , then  $f$  is iteratively materially implemented by mechanism  $\mathcal{M}$  and  $S^{\mathcal{M}}$  satisfies the interim best response property.*

The ‘material’ qualification will only be used in the necessity part of Theorem 2 where we shall invoke the above Corollary 2.2. There we shall use the fixed-point property of  $S^{\mathcal{M}}$ , stated earlier in (1), to derive the robust monotonicity condition. In the sufficiency part of the proof, a non-empty set  $S^{\mathcal{M}}$  is obtained in the augmented mechanism by virtue of ex post incentive compatibility. Similarly, for the direct implementation results, a non-empty set  $S^{\mathcal{M}}$  exists by ex post incentive compatibility and the ‘material’ qualification is not needed at all. The following implication of robust implementability will be used to establish robust monotonicity in Theorem 2.

**Lemma 1 (Truth-telling as Best Response)**

*If  $f$  is iteratively materially implemented by mechanism  $\mathcal{M}$  and  $S^{\mathcal{M}}$  satisfies the interim best response property, then for all  $i$  and  $\theta_{-i} \in \Theta_{-i}$ , there exists  $\nu_i \in \Delta(S_{-i}^{\mathcal{M}}(\theta_{-i}))$ ,*

$$u_i(f(\theta_i, \theta_{-i}), (\theta_i, \theta_{-i})) \geq \sum_{m_{-i}} \nu_i(m_{-i}) u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})) \quad (3)$$

for all  $m_i \in M_i$  and  $\theta_i \in \Theta_i$ .

**Proof.** Applying the definition of the interim best response property for  $i$  and the degenerate distribution putting probability 1 on  $\theta_{-i}$ , we have that there exists  $\nu_i \in \Delta(S_{-i}^{\mathcal{M}}(\theta_{-i}))$  such that

$$\emptyset \neq \arg \max_{m_i} \sum_{m_{-i}, \theta_{-i}} \nu_i(m_{-i}) u_i((m_i, m_{-i}), (\theta_i, \theta_{-i})) \subseteq S_i^{\mathcal{M}}(\theta_i) \text{ for all } \theta_i \in \Theta_i.$$

---

<sup>12</sup>Our terminology mirrors the language of modal logic where proposition  $A$  materially implies  $B$  whenever  $A$  is false, as well as when both  $A$  and  $B$  are true, see Hughes and Creswell (1996).

But by iterative material implementability,  $m \in S^{\mathcal{M}}(\theta) \Rightarrow g(m) = f(\theta)$ . So

$$u_i(f(\theta_i, \theta_{-i}), (\theta_i, \theta_{-i})) \geq \sum_{m_{-i}} \nu_i(m_{-i}) u_i((m_i, m_{-i}), (\theta_i, \theta_{-i})),$$

for all  $m_i \in M_i$  and  $\theta_i \in \Theta_i$ . ■

Lemma 1 shows how small the gap between the ex post and interim best response property is. It establishes that truth-telling is a best response against some beliefs over messages  $m_{-i}$  for any given payoff type profile  $\theta_{-i}$ .

## 5 Direct Implementation

We first characterize when robust direct implementation is possible. By "direct mechanism", we mean the mechanism where each agent simply reports his payoff type  $\theta_i$ , and so  $M_i = \Theta_i$  for all  $i$  and  $g(\theta) = f(\theta)$  for all  $\theta$ . Thus we assume that agents do not report their "higher order belief" types. "Truth-telling behavior" is the strategy of always truthfully reporting your payoff type.

### Definition 10 (Robust Direct Implementation)

*SCF  $f$  is robustly directly implementable if truth-telling is an equilibrium of the direct mechanism, and  $f$  is interim implementable on any type space  $\mathcal{T}$  by the direct mechanism.*

A *deception* is a profile  $\beta = (\beta_1, \dots, \beta_I)$ , where:

$$\beta_i : \Theta_i \rightarrow 2^{\Theta_i} / \emptyset,$$

with  $\theta_i \in \beta_i(\theta_i)$  for all  $i, \theta_i$ .

### Definition 11 (Acceptable / Unacceptable Deception)

*A deception is acceptable if  $\theta' \in \beta(\theta) \Rightarrow f(\theta') = f(\theta)$ . A deception is unacceptable if it is not acceptable.*

The inverse mapping of the deception  $\beta_i$  represents the set of true type profiles  $\theta_i$  which could lead to a deception  $\theta'_i$  and we write

$$\beta_i^{-1}(\theta'_i) = \{\theta_i : \theta'_i \in \beta_i(\theta_i)\}.$$

A deception is a message correspondence profile in the special case of a direct mechanism.

### Definition 12 (Ex Post Incentive Compatibility)

*Social choice function  $f$  satisfies ex post incentive compatibility if*

$$u_i(f(\theta_i, \theta_{-i}), (\theta_i, \theta_{-i})) \geq u_i(f(\theta'_i, \theta_{-i}), (\theta_i, \theta_{-i})),$$

for all  $i, \theta_i, \theta'_i$  and  $\theta_{-i}$ .



**Definition 13 (Direct Robust Monotonicity)**

Social choice function  $f$  satisfies direct robust monotonicity if, for every unacceptable deception  $\beta$ , there exist  $i, \theta_i, \theta'_i \in \beta_i(\theta_i)$  such that for all  $\theta'_{-i} \in \Theta_{-i}$  and  $\rho_i \in \Delta(\theta'_{-i}, \beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $\theta''_i \in \Theta_i$  such that

$$\sum_{\theta'_{-i} \in \Theta_{-i}} \rho_i(\theta'_{-i}, \theta_{-i}) u_i(f(\theta''_i, \theta'_{-i}), (\theta_i, \theta_{-i})) > \sum_{\theta'_{-i} \in \Theta_{-i}} \rho_i(\theta'_{-i}, \theta_{-i}) u_i(f(\theta'_i, \theta'_{-i}), (\theta_i, \theta_{-i})). \quad (4)$$

We will show direct robust monotonicity is a necessary and sufficient condition for implementation in the direct mechanism. In consequence, the designer can only offer those rewards to agent  $i$  which can be generated through the social choice function  $f$  by a report of agent  $i$ , say  $\theta''_i$ . The focus on rewards which can be generated through type reports in the direct mechanism implies that the incentive compatibility condition for claiming the reward  $y$  in the correct circumstances (if and only if the other agents deceive) is automatically satisfied, provided that the social choice function  $f$  satisfies ex post incentive compatibility.

**Lemma 2 (Robust Implementation as Fixed Point)**

$f$  satisfies direct robust monotonicity if and only if  $\beta \leq b(\beta) \Rightarrow \beta$  acceptable.

**Proof.** ( $\Rightarrow$ ) The proof is by contradiction. Thus suppose that  $f$  satisfies direct robust monotonicity and that the deception  $\beta$  is a fixed point under the mapping  $b$ , but  $\beta$  is unacceptable. In the direct mechanism, the set of messages is  $M_i = \Theta_i$  for all  $i$ . Then by hypothesis of direct robust monotonicity, there exists  $i, \theta_i, \theta'_i \in \beta_i(\theta_i)$  such that for all  $\theta'_{-i} \in \Theta_{-i}$  and  $\rho_i \in \Delta(\theta'_{-i}, \beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $\theta''_i \in \Theta_i$  satisfying the strict inequality (4). But this implies that  $\theta'_i \notin b_i(\beta)[\theta_i]$  which contradicts the fixed point property of  $\beta$ .

( $\Leftarrow$ ) The proof is again by contradiction. Thus suppose that the fixed point property  $\beta \leq b(\beta)$  indeed implies that  $\beta$  is acceptable, but that  $f$  does not satisfy direct robust monotonicity. In other words, let us suppose that there exists an unacceptable deception  $\beta$  for which we cannot find  $i, \theta_i$  and  $\theta'_i \in \beta_i(\theta_i)$  such that the inequality (4) can be satisfied. By hypothesis,  $\beta$  is unacceptable, and it follows that the premise of the hypothesis, namely the fixed point property cannot be satisfied by  $\beta$ . But this implies that there exists  $i$  and  $\theta'_i$  with  $\theta'_i \notin b_i(\beta)[\theta_i]$ . But the exclusion means that for every  $\rho_i \in \Delta(\Theta'_{-i} \times \Theta_{-i})$  such that  $\rho_i(\theta'_{-i}, \theta_{-i}) > 0 \Rightarrow \theta'_j \in \beta_j(\theta_j)$  for each  $j \neq i$ , there exists  $\theta''_i$  such that

$$\sum_{\theta'_{-i}, \theta_{-i}} \rho_i(\theta'_{-i}, \theta_{-i}) [u_i(f(\theta''_i, \theta'_{-i}), (\theta_i, \theta_{-i})) - u_i(f(\theta'_i, \theta'_{-i}), (\theta_i, \theta_{-i}))] > 0,$$

but this delivers the desired contradiction. ■

**Theorem 1 (Robust Direct Implementation)**

*Social choice function  $f$  is robustly directly implementable if and only if  $f$  satisfies EPIC and direct robust monotonicity.*

**Proof.** The existence of a truthtelling equilibrium on every type space is equivalent to EPIC (see Bergemann and Morris (2004), Proposition 2). With the existence of a truthtelling equilibrium, verifying robust implementation reduces to verifying iterative implementation by Proposition 1. Iterative implementation requires that the largest fixed point of  $b$  is acceptable. But Lemma 2 shows that this is guaranteed by direct robust monotonicity. ■

**6 Robust Implementation**

A "robust monotonicity" condition is key to our main result. In the direct mechanism, where agents other than  $i$  report themselves to be types  $\theta_{-i}$ , agent  $i$  can obtain outcomes  $f(\theta'_i, \theta_{-i})$  for any  $\theta'_i$ . But once we allow augmented mechanisms, we could conceivably offer agent  $i$  a larger set of lotteries if he reports deviant behavior of his opponents. We need to identify, for any given report  $\theta_{-i}$ , the set of lotteries with the property that whatever agent  $i$ 's actual type, he would never prefer such an allocation to what he would obtain under the social choice function if other agents were reporting truthfully. Thus:

$$Y_i(\theta_{-i}) = \{y \in Y : u_i(y, (\theta'_i, \theta_{-i})) \leq u_i(f(\theta'_i, \theta_{-i}), (\theta'_i, \theta_{-i})) \text{ for all } \theta'_i \in \Theta_i\}.$$

Henceforth, we refer to the set  $Y_i(\theta_{-i})$  as the *reward set* (for agent  $i$ ).

To understand the robust monotonicity condition, it is useful to first think about agents playing the direct mechanism. Suppose that it was common knowledge that in the direct mechanism, type  $\theta_i$  of player  $i$  will send a report  $\theta'_i \in \beta_i(\theta_i)$ . If  $\beta$  is acceptable, we would know that  $f$  was being implemented. But if  $\beta$  is unacceptable, we must find a type of some agent who is prepared to report that other players are misreporting. But for the "whistle-blower" who is going to report that we are in a bad equilibrium, we cannot know what he believes about the types of the other players, nor can we know what message he expects to hear except that it is a message consistent with the deception. Finally, the reward that he is offered must not mess up the truth-telling behavior in the good equilibrium. This gives the following condition:

**Definition 14 (Robust Monotonicity)**

*Social choice function  $f$  satisfies robust monotonicity if for every unacceptable deception  $\beta$ , there exist  $i$ ,  $\theta_i$ ,  $\theta'_i \in \beta_i(\theta_i)$  such that, for all  $\theta'_{-i} \in \Theta_{-i}$  and  $\psi_i \in \Delta(\beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $y$  such that:*

$$\sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i}) u_i(y, (\theta_i, \theta_{-i})) > \sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i}) u_i(f(\theta'_i, \theta'_{-i}), (\theta_i, \theta_{-i})), \quad (5)$$

while

$$u_i(f(\theta''_i, \theta'_{-i}), (\theta_i, \theta'_{-i})) \geq u_i(y, (\theta''_i, \theta'_{-i})), \quad \text{for all } \theta''_i \in \Theta_i. \quad (6)$$

Notice that this condition has a very similar form to the direct robust monotonicity condition. We simply allow a richer set of rewards in an augmented mechanism. In addition, in the augmented mechanism, agent  $i$  can propose an allocation conditional on the misreport  $\theta'_{-i}$  of the other agents. It therefore suffices that the reward  $y$  agent  $i$  proposes in the event of report  $\theta'_{-i}$  is successful for all possible distributions  $\psi_i(\theta_{-i})$  over the set of types  $\beta_{-i}^{-1}(\theta'_{-i})$  which could have reported  $\theta'_{-i}$  under the deception  $\beta_{-i}$ . In contrast, if we were to seek robust direct implementation, the report of agent  $i$  has to lead to rewards which work for all possible misreports  $\theta'_{-i}$  and underlying true type profiles  $\theta_{-i}$ . In consequence, the expectation had to be taken both with respect to the reports and the true types, i.e.  $\rho(\theta_{-i}, \theta'_{-i})$ .

If we compare the notion of robust monotonicity with the notion of Bayesian monotonicity for a given type space (see Definition 24), then three major differences appear. First, the notion of a deception is set-valued rather than point-valued. Second, the reward has to be successful against all possible distributions  $\psi_i(\theta_{-i})$  over true payoff profile rather than a single distribution, the posterior derived from the common prior. Third, the incentive constraints for the reward, (6), have to be satisfied ex post rather than interim. All three modifications directly stem from the robustness concern. The deception has to be set valued as in a rich type, a given payoff type  $\theta_i$  can now generate different misreports  $\theta'_i$  through distinct types  $t_i$  who all share the same true payoff type. Similarly, in a rich type space, there might be many distributions over payoff types,  $\psi_i(\theta_{-i})$ , which adopt a given misreport  $\theta'_i$ . Finally, the ex post incentive constraints regarding the reward  $y$  are necessary to maintain the truth-telling equilibrium in the direct revelation segment of the augmented mechanism. This is the robustness property of the ex post constraints in the direct mechanism developed in Bergemann and Morris (2004).

Robust monotonicity is strictly stronger than both Maskin monotonicity and ex post monotonicity (a necessary and almost sufficient condition for ex post implementation described in Bergemann and Morris (2005)). To get a sense of the strength of the condition, we can return to the examples of Section 3:

1. In the majority rule example, Maskin monotonicity and ex post monotonicity are both satisfied, but robust monotonicity fails.
2. In the coordination example, robust monotonicity (and thus Maskin and ex post monotonicity) are satisfied.

3. In the public goods example, robust monotonicity (and thus Maskin and ex post monotonicity) are satisfied if there is not too much interdependence of preferences. If there is too much interdependence of preferences, robust monotonicity fails but both Maskin and ex post monotonicity are satisfied.
4. In the private good example, ex post monotonicity holds, but Maskin monotonicity (and thus robust monotonicity) fails. However, in a perturbed version of this example, we again have robust monotonicity satisfied if only if there is not too much interdependence of preferences.

Although we will not use it extensively, notice that if  $M_i$  and  $\Theta_i$  are finite, then standard duality arguments imply the following alternative characterization: social choice function  $f$  satisfies *robust monotonicity* if for every unacceptable deception  $\beta$ , there exist  $i$ ,  $\theta_i$ ,  $\theta'_i \in \beta_i(\theta_i)$  such that, for all  $\theta'_{-i} \in \Theta_{-i}$ , there exists  $y \in Y_i(\theta'_{-i})$  such that:

$$u_i(y, (\theta_i, \theta_{-i})) > u_i(f(\theta'_i, \theta'_{-i}), (\theta_i, \theta_{-i})),$$

for all  $\theta_{-i}$  such that  $\theta'_{-i} \in \beta_{-i}(\theta_{-i})$ . Note that this characterization uses the requirement that  $Y$  is a lottery space.

Finally, we need an extremely weak economic condition to ensure that it is always possible to reward and punish each player independently of the others.

### Definition 15 (Bad Outcome)

The bad outcome property is satisfied if, for each  $i$ , there exists  $\underline{y}_i \in Y$  such that, for all  $\theta_i \in \Theta_i$ ,  $\psi_i \in \Delta(\Theta_{-i})$  and  $\theta'_{-i} \in \Theta_{-i}$ , there exists  $y \in Y_i(\theta'_{-i})$  such that

$$\sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(y, (\theta_i, \theta_{-i})) > \sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(\underline{y}_i, (\theta_i, \theta_{-i})).$$

This property says that there exists a bad outcome  $\underline{y}_i$  for each player  $i$  such that we can always offer him a lottery that makes him better off whatever his beliefs about other players' types and whatever reports other players are making. The bad outcome property, together with the use of lotteries, allows us to dispense with any no veto property which typically appear in the sufficient conditions. In addition, we can omit the usual cardinality assumption of  $I \geq 3$ .

### Theorem 2 (Robust Implementation)

1. If  $f$  is robustly implementable, then  $f$  satisfies EPIC and robust monotonicity;
2. if  $f$  satisfies EPIC, robust monotonicity and the bad outcome property, then  $f$  is robustly implementable.

**Proof. (1.)** We first prove that robust implementability implies EPIC and robust monotonicity. We do so by appealing to the necessary conditions for robust implementation in Corollary 2.

We first establish EPIC. By Lemma 1, there exists  $\nu_i \in \Delta(S_{-i}^{\mathcal{M}}(\theta_{-i}))$ ,

$$u_i(f(\theta_i, \theta_{-i}), (\theta_i, \theta_{-i})) \geq \sum_{m_{-i}} \nu_i(m_{-i}) u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})),$$

for all  $m_i \in M_i$  and  $\theta_i \in \Theta_i$ . If we choose  $m_i \in S_i^{\mathcal{M}}(\theta'_i)$ , iterative material implementation implies that  $g(m_i, m_{-i}) = f(\theta'_i, \theta_{-i})$  for all  $m_{-i} \in S_{-i}^{\mathcal{M}}(\theta_{-i})$ . So

$$u_i(f(\theta_i, \theta_{-i}), (\theta_i, \theta_{-i})) \geq u_i(f(\theta'_i, \theta_{-i}), (\theta_i, \theta_{-i})),$$

for all  $\theta'_i \in \Theta_i$ .

We next establish robust monotonicity. Fix an unacceptable deception  $\beta$  and suppose that  $f$  is iteratively materially implementable. There must exist a message correspondence profile  $S$  such that

$$b(S) \leq S,$$

and

$$S_i^{\mathcal{M}}(\theta'_i) \subseteq S_i(\theta_i), \quad (7)$$

for all  $i$ ,  $\theta_i$  and  $\theta'_i \in \beta_i(\theta_i)$ ; but

$$S_i^{\mathcal{M}}(\theta'_i) \not\subseteq b_i(S)[\theta_i], \quad (8)$$

for all  $i$ ,  $\theta_i$  and  $\theta'_i \in \beta_i(\theta_i)$ . The existence of such an  $S$  can be established constructively. Clearly  $\bar{S}$  satisfies (7). Iteratively apply the operator  $b$ . By iterative implementation, there exists  $k$  (perhaps transfinite) such that:

$$S \triangleq b^k(\bar{S}) \quad (9)$$

satisfies (8). Thus there exists  $k$  such that  $b^k(\bar{S})$  satisfies (7) and  $b^{k+1}(\bar{S})$  satisfies (8).

By (8), simply pick

$$\hat{m}_i \in S_i(\theta_i) \cap S_i^{\mathcal{M}}(\theta'_i) \quad \text{and} \quad \hat{m}_i \notin b_i(S)[\theta_i] \cap S_i^{\mathcal{M}}(\theta'_i).$$

Since message  $\hat{m}_i \notin b_i(S)[\theta_i]$ , we know that for every  $\lambda_i \in \Delta(M_{-i} \times \Theta_{-i})$  such that

$$\lambda_i(m_{-i}, \theta_{-i}) > 0 \Rightarrow m_j \in S_j(\theta_j) \quad \text{for all } j \neq i,$$

there exists  $m_i^*$  such that

$$\sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(m_i^*, m_{-i}), (\theta_i, \theta_{-i})) > \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(\hat{m}_i, m_{-i}), (\theta_i, \theta_{-i})). \quad (10)$$

Next we identify a particular belief  $\lambda_i(m_{-i}, \theta_{-i})$  for which the inequality (10) holds. By (3) in Lemma 1, there exists  $\nu_i \in \Delta(S_{-i}^{\mathcal{M}}(\theta'_{-i}))$  such that

$$\sum_{m_{-i}} \nu_i(m_{-i}) u_i(g(m_i, m_{-i}), (\theta''_i, \theta'_{-i})) \leq u_i(f(\theta''_i, \theta'_{-i}), (\theta''_i, \theta'_{-i})), \quad (11)$$

for all  $m_i \in M_i$  and  $\theta''_i \in \Theta_i$ . Thus for any  $\psi_i \in \Delta(\Theta_{-i})$ , we can set

$$\lambda_i(m_{-i}, \theta_{-i}) = \nu_i(m_{-i}) \psi_i(\theta_{-i}).$$

Applying the above claim (10), there exists  $m_i^*$  such that:

$$\sum_{\theta_{-i}, m_{-i}} \psi_i(\theta_{-i}) \nu_i(m_{-i}) u_i(g(m_i^*, m_{-i}), (\theta_i, \theta_{-i})) > \sum_{\theta_{-i}, m_{-i}} \psi_i(\theta_{-i}) \nu_i(m_{-i}) u_i(g(\widehat{m}_i, m_{-i}), (\theta_i, \theta_{-i})).$$

But  $\nu_i(m_{-i}) > 0 \Rightarrow (\widehat{m}_i, m_{-i}) \in S^{\mathcal{M}}(\theta')$ , so by iterative material implementation:

$$g(\widehat{m}_i, m_{-i}) = f(\theta').$$

We also observe that as we defined  $S$  to be the set obtained after the  $k$ -th iteration of the operator  $b$ , see (9), if  $\theta'_{-i} \in \beta_{-i}(\theta_{-i})$ , then  $\nu_i(m_{-i}) > 0 \Rightarrow m_{-i} \in S_{-i}(\theta_{-i})$ . Thus for every  $\psi_i \in \Delta(\beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $m_i^*$  such that

$$\sum_{\theta_{-i}, m_{-i}} \psi_i(\theta_{-i}) \nu_i(m_{-i}) u_i(g(m_i^*, m_{-i}), (\theta_i, \theta_{-i})) > \sum_{\theta_{-i}, m_{-i}} \psi_i(\theta_{-i}) \nu_i(m_{-i}) u_i(f(\theta'), (\theta_i, \theta_{-i})). \quad (12)$$

Now, the inequality (12) essentially establishes guarantees the reward inequality for robust monotonicity. We can complete the argument by letting  $y$  be the lottery with

$$y(z) \equiv \sum_{m_{-i}} g(m_i^*, m_{-i}) \nu_i(m_{-i}).$$

We now have established that for each  $\theta'_{-i} \in \beta_{-i}(\theta_{-i})$  and  $\psi_i \in \Delta(\beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $y$  such that (by (11))

$$u_i(y, (\theta''_i, \theta'_{-i})) \leq u_i(f(\theta''_i, \theta'_{-i}), (\theta''_i, \theta'_{-i})),$$

for all  $\theta''_i \in \Theta_i$ , and thus  $y \in Y_i(\theta'_{-i})$ .<sup>13</sup> And by (12) we then have:

$$\sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(y, (\theta_i, \theta_{-i})) > \sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(f(\theta'), (\theta_i, \theta_{-i})).$$

---

<sup>13</sup>Note that this step implies that even if we had restricted attention to mechanisms with deterministic outcomes, our robust monotonicity would only have established that there exists a lottery (not necessarily a deterministic outcome) sufficient to reward a whistle-blower.

(2.) We now prove that EPIC, robust monotonicity and the bad outcome property imply robust implementation. We do so by explicitly constructing the implementing mechanism.

Each agent  $i$  sends a message  $m_i = (m_i^1, m_i^2, m_i^3, m_i^4)$ , where  $m_i^1 \in \Theta_i$ ,  $m_i^2 \in \mathbb{Z}^+$ ,  $m_i^3 : \Theta_{-i} \rightarrow \Delta(Y)$  with  $m_i^3(\theta_{-i}) \in Y_i(\theta_{-i})$ ,  $m_i^4 \in Y$ . The outcome  $g(m)$  is determined by the following rules.

Rule 1: If  $m_i^2 = 1$  for all  $i$ , pick  $f(m^1)$ .

Rule 2: If there exists  $j \in I$  such that  $m_i^2 = 1$  for all  $i \neq j$  and  $m_j^2 > 1$ , then pick  $m_j^3(m_{-j}^1)$  with probability  $1 - \frac{1}{m_j^2+1}$  and  $\underline{y}_i$  with probability  $\frac{1}{m_j^2+1}$ .

Rule 3: In all other cases, for each  $i$ , with probability  $\frac{1}{I} \left(1 - \frac{1}{m_i^2+1}\right)$  pick  $m_i^4$ , and with probability  $\frac{1}{I} \left(\frac{1}{m_i^2+1}\right)$  pick  $\underline{y}_i$ .

We first show that it is never a best reply for type  $\theta_i$  to send a message with  $m_i^2 > 1$  (i.e.,  $m_i \in b_i(\bar{S}) \Rightarrow m_i^2 = 1$ ). Suppose that  $\theta_i$  has conjecture  $\lambda_i \in \Delta(M_{-i} \times \Theta_{-i})$ . We can partition the messages of other players as follows:

$$M_{-i}^*(\theta_{-i}) = \{m_{-i} : m_j^2 = 1 \text{ for all } j \neq i \text{ and } m_{-i}^1 = \theta_{-i}\},$$

and

$$\widehat{M}_{-i} = \{m_{-i} : m_j^2 > 1 \text{ for some } j \neq i\}.$$

By the bad outcome property, we know that there exists  $m_i^4 \in Y$  such that, if

$$\sum_{m_{-i} \in \widehat{M}_{-i}, \theta_{-i} \in \Theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) > 0,$$

then

$$\sum_{m_{-i} \in \widehat{M}_{-i}, \theta_{-i} \in \Theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(m_i^4, \theta) > \sum_{m_{-i} \in \widehat{M}_{-i}, \theta_{-i} \in \Theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(\underline{y}_i, \theta).$$

And we know that there exists  $m_i^3$  such that, if

$$\sum_{m_{-i} \in M_{-i}^*(\theta'_{-i}), \theta_{-i} \in \Theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) > 0,$$

then

$$\sum_{m_{-i} \in M_{-i}^*(\theta'_{-i}), \theta_{-i} \in \Theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(m_i^3(\theta'_{-i}), \theta) > \sum_{m_{-i} \in M_{-i}^*(\theta'_{-i}), \theta_{-i} \in \Theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(\underline{y}_i, \theta).$$

Thus if  $(m_i^1, m_i^2, m_i^3, m_i^4)$  with  $m_i^2 > 1$  were a best response, then  $(m_i^1, m_i^2 + 1, m_i^3, m_i^4)$  would be an even better response, contradiction.

Now fix any  $S$  with  $m_i \in S_i(\theta_i) \Rightarrow m_i^2 = 1$ . Let

$$\beta_i(\theta_i) = \{\theta'_i : (\theta'_i, 1, m_i^3, m_i^4) \in S_i(\theta_i) \text{ for some } (m_i^3, m_i^4)\}.$$

First observe that EPIC implies that  $\theta_i \in \beta_i(\theta_i)$ . We will argue that if  $\beta$  is not acceptable, then  $b(S) \neq S$ . By robust monotonicity, we know that there exists  $i$ ,  $\theta_i$ ,  $\theta'_i \in \beta_i(\theta_i)$  such that, for all  $\theta'_{-i} \in \Theta_{-i}$  and  $\psi_i \in \Delta(\beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $y \in Y_i(\theta'_{-i})$  such that

$$\sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i}) u_i(y, (\theta_i, \theta_{-i})) > \sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i}) u_i(f(\theta'_i, \theta'_{-i}), (\theta_i, \theta_{-i})).$$

But now for any conjecture  $\lambda_i \in \Delta\left(\left\{(m_{-i}, \theta_{-i}) : m_j^2 = 1 \text{ for all } j \neq i\right\}\right)$ , there exists  $m_i^3$  (with  $m_i^3(\theta_{-i}) \in Y_i(\theta_{-i})$ ) such that

$$\sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(m_i^3(m_{-i}^1), \theta) > \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(f(\theta'_i, m_{-i}^1), (\theta_i, \theta_{-i})).$$

Thus message  $(\theta'_i, 1, m_i^3, m_i^4)$  is never a best response for type  $\theta_i$ .

We conclude that if

$$\beta_i(\theta_i) = \{\theta'_i : (\theta'_i, 1, m_i^3, m_i^4) \in S_i^{\mathcal{M}}(\theta_i) \text{ for some } (m_i^3, m_i^4)\},$$

then  $\beta$  is acceptable. Thus  $f$  is iteratively materially implemented.

Finally observe that  $S^{\mathcal{M}}$  must satisfy the ex post best response property, with type  $\theta_i$  sending a message of the form  $(\theta'_i, 1, m_i^3, m_i^4)$ , so robust implementation is possible by Corollary 2. ■

The proof directly uses the link between iterative and robust implementation for the necessity as well as the sufficiency part. We briefly sketch the idea of the necessity part of the proof. If  $f$  is robustly implementable, then it is iteratively implementable by Corollary 2. From iterative implementability, we then want to show that  $f$  satisfies strict robust monotonicity. We then consider a given *and* unacceptable deception  $\beta$ . We start the process of iterative elimination and stop it at a specific round, denoted by  $k$ . This round  $k$  is the first round at which we can find an agent  $i$ , a true type profile  $\theta_i$  and a report  $\theta'_i \in \beta_i(\theta_i)$ , such that a message, denoted by  $\hat{m}_i$ , which will survive the process of iterated elimination for type  $\theta'_i$ , fails to survive the  $k$ -th round of elimination for type  $\theta_i$ . We then show that the elimination of message  $\hat{m}_i$  at round  $k$  implies that the social choice function  $f$  satisfies strict robust monotonicity with respect to the deception  $\beta$ . Briefly, if  $\hat{m}_i$  survives the process of elimination for type  $\theta'_i$ , the message  $\hat{m}_i$  acts in the mechanism so as to report a payoff type  $\theta'_i$ . If it is eliminated at round  $k$  for payoff type  $\theta_i$ , then this means that for any belief agent  $i$  has over the remaining agents, there exists a message  $m_i^*$  which leads to an allocation through  $g$  which is strictly preferred by agent  $i$  when he has a payoff type  $\theta_i$ . The significance of



round  $k$  being the first round for which such an elimination relative to the deception  $\beta$  occurs, is that at this round, there do not yet exist any restrictions about message and payoff type profile regarding the other players deception. The fact then that  $\hat{m}_i$  can be eliminated allows us to use full strength of the elimination argument to establish robust monotonicity. In the context of the proof it is interesting to note that the key step from iterative elimination to robust monotonicity is an argument which involves the early stages of the elimination process rather than the limit of iteration process.

The results of Theorem 2 rely both on allowing lotteries and the bad outcome property. In Section 9.1, both assumptions are discussed and a simple example satisfying EPIC and robust monotonicity but not robustly implementable without lotteries is described.

The novel difficulty in obtaining the necessary results for implementation arise from the robustness requirement. If  $f$  is robustly implemented, the mechanism which achieves implementation could be badly behaved with respect to the existence of best responses against all possible beliefs about action and type profiles of the other agents. This difficulty did not appear in the direct robust implementation as direct implementation guaranteed ex post incentive compatibility and inside the direct mechanism the (non-empty) existence of  $S^M$ . In Section 9.2 we reconsider the robust implementation result by restricting attention to *nice mechanisms*, mechanisms in which best responses always exist for all beliefs over payoff type and message profiles.

Finally, in the Appendix we provide a direct proof that interim monotonicity on all type spaces is equivalent to robust monotonicity. The result, contained in Proposition 6, explicitly constructs a type space to show that if robust monotonicity fails then there also exists a type space  $\mathcal{T}$  for which Bayesian monotonicity fails.

## 7 Virtual Implementation

In this section, we extend our robustness analysis from interim to virtual implementation. This section is more limited in scope than the previous sections. The notion of virtual Bayesian implementation is widely considered to be much more permissive than interim implementation (see Abreu and Matsushima (1992b), Duggan (1997) and Serrano and Vohra (2005)).

We shall define a robust version of virtual implementation and then give a simple condition, type indistinguishability, to obtain an impossibility result for robust virtual implementation. The impossibility result will reappear in Section 8. There we consider an environment in which the payoff types can be linearly aggregated with respect to the impact of the types on the utility of each agent. In this linear environment, we will be able to derive an exact and sharp bound for the possibility of robust implementation. Conversely, if robust implementation fails the bound, then

robust virtual implementation will also be impossible. With respect to virtual implementation, a message of the robustness analysis consequently will be that the difference between interim and virtual implementation shrinks substantially once we impose the robustness requirement on implementation.

**Definition 16 (Virtual Implementation)**

Social choice function  $f$  is  $\varepsilon$  implementable by mechanism  $\mathcal{M}$  on type space  $\mathcal{T}$ , if there exists an equilibrium of the game  $(\mathcal{T}, \mathcal{M})$  and every equilibrium  $\sigma$  of the game  $(\mathcal{T}, \mathcal{M})$  satisfies

$$\sum_{m \in M} g\left(f\left(\widehat{\theta}(t)\right) | m\right) \sigma(m|t) \geq 1 - \varepsilon.$$

Social choice function  $f$  is virtually implementable if it is  $\varepsilon$ -implementable for all  $\varepsilon > 0$ .

We denote by  $g(y|m)$  the probability that the outcome  $y$  is realized conditional on receiving message  $m$ . We next present the robust version of virtual implementation.

**Definition 17 (Virtual Robust Implementation)**

Social choice function  $f$  is  $\varepsilon$ -robustly implementable by mechanism  $\mathcal{M}$  if, on every type space  $\mathcal{T}$ , there exists an equilibrium of the game  $(\mathcal{T}, \mathcal{M})$  and every equilibrium  $\sigma$  of the game  $(\mathcal{T}, \mathcal{M})$  satisfies

$$\sum_{m \in M} g\left(f\left(\widehat{\theta}(t)\right) | m\right) \sigma(m|t) \geq 1 - \varepsilon.$$

Social choice function  $f$  is virtually robustly implementable if it is  $\varepsilon$ -robustly implementable for all  $\varepsilon > 0$ .

In this section, we restrict attention to mechanisms that satisfy a best response property.

**Definition 18 (Nice Mechanism)**

Mechanism  $\mathcal{M}$  is nice if for all  $i$  and  $\theta_i \in \Theta_i$ ,

$$S^{\mathcal{M}}(\theta_i) \neq \emptyset,$$

$$\text{and } \arg \max_{m_i} \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(m_i, m_{-i}), (\theta_i, \theta_{-i})) \neq \emptyset,$$

for all  $\lambda_i \in \Delta(M_{-i} \times \Theta_{-i})$ .

Thus rationalizable sets are non-empty and best responses exist to all conjectures. This restriction is related to restrictions studied in the complete information literature, such as the bounded

mechanism condition of Jackson (1992), but we have not yet identified the exact relation. Note that this property is automatically satisfied if message spaces are finite. Next we present a sufficient condition to obtain an impossibility result for robust virtual implementation.<sup>14</sup>

**Definition 19 (Indistinguishable Payoff Types)**

Payoff types are indistinguishable if, for all  $i$ , there exists  $\psi_i : \theta_i \rightarrow \Delta(\Theta_{-i})$  such that for all  $\theta_i, \theta'_i \in \Theta_i$  and  $y, y' \in Y$ ,

$$\begin{aligned} \sum_{\theta_{-i}} u_i(y, (\theta_i, \theta_{-i})) \psi_i(\theta_{-i}|\theta_i) &\geq \sum_{\theta_{-i}} u_i(y', (\theta_i, \theta_{-i})) \psi_i(\theta_{-i}|\theta_i) \\ \Leftrightarrow \sum_{\theta_{-i}} u_i(y, (\theta'_i, \theta_{-i})) \psi_i(\theta_{-i}|\theta'_i) &\geq \sum_{\theta_{-i}} u_i(y', (\theta'_i, \theta_{-i})) \psi_i(\theta_{-i}|\theta'_i). \end{aligned}$$

We note that any two types,  $\theta_i$  and  $\theta'_i$ , have to agree in their ranking of the alternatives only for some specific posterior distribution,  $\psi_i(\theta_{-i}|\theta_i)$  and  $\psi_i(\theta_{-i}|\theta'_i)$ , but not for all posterior distributions. On the other hand, the ranking has to be constant for any arbitrary pair,  $\theta_i, \theta'_i \in \Theta_i$ .

**Proposition 3 (Failure of Virtual Robust Implementation)**

If types are indistinguishable, then  $f$  is virtually robustly implementable by a nice mechanism if and only if it is constant.

**Proof.** Clearly  $f$  is virtually robustly implementable if it is constant. Suppose that types are not indistinguishable. Then there exists  $v_i : Y \rightarrow \mathbb{R}$  such that type  $\theta_i$  with beliefs  $\psi_i(\theta_{-i})$  prefers  $y$  to  $y'$  if and only if  $v_i(y) \geq v_i(y')$ . We argue by induction that for all  $i$  and  $k$ , there exists  $m_i^k$  such that  $m_i^k \in S_i^k(\theta_i)$  for all  $\theta_i \in \Theta_i$ . It is clearly true for  $k = 0$  and to begin the inductive step suppose it is true for  $k$ . Fix any  $i$  and fix

$$m_i^{k+1} \in \arg \max_{m_i} v_i \left( g \left( m_i, m_{-i}^k \right) \right),$$

for some  $m_{-i}^k \in S_{-i}^k$ . Now  $m_i^{k+1} \in S_i^{k+1}(\theta_i)$  for all  $\theta_i$ . Thus for each  $i$  there exists  $m_i^* \in S_i^{\mathcal{M}}(\theta_i)$  for all  $\theta_i$ . Thus for any mechanism  $\mathcal{M}$ , there exists  $y \in Y$  such that on any type space, outcome  $y$  is always realized. Thus virtual robust implementation requires that for every  $\varepsilon > 0$ , there exists  $y \in Y$  such that  $\|f(\theta) - y\| \leq \varepsilon$ . This in turn requires that  $f$  is constant. ■

The same argument will imply that if types are indistinguishable, then a social choice function cannot be implemented by iterated deletion of "weakly ex post dominated" messages, as in Chung and Ely (2001).

---

<sup>14</sup>This condition can be further weakened with an iterative indistinguishability condition that is the robust analogue of the measurability condition in Abreu and Matsushima (1992b).

## 8 Linear Aggregation of Interdependent Types

In this section, we present a class of environments with interdependent preferences, for which we can derive precise implementation results. The environment is defined to be the class of interdependent preferences where the payoff types can be aggregated linearly for the utility representation. This class of preferences allows us to clearly illustrate the link between iterative and robust implementation. With the linear aggregation of the payoff types, we obtain two very distinct implementation results, separated by a sharp bound on the size of the interdependence in the valuations. If there is not too much interdependence, then robust implementation is possible. Conversely, if there is too much interdependence, then robust implementation will be impossible, but more surprisingly, even robust virtual implementation is impossible.

Thus we consider preferences  $u_i(y, \theta)$  which permit linear aggregation of the payoff types as follows:

$$u_i(y, \theta) \triangleq v_i\left(y, \theta_i + \gamma \sum_{j \neq i} \theta_j\right),$$

with  $\gamma \in \mathbb{R}$  and  $\Theta_i = [0, 1]$  for all  $i$ . We define the value of the linear aggregator for agent  $i$  as  $\mu_i$ :

$$\mu_i : \mathbb{R}^I \rightarrow \mathbb{R},$$

and accordingly can write the preferences of agent  $i$  as a function of the allocation and the linear aggregator  $\mu_i$ , or  $v_i(y, \mu_i)$ .

With the linear aggregation of the payoff types, we obtain two very distinct implementation results, separated by a sharp bound on the size of the interdependence in the valuations. We begin the impossibility result under “too much” interdependence.

### Theorem 3 (Virtual Implementation)

*If  $\gamma > \frac{1}{I-1}$ , then types are indistinguishable and virtual robust implementation is impossible.*

**Proof.** Suppose  $\gamma > \frac{1}{I-1}$ . Let type  $\theta_i$  always put probability 1 on a profile  $\theta_{-i}$  with

$$\sum_{j \neq i} \theta_j = \frac{1}{\gamma} \left( \frac{1}{2} - \theta_i \right).$$

Then

$$\begin{aligned} \theta_i + \gamma \sum_{j \neq i} \theta_j &= \theta_i + \gamma \left( \frac{1}{\gamma} \left( \frac{1}{2} - \theta_i \right) \right) \\ &= \theta_i + \frac{1}{2} - \theta_i \\ &= \frac{1}{2}. \end{aligned}$$

This implies that types are indistinguishable. While Proposition 3 is written for discrete types, it extends straightforwardly to compact type spaces. ■

On the other hand, we have positive implementation results when there is not too much interdependence in preferences. Before we can state the positive implementation results, we need to make the usual assumptions regarding single crossing to guarantee ex post incentive compatibility in the direct mechanism.

**Definition 20 (Single Crossing)**

$v_i$  satisfies single crossing if  $v_i(y', \mu) = v_i(y, \mu)$  and  $v_i(y', \mu') > v_i(y, \mu')$  imply  $v_i(y', \tilde{\mu}) > v_i(y, \tilde{\mu})$  for all  $\tilde{\mu} > \mu$  (if  $\mu' > \mu$ ) and for all  $\tilde{\mu} < \mu$  (if  $\mu' < \mu$ ).

This condition requires that for any pair of lotteries  $y$  and  $y'$ , there is at most one value of the aggregate of individual types,  $\mu$ , where preferences over the lotteries switch. We state the single crossing condition in terms of the linear aggregator  $\mu$  as we do not make any ordering assumption about the set of allocations  $Y$ .

**Definition 21 (Non-degeneracy)**

$v_i$  is nondegenerate if for all  $y$ ,  $\mu \neq \mu'$ , there exists  $y'$  such that  $v_i(y', \mu) = v_i(y, \mu)$  and  $v_i(y', \mu') > v_i(y, \mu')$ .

With robust implementation, we need to strengthen the ex post incentive constraints in the direct mechanism to be strict inequalities.

**Definition 22 (Strict Ex Post Incentive Compatibility)**

Social choice function  $f$  satisfies strict ex post incentive compatibility if  $f(\theta_i, \theta'_{-i}) \neq f(\theta'_i, \theta'_{-i})$  for some  $\theta'_{-i}$  implies

$$u_i(f(\theta_i, \theta_{-i}), (\theta_i, \theta_{-i})) > u_i(f(\theta'_i, \theta_{-i}), (\theta_i, \theta_{-i}))$$

for all  $i$  and  $\theta_{-i}$ .

We can now establish the converse to Theorem 3.

**Theorem 4 (Robust Implementation)**

If  $\gamma < \frac{1}{I-1}$ , strict EPIC holds and each  $v_i$  satisfies single crossing and non-degeneracy, then robust monotonicity is satisfied.

**Proof.** Consider the deception  $\beta$ . Let

$$\Delta = \sup_{i, \theta_i, \theta'_i \in \beta_i(\theta_i)} |\theta'_i - \theta_i|.$$

If  $\Delta = 0$ , then the deception is acceptable and we are done. If  $\Delta > 0$ , fix any  $0 < \varepsilon < 1 - \gamma(I - 1)$  and  $(i, \theta_i, \theta'_i)$  with  $|\theta'_i - \theta_i| > \Delta(1 - \varepsilon)$ . Fix any  $\theta'_{-i}$ . By non-degeneracy of  $v_i$ , there exists  $y'$  such that  $u_i(y', \theta') = u_i(f(\theta'), \theta')$  and  $u_i(y', (\theta_i, \theta'_{-i})) > u_i(f(\theta'), (\theta_i, \theta'_{-i}))$ . Let  $y^n$  be a lottery putting probability  $\frac{1}{n}$  on  $y'$  and probability  $1 - \frac{1}{n}$  on  $f(\theta')$ . Now  $y^n \rightarrow f(\theta')$ ,  $u_i(y^n, \theta') = u_i(f(\theta'), \theta')$  and  $u_i(y^n, (\theta_i, \theta'_{-i})) > u_i(f(\theta'), (\theta_i, \theta'_{-i}))$  for all  $n$ .

We first establish that for sufficiently large  $n$ ,  $y^n \in Y_i^*(\theta'_{-i})$ . If this wasn't true, then for every  $n$ , there would exist  $\tilde{\theta}_i$  such that

$$u_i\left(y^n, \left(\tilde{\theta}_i, \theta'_{-i}\right)\right) > u_i\left(f\left(\tilde{\theta}_i, \theta'_{-i}\right), \left(\tilde{\theta}_i, \theta'_{-i}\right)\right).$$

By continuity, we must have

$$u_i\left(f\left(\theta'\right), \left(\tilde{\theta}_i, \theta'_{-i}\right)\right) \geq u_i\left(f\left(\tilde{\theta}_i, \theta'_{-i}\right), \left(\tilde{\theta}_i, \theta'_{-i}\right)\right),$$

for some  $\tilde{\theta}_i$ , contradicting strict EPIC.

Now suppose that  $\theta'_{-i} \in \beta_{-i}(\theta_{-i})$ . Observe that

$$\left(\theta_i + \gamma \sum_{j \neq i} \theta_j\right) - \left(\theta'_i + \gamma \sum_{j \neq i} \theta'_j\right) = (\theta_i - \theta'_i) + \gamma \sum_{j \neq i} (\theta'_j - \theta_j).$$

But by hypothesis:

$$|\theta_i - \theta'_i| > \Delta(1 - \varepsilon) > \gamma(I - 1)\Delta,$$

while

$$\begin{aligned} \gamma \left| \sum_{j \neq i} (\theta'_j - \theta_j) \right| &\leq \gamma \sum_{j \neq i} |\theta'_j - \theta_j| \\ &\leq \gamma(I - 1)\Delta. \end{aligned}$$

So the sign of

$$\left(\theta_i + \gamma \sum_{j \neq i} \theta_j\right) - \left(\theta'_i + \gamma \sum_{j \neq i} \theta'_j\right),$$

equals the sign of

$$\theta_i - \theta'_i.$$

But now  $u_i(y^n, \theta') = u_i(f(\theta'), \theta')$  and  $u_i(y^n, (\theta_i, \theta'_{-i})) > u_i(f(\theta'), (\theta_i, \theta'_{-i}))$  implies, by single crossing of  $v_i$ ,  $u_i(y^n, (\theta_i, \theta_{-i})) > u_i(f(\theta'), (\theta_i, \theta_{-i}))$  for all  $\theta_{-i}$  such that  $\theta'_{-i} \in \beta_{-i}(\theta_{-i})$ . This establishes robust monotonicity. ■

The proof of Theorem 4 use the following four types,  $\theta_i, \tilde{\theta}_i, \theta_i^n, \theta_i'$  of agent  $i$ , ordered on the real line:

$$\begin{array}{ccccccc} \theta_i & & \tilde{\theta}_i & & \theta_i^n & & \theta_i' \\ & & \longleftarrow & & \longrightarrow & & \\ & & & & & & \mathbb{R} \end{array}$$

The current implementation results were obtained within an environment of linear and symmetric aggregation:

$$\mu_i = \theta_i + \gamma \sum_{j \neq i} \theta_j.$$

We can generalize the implementation results in this section to a general, not necessarily symmetric model of linear aggregation:

$$\mu_i = \theta_i + \sum_{j \neq i} \gamma_{ij} \theta_j.$$

The bound on the interdependence then comes through the eigenvalue of the interaction matrix  $\gamma = \{\gamma_{ij}\}$ . We suspect that an appropriate linearization argument would extend the technique to general non-linear environments.

Finally, if we impose some additional structure on the payoffs as a function of the reports in the direct mechanism, then we can use the same argument as in Theorem 4 to obtain direct robust monotonicity results.

**Definition 23 (Strictly Concave Deviations)**

For any  $\theta$ ,  $u_i \left( f \left( \tilde{\theta}_i, \theta_{-i} \right), \left( \theta_i, \theta_{-i} \right) \right)$  is strictly concave in  $\tilde{\theta}_i$ .

**Proposition 4 (Direct Implementation)**

If  $\gamma < \frac{1}{I-1}$ , strict EPIC and strictly concave deviations hold and each  $v_i$  satisfies single crossing, non-degeneracy, then robust direct monotonicity is satisfied.

**Proof.** The argument in the previous theorem goes through where we instead to define  $y^n = f(\theta_i^n, \theta_{-i}')$  with  $\theta_i^n \rightarrow \theta_i'$  and the sign of  $\theta_i^n - \theta_i'$  equal to the sign of  $\theta_i - \theta_i'$ . Since only the direct mechanism messages are used, the proof then establishes direct robust monotonicity rather than robust monotonicity. ■

## 9 Extensions, Variations and Discussion

### 9.1 Lotteries, Pure Strategies and Bayesian Implementation

In this section, we discuss how our main Theorem 2 is related to the classic literature on Bayesian implementation developed by Postlewaite and Schmeidler (1986), Palfrey and Srivastava (1989)

and Jackson (1991). These authors asked whether it was possible to implement a social choice function in equilibrium on a fixed type space  $\mathcal{T}$ .<sup>15</sup> These authors analyzed the classic problem where attention was restricted to pure strategy equilibria and deterministic mechanisms. Thus the social choice function is a mapping  $f : \Theta \rightarrow Z$ .

Having fixed a type space, the natural notion of a pure strategy deception on the fixed type space is a collection  $\alpha = (\alpha_1, \dots, \alpha_I)$ , with each  $\alpha_i : T_i \rightarrow T_i$ . Thus  $\alpha : T \rightarrow T$  is defined by  $\alpha(t) = (\alpha_i(t_i))_{i=1}^I$ . The key monotonicity notion, translated into our language, was then the following:

**Definition 24 (Bayesian Monotonicity)**

*Social choice function  $f$  satisfies Bayesian monotonicity on type space  $\mathcal{T}$  if, for every deception  $\alpha$  with  $f(\hat{\theta}(t)) \neq f(\hat{\theta}(\alpha(t)))$  for some  $t$ , there exists  $i$ ,  $t_i$  and  $h : T \rightarrow Z$  such that*

$$\sum_{t_{-i} \in T_{-i}} u_i \left( h(\alpha(t)), \hat{\theta}(t) \right) \hat{\pi}_i(t_{-i}) [t_i] > \sum_{t_{-i} \in T_{-i}} u_i \left( f(\hat{\theta}(\alpha(t))), \hat{\theta}(t) \right) \hat{\pi}_i(t_{-i}) [t_i], \quad (13)$$

and

$$\begin{aligned} & \sum_{t_{-i} \in T_{-i}} u_i \left( f(\hat{\theta}(t'_i, t_{-i})), \hat{\theta}(t'_i, t_{-i}) \right) \hat{\pi}_i(t_{-i}) [t'_i] \\ & \geq \sum_{t_{-i} \in T_{-i}} u_i \left( h(\alpha_i(t_i), t_{-i}), \hat{\theta}(t'_i, t_{-i}) \right) \hat{\pi}_i(t_{-i}) [t'_i], \quad \forall t'_i. \end{aligned} \quad (14)$$

Jackson shows that this condition is necessary for Bayesian implementation, and that a slight strengthening, Bayesian monotonicity no veto, is sufficient.

In the Appendix (Section 11.3), we show that our robust monotonicity condition is equivalent to the requirement that Bayesian monotonicity is satisfied on all type spaces. Thus we have an alternative way of showing the necessity of robust monotonicity for robust implementation. Figure 1 gave a graphical representation of how the results fit together.

But note that this line of argument would establish the necessity of robust implementation if the planner is restricted to deterministic mechanisms (a disadvantage) but he can assume that players follow pure strategies (an advantage). How do these assumptions matter?

First, observe that the advantage of restricting attention to pure strategies goes away completely when we require implementation on all type spaces: if there is a mixed strategy equilibrium that results in a socially sub-optimal outcome on some type space, we can immediately construct a larger type space (purifying the original equilibrium) where the socially sub-optimal outcome is played in

---

<sup>15</sup>They allowed for more general social choice sets, but we restrict attention to functions for our comparison.



a pure strategy equilibrium. Thus our robust analysis conveniently removes that unfortunate gap between pure and mixed strategy implementation that has plagued the implementation literature.

We use the extension to stochastic mechanisms in just two places. Ex post incentive compatibility and robust monotonicity would remain necessary conditions even if we restricted attention to deterministic mechanisms (the arguments would be unchanged). But, as we note in Footnote 13, even if lotteries were not used in the implementing mechanism, the implied robust monotonicity condition would involve lotteries (as rewards for whistle-blowers). But if lotteries were not allowed, our sufficiency argument would then require a slightly strengthened version of the robust monotonicity condition, with the lottery  $y$  replaced by a deterministic outcome. Our sufficiency argument also uses lotteries under Rules 1 and 2. As in a recent paper by Benoit and Ok (2004) on complete information implementation, we use lotteries to significantly weaken the sufficient conditions, so that we require only the "bad outcome" property in addition to EPIC and robust monotonicity. If we did not allow lotteries in this part of the argument, we would require a much stronger economic condition in the spirit of Jackson's "Bayesian monotonicity no veto" condition. We have developed combined robust monotonicity and economic conditions (not reported here) sufficient for interim implementation on all full support types spaces. However, an additional complication is that, without lotteries in the implementing mechanism, we cannot establish sufficiency on type spaces where agents have disjoint supports.

It is possible to construct a simple example where EPIC and robust monotonicity are not sufficient for robust monotonicity without lotteries by taking the coordination example of Section 3.2 but removing the outcomes  $e$  and  $f$ . As we show in the Appendix (Section 11.4), robust implementation is then not possible in this example despite the fact that the social choice function selects a unique strictly Pareto-dominant outcome at every type profile.

## 9.2 Nice Mechanisms

In our analysis of robust implementation, we deliberately allowed for very badly behaved infinite mechanisms in order to make a tight connection with the existing literature and to get tight results. Many authors have argued that "integer game" constructions, like that we use in Theorem 2, should not be taken seriously (see, e.g., Abreu and Matsushima (1992a) and Jackson (1992)). In our analysis of virtual robust implementation in Section 7, we restricted attention to "nice" mechanisms with best responses always well defined. Much of our analysis of the relation between iterative and robust implementation, and the characterization of robust implementation, would be much simpler with the restriction to nice mechanisms. In fact with nice mechanisms we obtain the following stronger necessary conditions for robust implementation.

**Definition 25 (Strict Robust Monotonicity)**

Social choice function  $f$  satisfies strict robust monotonicity if for every unacceptable deception  $\beta$ , there exist  $i$ ,  $\theta_i$ ,  $\theta'_i \in \beta_i(\theta_i)$  such that, for all  $\theta'_{-i} \in \Theta_{-i}$  and  $\psi_i \in \Delta(\beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $y$  such that

$$\sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i}) u_i(y, (\theta_i, \theta_{-i})) > \sum_{\theta'_{-i} \in \Theta_{-i}} \psi_i(\theta'_{-i}) u_i(f(\theta'_i, \theta'_{-i}), (\theta_i, \theta_{-i})), \quad (15)$$

while

$$u_i(f(\theta''_i, \theta_{-i}), (\theta'_i, \theta'_{-i})) > u_i(y, (\theta''_i, \theta'_{-i})), \quad \forall \theta''_i. \quad (16)$$

The proof of the subsequent robust implementation result closely follows the proof of necessity part of Theorem 2. But the existence of best responses in the definition of a nice mechanism allows the necessary conditions to be strengthened to their strict versions.

**Proposition 5 (Robust Implementation)**

If  $f$  is robustly implementable by a nice mechanism, then  $f$  satisfies strict EPIC and strict robust monotonicity.

**Proof.** See Appendix. ■

With the restriction to nice mechanisms, the relationship between iterative deletion and robust implementation emerges more directly. Finally, we should mention that we do not have general sufficient conditions for robust implementation by nice mechanisms, just instances where robust implementation is possible in the direct mechanism and examples where it is possible in nice augmented mechanisms (e.g., the coordination example of Section 3.2).

**9.3 Extensions**

The previous sections examined the importance of our assumptions about lotteries over outcomes and restrictions on mechanisms. We also restricted attention in our main analysis to the case of discrete but perhaps infinite payoff types  $\Theta_i$  and types  $T_i$ , although our examples and linear aggregation results dealt with compact  $\Theta_i$ .

Many of our results would extend easily to more general  $\Theta_i$  and  $T_i$ . This is true of the direct implementation analysis, the necessary conditions for robust implementation and the virtual implementation analysis. The sufficiency for robust implementation might be more delicate.

**10 Conclusion**

This paper examined the robustness of the classic implementation problem. We formalized robustness by requiring that the implementation problem remains solvable as we gradually relax common

knowledge among the agents and the designer. The weakening of common knowledge was achieved by considering large type spaces in which the private information of the individual agents becomes more prominent.

In contrast to our earlier results on truthful implementation, Bergemann and Morris (2004), robust implementation is in general a more demanding notion of implementation than ex post implementation. It remains an open question whether a systematic relationships between ex post, interim and robust implementation do arise in specific environments such as single crossing or supermodular environments. The analysis of the environment with interdependent values and linear aggregation in Section 8 clearly suggests that a systematic relationship can be established for many interesting environments. We also extended the robustness argument to the notion of virtual implementation. While your analysis here was preliminary, it clearly offered evidence that the distance between interim and virtual implementation may shrink considerably after imposing robustness on the implementation concept.

## 11 Appendix

### 11.1 Virtual Implementation in the Single Unit Auction

We complete the iterative implementation argument for the single unit auction example of Section 3.4. We study the following symmetric  $\varepsilon$ -efficient allocation rule is the following:

$$x_i^*(\theta) = \varepsilon \frac{1}{I} \theta_i + (1 - \varepsilon) x_i^{**}(\theta).$$

The corresponding essentially unique ex post transfer rule is:

$$y_i(\theta) = \varepsilon \frac{1}{I} \gamma \left( \sum_{j \neq i} \theta_j \right) \theta_i + \frac{1}{2I} \varepsilon (\theta_i)^2 + (1 - \varepsilon) \left( \max_{j \neq i} \theta_j + \gamma \sum_{j \neq i} \theta_j \right) x_i^{**}(\theta).$$

Thus if the true type profile is  $\theta$  and agents report themselves to be type profile  $\theta'$ , agent  $i$ 's expected utility is

$$\left( \theta_i + \gamma \sum_{j \neq i} \theta_j \right) x_i^*(\theta') - y_i(\theta'),$$

or

$$\begin{aligned} & \left( \theta_i + \gamma \sum_{j \neq i} \theta_j \right) \left( \varepsilon \frac{1}{I} \theta'_i + (1 - \varepsilon) x_i^{**}(\theta') \right) \\ & - \varepsilon \frac{1}{I} \gamma \left( \sum_{j \neq i} \theta'_j \right) \theta'_i - \frac{1}{2I} \varepsilon (\theta'_i)^2 - (1 - \varepsilon) \left( \max_{j \neq i} \theta'_j + \gamma \sum_{j \neq i} \theta'_j \right) x_i^{**}(\theta') \end{aligned}$$

or

$$\varepsilon \frac{1}{I} \theta'_i \left( \theta_i - \frac{1}{2} \theta'_i + \gamma \sum_{j \neq i} (\theta_j - \theta'_j) \right) + (1 - \varepsilon) \left( \left( \theta_i - \max_{j \neq i} \theta'_j \right) + \gamma \sum_{j \neq i} (\theta_j - \theta'_j) \right) x_i^{**}(\theta').$$

The payoff gain to agent  $i$  of reporting himself to be type  $\theta_i + x$  when his true type is  $\theta_i$  is

$$\begin{aligned} & \varepsilon \frac{1}{I} \left( (\theta_i + x) \left( \theta_i - \frac{1}{2} (\theta_i + x) + \gamma \sum_{j \neq i} (\theta_j - \theta'_j) \right) - \theta_i \left( \frac{1}{2} \theta_i + \gamma \sum_{j \neq i} (\theta_j - \theta'_j) \right) \right) \\ & + (1 - \varepsilon) \left( \begin{aligned} & \left( \left( \theta_i - \max_{j \neq i} \theta'_j \right) + \gamma \sum_{j \neq i} (\theta_j - \theta'_j) \right) x_i^{**}(\theta_i + x, \theta'_{-i}) \\ & - \left( \left( \theta_i - \max_{j \neq i} \theta'_j \right) + \gamma \sum_{j \neq i} (\theta_j - \theta'_j) \right) x_i^{**}(\theta_i, \theta'_{-i}) \end{aligned} \right), \end{aligned}$$

which is equal to

$$\begin{aligned} & \varepsilon \frac{1}{I} \left( x \gamma \sum_{j \neq i} (\theta_j - \theta'_j) - \frac{1}{2} x^2 \right) \\ & + (1 - \varepsilon) \left( \left( \theta_i - \max_{j \neq i} \theta'_j \right) + \gamma \sum_{j \neq i} (\theta_j - \theta'_j) \right) (x_i^{**}(\theta_i + x, \theta'_{-i}) - x_i^{**}(\theta_i, \theta'_{-i})). \end{aligned}$$

Now, the first term is maximized by setting

$$x = \gamma \sum_{j \neq i} (\theta_j - \theta'_j).$$

and the second term is maximized by choosing  $x > \max_{j \neq i} \theta'_j - \theta_i$  if

$$\theta_i > \max_{j \neq i} \theta'_j - \gamma \sum_{j \neq i} (\theta_j - \theta'_j);$$

anything if

$$\theta_i = \max_{j \neq i} \theta'_j - \gamma \sum_{j \neq i} (\theta_j - \theta'_j);$$

and  $x < \max_{j \neq i} \theta'_j - \theta_i$  if

$$\theta_i < \max_{j \neq i} \theta'_j - \gamma \sum_{j \neq i} (\theta_j - \theta'_j).$$

Thus the whole expression is maximized setting

$$x = \gamma \sum_{j \neq i} (\theta_j - \theta'_j).$$

This is exactly the same best response property as we obtained in the public good game. Therefore we get robust implementation in the direct mechanism if  $\gamma < \frac{1}{I-1}$  as in the example of the public good in Section 3.3.

## 11.2 A Badly Behaved Mechanism

The example illustrates the gap between the necessary and sufficient conditions in Proposition 2. Specifically, it shows that there can be an equilibrium for every type space  $\mathcal{T}$  in a mechanism, yet  $S^{\mathcal{M}}$  does not satisfy the ex post best response property.

In the example, there are two agents and there is complete information, so each agent has a unique type. There are a finite number of outcomes  $Z = \{a, b, c\}$ . The payoffs are given by the

following table:

	$a$	$b$	$c$
agent 1	-1	0	+1
agent 2	0	0	0

The planner's choice (in the unique payoff state) is  $a$ . Thus it is trivial to robustly implement the social choice function. But suppose that the planner chooses the following (strange) mechanism:  $M_1 = \{1, 2, 3, \dots\}$ ,  $M_2 = \{1, 2\}$  and

$$g(m_1, m_2) = \begin{cases} a, & \text{if } m_1 = 1 \\ b, & \text{if } m_1 > 1 \text{ and } m_2 = 1 \\ \left[\frac{1}{m_1}, b; \left(1 - \frac{1}{m_1}\right), c\right], & \text{if } m_1 > 1 \text{ and } m_2 = 2 \end{cases}$$

where  $\left[\frac{1}{m_1}, b; \left(1 - \frac{1}{m_1}\right), c\right]$  is the lottery putting probability  $\frac{1}{m_1}$  on  $b$  and probability  $\left(1 - \frac{1}{m_1}\right)$  on  $c$ . Thus  $g(m_1, m_2)$  can be represented by the following table:

$g$	1	2
1	$a$	$a$
2	$b$	$\left(\frac{1}{2}, b; \frac{1}{2}, c\right)$
3	$b$	$\left(\frac{1}{3}, b; \frac{2}{3}, c\right)$
$\vdots$	$\vdots$	$\vdots$
$k$	$b$	$\left(\frac{1}{k}, b; 1 - \frac{1}{k}, c\right)$
$\vdots$	$\vdots$	$\vdots$

Thus the agents are playing the following complete information game:

$m_1/m_2$	1	2
1	0, 0	0, 0
2	-1, 0	$\frac{1}{2}, 0$
3	-1, 0	$\frac{2}{3}, 0$
$\vdots$	$\vdots$	$\vdots$
$k$	-1, 0	$1 - \frac{1}{k}, 0$
$\vdots$	$\vdots$	$\vdots$

Now on any type space, there is always an equilibrium where player 1 chooses action 1 and player 2 chooses action 1, and outcome  $a$  is chosen. Moreover, on any type space, in any equilibrium, outcome  $a$  is always chosen: if player 1 ever has a best response not to play 1 then he has no best

response. So he always plays 1 in equilibrium. Thus the trivial social choice function is robustly implemented by this mechanism.

While only message 1 survives iterated deletion of never best responses for player 1, both messages survive iterated deletion of never best responses for player 2. Thus we have  $S_1^{\mathcal{M}} = \{1\}$  and  $S_2^{\mathcal{M}} = \{1, 2\}$ . Note that  $S^{\mathcal{M}}$  satisfies the interim best response property, see Definition 6, but not the ex post best response property, see Definition 5. For we observe that

$$u_1(g(1, 2)) = u_1(a) = 0 < \frac{1}{2} = u_1(g(2, 2)),$$

violating the ex post best response property.

The insight of the example is that the quantifier “for every type space  $\mathcal{T}$ ” does not necessarily guarantee that all actions which will be chosen with positive probability in some equilibrium and for some type space, will also be chosen with probability one in some equilibrium for some type space. For this reason, the quantifier “for every type space  $\mathcal{T}$ ” does not allow us to establish a local, i.e. ex post best response property of every action in  $S^{\mathcal{M}}$ .

### 11.3 Bayesian Monotonicity

The next proposition establishes the equivalence between robust monotonicity and Bayesian monotonicity on every type space by means of a constructive proof (via a specific type space). The constructive element is the identification of a type space on which Bayesian monotonicity is guaranteed to fail if robust monotonicity fails. It is worthwhile to note that the specific type space is much smaller than the universal type.

In some sense, the notion of robustness is more subtle in the context of full rather than partial implementation. With partial implementation, i.e. truth-telling in the direct mechanism, the universal type space is by definition the most difficult type space to obtain truth-telling. In the universal type space, every agent has the maximal number of possible misreports and hence the designer faces the maximal number of incentive constraints. In the context of full implementation, the trade-off is ambiguous. As a larger type space contains by definition more types, it offers every agent more possibilities to misreport. But then, just as a larger type space made truth-telling more difficult to obtain, the other equilibria might also cease to exist after the introduction of additional types. This second part offers the possibility that larger type spaces facilitate rather than complicate the full implementation problem.

#### Proposition 6 (Equivalence)

*Social choice function  $f$  satisfies Bayesian monotonicity on every type space if and only if it satisfies robust monotonicity.*

**Proof.** ( $\Rightarrow$ ) We will show that if robust monotonicity fails, we can construct a type space where Bayesian monotonicity fails. The argument will be constructive.

Fix an unacceptable deception  $\beta$ . Suppose that robust monotonicity fails. Then for each  $i$ ,  $\theta_i$ ,  $\theta'_i \in \beta_i(\theta_i)$ , there exist

$$\theta_{-i}[\theta_i, \theta'_i] \in \Theta_{-i} \quad \text{and} \quad \psi_i[\theta_i, \theta'_i] \in \Delta(\beta_{-i}^{-1}(\theta_{-i}[\theta_i, \theta'_i])) \quad (17)$$

such that:

$$u_i(f(\theta''_i, \theta_{-i}[\theta_i, \theta'_i]), (\theta''_i, \theta_{-i}[\theta_i, \theta'_i])) \geq u_i(y, (\theta''_i, \theta_{-i}[\theta_i, \theta'_i])), \quad \forall \theta''_i \in \Theta_i \quad (18)$$

implies

$$\sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i})[\theta_i, \theta'_i] u_i(f(\theta'_i, \theta_{-i}[\theta_i, \theta'_i]), (\theta_i, \theta_{-i})) \geq \sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i})[\theta_i, \theta'_i] u_i(y, (\theta_i, \theta_{-i})). \quad (19)$$

Now we construct a type space around  $\theta_i, \theta'_i$  and  $\psi_i[\theta_i, \theta'_i]$  given by (17) for which Bayesian monotonicity fails. First, agent  $i$  has a set of "deception" types  $T_i^1$  which are isomorphic to  $\Xi_i = \{(\theta_i, \theta'_i) : \theta_i \in \Theta_i \text{ and } \theta'_i \in \beta_i(\theta_i)\}$ ; thus there exists a bijection  $\xi_i^1 : T_i^1 \rightarrow \Xi_i$ . The type responding to  $(\theta_i, \theta'_i)$  has payoff type  $\theta_i$  and believes that the other agents are of type:

$$\left( [\xi_j^1]^{-1}(\theta_j, \theta_{ij}[\theta_i, \theta'_i]) \right)_{j \neq i}$$

with probability  $\psi_i(\theta_{-i})[\theta_i, \theta'_i]$ . Second, agent  $i$  has a set of "pseudo-complete information types"  $T_i^2$ , which are isomorphic to  $\Theta$ ; thus there exists a bijection  $\xi_i^2 : T_i^2 \rightarrow \Theta$ . The type corresponding to  $\theta$  has payoff type  $\theta_i$  and he is convinced that each other agent  $j$  is type  $[\xi_j^1]^{-1}(\theta_j, \theta_j)$ .

Slightly more formally, we have

$$T_i = T_i^1 \cup T_i^2.$$

If  $t_i \in T_i^1$  and  $\xi_i^1(t_i) = (\theta_i, \theta'_i)$ , then

$$\widehat{\theta}_i(t_i) = \theta_i;$$

if  $t_i \in T_i^2$  and  $\xi_i^2(t_i) = \theta$ , then

$$\widehat{\theta}_i(t_i) = \theta_i.$$

If  $t_i \in T_i^1$  and  $\xi_i^1(t_i) = (\theta_i, \theta'_i)$ , then

$$\pi_i^*(t_{-i})[t_i] = \begin{cases} \psi_i(\theta_{-i})[\theta_i, \theta'_i], & \text{if } t_{-i} \in T_{-i}^1 \text{ and } \theta_{-i} = \left( [\xi_j^1]^{-1}(\theta_j, \theta_{ij}[\theta_i, \theta'_i]) \right)_{j \neq i} \\ 0, & \text{if otherwise} \end{cases}$$



If  $t_i \in T_i^2$  and  $\xi_i^2(t_i) = \theta$ , then

$$\pi_i^*(t_{-i})[t_i] = \begin{cases} 1, & \text{if } t_{-i} \in T_{-i}^1 \text{ and } \theta_{-i} = \left( [\xi_j^1]^{-1}(\theta_j, \theta_{ij}[\theta_i, \theta'_i]) \right)_{j \neq i} \\ 0, & \text{if } \text{otherwise} \end{cases}$$

Now consider the Bayesian deception on this type space where each type  $[\xi_i^1]^{-1}(\theta_i, \theta'_i)$  reports himself to be type  $[\xi_i^1]^{-1}(\theta'_i, \theta'_i)$ , and all other types report their types truthfully. Thus

$$\alpha_i(t_i) = \begin{cases} [\xi_i^1]^{-1}(\theta'_i, \theta'_i), & \text{if } t_i = [\xi_i^1]^{-1}(\theta_i, \theta'_i) \\ t_i, & \text{if } \text{otherwise} \end{cases}.$$

Since  $\beta$  was unacceptable, we must have that  $f(\widehat{\theta}(t)) \neq f(\widehat{\theta}(\alpha(t)))$  for some  $t$ . Thus the Bayesian monotonicity condition (Definition 24) for this type space requires that there exist  $i$ ,  $t_i$  and  $h : T \rightarrow Z$  such that

$$\sum_{t_{-i} \in T_{-i}} u_i(h(\alpha(t)), \widehat{\theta}(t)) \widehat{\pi}_i(t_{-i})[t_i] > \sum_{t_{-i} \in T_{-i}} u_i(f(\widehat{\theta}(\alpha(t))), \widehat{\theta}(t)) \widehat{\pi}_i(t_{-i})[t_i], \quad (20)$$

and

$$\begin{aligned} & \sum_{t_{-i} \in T_{-i}} u_i(f(\widehat{\theta}(t''_i, t_{-i})), \widehat{\theta}(t''_i, t_{-i})) \widehat{\pi}_i(t_{-i})[t''_i] \\ & \geq \sum_{t_{-i} \in T_{-i}} u_i(h(\alpha_i(t_i), t_{-i}), \widehat{\theta}(t''_i, t_{-i})) \widehat{\pi}_i(t_{-i})[t''_i], \quad \forall t''_i. \end{aligned} \quad (21)$$

The  $t_i$  cannot be an element of  $T_i^2$ , because such a type does not expect any deviation from truth-telling under the deception. So it must be an element of  $T_i^1$ , with  $\xi_i^1(t_i) = (\theta_i, \theta'_i)$ . Now condition (20) becomes

$$\begin{aligned} & \sum_{\theta_{-i} \in \Theta_{-i}} u_i \left( h \left( [\xi_i^1]^{-1}(\theta'_i, \theta'_i), \left( \left( [\xi_j^1]^{-1}(\theta_{ij}[\theta_i, \theta'_i], \theta_{ij}[\theta_i, \theta'_i]) \right)_{j \neq i} \right), (\theta_i, \theta_{-i}) \right) \psi_i(\theta_{-i})[\theta_i, \theta'_i] \right) \\ & > \sum_{\theta_{-i} \in \Theta_{-i}} u_i(f(\theta'_i, \theta_{-i}[\theta_i, \theta'_i]), (\theta_i, \theta_{-i})) \psi_i(\theta_{-i})[\theta_i, \theta'_i]. \end{aligned} \quad (22)$$

But letting  $t''_i$  in condition (21) be in  $T_i^2$  with  $\xi_i^2(t''_i) = (\theta''_i, \theta_{-i}[\theta_i, \theta'_i])$ , we have

$$\begin{aligned} & u_i(f(\theta''_i, \theta_{-i}[\theta_i, \theta'_i]), (\theta''_i, \theta_{-i}[\theta_i, \theta'_i])) \\ & \geq u_i \left( h \left( [\xi_i^1]^{-1}(\theta'_i, \theta'_i), \left( \left( [\xi_j^1]^{-1}(\theta_{ij}[\theta_i, \theta'_i], \theta_{ij}[\theta_i, \theta'_i]) \right)_{j \neq i} \right), (\theta''_i, \theta_{-i}[\theta_i, \theta'_i]) \right) \right) \end{aligned} \quad (23)$$

for all  $\theta''_i$ . Setting

$$z = h \left( [\xi_i^1]^{-1}(\theta'_i, \theta'_i), \left( \left( [\xi_j^1]^{-1}(\theta_{ij}[\theta_i, \theta'_i], \theta_{ij}[\theta_i, \theta'_i]) \right)_{j \neq i} \right) \right),$$

condition (22) becomes

$$\begin{aligned} & \sum_{\theta_{-i} \in \Theta_{-i}} u_i(z, (\theta_i, \theta_{-i})) \psi_i(\theta_{-i}) [\theta_i, \theta'_i] \\ & > \sum_{\theta_{-i} \in \Theta_{-i}} u_i(f(\theta'_i, \theta_{-i} [\theta_i, \theta'_i]), (\theta_i, \theta_{-i})) \psi_i(\theta_{-i}) [\theta_i, \theta'_i]. \end{aligned}$$

while condition (23) requires  $z \in Y_i(\theta_{-i} [\theta_i, \theta'_i])$ . But these latter claims contradict our initial assumption that robust monotonicity fails (i.e., (18)). Thus Bayesian monotonicity fails for this type space and the claim is proved.

( $\Leftarrow$ ) Suppose  $f$  satisfies robust monotonicity. Fix any type space  $\mathcal{T}$  and any deception  $\alpha$  with  $f(\widehat{\theta}(t)) \neq f(\widehat{\theta}(\alpha(t)))$  for some  $t$ . Define  $\beta$  by

$$\beta_i(\theta_i) = \left\{ \theta'_i : \exists t_i \text{ such that } \widehat{\theta}_i(t_i) = \theta_i \text{ and } \widehat{\theta}_i(\alpha_i(t_i)) = \theta'_i \right\}.$$

Deception  $\beta$  is unacceptable, so by robust monotonicity, there exist  $i, \theta_i, \theta'_i \in \beta_i(\theta_i)$  such that for every  $\theta'_{-i} \in \Theta_{-i}$  and  $\psi_i \in \Delta(\beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $y[\theta'_{-i}, \psi_i] \in Y_i(\theta'_{-i})$  such that

$$\sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i}) u_i(y[\theta'_{-i}, \psi_i], (\theta_i, \theta_{-i})) > \sum_{\theta_{-i} \in \Theta_{-i}} \psi_i(\theta_{-i}) u_i(f(\theta'_i, \theta'_{-i}), (\theta_i, \theta_{-i})). \quad (24)$$

Now choose any  $t_i$  such that  $\widehat{\theta}_i(t_i) = \theta_i$  and  $\widehat{\theta}_i(\alpha_i(t_i)) = \theta'_i$ . For every (mis-)report  $\theta'_{-i}$ , we now derive a distribution over payoff types  $\theta_{-i}$  which represents the likelihood that the report  $\theta'_{-i}$  comes from the true payoff type profile  $\theta_{-i}$ , given the type space  $\mathcal{T}$ . For each  $\theta'_{-i}$ , define  $\psi_i[\theta'_{-i}] \in \Delta(\Theta_{-i})$  by

$$\psi_i(\theta_{-i}) [\theta'_{-i}] \triangleq \frac{\sum_{\{t_{-i}: \widehat{\theta}_j(\alpha_j(t_j)) = \theta'_j \text{ and } \widehat{\theta}_j(t_j) = \theta_j, \forall j \neq i\}} \widehat{\pi}_i(t_{-i}) [t_i]}{\sum_{\{t_{-i}: \widehat{\theta}_j(\alpha_j(t_j)) = \theta'_j, \forall j \neq i\}} \widehat{\pi}_i(t_{-i}) [t_i]}. \quad (25)$$

Now let  $h$  satisfy

$$h(t'_i, t_{-i}) \triangleq \begin{cases} y[\widehat{\theta}_{-i}(t_{-i}), \psi_i[\widehat{\theta}_{-i}(t_{-i})]] & \text{if } t'_i = \alpha_i(t_i) \\ f(\widehat{\theta}(t'_i, t_{-i})) & \text{if otherwise} \end{cases}. \quad (26)$$

To establish Bayesian monotonicity, it is enough to show that the two inequalities of Bayesian monotonicity are satisfied, or:

$$\sum_{t_{-i}} u_i(h(\alpha(t)), \widehat{\theta}(t)) \widehat{\pi}_i(t_{-i}) [t_i] > \sum_{t_{-i}} u_i(f(\widehat{\theta}(\alpha(t))), \widehat{\theta}(t)) \widehat{\pi}_i(t_{-i}) [t_i], \quad (27)$$

and

$$\begin{aligned} & \sum_{t_{-i}} u_i(f(\widehat{\theta}(t'_i, t_{-i})), \widehat{\theta}(t'_i, t_{-i})) \widehat{\pi}_i(t_{-i}) [t'_i] \\ & \geq \sum_{t_{-i}} u_i(h(\alpha_i(t_i), t_{-i}), \widehat{\theta}(t'_i, t_{-i})) \widehat{\pi}_i(t_{-i}) [t_i], \quad \forall t'_i. \end{aligned} \quad (28)$$

By inserting the posterior beliefs  $\psi_i$  and the rewards  $h(t'_i, t_{-i})$ , as defined above in (25) and (26) respectively, we can rewrite the two sides of the inequality (27) as follows:

$$\begin{aligned} & \sum_{t_{-i}} u_i \left( h(\alpha(t)), \widehat{\theta}(t) \right) \widehat{\pi}_i(t_{-i}) [t_i] \\ = & \sum_{\theta'_{-i}} \left( \sum_{\{t_{-i}: \widehat{\theta}_j(\alpha_j(t_j)) = \theta'_j, \forall j \neq i\}} \widehat{\pi}_i(t_{-i}) [t_i] \right) \sum_{\theta_{-i}} \psi_i(\theta_{-i}) [\theta'_{-i}] u_i(y[\theta'_{-i}, \psi_i[\theta'_{-i}]], \theta) \end{aligned}$$

and

$$\begin{aligned} & \sum_{t_{-i}} u_i \left( f(\widehat{\theta}(\alpha(t))), \widehat{\theta}(t) \right) \widehat{\pi}_i(t_i) [t_{-i}] \\ = & \sum_{\theta'_{-i}} \left( \sum_{\{t_{-i}: \widehat{\theta}_j(\alpha_j(t_j)) = \theta'_j, \forall j \neq i\}} \widehat{\pi}_i(t_{-i}) [t_i] \right) \sum_{\theta_{-i}} \psi_i(\theta_{-i}) [\theta'_{-i}] u_i(f(\theta'), \theta) \end{aligned}$$

so (27) follows from (24). Also

$$\begin{aligned} & \sum_{t_{-i}} u_i \left( h(\alpha_i(t_i), t_{-i}), \widehat{\theta}(t'_i, t_{-i}) \right) \widehat{\pi}_i(t_{-i}) [t'_i] \\ = & \begin{cases} \sum_{t_{-i}} u_i \left( y \left[ \widehat{\theta}_{-i}(t_{-i}), \psi_i \left[ \widehat{\theta}_{-i}(t_{-i}) \right] \right], \widehat{\theta}(t'_i, t_{-i}) \right) \widehat{\pi}_i(t_{-i}) [t'_i] & \text{if } t'_i = \alpha_i(t_i) \\ \sum_{t_{-i}} u_i \left( f \left( \widehat{\theta}(t'_i, t_{-i}), \widehat{\theta}(t'_i, t_{-i}) \right) \right) \widehat{\pi}_i(t_{-i}) [t'_i], & \text{if } t'_i \neq \alpha_i(t_i) \end{cases} \end{aligned}$$

Now  $y \left[ \widehat{\theta}_{-i}(t_{-i}), \psi_i \left[ \widehat{\theta}_{-i}(t_{-i}) \right] \right] \in Y_i \left( \widehat{\theta}_{-i}(t_{-i}) \right)$  implies (28). ■

The proof may appear rather intricate in its details. We next give a brief outline of the basic steps to show how interim implies robust monotonicity. The proof proceeds by contrapositive. We start with an unacceptable deception  $\beta$  which by hypothesis fails robust monotonicity and hence satisfies the inequalities (18) and (19). For the given deception  $\beta$ , we then create a type space, consisting of two components for every agent  $i$ . The first component for agent  $i$  is created by the set of pairs of payoff types  $(\theta_i, \theta'_i)$ , where the first entry is the true payoff type and the second entry is a feasible deception (under  $\beta$ ), or  $\theta'_i \in \beta_i(\theta_i)$ . For this reason, we refer to these types as “deception types.” For every such pair  $(\theta_i, \theta'_i)$  there exists at least one particular payoff profile  $\theta'_{-i}$  which acts as a misreport. Under the deception  $\beta$ , this payoff profile  $\theta'_{-i}$  could have been reported by all true payoff profiles which are in the support of  $\psi_i$ . Consequently, the belief component of type  $(\theta_i, \theta'_i)$  is given by simply adopting  $\psi_i(\cdot | \theta_i, \theta'_i)$ . The second component consists of “pseudo complete information types”, described by  $t_i = \theta \in \Theta$ . Each such type has a belief that assigns probability one to the event that the true payoff profile is given by  $\theta$  and that all other agents report the deception type  $(\theta_j, \theta_j)$ , and hence the “pseudo” in the labelling.

Given this type space  $T_i$ , we then consider a particular deception  $\alpha_i : T_i \rightarrow T_i$ . The deception  $\alpha_i$  is localized around the “deception types” and the “pseudo complete information types” report truthfully. The deception  $\alpha_i$  consists of agent  $i$  always reporting his deception type rather than his true type, or  $\alpha_i(\theta_i, \theta'_i) = (\theta'_i, \theta'_i)$ . We then verify whether  $f$  is interim monotone under  $\alpha$ . The existence of the pseudo complete information types  $\theta$  forces the interim incentive compatibility conditions to reduce to ex post incentive compatibility conditions. This guarantees the hypothesis in the robust monotonicity notion, namely inequality (18), and thus leads to the conclusion in form of the inequalities (19). But then we obtain a contradiction to the reward condition of interim monotonicity, unless the hypothesis for the interim monotonicity condition, namely  $f \neq f \circ \alpha$ , is not satisfied, i.e.  $f = f \circ \alpha$  holds, but of course this implies that  $\beta$  is acceptable.

### 11.4 Coordination Example 2

The next example is the pure coordination game, which we first considered in Section 3.2, without the additional allocations,  $z$  and  $z'$ . It illustrates the importance of lotteries for robust implementation. The example will satisfy EPIC and robust monotonicity, yet it cannot be robustly implemented without the use of lotteries. On the other hand if lotteries are allowed then the lottery which selects each of the four possible outcomes with equal probability constitutes a bad outcome, and hence the sufficient conditions for robust implementation would be satisfied with lotteries.

The examples has two agents,  $i = 1, 2$  and each agent  $i$  has two possible types,  $\theta_i$  and  $\theta'_i$ . There are four possible outcomes:  $Z = \{a, b, c, d\}$ . Agents' payoffs are given by:

<b>a</b>	$\theta_2$	$\theta'_2$
$\theta_1$	3, 3	0, 0
$\theta'_1$	0, 0	1, 1

<b>b</b>	$\theta_2$	$\theta'_2$
$\theta_1$	0, 0	3, 3
$\theta'_1$	1, 1	0, 0

<b>c</b>	$\theta_2$	$\theta'_2$
$\theta_1$	0, 0	1, 1
$\theta'_1$	3, 3	0, 0

<b>d</b>	$\theta_2$	$\theta'_2$
$\theta_1$	1, 1	0, 0
$\theta'_1$	0, 0	3, 3

The social choice function  $f$  selects the efficient outcome in every state:

<b>f</b>	$\theta_2$	$\theta'_2$
$\theta_1$	$a$	$b$
$\theta'_1$	$c$	$d$

As in the example in Section 3.2, the social choice function is strictly ex post incentive compatible but there is another equilibrium in the "direct mechanism" where each agent misreports his type, and each agent gets a payoff of 1.

Robust monotonicity is clearly satisfied even if the rewards  $y$  are restricted to be the deterministic allocations  $Z$ . We will show that robust implementation is not possible even in an infinite

mechanism if we restrict attention to deterministic mechanisms. Fix a mechanism  $\mathcal{M}$ . Let

$$S_i^*(\theta_i) = \{m_i : g(m_i, m_j) = f(\theta_i, \theta_j) \text{ for some } m_j, \theta_j\},$$

be the set of messages for agent  $i$  which would select the allocation recommended by the social choice function for some  $m_j, \theta_j$ . Let  $S_i^k(\theta_i) = b_i(b^{k-1}(\bar{S}))[\theta_i]$  (using transfinite induction if necessary). We now show by induction that,  $S_i^*(\theta_i) \subseteq S_i^k(\theta_i)$  for all  $k$  using the structure of the payoffs. Suppose that this is true for  $k$ . Then for any  $m_i \in S_i^*(\theta_i) \subseteq S_i^k(\theta_i)$ , there exists  $m_j \in S_j^*(\theta_j) \subseteq S_j^k(\theta_j)$  such that  $g(m_i, m_j) = f(\theta_i, \theta_j)$ . Thus there does not exist  $\nu_i \in \Delta(M_i)$  such that

$$\sum_{m'_i} \nu_i(m'_i) u_i(g(m'_i, m_j), (\theta_i, \theta_j)) > u_i(g(m_i, m_j), (\theta_i, \theta_j)) = 3.$$

So  $m_i \in S_i^{k+1}(\theta_i)$ .

Thus we must have that  $(m_1, m_2) \in S_1^*(\theta_1) \times S_2^*(\theta_2)$  implies  $g(m_1, m_2) = f(\theta_1, \theta_2)$ . Let  $m_i^*(\cdot)$  be any selection from  $S_i^*(\cdot)$ . Now let  $k^*$  be the lowest  $k$  such that, for some  $i$ ,

$$m_i^*(\theta'_i) \notin S_i^k(\theta_i).$$

Without loss of generality, let  $i = 1$ . Note  $m_2^*(\theta'_2) \in S_2^{k-1}(\theta_2)$  by definition of  $k^*$ . If agent 1 was type  $\theta_1$  and was sure his opponent were type  $\theta_2$  and choosing action  $m_2^*(\theta'_2)$ , we know that he could guarantee himself a payoff of 1 by choosing  $m_1^*(\theta'_1)$ . Since  $m_1^*(\theta'_1)$  is deleted for type  $\theta_1$  at round  $k^*$ , we know that there exists  $\nu_1 \in \Delta(M_1)$  such that

$$\sum_{m'_1} \nu_1(m'_1) g_1(m'_1, m_2^*(\theta'_2)) > 1,$$

and thus there exists  $m'_1$  such that  $g_1(m'_1, m_2^*(\theta'_2)) = f(\theta_1, \theta_2)$ . This implies that  $m_2^*(\theta'_2) \in S_2^*(\theta_2)$ , a contradiction.

The example uses the fact that the social choice function always selects an outcome that is strictly Pareto-optimal and - paradoxically - it is this feature which inhibits iterative implementation in the current example. Borgers (1995) proves the impossibility of complete information implementation of non-dictatorial social choice functions in iteratively undominated strategies when the set of feasible preference profiles includes such unanimous preference profiles and the argument here is reminiscent of Borgers' argument.

## 11.5 Nice Mechanisms and Strict Robust Monotonicity

**Proof of Proposition 5.** The restriction to nice mechanisms ensures that  $S^{\mathcal{M}}$  is non-empty. It follows that if mechanism  $\mathcal{M}$  iteratively implements  $f$ , then, for each  $i$ , there exists  $m_i^* : \Theta_i \rightarrow M_i$

such that

$$g(m^*(\theta)) = f(\theta) \text{ and } m^*(\theta) \in S^{\mathcal{M}}(\theta),$$

we can simply let  $m_i^*(\theta_i)$  be any element of  $S_i^{\mathcal{M}}(\theta_i)$ .

We first establish strict EPIC. Suppose strict EPIC fails. Then there exists  $i$ ,  $\theta$  and  $\theta'_i$  such that  $f(\theta'_i, \theta_{-i}) \neq f(\theta_i, \theta_{-i})$  and

$$u_i(f(\theta'_i, \theta_{-i}), \theta) \geq u_i(f(\theta), \theta).$$

Now, for any message  $m_i$  with

$$m_i \in \arg \max_{m'_i} u_i(g(m'_i, m_{-i}^*(\theta_{-i})), (\theta_i, \theta_{-i})),$$

since  $m_{-i}^*(\theta_{-i}) \in S_i^\infty(\theta_{-i})$ , we must have  $m_i \in S_i^\infty(\theta_i)$  and thus  $g(m_i, m_{-i}^*(\theta'_{-i})) = f(\theta_i, \theta'_{-i})$  for all  $\theta'_{-i}$ . Thus let

$$m_i^*(\theta'_i) \in \arg \max_{m'_i} u_i(g(m'_i, m_{-i}^*(\theta_{-i})), (\theta_i, \theta_{-i})),$$

and  $f(\theta'_i, \theta'_{-i}) = g(m_i^*(\theta'_i), m_{-i}^*(\theta'_{-i})) = f(\theta_i, \theta'_{-i})$  for all  $\theta'_{-i}$ , a contradiction.

Now we establish strict robust monotonicity. Fix an unacceptable deception  $\beta$ . Let  $\widehat{k}$  be the largest  $k$  such that for every  $i$ ,  $\theta_i$  and  $\theta'_i \in \beta_i(\theta_i)$ ,

$$S_i^\infty(\theta'_i) \subseteq S_i^k(\theta_i).$$

We know that such a  $\widehat{k}$  exists because  $S_i^0(\theta_i) \cap S_i^\infty(\theta'_i) = S_i^\infty(\theta'_i)$  and, since  $\mathcal{M}$  iteratively implements  $f$ , we must have  $S_i^\infty(\theta_i) \cap S_i^\infty(\theta'_i) = \emptyset$ .

Now we know that there exists  $i$  and  $\theta'_i \in \beta_i(\theta_i)$  such that

$$S_i^{\widehat{k}+1}(\theta_i) \cap S_i^\infty(\theta'_i) \neq S_i^\infty(\theta'_i).$$

Let

$$\widehat{m}_i \in S_i^{\widehat{k}}(\theta_i) \cap S_i^\infty(\theta'_i),$$

and

$$\widehat{m}_i \notin S_i^{\widehat{k}+1}(\theta_i) \cap S_i^\infty(\theta'_i).$$

Since message  $\widehat{m}_i$  gets deleted for  $\theta_i$  at round  $\widehat{k} + 1$ , we know that for every  $\lambda_i \in \Delta(M_{-i} \times \Theta_{-i})$  such that

$$\lambda_i(m_{-i}, \theta_{-i}) > 0 \Rightarrow m_j \in S_j^{\widehat{k}}(\theta_j) \text{ for all } j \neq i,$$

there exists  $m_i^*$  such that

$$\sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(m_i^*, m_{-i}), (\theta_i, \theta_{-i})) > \sum_{m_{-i}, \theta_{-i}} \lambda_i(m_{-i}, \theta_{-i}) u_i(g(\hat{m}_i, m_{-i}), (\theta_i, \theta_{-i})).$$

Let

$$\hat{m}_j \in S_j^\infty(\theta'_j)$$

for all  $j \neq i$ . Now the above claim remains true if we restrict attention to distributions  $\lambda_i$  putting probability 1 on  $\hat{m}_{-i}$ . Thus for every  $\psi_i \in \Delta(\Theta_{-i})$  such that

$$\psi_i(\theta_{-i}) > 0 \Rightarrow \hat{m}_j \in S_j^{\hat{k}}(\theta_j) \text{ for all } j \neq i,$$

there exists  $m_i^*$  such that

$$\sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(g(m_i^*, \hat{m}_{-i}), (\theta_i, \theta_{-i})) > \sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(g(\hat{m}_i, \hat{m}_{-i}), (\theta_i, \theta_{-i})).$$

But  $\hat{m} \in S^\infty(\theta')$ , so (since  $\mathcal{M}$  iteratively implements  $f$ ),  $g(\hat{m}_i, \hat{m}_{-i}) = f(\theta')$ . Also observe that if  $\theta'_{-i} \in \beta_{-i}(\theta_{-i})$ , then  $\hat{m}_{-i} \in S_{-i}^{\hat{k}}(\theta_{-i})$ . Thus for every  $\psi_i \in \Delta(\beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $m_i^*$  such that

$$\sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(g(m_i^*, \hat{m}_{-i}), (\theta_i, \theta_{-i})) > \sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(f(\theta'), (\theta_i, \theta_{-i})),$$

which establishes the reward inequality, (15), of strict robust monotonicity.

Now suppose the incentive inequality, (16), are not satisfied strictly, and hence:

$$u_i(g(m_i^*, \hat{m}_{-i}), (\tilde{\theta}_i, \theta'_{-i})) \geq u_i(f(\tilde{\theta}_i, \theta'_{-i}), (\tilde{\theta}_i, \theta'_{-i})).$$

Now, for any

$$m_i \in \arg \max_{m'_i} u_i(g(m'_i, \hat{m}_{-i}), (\tilde{\theta}_i, \theta'_{-i})), \quad (29)$$

since  $\hat{m}_{-i} \in S_{-i}^\infty(\theta'_{-i})$ , we must have  $m_i \in S_i^\infty(\tilde{\theta}_i)$  and thus  $g(m_i, \hat{m}_{-i}) = f(\theta_i, \theta'_{-i})$ . Thus from (29) we also know that  $m_i^*$  achieves the maximum:

$$m_i^* \in \arg \max_{m'_i} u_i(g(m'_i, \hat{m}_{-i}), (\tilde{\theta}_i, \theta'_{-i}))$$

and, for all  $\tilde{\theta}_i$ , if

$$u_i(g(m_i^*, \hat{m}_{-i}), (\tilde{\theta}_i, \theta'_{-i})) \geq u_i(f(\tilde{\theta}_i, \theta'_{-i}), (\tilde{\theta}_i, \theta'_{-i})),$$

then  $g(m_i^*, \hat{m}_{-i}) = f(\tilde{\theta}_i, \theta'_{-i})$ .

Now setting  $y \equiv g(m_i^*, \hat{m}_{-i})$ , we have established that for each  $\theta'_{-i} \in \beta_{-i}(\theta_{-i})$  and  $\psi_i \in \Delta(\beta_{-i}^{-1}(\theta'_{-i}))$ , there exists  $y$  such that  $y \in Y_i^*(\theta'_{-i})$  and

$$\sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(g(m_i^*, \hat{m}_{-i}), (\theta_i, \theta_{-i})) > \sum_{\theta_{-i}} \psi_i(\theta_{-i}) u_i(f(\theta'), (\theta_i, \theta_{-i})),$$

which concludes the proof. ■



## 12 Notation

$\alpha_i(\theta_i)$  deception

$b_i(S)[\theta_i]$  never best response operator

$\beta_i(\theta_i)$  deception correspondence

$\underline{\beta}_i^k, \overline{\beta}_i^k$  lower and upper bound for best responses

$f$  social choice function

$g$  outcome function

$\gamma$  interdependence parameter

$I$  number of agents

$\lambda_i(\theta_{-i}, m_{-i})$  belief of agent  $i$

$\mu_i$  linear aggregator of payoff types

$\nu_i(m_{-i})$  belief over messages

$\psi_i(\theta_{-i})$  belief over payoff types

$\psi_i[\cdot]$  specific belief in Bayesian monotonicity

$\pi_i(t_{-i})[t_i]$  conditioning event

$\rho(\theta'_{-i}, \theta_{-i})$  belief over reports and types (in direct mechanism)

$S_i, S$  message correspondence

$S_i^k$  strategies surviving  $k$  - th round of elimination

$t_i, t$  type (profile)

$\theta_i, \theta$  payoff type

$\theta_{-i}[\cdot, \cdot]$  payoff type in Bayesian monotonicity

$u_i$  utility

$v_i$  utility with linear aggregator

$y$  allocation

$y[\cdot, \cdot]$  reward in Bayesian monotonicity

$y_i$  monetary transfer

$Y = \Delta(Z)$  extended outcome space

$Y_i(\theta_{-i})$  incentive compatible reward set

$Y_i^*(\theta_{-i})$  strictly incentive compatible reward set

$\zeta$  message

$Z$  outcome space

## References

- ABREU, D., AND H. MATSUSHIMA (1992a): “Virtual Implementation in Iterative Undominated Strategies: Complete Information,” *Econometrica*, 60, 993–1008.
- (1992b): “Virtual Implementation In Iteratively Undominated Strategies: Incomplete Information,” Discussion paper, Princeton University and University of Tokyo.
- BATTIGALLI, P. (1999): “Rationalizability in Incomplete Information Games,” Discussion Paper ECO 99/17, European University Institute.
- BATTIGALLI, P., AND M. SINISCALCHI (2003): “Rationalization and Incomplete Information,” *Advances in Theoretical Economics*, 3, Article 3.
- BENOIT, J., AND E. OK (2004): “Nash Implementation Without No Veto,” Discussion paper, New York University.
- BERGEMANN, D., AND S. MORRIS (2004): “Robust Mechanism Design,” *Econometrica*, forthcoming.
- (2005): “Ex Post Implementation,” Discussion Paper 1502, Cowles Foundation for Research in Economics, Yale University.
- BORGERS, T. (1995): “A Note on Implementation and Strong Dominance,” in *Social Choice, Welfare and Ethics*, ed. by W. Barnett, H. Moulin, M. Salles, and N. Schofield. Cambridge University Press, Cambridge.
- BRANDENBURGER, A., AND E. DEKEL (1987): “Rationalizability and Correlated Equilibria,” *Econometrica*, 55, 1391–1402.
- CHUNG, K.-S., AND J. C. ELY (2001): “Efficient and Dominance Solvable Auctions with Interdependent Valuations,” Discussion paper, Northwestern University.
- DEKEL, E., D. FUDENBERG, AND S. MORRIS (2005): “Interim Rationalizability,” Discussion paper, Tel-Aviv University, Harvard University and Yale University.
- DUGGAN, J. (1997): “Virtual Bayesian Implementation,” *Econometrica*, 65, 1175–1199.
- HUGHES, G., AND M. CRESWELL (1996): *A New Introduction Into Modal Logic*. Routledge, London.
- JACKSON, M. (1991): “Bayesian Implementation,” *Econometrica*, 59, 461–477.

- JACKSON, M. (1992): “Implementation in Undominated Strategies: A Look at Bounded Mechanisms,” *Review of Economic Studies*, 59, 757–775.
- LIPMAN, B. (1994): “A Note on the Implications of Common Knowledge of Rationality,” *Games and Economic Behavior*, 6, 114–129.
- PALFREY, T., AND S. SRIVASTAVA (1989): “Mechanism Design with Incomplete Information: A Solution to the Implementation Problem,” *Journal of Political Economy*, 97, 668–691.
- POSTLEWAITE, A., AND D. SCHMEIDLER (1986): “Implementation in Differential Information Economies,” *Journal of Economic Theory*, 39, 14–33.
- SERRANO, R., AND R. VOHRA (2005): “A Characterization of Virtual Bayesian Implementation,” *Games and Economic Behavior*, 50, 312–331.
- WILSON, R. (1987): “Game-Theoretic Analyses of Trading Processes,” in *Advances in Economic Theory: Fifth World Congress*, ed. by T. Bewley, pp. 33–70, Cambridge. Cambridge University Press.