

# Mechanism design without commitment<sup>\* †</sup>

Hannu Vartiainen<sup>‡</sup>  
HECER and University of Helsinki

Jan 2014

## Abstract

This paper identifies mechanisms that are implementable even when the planner cannot commit to the rules of the mechanism. The standard approach is to require mechanism to be robust against redesign. This often leads nonexistence of acceptable mechanisms. The novelty of this paper to require robustness against redesigns that are themselves robust against redesigns that are themselves robust against... . That is, we allow the planner to costlessly redesign the mechanism any number of times, and identify redesign strategies that are both optimal and dynamically consistent. A mechanism design strategy that credibly implements a direct mechanism after all histories is shown to exist. The framework is applied to bilateral bargaining situations. We demonstrate that a welfare maximizing second best mechanism can be implemented even without commitment.

*Keywords:* mechanisms, commitment, consistency, optimality, bilateral bargaining.

*JEL:* C72, D44, D78.

## 1 Introduction

Mechanism design theory provides powerful tools for the planner to implement desired outcomes in collective choice situations with incomplete information. However, the theory relies on an assumption that is both limiting and, at times, unreasonable: that the planner herself can commit to the mechanism. This assumption is crucial since the incentive compatibility of the mechanism requires that it is played as planned. What exacerbates the problem is that the optimality properties of the mechanism may change when information is being revealed during the play. Hence, given ex post information, another continuation mechanism may begin to dominate the original mechanism and the planner is tempted to change the rules of the game.

---

<sup>\*</sup>Preliminary and incomplete.

<sup>†</sup>I thank Klaus Kultti, Juuso Välimäki, Pauli Murto, and Hannu Salonen for useful discussions and good comments.

<sup>‡</sup>E-mail: hannu.vartiainen@helsinki.fi. Address: HECER, Arkadiankatu 7, 00014 University of Helsinki, Finland.

In the literature on commitment in mechanism design, the usual approach is to appeal to the *incentive compatibility principle* (Myerson 1991, 1979) by assuming that parties anticipate how the mechanism will be redesigned. As in Neeman and Pavlov (2012), the foreseen renegotiation can then be incorporated into the original mechanism, and the attention can be limited to mechanisms that are robust against ex post renegotiation.<sup>1</sup>

Renegotiation-proofness is not, however, an entirely satisfactory concept. The problem is that it is often too restrictive. For example, in private valuations environment it admits only ex post efficient mechanisms (Neeman and Pavlov, 2012). Hence, not all set-ups support a renegotiation-proof mechanism, as the next canonical example demonstrates.

Consider the case of bilateral bargaining. There is a single indivisible good, a buyer, and a seller. Agents' privately known valuations are independently drawn from an interval. By the remarkable result of Myerson and Satterthwaite (1983), there is no incentive compatible, individually rational, and budget balanced mechanism that allocates the good to the agent with the highest valuation. Thus any feasible mechanism occasionally implements the inefficient no-trade outcome. But then the agents are tempted to renegotiate the mechanism rather than follow its instructions whenever no-trade outcome should materialize.

Renegotiation-proofness may thus be thought as a sufficient but not necessary condition for mechanisms that are implementable without commitment. The conceptual problem with renegotiation-proofness is that it permits all blocking mechanisms, even those that are not themselves credible. The natural way to restrict redesigns is to ask also the new mechanisms to be robust against renegotiation, when exposed to the same criterion as the original mechanism. This approach - the theme of this paper - provides a consistent way to close the gap between the necessary and sufficient conditions for mechanisms without commitment.

This paper develops a framework to identify implementable mechanisms when the planner cannot commit to the mechanism. Instead, she is permitted to redesign the mechanism *any* number of times without a cost. The key idea is to require robustness against redesigns that are themselves robust against redesigns that are themselves robust against... . The framework is portable to any mechanism design scenario. The structural assumptions that guarantee the existence of the solution are that the agents' type sets are finite and that their preferences exhibit value distinction (no pure belief types). We put no restrictions (except continuity) on the preferences of the planner.

As the starting point we take the observation that potential redesigns take place in sequential order and, hence, can be thought as a strategy. We identify redesign strategies that are *dynamically consistent*. By appealing to the incentive compatibility principle, our research strategy is to reduce, after each history, the continuation equilibria to a single direct incentive compatible mechanism.<sup>2</sup> In order to do this, we separate the two tasks of

---

<sup>1</sup>See also Forges (1995) and Dewatripont (1989). Other contributions on mechanism design without commitment include Segal and Whinston (2002), Freixas et al. (1985), McAfee and Vincent (1997), Baliga and Sjöström (1997), Bester and Strautz (2001), Skreta (2006, 2011), and Vartiainen (2012).

<sup>2</sup>Incentive compatibility principle: any equilibrium of the mechanism selection can be represented as a direct

a mechanism: information processing and implementation. An information processing device generates a public signal on the basis of the agents' reports, and simulates the information flow in the continuation game.<sup>3</sup> An implementation device then reflects what outcomes are implemented on the basis of revealed information. That is, after communication has been taken place via an information processing device, the planner reconsiders whether to implement the outcome suggested by the implementation device, or to design a new mechanism given the posterior information. Hence she cannot commit to the implementation device. However, no restrictions are put on how she coordinates communication between the parties through the information processing device.

The central question is what conditions should we put on the sequences of direct mechanism that reflect dynamically consistent redesign strategy. In the bilateral bargaining example above, the conditions should embody the intuition that a feasible mechanism is not renegotiated *ex post*, after the outcome has been revealed, to a new mechanism that is *itself* not subject to renegotiation, and so forth. More generally, after each history, the designer must be able to commit to the mechanism that the strategy assigns to her, given the counterfactual of not doing so.

The planner's mechanism selection strategy must be specified for all histories, compactly summarized by sequences of beliefs. Our solution concept guarantees that, after each history, the chosen mechanism gives the agents the incentives to play truthfully the information processing device and planner the incentives to obediently follow the implementation device. The two conditions that are necessary and sufficient for the mechanism design strategy to meet these desiderata are *optimality* and *consistency*. The former implies that, after all histories, the prescribed mechanism must maximize the planner's preferences among all the mechanisms that are feasible. This condition is dubbed as *Bellman optimality*. The latter condition requires that the mechanism prescribed by the strategy today must not be in conflict with the mechanism prescribed to her in the future.

Our main result is that a Bellman optimal and consistent mechanism design strategy always exists. The proof, which relies on a fixed point argument, uses history dependent mechanism design strategies. Indeed, there may be no history *independent* design strategy that meets the two desiderata.

Our approach highlights the central aspect of the mechanism design problem when the mechanism can be redesigned or renegotiated: it is not only the a priori incentives to reveal information that matter for the design but also how information flows within the mechanism are managed. Information that is revealed along the play may adversely affect the incentives at later stages (in Freixas et al. 1985, this property is called the "ratchet effect") which, given farsighted agents, affects the incentives already at the information revelation stage. How the information processing device should optimally be designed is the central - but difficult - question. The information processing device must be informative enough to allow implementing the desired outcome. But this still leaves

---

single stage mechanism that is truthfully played and obediently implemented.

<sup>3</sup>Assuming public signals restricts away private communication. This is a simplification. See Skreta (2006, 2010) for analyses of mechanism design without commitment but with private communication.

much freedom for the designer, and effective solutions often exist. We demonstrate the power of designing information processing devices in the bilateral bargaining context.

Our second result studies mechanisms that can be implemented without commitment in the canonical bargaining set up of Myerson and Satterthwaite (1983). The central question we ask whether the commitment inability rule out the possibility to implement the second best mechanism (Pareto-optimal in the class of incentive compatible, individually rational, and budget balanced mechanisms)? Our answer to this question is the affirmative: there is a Bellman optimal and consistent mechanism design strategy that implements the *incentive efficient* bargaining mechanism even if the agents do not have any external ways to commit to the inefficient no-trade outcomes. The driving force behind this result is that, by managing what information is being revealed during the bargaining process, the planner can induce a situation ex post where the buyer and the seller can commit not to continue bargaining any further even if they know that mutually beneficial transactions would still be possible. Interestingly, this calls for an information structure that is not as coarse as possible nor as fine as possible, but rather something in the middle. Specifically, the information structure that permits this is the one in which the agents conceive it possible that the agents' valuations are equally high, the agents valuations are equally low, or the buyer has the high valuation and the seller the low valuation. We demonstrate that, under such occurrences, the bargainers cannot reliably execute trade as it would require no trade in both the cases where the valuations are equal which cannot be committed to.

The novelty of our approach is that renegotiated mechanism is subjected to the same criticism than the original mechanism but otherwise possible mechanism/communication structures are not restricted. The key difference to Neeman and Pavlov (2012) and Forges (1995) is that they only focus on one-step counterfactuals whereas we account for the infinite hierarchy of counterfactuals. As a consequence, their solutions have more cutting power but suffer from existence problems.

Bester and Strausz (2001) study the one-agent scenario where the principal cannot commit to a certain action after the agent has communicated his type. Their main achievement is in showing that implementable outcomes can still be characterized via a version of the revelation principle. This result, however, heavily relies on the restricted form of the commitment problem. The principal can commit not to employ another mechanism once the agent has communicated his information. In particular, she can commit not to add another layer of mechanism on top of the old one. In contrast, we allow the planner to change the mechanism without restrictions.

Commitment is a critical question in the context of bargaining. The famous Coase Theorem asserts that, in the absence of commitment, the uninformed seller cannot commit to selling the good above her own reservation valuation. A mechanism design version of this theorem is provided by Ausubel and Deneckere (1989). McAfee and Vincent (1997) focus on a related question of designing an auction in a multi-agent environment when the seller cannot commit to the reserve price. They obtain a version of the Coase Theorem: when the opportunity cost of waiting vanishes, the seller is forced to sell without a reserve price. Skreta (2006, 2011) studies more auction design when the seller

has more flexibility in changing the rules of the game. Allowing remarkably rich strategy set for the seller, she is able to characterize the equilibrium mechanism. Her analysis relies on the assumption that redesigning the game is costly for the seller. Vartiainen (2011) approaches auction design without commitment from another angle. No waiting or other redesign costs are assumed. Applying the same solution as this paper, the key assumption in Vartiainen (2011) is that the information processing device prevents private communication between the seller and any individual agents. It is shown that the unique mechanism that is implementable by using a stationary mechanism design strategy implements the English auction in all cases.

General analyses of mechanism design without commitment include Holmström and Myerson (1983), Green and Laffont (1985), Baliga et al. (1997), and Lagunoff (1992). None of these does, however, address the main question of this paper: how to design mechanism when the planner can change the rules of the game as many times she wishes. The focus of Holmström and Myerson (1983) is in the question of ex ante committing to a particular rule. Their criterion "durability" excludes mechanisms that are not robust against a subset of types revealing that they belong to this set by designing a new mechanism for the types in this set. The posterior implementability concept of Green and Laffont (1985) demands that the incentives of the agents must not be sensitive to them understanding which outcome becomes implemented. As in this paper, Baliga et al. (1997) study mechanism design when planner is also a player. However, their focus is in Nash implementation which renders the informational processing property of the mechanism quite different. Lagunoff (1992) studies repeated redesign of complete information mechanism. He aim is to show that, under rather mild conditions, the any outcome that can potentially become implemented is Pareto optimal.

This paper is organized as follows. Section 2 specifies the set up and introduces the solution concept. Section 3 proves the existence of the solution. Section 4 applies the solution to the bilateral bargaining set up, and Section 5 provides concluding discussion.

## 2 Set up

**Preferences** There is a set  $\{1, \dots, n\}$  of agents, a planner, and a *finite* set of physical outcomes  $X$ . Agent  $i$ 's privately known type  $\theta_i$  is drawn from a *finite* set  $\Theta_i$ . Write  $\Theta = \times_{i \in N} \Theta_i$  with a typical element  $\theta = (\theta_i)_{i=1}^n$ , and  $\Theta_{-i} = \times_{j \neq i} \Theta_j$  with a typical element  $\theta_{-i} = (\theta_j)_{j \neq i}$ .<sup>4</sup> Denote the set of probability distributions on a (countable) set  $A$  by  $\Delta A$ . Denote a typical element of  $\Delta \Theta$  by  $p$  and by  $p_i(\cdot : \theta_i)$  the conditional distribution over  $\Theta_{-i}$  given  $p$  and the agent  $i$ 's type  $\theta_i$ . The support of the probability distribution  $p$  is denoted by  $\text{supp}(p)$ .<sup>5</sup>

Agent  $i$ 's vNM utility functions  $u$  exhibit *private valuations* and are of form  $u_i : X \times \Theta_i \rightarrow \mathbb{R}$ . We also assume that types satisfy *value distinction*: for all  $i \in \{1, \dots, n\}$ , for all  $\theta_i, \theta'_i \in \Theta_i$  there are lotteries  $\ell, \ell' \in \Delta X$  such that  $u_i(\ell, \theta_i) > u_i(\ell', \theta_i)$  and  $u_i(\ell, \theta'_i) < u_i(\ell', \theta'_i)$ . This assumption precludes pure belief types.

<sup>4</sup>That is,  $p_i(\theta_i) = \sum_{\theta_{-i}} p(\theta_i, \theta_{-i})$ .

<sup>5</sup> $\text{supp}(p) = \{\theta : p(\theta) > 0\}$ .

The agents and the planner want to maximize their expected payoff. As the outcome of the game may depend on the types of the agents', expectations are defined with respect to the *outcome function*  $f : \Theta \rightarrow \Delta X$  that specifies a probability distribution over outcomes for each type profile. Denote by

$$F = \{f : \Theta \rightarrow \Delta X\}$$

the set of all outcome functions. Endowed with the uniform metric,  $F$  is a compact metric space.

Given a probability distribution over the agents' types  $p \in \Delta\Theta$  and an outcome function  $f : \Theta \rightarrow \Delta X$ , agent  $i$ 's expected payoff is

$$\sum_{\theta_{-i}} \sum_x p(\theta_{-i} : \theta_i) u_i(x, \theta_i) f(x : \theta).$$

Planner's preferences are captured by a Bernoulli utility function  $v : X \times \Theta \rightarrow \mathbb{R}_+$  such that her expected payoff of  $f$  under  $p$  is given by

$$\sum_{\theta} \sum_x p(\theta) v(x, \theta) f(x : \theta).$$

Denote by  $v(f, p)$  the expected value of an outcome function  $f$  under  $p$ , and by  $v(x, p)$  the expected value of the outcome  $x$  under  $p$ .

**Mechanism** To implement an outcome function  $f$  the planner must elicit information from the agents by using a *mechanism*. A mechanism does two things: processes information and implements an outcome. To study commitment problems of the planner, we separate these tasks. A mechanism is a composite function

$$g \circ r : \Theta \rightarrow \Delta X,$$

consisting of an information processing device  $r$  and an implementation device  $g$  such that

$$r : \Theta \rightarrow \Delta S \quad \text{and} \quad g : S \rightarrow X,$$

where  $\Delta S$  is the set of probability distributions over a set  $S$  which we assume to be *countably infinite*.

A composite mechanism works as follows. After receiving the agents' messages  $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_n)$ , the information processing device  $r$  generates a (possibly random) public signal  $s \in S$  such that  $r(s : \hat{\theta}) > 0$ . This signal is used by the outcome function  $g$  to implement the outcome  $g(s) \in X$ . The signal  $s$  is the only information anyone - including the planner - obtains from  $r$ .<sup>6</sup>

---

<sup>6</sup>That the implementation device  $g$  is deterministic reflects the idea that the designer cannot make partial commitment, e.g. in the probabilistic sense, concerning implementation before the outcome is actually implemented. However, allowing random implementation device would not affect our results.

Because the set  $S$  is infinite it will be convenient - and without loss of generality - to assume that  $g$  is given and has the property that

$$g^{-1}(x) = \{s \in S : g(s) = x\}$$

contains infinitely many elements for all  $x \in X$ . Then also  $g(S) = X$ . Given this specification of  $g$ , the mechanism selection problem of the planner reduces to one of choosing  $r$ .

A mechanism

Many composite mechanisms induce the same outcome function. A particular example of a composite mechanism is the *direct* mechanism where  $g$  is a one-to-one function. This mechanism reveals the least amount of information necessary to implement the outcome specified by the outcome function  $f$ . In the other extreme there is the *fully revealing* mechanism that has the property that  $r$  is one-to-one. Under such mechanism, the agents fully reveal their types to the designer who then takes an action. It is clear that a fully revealing mechanism is likely to suffer from the planner's commitment problems.. Once the planner becomes informed of all the relevant information, she often is no longer interested in implementing the planned outcome. However, as we demonstrate in Section 5, committing to a mechanism may mean that some information should be induced - the direct mechanism is, in general, not the right mechanism either.

**Redesigning game** To fix the ideas, we now extend the single stage mechanism selection game by allowing the planner to redesign the mechanism repeatedly, once an outcome has been realized. Consider the following multistage mechanism design game. At each stage  $t = 0, 1, 2, \dots$ , the planner with preferences  $v$  announces a mechanism  $g \circ r^t$ , where  $r^t : \Theta \rightarrow \Delta S$ .<sup>7</sup> Given  $r^t$ , the agents choose an action profile  $\hat{\theta}^t \in \Theta$  which produces a randomized output  $r^t(\cdot : \hat{\theta}^t) \in \Delta S$ . The agents and the planner update their beliefs based on the observed signal  $s$ . Under the derived belief, the planner either implements  $g(s)$  or moves the game to period  $t + 1$ , in which case the game repeats itself.<sup>8</sup>

This is a proper extensive form game with incomplete information and imperfect monitoring. In any sequential equilibrium, after each history, (i) the players update their beliefs using the Bayes' rule when possible, (ii) maximize their expected payoffs given the strategies of the other players and their beliefs. In particular, the planner updates her beliefs in each period after observing the signal that this period's information processing device generates.

Planner's off-equilibrium beliefs are passive: whenever she observes a zero-probability signal, she believes that it is a result of a mistake and adopts the beliefs that she would obtain under some on equilibrium signal.

The game exhibits many sequential equilibria (including the "babbling" one). Our question is to select the natural one. As in the standard Bayesian mechanism design

<sup>7</sup> Allowing messages spaces larger than  $\times_i \Theta_i$  would not affect the results.

<sup>8</sup> Nothing would change if we allow the planner to discount her payoff by factor  $\delta \in (0, 1]$ . If  $\delta = 1$ , then the payoff of all parties from infinitely long delay in implementing an outcome is 0.

literature, we give the planner a leading role in choosing the continuation equilibrium. That is, we consider a particular type of equilibria where, at each public history of the game, the continuation equilibrium is the best among ones that the planner does not want to change any further. In a dynamic set up, this entails specifying a dynamic consistent equilibrium selection strategy. In particular, we need to know how the planner chooses the continuation equilibrium after *each* history of play. In the next section, we develop a solution that fullfills these aims.

### 3 The Solution

The planner's problem is that she cannot commit to the implementation device  $g$  at the ex post stage of the mechanism. Rather, once the signal  $s$  has been produced by the information processing device  $r$  she may be tempted to design a new mechanism under her post-signal belief. In this section, we develop a solution that identifies, under each history, planner's optimal equilibrium in the mechanism design game subject to the constraint that she will not change the equilibrium in the future, i.e. optimal mechanisms that the planner can commit to. Such a mechanism selection strategy is solved in two nested parts. First we specify mechanisms that the agents can commit to under the hypothesis that the planner can. Then we identify conditions under which the planner indeed can commit to the mechanism given that the agents play truthfully. This requires defining which mechanism the planner would implement under possible ex post beliefs.

**Agents' incentives** In order to study mechanisms that are consistent with the agents' incentives, let us assume that, at any given stage of the game, the planner can commit to implement the current mechanism as planned. What matters to the agents incentives is the outcome function associated to the mechanisms. Given  $p$ , mechanism  $g \circ r$  induces an outcome function  $f$  if  $f = g \circ r : \Theta \rightarrow \Delta X$ . That is, for any  $x \in X$ ,

$$f(x : \theta) = \sum_{s \in g^{-1}(x)} r(s : \theta), \quad \text{for all } \theta \in \Theta.$$

Type  $\theta_i$ 's *interim* payoff from a mechanism  $g \circ r$  under prior beliefs  $p$  when he reports  $\hat{\theta}_i$  to the mechanism is

$$\sum_{\theta_{-i}} \sum_x p(\theta_{-i} : \theta_i) u_i(g(s), \theta_i) r(s : \theta_{-i}, \hat{\theta}_i).$$

By the *revelation principle* (Myerson, 1979), an implementable mechanism  $g \circ r$  must be *incentive compatible* (IC):

$$\sum_{\theta_{-i}} \sum_s p(\theta) u_i(g(s), \theta_i) [r(s : \theta) - r(s : \theta_{-i}, \theta'_i)] \geq 0, \quad \text{for all } \theta_i, \theta'_i \in \Theta_i, \text{ for all } i = 1, \dots, n. \quad (1)$$

Conversely, if a mechanism at stage  $t$  is incentive compatible, and the planner can commit to implement any of its recommendations, then the agents are willing to play it truthfully.

However, we also need a more reliable condition that specifies not only when the mechanism can be implemented but also when it can be guaranteed to become implemented. To this end, we say that  $r$  is *robustly incentive compatible* if  $r(\text{supp}(p)) = r(\Theta)$  and, for any  $i$  and for any  $\theta_i \in \text{supp}(p_i)$ ,  $r(\cdot : \theta_{-i}, \theta_i) \neq r(\cdot : \theta_{-i}, \hat{\theta}_i)$  for some  $\theta_{-i}$  implies

$$\sum_{\theta_{-i}} \sum_s p(\theta) u_i(g(s), \theta_i) \left[ r(s : \theta) - r(s : \theta_{-i}, \hat{\theta}_i) \right] > 0, \quad \text{for all } \hat{\theta}_i \in \Theta_i.$$

That is, robustness requires that whenever the functioning of the mechanism requires the knowledge of the type of  $i$ , he should be given strict incentives to reveal it truthfully. Moreover, robustness restricts the mechanism by not allowing off-equilibrium signals  $s$ . Under such signals, it is a matter of off-equilibrium beliefs and continuation play whether the suggested outcome can be implemented. With robustly incentive compatible mechanism, such concerns are not warranted.

Denote by

$$\begin{aligned} IC(p) &= \{r \in \Phi : g \circ r \text{ is incentive compatible under } p\} \\ \widehat{IC}(p) &= \{r \in \Phi : g \circ r \text{ is robustly incentive compatible under } p\} \end{aligned}$$

the set of information processing devices that induce incentive compatible and strictly incentive compatible mechanisms under  $p$ .

Truthful announcements form a Bayes-Nash equilibrium in an incentive compatible mechanism  $r$  if the planner can *commit* to follow  $g$  after  $r$  has produced a signal  $s$ . Thus a mechanism maximizing the planner's payoff in  $IC(p)$  can be interpreted as the her full commitment benchmark. However, the many composite mechanisms that generate the same outcome function are not equivalent in terms of the planner's incentives. Different mechanisms reveal different amount of information to the planner and hence may affect her strategic possibilities ex post.

**Planner's incentives** The planner can condition her design strategy on the past design history. The stage  $t$  public history is summarized by a sequence

$$h = ((r^0, s^0), (r^1, s^1), \dots, (r^t, s^t)),$$

where  $g \circ r^k$  is the mechanism selected by the planner at stage  $k$ , and  $s^k$  is the output generated by the mechanism given the actions of the agents. Denote by  $H$  the set of all finite public histories and by  $\emptyset$  the initial history. Denote by  $h$  a typical element of  $H$  and by  $(h, (r, s))$  the concatenation of  $h$  and the public outcome  $(r, s)$  at the next stage. Then also  $(h, (r, s))$  is a public history.

Denote by  $p|_h$  planner's belief at history  $h$  (and hence by  $p|_{\emptyset}$  the initial belief). Assuming that the agents report their types truthfully, the signal  $s$  generated by a

mechanism  $r$  induces a posterior distribution  $p_{|h,(r,s)} \in \Delta\Theta$  such that

$$p_{|h,(r,s)}(\theta) = \frac{p_{|h}(\theta)r(s:\theta)}{\sum_{\theta' \in \Theta} p_{|h}(\theta')r(s:\theta')}, \quad \text{whenever } s \in r(\text{supp}(p_{|h})).$$

When  $s \notin r(\text{supp}(p_{|h}))$ , no restrictions are put on the posterior belief  $p_{|h,(r,s)}(\cdot)$ .

Any implementable mechanism  $g \circ r$  must be robust against the planner's temptation to redesign it ex post. That is, of replacing the outcome  $g(s)$  with *another* mechanism in  $\Phi$  that is preferred to the outcome  $g(s)$  under the posterior belief generated by the signal  $s$  the information processing device  $r$ . Our task is to identify conditions under which she will not do that.

Let the designer's mechanism design strategy be captured by a *choice rule*  $\sigma$  that specifies her mechanism choice for each history  $h \in H$ . Since the planner can only utilize mechanisms that she can commit to, the choice rule has to satisfy

$$\sigma : H \rightarrow \Phi \cup X \quad \text{such that} \quad \sigma(h) \in IC(p_{|h}) \cup X, \quad \text{for all } h \in H. \quad (2)$$

If  $\sigma(h) \in X$ , then the planner implements the outcome  $\sigma(h)$  at history  $h$ . If  $\sigma(h) \in IC(p_{|h})$ , then the planner continues the game by designing a new mechanism given the belief  $p_{|h}$ . The *function*  $\sigma(\cdot)$  represents the dynamic mechanism selection strategy of the seller, conditioned on histories. Note that the strategy is defined for *all* histories, including the off-equilibrium ones.

We now identify properties that the strategy  $\sigma$  should satisfy. We argue that the sequential rationality of the planner, and the players' knowledge of this, requires that  $\sigma$  reflect internal consistency and optimization. First we describe the set of mechanisms that the planner can commit to today given that  $\sigma$  is followed in the future. Denote by  $C^\sigma(h)$  planner's *maximal choice set* at history  $h$ , given  $\sigma$ . That is, the set of incentive compatible mechanisms that are *not* subject to redesign under the hypothesis that  $\sigma$  is followed ex post:

$$C^\sigma(h) = \{r \in IC(p_{|h}) : \sigma(h, (r, s)) = g(s), \text{ for all } s \in r(\Theta)\}.$$

Similarly, denote by  $\widehat{C}^\sigma(h)$  planner's *robust choice set* at history  $h$ , given  $\sigma$ . The set of robustly incentive compatible mechanisms that are not subject to redesign under the hypothesis that  $\sigma$  is followed ex post

$$\widehat{C}^\sigma(h) = \left\{ r \in \widehat{IC}(p_{|h}) : \sigma(h, (r, s)) = g(s), \text{ for all } s \in r(\Theta) \right\}.$$

Clearly  $\widehat{C}^\sigma(h) \subseteq C^\sigma(h)$ , reflecting the idea that  $\widehat{C}^\sigma(h)$  contains only elements that the planner can *confidently* implement, as they give strict incentives to the agents whenever incentives are needed, and they do not depend on off-equilibrium beliefs.

Choice sets  $C^\sigma(h)$  and  $\widehat{C}^\sigma(h)$  are defined with respect to the assumed  $\sigma$ . We now formally specify conditions that sequential rationality imposes on the choice rule  $\sigma$  itself. The first condition requires consistency in the sense that employing  $\sigma$  *ex ante* should not contradict  $\sigma$  being employed *ex post*.

**Definition 1 (Consistency)** *Choice rule  $\sigma$  is consistent if  $\sigma(h) \in C^\sigma(h) \cup X$ , for all  $h \in H$ .*

Our second condition reflects optimality.

**Definition 2 (Optimality)** *Choice rule  $\sigma$  is Bellman optimal if, for all  $h \in H$ ,*

1.  $\sigma(h) \in C^\sigma(h)$  implies

- $\delta v(\sigma(h), p|_h) \geq \delta v(g \circ r, p|_h)$ , for all  $r \in C^\sigma(h)$ ,
- $\delta v(\sigma(h), p|_h) \geq v(x, p|_h)$ , for all  $x \in X$ ,

2.  $\sigma(h) \in X$  implies

- $v(\sigma(h), p|_h) \geq \delta v(g \circ r, p|_h)$ , for all  $r \in \widehat{C}^\sigma(h)$ ,
- $v(\sigma(h), p|_h) \geq v(x, p|_h)$ , for all  $x \in X$ .

In words, when the planner implements a mechanism, this mechanism must be optimal in the class of mechanisms that she can commit to. Moreover, the mechanism must induce at least as high payoff as implementing an outcome immediately. If the planner chooses to implement a pure outcome, then this outcome must induce at least as high payoff as any feasible mechanism, or any other immediately implemented pure outcome.

Bellman optimality can be viewed as a one-time deviation restriction to the planner's design strategy: after any history  $h$ , she will not profit from a one-shot deviation to  $\sigma$  given that she will later follow  $\sigma$ . This entails that whenever the planner chooses a mechanism, that mechanism must be at least as good as any other feasible mechanisms, and it must be at least as profitable as implementing an outcome immediately. Conversely, without Bellman optimality,  $\sigma$  could not be convincingly committed to since the planner is able to make a reliable deviation to it after some history.

To see that a choice rule  $\sigma$  meeting the three desiderata reflects a planner's equilibrium selection strategy in the redesign game specified in the previous section, note that

- starting from any stage, a continuation (sequential) equilibrium of the redesign game can, by the revelation principle and the assumption that  $X$  is finite and  $S$  infinite be replicated by a single stage mechanism that the agents play truthfully,
- hence a scheme that specifies the planner's choice of the continuation equilibrium in class of equilibria that she can commit to can be simulated by a choice rule  $\sigma$  that assumes that each continuation equilibrium is of length 1.

Hence, by the revelation principle, and under the hypothesis that the designer can commit to the choice rule  $\sigma$  :

- A mechanism  $\phi$  is truthfully playable if  $\phi \in C^\sigma(h)$ , since then it will *not* be redesigned ex post.

- A mechanism  $\phi$  is *not* truthfully playable if  $\phi \notin C^\sigma(h)$ , since then it will be redesigned ex post.

## 4 Existence

We now state the main result of the paper.

**Theorem 1** *A Bellman optima and consistent mechanism selection strategy  $\sigma$  exists.*

The remainder of this section is devoted to proving the result. To construct a mechanism design strategy  $\sigma$  that meets the two desiderata, we first separate beliefs into two classes, ones under which an outcome can be implemented and ones under which it cannot be. First, denote the indirect utility function of the planner by

$$\max_{x \in X} v(x, p) = \bar{v}(p)$$

Consistently with the previous notation, denote by  $p_{|(r,s)}$  the posterior belief after signal  $s$  of an information processing device  $r$  starting from any given prior distribution  $p$  :

$$p_{|(r,s)}(\theta) = \frac{p(\theta)r(s : \theta)}{\sum_{\theta' \in \Theta} p(\theta')r(s : \theta')}, \quad \text{whenever } s \in r(\text{supp}(p)).$$

We will prove the existence via a series of subresults. The proof will rely on a fixed point argument which requires certain continuity properties from the correspondence of the set of incentive compatible mechanism, as a function of the prior distribution  $p$ . Note first that, by the closed graph theorem, the set of incentive compatible outcome functions  $g \circ r$  is upper hemicontinuous as a correspondence of the prior distribution  $p$ . However, given that  $g$  is defined on a countably infinite set  $S$ , the set of information processing devices  $r$  that induce incentive compatible outcome functions does not share these nice continuity properties. In particular,  $IC(\cdot)$  need not be . To restore the continuity properties of the mechanism correspondence, we need to focus on a restricted class of mechanisms - that only use a finite set of signals  $s$ . For any  $T \subseteq S$ , denote the set of incentive compatible mechanisms  $r : \Theta \rightarrow \Delta T$  under  $p$  by

$$IC^T(p) = \{r : g \circ r \text{ is incentive compatible under } p \text{ and } r : \Theta \rightarrow \Delta T\}.$$

Further, denote  $\widehat{IC}^T(p) = IC^T(p) \cap \widehat{IC}(p)$ .

**Lemma 1**  *$IC^T(\cdot)$  is upper hemicontinuous for any finite  $T$ .*

Lower hemicontinuity  $\widehat{IC}^T(p)$  or  $IC^T(p)$  cannot be guaranteed, however. We state that they meet the following weaker continuity notion.

**Lemma 2** *Let  $r \in \widehat{IC}^{r(\text{supp}(p))}(p)$ . For any  $p'$  sufficiently close to  $p$ , there is an  $r'$  such that  $r' \in IC^{r(\text{supp}(p))}(p')$ .*

**Proof.** Let  $r'$  agree with  $r$  on  $\times_i \text{supp}(p_i)$ . To specify  $r'$  outside  $\times_i \text{supp}(p_i)$ , construct a communication strategy of the agents where each  $\theta_i \in \text{supp}(p_i)$  reports his type truthfully truthfully, and each  $\theta_i \in \text{supp}(p'_i) \setminus \text{supp}(p_i)$  reports their best responses in  $\text{supp}(p_i)$  to all agents strategies. Use then the standard revelation argument to specify  $r'$  for type profiles  $\theta \in \times_i \text{supp}(p_i) \setminus \times_i \text{supp}(p_i)$ .

Our claim is that  $r' \in IC^{r(\text{supp}(p))}(p')$ . It suffices to show that  $r'$  is incentive compatible on  $\times_i \text{supp}(p_i)$ , for  $p'$  sufficiently close to  $p$ . This follows from the assumption that  $r \in \widehat{IC}^{r(\text{supp}(p))}(p)$  :

- If  $r(\cdot : \theta_{-i}, \theta_i) \neq r(\cdot : \theta_{-i}, \theta'_i)$  for some  $\theta_{-i}$ , then the incentive constraint for  $\theta_i \in \text{supp}(p_i)$  and  $\theta'_i \in \Theta_i$  holds strictly under  $p$ , and hence it holds also for  $p'$  sufficiently close to  $p$ .

- If  $r(\cdot : \theta_{-i}, \theta_i) = r(\cdot : \theta_{-i}, \theta'_i)$  for all  $\theta_{-i}$ , then the incentive constraint for  $\theta_i \in \text{supp}(p_i)$  and  $\theta'_i \in \Theta_i$  holds as equality for for all  $p'$ . ■

The previous lemmata shows that, when restricted to finite set of signals, the correspondences of incentive compatible information processing devices have nice continuity properties. The next lemma establishes that it is without loss of generality to focus on information processing devices  $r$  that do have finite support. That is, we show that any  $r$  that associates positive probability to infinitely many signals can be simulated by a one that only sends signals from a finite set, without affecting the incentive properties of the mechanism, or the posteriors that the device may generate.

For any  $r$  and  $g$ , use the notation  $r(T : \theta) = \sum_{s \in T} r(s : \theta)$  for any  $T \subseteq S$ ,  $r(\Theta) = \{s : r(s : \theta) > 0, \text{ for some } \theta \in \Theta\}$ , and  $g^{-1}(x) = \{s : x = g(s)\}$ .

**Lemma 3** Fix  $p$ . For any  $r : \Theta \rightarrow \Delta S$  there is  $r' : \Theta \rightarrow \Delta S$  such that

- (i)  $\{p_{|(r',s)}\}_{s \in g^{-1}(x)} \subseteq \{p_{|(r,s)}\}_{s \in g^{-1}(x)}$ ,
- (ii)  $r(g^{-1}(x) : \theta) = r'(g^{-1}(x) : \theta)$  for all  $x$ , and for all  $\theta$ ,
- (iii)  $r(\Theta)$  contains at most  $|X|(|\Theta| + 1) + 1$  elements.

**Proof.** We construct  $r'$  that meets the desiderata (i) - (ii) with respect to any given  $x$ , i.e. (i)  $\{p_{|(r',s)}\}_{s \in g^{-1}(x)} \subseteq \{p_{|(r,s)}\}_{s \in g^{-1}(x)}$ , (ii)  $r(g^{-1}(x) : \theta) = r'(g^{-1}(x) : \theta)$  for all  $\theta$ . Moreover,  $r'$  will have the property that  $\{s : r'(s : \theta) > 0, \theta \in \Theta, g(s) = x\}$  contains at most  $|\Theta| + 1$  elements. Applying the argument to all elements of  $X$  implies (iii), and proves the lemma.

Given a prior  $p \in \Delta \Theta$  denote by  $p_x \in \Delta \Theta$  the conditional probability distribution in the event  $g^{-1}(x)$ , i.e.

$$p_x(\theta) = \frac{p(\theta)r(g^{-1}(x) : \theta)}{\sum_{\theta'} p(\theta')r(g^{-1}(x) : \theta')}, \quad \text{for all } \theta \in \Theta, \quad (3)$$

Given signal  $s \in g^{-1}(x)$ , the posterior  $p_{|(r,s)}$  is defined by

$$p_{|(r,s)}(\theta) = \frac{p(\theta)r(s : \theta)}{\sum_{\theta'} p(\theta')r(s : \theta')}, \quad \text{for all } \theta \in \Theta.$$

Then we may write the distribution as specified by (3) in the form

$$p_x(\theta) = \sum_{s \in g^{-1}(x)} p_{|(r,s)}(\theta) \lambda_s, \quad \text{for all } \theta \in \Theta$$

where

$$\lambda_s = \frac{\sum_{\theta' \in \Theta} p(\theta') r(s : \theta')}{\sum_{\theta' \in \Theta} p(\theta') r(g^{-1}(x) : \theta')}, \quad \text{for all } s \in g^{-1}(x).$$

Thus  $p_x$  is a convex combination of the vectors  $\{p_{|(r,s)}\}_{s \in g^{-1}(x)} \subset \Delta\Theta$ . It then follows, by Carathéodory's Theorem, that there is  $\bar{S} \subseteq g^{-1}(x)$ , consisting of  $|\Theta| + 1$  or fewer elements, such that  $p_x$  is a convex combination of  $\{p_{|(r,s)}\}_{s \in \bar{S}}$ . That is, there is a vector  $(\bar{\lambda}_s)_{s \in \bar{S}}$  of nonnegative weights summing to one such that

$$p_x(\theta) = \sum_{s \in \bar{S}} p_{|(r,s)}(\theta) \bar{\lambda}_s, \quad \text{for all } \theta \in \Theta. \quad (4)$$

Construct  $r'$  such that, for all  $\theta \in \Theta$ ,

$$r'(s : \theta) = \begin{cases} r(s : \theta), & \text{for all } s \in S \setminus g^{-1}(x), \\ \frac{r(g^{-1}(x) : \theta) p_{|(r,s)}(\theta) \bar{\lambda}_s}{p_x(\theta)}, & \text{for all } s \in \bar{S}, \\ 0, & \text{for all } s \in g^{-1}(x) \setminus \bar{S}. \end{cases}$$

To verify that  $r'$  meets (i), it suffices that

$$p_{|(r',s)} = p_{|(r,s)}, \quad \text{for all } s \in \bar{S}.$$

To see this, note that, for all  $s \in \bar{S}$ ,

$$\begin{aligned} p_{|(r',s)}(\theta) &= \frac{p(\theta) r'(s : \theta)}{\sum_{\theta'} p(\theta') r'(s : \theta')} \\ &= \frac{p_x(\theta) r'(s : \theta) / r(g^{-1}(x) : \theta)}{\sum_{\theta'} p_x(\theta') r'(s : \theta') / r(g^{-1}(x) : \theta')} \\ &= \frac{p_{|(r,s)}(\theta) \bar{\lambda}_s}{\sum_{\theta'} p_{|(r,s)}(\theta') \bar{\lambda}_s} \\ &= p_{|(r,s)}(\theta), \quad \text{for all } \theta \in \Theta, \end{aligned}$$

where the first equality is the definition of  $p_{|(r',s)}(\theta)$ , the second follows by (3), the third by the definition of  $r'(s : \theta)$  on  $\bar{S}$ , and the final equality from  $\sum_{\theta'} p_{|(r,s)}(\theta') = 1$ .

To see that  $r'$  meets (ii), note that

$$\begin{aligned} r'(g^{-1}(x) : \theta) &= \sum_{s \in g^{-1}(x)} r'(s : \theta) \\ &= r(g^{-1}(x) : \theta) \frac{\sum_{s \in \bar{S}} p_{|(r,s)}(\theta) \bar{\lambda}_s}{p_x(\theta)} \\ &= r(g^{-1}(x) : \theta), \end{aligned}$$

where the second equality follows from  $r'(s : \theta) = r(s : \theta)$  for all  $s \notin g^{-1}(x)$  and the third by (4). Moreover,  $\{s : r'(s : \Theta) > 0, g(s) = x\} = \bar{S}$ , which consists of  $|\Theta| + 1$  elements, and hence  $r'$  meets the desideratum (iii). ■

Now we invoke an iterative procedure to identify distributions  $p$  under which the planner can commit to implement a pure outcome, and distributions under which she can commit to implement a mechanism. First we need two operators that specify primitive conditions under which the planner can and can not implement a pure outcome. Let  $p \in \Delta\Theta$  and  $D \subseteq \Delta\Theta$ . Define

$$\begin{aligned} M(p, D) &= \left\{ r \in IC(p) : \begin{array}{l} \delta v(g \circ r, p) \geq \bar{v}(p) \text{ and, for all } s \in r(\Theta), \\ v(g(s), p|_{(r,s)}) \geq \delta \bar{v}(p|_{(r,s)}) \text{ and } p|_{(r,s)} \in D \end{array} \right\}, \\ \widehat{M}(p, D) &= \left\{ r \in \widehat{IC}(p) : \begin{array}{l} \delta v(g \circ r, p) > \bar{v}(p) \text{ and, for all } s \in r(\Theta), \\ v(g(s), p|_{(r,s)}) > \delta \bar{v}(p|_{(r,s)}) \text{ and } p|_{(r,s)} \in D \end{array} \right\}. \end{aligned}$$

That is,  $M(p, D)$  constitutes the prior distributions under which it is weakly profitable for the planner to design a mechanism under the constraint the mechanism must induce a posterior in  $D$  and the planner can commit to the outcome of the mechanism rather than implementing any constant mechanism that takes one period to implement.  $\widehat{M}(p, D)$  is a more robust version of the same concept.

The following lemmata establishes two important continuity properties of these concepts.

**Lemma 4** *Let  $D$  be a open set in  $\Delta\Theta$ . Then there is a open set  $B \subseteq \Delta\Theta$  such that  $\{p : \widehat{M}(p, D) \neq \emptyset\} \subseteq B \subseteq \{p : M(p, D) \neq \emptyset\}$ .*

**Proof.** It suffices to show that each point in  $\{p : \widehat{M}(p, D) \neq \emptyset\}$  has an open neighborhood that is contained in  $\{p : M(p, D) \neq \emptyset\}$ . Let  $\widehat{M}(p, D) \neq \emptyset$ . By assumption, there is  $r \in \widehat{IC}^{r(\text{supp}(p))}(p)$  such that  $\delta v(g \circ r, p) > \bar{v}(p)$  and  $p|_{(r,s)} \in D$ , for all  $s \in r(\Theta)$ . By Lemma 3, we may assume that  $r(\Theta)$ , and hence  $r(\text{supp}(p))$ , contains finitely many elements. By Lemma 2 there is, for any  $p'$  sufficiently close to  $p$ , an  $r'$  such that  $r' \in IC^{r(\text{supp}(p))}(p')$ . Continuity of  $v$  implies that  $\delta v(g \circ r, p') > \bar{v}(p')$ , for  $p'$  sufficiently close to  $p$ . Moreover, since  $D$  is an open set and  $r(\text{supp}(p))$  contains finitely many elements,  $p|_{(r,s)} \in D$  implies  $p'|_{(r',s)} \in D$ , for all  $s \in r'(\text{supp}(p')) \subseteq r(\text{supp}(p))$ , for  $p'$  sufficiently close to  $p$ . Thus  $r' \in M(p', D)$ . Since this holds for any  $p'$  sufficiently close to  $p$ , we conclude that  $p$  has an open neighborhood that is contained in  $\{q : M(q, D) \neq \emptyset\}$ . ■

**Lemma 5** *Let  $D$  be a open set in  $\Delta\Theta$ . Then there is a open set  $G \subseteq \Delta\Theta$  such that  $\{p : M(p, \Delta\Theta \setminus D) = \emptyset\} = G$ .*

**Proof.** It suffices to show that  $\Delta\Theta \setminus G$  is closed. Let  $p$  be a limit point of  $\Delta\Theta \setminus G$ . Then there is a sequence  $\{p^t\}$  in  $\Delta\Theta \setminus G$  converging to  $p$ . Thus there also a sequence

$\{r^t\}$  such that, for all  $t$ ,  $r^t \in IC(p^t)$ ,  $v(g \circ r^t, p^t) \geq \bar{v}(p^t)$ , and  $p_{|(r^t, s)}^t \in \Delta\Theta \setminus D$ , for all  $s \in r^t(\text{supp}(p^t))$ . By Lemma 3, we may let the support of  $r^t$  be contained in a finite set  $\bar{S}$ . Since  $\{r^t\}$  lies in a compact space, it has a convergent subsequence and a limit point  $r^*$ . Since  $\bar{S} \times X$  is a finite set, this subsequence has a subsequence, also denoted by  $\{r^t\}$ . Since  $\Delta\Theta \setminus D$  is a closed set and  $p_{|(r^t, s)}^t \in \Delta\Theta \setminus D$  for all  $t$ , the sequence  $\{p_{|(r^t, s)}^t\}$  has a convergent subsequence and a limit point  $p_{|(r^*, s)} \in \Delta\Theta \setminus D$ , for all  $s \in \bar{S}$ . By the upper hemicontinuity of the  $IC^{\bar{S}}$  correspondence,  $r^* \in IC^{\bar{S}}(p)$ . By the continuity of the  $v$  function,  $v(g \circ r^*, p) \geq \bar{v}(p)$ . Thus  $p \in \Delta\Theta \setminus G$ . ■

Note that

$$\{p : M(p, D) \neq \emptyset\} \subseteq \{p : M(p, D') \neq \emptyset\}, \quad \text{for all } D \subseteq D' \subseteq \Delta\Theta, \quad (5)$$

and

$$\{p : M(p, D) \neq \emptyset\} = \Delta\Theta \setminus \{p : M(p, D) = \emptyset\}, \quad \text{for all } D \subseteq \Delta\Theta. \quad (6)$$

Now we use Lemmata 4 and 5 to construct two sequences of sets. The construction is recursive. Let  $B^0 = \emptyset$ . Take any  $k = 1, 2, \dots$ . Given sequences  $\{G_j\}_{j=0}^{k-1}$  and  $\{B_j\}_{j=0}^k$  of *open* subsets of  $\Delta\Theta$ , choose an *open* set  $G_k$  such that

$$G_k = \{p : M(p, \Delta\Theta \setminus \cup_{j \leq k} B_j) = \emptyset\} \quad (7)$$

and an *open* set  $B_{k+1}$  such that

$$\left\{p : \widehat{M}(p, \cup_{j \leq k} G_j) \neq \emptyset\right\} \subseteq B_{k+1} \subseteq \{p : M(p, \cup_{j \leq k} G_j) \neq \emptyset\}. \quad (8)$$

Since  $\cup_{j \leq k} B_j$  and  $\cup_{j \leq k-1} G_j$  are open sets by assumption, the desired  $G_k$  and  $B_{k+1}$  exist, by Lemmata 5 and 5.

The sequences  $\{G_k\}_{k=0}^{\infty}$  and  $\{B_k\}_{k=0}^{\infty}$  meeting (7) and (8) at each step  $k$  are now constructed. By (5) and (6)

$$G_0 \subseteq \dots \subseteq G_k \subseteq \dots \quad (9)$$

Define

$$B^* = \cup_k B_k \quad \text{and} \quad G^* = \cup_k G_k.$$

**Lemma 6**  $G^*$  and  $B^*$  are open and disjoint subsets of  $\Delta\Theta$ .

**Proof.** Since a union of a collections of open sets is open,  $G^*$  and  $B^*$  are open sets. To see that they are disjoint, it suffices that

$$(\cup_{j \leq k} G^j) \cap (\cup_{j \leq k} B^j) = \emptyset, \quad \text{for all } k = 0, 1, \dots \quad (10)$$

The proof is by induction. As  $B^0 = \emptyset$ , clearly  $G^0 \cap B^0 = \emptyset$ . Suppose that (10) holds up to step  $k = 0, \dots$ . We show that it also holds for step  $k + 1$ . We have

$$\begin{aligned} B_{k+1} &\subseteq \{p : M(p, \cup_{j \leq k} G_j) \neq \emptyset\} \\ &\subseteq \{p : M(p, \Delta\Theta \setminus \cup_{j \leq k} B_j) \neq \emptyset\} \\ &= \Delta\Theta \setminus \{p : M(p, \Delta\Theta \setminus \cup_{j \leq k} B_j) = \emptyset\} \\ &= \Delta\Theta \setminus \cup_{j \leq k} G_j, \end{aligned}$$

where the first set inclusion follows from (8), the second from (5), the first equality from (6), and second equality from the definition of  $G_k$ . Since, by the hypothesis,  $\cup_{j \leq k} B^j \subseteq \Delta\Theta \setminus \cup_{j \leq k} G^j$ , also  $\cup_{j \leq k+1} B^j \subseteq \Delta\Theta \setminus \cup_{j \leq k} G^j$ . Thus

$$(\cup_{j \leq k} G^j) \cap (\cup_{j \leq k+1} B^j) = \emptyset. \quad (11)$$

Similarly,

$$\begin{aligned} G_{k+1} &= \{p : M(p, \Delta\Theta \setminus \cup_{j \leq k+1} B_j) = \emptyset\} \\ &= \Delta\Theta \setminus \{p : M(p, \Delta\Theta \setminus \cup_{j \leq k+1} B_j) \neq \emptyset\} \\ &\subseteq \Delta\Theta \setminus \{p : M(p, \cup_{j \leq k} G^j) \neq \emptyset\} \\ &= \Delta\Theta \setminus \cup_{j \leq k} \{p : M(p, \cup_{i \leq j} G^i) \neq \emptyset\} \\ &\subseteq \Delta\Theta \setminus \cup_{j \leq k+1} B_j, \end{aligned}$$

where the first equality follows from (7), the second from (6), the first set inclusion from (11) and (5), and the final set inclusion from (8), applied to all  $j = 0, \dots, k$ . Thus, by (9),

$$(\cup_{j \leq k+1} G^j) \cap (\cup_{j \leq k+1} B^j) = \emptyset.$$

■

Since  $B^*$  is an open set, we conclude:

**Corollary 1**  $\Delta\Theta \setminus B^*$  is a compact set.

By construction, the sets  $G^*$  and  $B^*$  have the fixed point property that

$$\begin{aligned} G^* &\subseteq \{p : M(p, \Delta\Theta \setminus B^*) = \emptyset\}, \\ B^* &\subseteq \{p : M(p, G^*) \neq \emptyset\}. \end{aligned}$$

Moreover,

$$\Delta\Theta \setminus (G^* \cup B^*) \subseteq \{p : \widehat{M}(p, G^*) = \emptyset\} \cap \{p : M(p, \Delta\Theta \setminus B^*) \neq \emptyset\}.$$

Intuitively, the role of the two sets,  $G^*$  and  $B^*$ , will be the following. Belief  $p \in G^*$  is "good" in a sense that under  $p$  the planner can commit to implement a pure outcome. Belief  $p \in B^*$  is "bad" in a sense that then the planner cannot commit to implement a pure outcome, precisely since she can commit to implement pure outcomes under posteriors in  $G^*$ . In  $\Delta\Theta \setminus (G^* \cup B^*)$ , neither of the previous properties hold. Instead, under any  $p \in \Delta\Theta \setminus (G^* \cup B^*)$ , a maximal pure outcome is not strictly dominated by a robust mechanism that ends up in posteriors in  $G^*$ . Moreover, under any  $p \in \Delta\Theta \setminus (G^* \cup B^*)$ , even a maximal pure outcome is weakly dominated by a mechanism that ends up in posteriors in  $\Delta\Theta \setminus B^*$ . Our mechanism design strategy  $\sigma$  will use these properties of beliefs in these sets.

Now we construct a strategy  $\sigma : H \rightarrow \Phi \cup X$  that meets our two desiderata. Denote, for any  $p \in \Delta\Theta$  and  $D \subseteq \Delta\Theta$ ,

$$\begin{aligned} K(p, D) &= \left\{ r \in IC(p) : \begin{array}{l} \text{for all } s \in r(\text{supp}(p)), \\ v(g(s), p_{|(r,s)}) \geq \delta\bar{v}(p_{|(r,s)}) \text{ and } p_{|(r,s)} \in D \end{array} \right\}, \\ \widehat{K}(p, D) &= \left\{ r \in \widehat{IC}(p) : \begin{array}{l} \text{for all } s \in r(\text{supp}(p)), \\ v(g(s), p_{|(r,s)}) > \delta\bar{v}(p_{|(g,r),s}) \text{ and } p_{|(r,s)} \in D \end{array} \right\}. \end{aligned}$$

Denote by  $\bar{r} : \Delta\Theta \rightarrow \Phi$  a rule that associates each  $p$  a mechanism that maximizes the planner's payoff in  $K(p, \Delta\Theta \setminus B^*)$ :

$$\bar{r}(p) \in \arg \max_{r \in K(p, \Delta\Theta \setminus B^*)} v(g \circ r, p). \quad (12)$$

The following lemma establishes that  $K(p, \Delta\Theta \setminus B^*)$  always contains a maximizer.

**Lemma 7**  $\bar{r}(p)$  exists whenever  $K(p, \Delta\Theta \setminus B^*)$  is nonempty.

**Proof.** Let  $K(p, \Delta\Theta \setminus B^*)$  be nonempty. By Lemma 3 it is without loss of generality to assume that  $r \in IC^T(p)$  where  $T$  is a finite set. By construction,  $IC^T(p)$  is a compact set and, by Corollary 1, so is  $\Delta\Theta \setminus B^*$ . Mapping  $r \mapsto p_{|(r,s)}$  and, hence, mapping  $r \mapsto v(g(s), p_{|(r,s)})$  are continuous in  $r$ , for any  $s \in T$ . Thus program (12) is about maximizing a continuous function in a compact space. ■

We now construct a mechanism design strategy that meets the two desiderata by using the notions of  $G^*$ ,  $B^*$ , and  $\bar{r}$ . First, partition the set  $H$  of histories into distinct "phases"  $H^0$  and  $H^1$  as follows. Let  $\emptyset \in H^0$ . For any  $h' = (h, (r, s)) \in H$ , let

$$(h, (r, s)) \in \begin{cases} H^1, & \text{if } p_{|h'} \in \Delta\Theta \setminus (G^* \cup B^*), \ v(g(s), p_{|h'}) \geq \delta\bar{v}(p_{|h'}), \text{ and } h \in H^0, \\ H^0, & \text{otherwise.} \end{cases} \quad (13)$$

Construct a *mechanism design strategy*  $\sigma$  that is measurable with respect to the partition  $\{H^0, H^1\}$  such that  $\sigma(\emptyset) = \bar{r}(p)$  and such that, for any  $h' = (h, (r, s)) \in H$ ,

$$\sigma(h') = \begin{cases} \bar{r}(p_{|h'}), & \text{if } p_{|h'} \in B^*, \\ g(s), & \text{if } p_{|h'} \in G^* \text{ and } v(g(s), p_{|h'}) \geq \delta\bar{v}(p_{|h'}), \\ \bar{r}(p_{|h'}), & \text{if } p_{|h'} \in G^* \text{ and } v(g(s), p_{|h'}) < \delta\bar{v}(p_{|h'}), \\ g(s), & \text{if } p_{|h'} \in \Delta\Theta \setminus (G^* \cup B^*), \ v(g(s), p_{|h'}) \geq \delta\bar{v}(p_{|h'}), \text{ and } h \in H^0, \\ \bar{r}(p_{|h'}), & \text{if } p_{|h'} \in \Delta\Theta \setminus (G^* \cup B^*) \text{ and either } v(g(s), p_{|h'}) < \delta\bar{v}(p_{|h'}) \text{ or } h \in H^1. \end{cases} \quad (14)$$

That is, the mechanism design strategy depends on the phase, the current belief of the planner, and the status quo outcome. By Lemma 6,  $G^*$  and  $B^*$  are disjoint and  $\sigma$  is well defined. We shall now show that the constructed strategy meets the two desiderata.

**Proposition 1** *The mechanism design strategy  $\sigma$  constructed in (13) and (14) is Bellman optimal and consistent.*

**Proof.** Recall that

$$M(p, Z) = \{r \in K(p, Z) : \delta v(r, p) \geq \bar{v}(p)\} \quad (15)$$

and that

$$\widehat{K}(p, Z) \subseteq K(p, Z). \quad (16)$$

First we describe the relevant choice sets for each history  $h' = (h, (r, s))$ . There are 5 distinct cases:

1. If  $p|_{h'} \in B^*$ , then  $C^\sigma(h') = K(p|_{h'}, \Delta\Theta \setminus B^*)$ .
2. If  $p|_{h'} \in G^*$  and  $v(g(s), p|_{h'}) \geq \delta \bar{v}(p|_{h'})$ , then  $\widehat{C}^\sigma(h') = \widehat{K}(p|_{h'}, \Delta\Theta \setminus B^*)$ .
3. If  $p|_{h'} \in G^*$  and  $v(g(s), p|_{h'}) < \delta \bar{v}(p|_{h'})$ , then  $C^\sigma(h') = K(p|_{h'}, \Delta\Theta \setminus B^*)$ .
4. If  $p|_{h'} \in \Delta\Theta \setminus (G^* \cup B^*)$ ,  $v(g(s), p|_{h'}) \geq \delta \bar{v}(p|_{h'})$ , and  $h \in H^0$ , then  $\widehat{C}^\sigma(h') = \widehat{K}(p|_{h'}, G^*)$ .
5. If  $p|_{h'} \in \Delta\Theta \setminus (G^* \cup B^*)$  and either  $v(g(s), p|_{h'}) < \delta \bar{v}(p|_{h'})$  or  $h \in H^1$ , then  $C^\sigma(h') = K(p|_{h'}, \Delta\Theta \setminus B^*)$ .

To conclude that  $\sigma$  is consistent:

- In Cases 1 and 5,  $\emptyset \neq M(p|_{h'}, \Delta\Theta \setminus B^*) \subseteq K(p|_{h'}, \Delta\Theta \setminus B^*)$ , by the construction of  $B^*$  and  $\Delta\Theta \setminus G^*$ , respectively. Since  $\sigma(h') = \bar{r}(p|_{h'}) \in K(p|_{h'}, \Delta\Theta \setminus B^*)$ , consistency is implied.
- In Case 3,  $K(p|_{h'}, \Delta\Theta \setminus B^*)$  is not empty since a constant mechanism that induces the payoff  $\delta \bar{v}(p|_{h'})$  belongs to  $K(p|_{h'}, \Delta\Theta \setminus B^*)$ . Since  $\sigma(h') = \bar{r}(p|_{h'}) \in K(p|_{h'}, \Delta\Theta \setminus B^*)$ , consistency is implied.
- In Cases 2 and 4,  $\sigma(h') = g(s) \in X$ , and consistency is implied.

To verify that  $\sigma$  is Bellman optimal:

- When  $\sigma(h') \in C^\sigma(h')$ , we check that  $\sigma$  meets part 1 of the condition.
  - In Cases 1, 3 and 5, by the definition of  $\bar{r}$ , no element of  $M(p|_{h'}, \Delta\Theta \setminus B^*)$ , and hence of  $K(p|_{h'}, \Delta\Theta \setminus B^*)$ , strictly payoff dominates  $\sigma(h') = \bar{r}(p|_{h'})$ .
  - In Cases 1 and 5, since  $\bar{r}(p|_{h'}) \in M(p|_{h'}, \Delta\Theta \setminus B^*)$ , also  $\delta v(\bar{r}(p|_{h'}), p|_{h'}) \geq \bar{v}(p|_{h'}) \geq v(x, p|_{h'})$ , for all  $x \in X$ .
  - In Case 3, since a constant mechanism that induces the payoff  $\delta \bar{v}(p|_{h'})$  belongs to  $K(p|_{h'}, \Delta\Theta \setminus B^*)$  we have, by the definition of  $\bar{r}$ ,  $\delta v(\bar{r}(p|_{h'}), p|_{h'}) \geq \delta \bar{v}(p|_{h'}) \geq \bar{v}(p|_{h'}) \geq v(x, p|_{h'})$ , for all  $x \in X$ .
- When  $\sigma(h') = g(s) \in X$ , we check that  $\sigma$  meets part 2 of the condition.

- In Case 2, by the construction of  $G^*$ , no element of  $M(p|_{h'}, \Delta\Theta \setminus B^*)$  and by (15) of  $K(p|_{h'}, \Delta\Theta \setminus B^*)$  and by (16) of  $\widehat{K}(p|_{h'}, \Delta\Theta \setminus B^*)$ , weakly payoff dominates  $\sigma(h') = g(s)$ .
- In Case 4,  $\widehat{M}(p|_{h'}, G^*)$  is empty by the construction of  $B^*$ , and hence no element of  $\widehat{K}(p|_{h'}, G^*)$  strictly payoff dominates  $\sigma(h') = g(s)$ .

■

## 4.1 Participation constraints

In games of incomplete commitment, it is natural to permit agents to exit the game. However, modeling choices need to be made as regards to when this will be possible. There are two primary possibilities: (i) Once the agents enter a mechanism, they commit to it until the planner changes its rules. At this point, they choose whether or not enter the new mechanism. (ii) Agents can exit the mechanism at any time.

As is clear from the analysis above, what is sufficient for the existence is the continuity of the relevant set of mechanisms. Previously, this was guaranteed by the value distinction assumption. With participation constraints, this condition to be strengthened.

**Interim individual rationality** The first alternative leads to the standard interim participation constraint. Normalizing the value of the outside option of a player to zero, a mechanism  $r$  is (interim) *individually rational* if

$$\sum_{\theta_{-i}} \sum_s p(\theta) u_i(g(s), \theta_i) r(s : \theta) \geq 0, \quad \text{for all } \theta_i \in \Theta_i, \text{ for all } i = 1, \dots, n.$$

We say that  $r$  is *robustly interim incentive compatible* if, for any  $i$  and for any  $\theta_i \in \text{supp}(p_i)$ ,  $r(\cdot : \theta_{-i}, \theta_i) \neq r(\cdot : \theta_{-i}, \hat{\theta}_i)$  for some  $\theta_{-i}$  implies

$$\sum_{\theta_{-i}} \sum_s p(\theta) u_i(g(s), \theta_i) r(s : \theta) > 0.$$

To recover the existence of the previous section we only need to replace the  $IC$  and  $\widehat{IC}$  correspondences with the correspondence of incentive compatible, individually rational information processing devices and that of robustly incentive compatible, robustly individually rational information processing devices, respectively. Replicating the steps in the existence proof with these correspondences is routine.

**Other notions of participation** That is, whenever the functioning of  $r$  relies on  $\theta_i$  revealing his type,  $\theta_i$  should have strict incentive to participate. and it is *ex post*

individually rational if<sup>9</sup>

$$u_i(g(s), \theta_i) \geq 0, \quad \text{for all } s \in r(\theta), \quad \text{for all } \theta \in \Theta, \quad \text{for all } i = 1, \dots, n.$$

Assume that the set of feasible outcomes

$$X^{IR} = \{x : u_i(x, \theta_i) \geq 0, \text{ for all } i, \text{ for all } \theta \in \Theta, \text{ for all } i = 1, \dots, n\}$$

Hence,  $X^{IR}$  comprises all outcomes that can be implemented without violating the participation constraints.

The problem is that incentive compatibility and ex post individual rationality are not independent: an agent might exercise the veto right after off-equilibrium histories. The following simple extension to incentive compatibility resolves the problem by allowing  $i$  to veto the outcome even after his untruthful announcements.<sup>10</sup> Denote by

$$\tilde{u}_i(x, \theta_i) := \max\{u_i(x, \theta_i), 0\}$$

Given  $p$ , a mechanism  $r$  is *veto-incentive compatible* if

$$\sum_{\theta_{-i}} p(\theta) \left[ \sum_s \tilde{u}_i(g(s), \theta_i) r(s : \theta) - \sum_s \tilde{u}_i(g(s), \theta_i) r(s : \theta_{-i}, \theta'_i) \right] \geq 0, \quad (17)$$

for all  $\theta_i, \theta'_i \in \Theta_i$ , for all  $i \in N$ .

Veto-incentive compatibility requires that truthful reporting forms a Bayes-Nash equilibrium even if vetoing is possible after an untruthful announcement. Any implementable mechanism must thus be veto-incentive compatible. For any  $p$ , denote the set of veto-incentive compatible mechanisms by  $VIC(p)$ . It is easy to see that any veto-incentive compatible mechanism is incentive compatible and ex post individually rational (but not vice versa).<sup>11</sup>

## 5 Application: Bilateral Bargaining

Since Myerson and Satterthwaite (1983), it has been well known that committing to bilateral bargaining mechanisms is difficult. Consider a situation where two agents, a buyer (agent 1) and a seller (agent 2), are about to trade a good. Agents' valuations  $\theta_1$  and  $\theta_2$  are drawn from the finite set  $\Theta_1 = \Theta_2 = \{0, \frac{1}{K}, \dots, \frac{K-1}{K}, 1\}$ , for some  $K \in \mathbb{N}$ .

Our focus is on *budget balanced* mechanisms. The set of possible outcomes is  $X = \{0, 1\} \times \mathbb{R}$  with a typical element  $(a, m)$  where  $a = 1$  if the good is transferred from the

<sup>9</sup> *Interim* individual rationality requires that participation be weakly profitable before the output has been realized. Ex post constraint has been analysed e.g. by Forges (1993, 1998) and Gresik (1991, 1996).

<sup>10</sup> Veto-incentive compatibility is due to Forges (1998), and is closely related to IC\* of Matthews and Postlewaite (1989).

<sup>11</sup> Choose  $\theta_i = \theta'_i$  in (17). We only need EXP-IR and IC in the remainder of the paper.

seller to the buyer and  $a = 0$  if not, and  $m$  is a monetary transfer from the buyer to the seller. Given valuations  $\theta_1, \theta_2$ , the payoffs of the agents from the outcome  $(a, m)$  are

$$\begin{aligned} u_1(a, m, \theta_1) &= a\theta_1 - m, \\ u_2(a, m, \theta_2) &= m - a\theta_2. \end{aligned}$$

Let the agents' types be independently distributed according to  $p_1 \in \Delta\Theta_1$  and  $p_2 \in \Delta\Theta_2$ . Assume that  $p_1$  and  $p_2$  have a full support.

An outcome function associated to the problem is a mapping  $(a, m) : \Theta \rightarrow \{0, 1\} \times \mathbb{R}$ . A mechanism is *ex post efficient* if  $a(\theta_1, \theta_2) = 1$  whenever  $\theta_1 \geq \theta_2$  and  $a = 0$  otherwise. A mechanism is *inefficient* if it is not *ex post efficient*.

Under prior distribution  $p$ , denote the expected payoffs of the agents 1 and 2 when they have valuations  $\theta_1$  and  $\theta_2$  and report  $\theta'_1$  and  $\theta'_2$ , respectively, by

$$\begin{aligned} \sum_{\theta_2} p(\theta_2 : \theta_1) [a(\theta'_1, \theta_2)\theta_1 - m(\theta'_1, \theta_2)], \\ \sum_{\theta_1} p(\theta_1 : \theta_2) [m(\theta_1, \theta'_2) - a(\theta_1, \theta'_2)\theta_2]. \end{aligned}$$

A direct mechanism is *incentive compatible* if

$$\begin{aligned} \sum_{\theta_2} p(\theta_2 : \theta_1) [a(\theta)\theta_1 - m(\theta)] &\geq \sum_{\theta_2} p(\theta_2 : \theta_1) [a(\theta'_1, \theta_2)\theta_1 - m(\theta'_1, \theta_2)], \text{ for all } \theta_1, \theta'_1 \in \Theta_1, \\ \sum_{\theta_1} p(\theta_1 : \theta_2) [m(\theta) - a(\theta)\theta_2] &\geq \sum_{\theta_1} p(\theta_1 : \theta_2) [m(\theta_1, \theta'_2) - a(\theta_1, \theta'_2)\theta_2], \text{ for all } \theta_2, \theta'_2 \in \Theta_2, \end{aligned}$$

and it is (interim) *individually rational* if

$$\begin{aligned} \sum_{\theta_2} p(\theta_2 : \theta_1) [a(\theta)\theta_1 - m(\theta)] &\geq 0, \text{ for all } \theta_2 \in \Theta_2, \\ \sum_{\theta_1} p(\theta_1 : \theta_2) [m(\theta) - a(\theta)\theta_2] &\geq 0, \text{ for all } \theta_1 \in \Theta_1 \end{aligned}$$

A mechanism  $(a, m)$  is *incentive efficient* if there is no other *incentive compatible*, *individually rational*, and *budget balanced* mechanism that generates higher expected payoffs to both the agents.

Let us interpret the planner as an impartial mediator who maximizes the joint surplus of the agents: for all  $p \in \Delta(\Theta_1 \times \Theta_2)$ ,

$$v((a, m), p) = \sum_{(\theta_1, \theta_2)} p(\theta_1, \theta_2) a(\theta_1, \theta_2) (\theta_1 - \theta_2)$$

Given this objective function, the planner has always an incentive not to stop with no-trade if there is still scope for further mutually beneficial trade.

The classic result due to Myerson and Satterthwaite (1983) says that when  $\theta_1$  and  $\theta_2$  are independently distributed on an interval and their absolutely continuous distributions

overlap, then the incentive and participation constraints prevent full efficiency: any incentive compatible, individually rational, and budget balanced mechanism implements an inefficient outcome with strictly positive probability. In particular, *any incentive efficient mechanism is inefficient*. This inefficiency raises the question of renegotiation. Would the parties stop bargaining once they know that all mutually beneficial transactions are not exhausted?

The aim of this section is to show that the agents' inability to commit to the mechanism does not prevent them implementing an incentive efficient contract, i.e. there is a consistent and Bellman optimal mechanism design rule that allows committing even to the inefficient outcome. Towards this end, we need to construct an information processing device under which consistent renegotiation is not feasible.

Before constructing the mechanism selection strategy that meets our desiderata, we need to extend the classic characterization results of Myerson and Satterthwaite (1983) to our discrete set up, as in the original context the set of valuations is a continuum (an interval). This is not a completely innocent modification of the model since the original Myerson-Satterthwaite (1993) result relies on an envelope argument, and hence requires the set of types to be connected.

Let  $\theta_1$  and  $\theta_2$  be independently distributed with distribution functions  $p_1$  and  $p_2$ . Given  $p_i$ , denote the cumulative distribution by

$$P_i(\theta_i) = \sum_{t \leq \theta_i} p_i(t), \quad \text{for } i = 1, 2,$$

and, for any  $\gamma \in [0, 1]$ ,

$$c_1(\theta_1, \gamma) = \theta_1 - \gamma \frac{1 - P_1(\theta_1)}{p_1(\theta_1)}, \quad \text{for all } \theta_1 \in T,$$

$$c_2(\theta_2, \gamma) = \theta_2 + \gamma \frac{P_2(\theta_2)}{p_2(\theta_2)}, \quad \text{for all } \theta_2 \in T.$$

We say that the two distribution functions  $p_1$  and  $p_2$  are *regular* if  $c_1(\cdot, 1)$  and  $c_2(\cdot, 1)$  are increasing.

We now establish a finite version of the classic result of Myerson and Satterthwaite (1983). The proof of the proposition is relegated to the appendix.

**Proposition 2** *Let  $\theta_1$  and  $\theta_2$  be independently distributed with regular distribution functions  $p_1$  and  $p_2$ , respectively. Then there is an incentive efficient direct mechanism  $(a^\gamma, m^\gamma)$  such that, for some  $\gamma \in (0, 1]$ ,*

$$\text{if } c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) \geq 0, \quad \text{then } a^\gamma(\theta_1, \theta_2) = 1 \quad (18)$$

$$\text{if } c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) < 0, \quad \text{then } a^\gamma(\theta_1, \theta_2) = 0. \quad (19)$$

From this result it is clear that, with sufficiently fine grid in  $\Theta_1 = \Theta_2$ , the incentive efficient direct mechanism  $(a^\gamma, m^\gamma)$  will be inefficient: an inefficient no-trade outcome

will materialize whenever

$$\gamma \frac{1 - P_1(\theta_1)}{p_1(\theta_1)} + \gamma \frac{P_2(\theta_2)}{p_2(\theta_2)} > \theta_1 - \theta_2 > 0.$$

We make two observations on the incentive efficient mechanism. These properties will be used to construct a mechanism on which the planner can commit to.

**Remark 1** *Let  $\theta_1$  and  $\theta_2$  be independently distributed with regular distribution functions  $p_1$  and  $p_2$ , respectively. Let  $(a^\gamma, m^\gamma)$  be an incentive efficient direct mechanism as defined in (18)-(19). Then, for any  $(\theta_1, \theta_2) \in \Theta_1 \times \Theta_2$ ,*

$$a^\gamma(\theta_1, \theta_2) = 0 \quad \text{implies} \quad \begin{cases} a^\gamma(\theta'_1, \theta_2) = 0, & \text{for all } \theta'_1 \leq \theta_1, \\ a^\gamma(\theta_1, \theta'_2) = 0, & \text{for all } \theta'_2 \geq \theta_2. \end{cases}$$

*In particular,  $\theta_1 > \theta_2$  and  $a(\theta_1, \theta_2) = 0$  imply  $a(\theta_1, \theta_1) = a(\theta_2, \theta_2) = 0$ .*

Our aim is to construct a mechanism that allows the parties to commit not to continue negotiation even when trade does not take place. To this end, the information processing device of the mechanism must be designed in such a way that the prescribed outcome can be committed to under the posterior information. Since the information structure with respect to the outcome function  $(a, m)$  is measurable is at most as coarse than that of  $r$ , we need to verify that that the outcome of the optimal mechanism does itself reveal unintended information. For our purposes, it suffices that there is an efficient mechanism that prescribes zero monetary transfer when trade does not take place. The no-trade outcome then only reveals that the types of the agents  $(\theta_1, \theta_2)$  satisfy (19).

This guarantees that, when trade does not take place, only this information is revealed. Gresik (1991) establishes the existence of such transfers in the continuous type sets case. For completeness, we construct such schemes in the current case when the types sets are finite. The proof of the following lemma appears in the appendix.

**Lemma 8** *Let  $\theta_1$  and  $\theta_2$  be independently distributed with regular distribution functions  $p_1$  and  $p_2$ , respectively. Then there is an incentive efficient direct mechanism  $(a^\gamma, m^\gamma)$  as defined in (18)-(19) such that the transfer rule  $m^\gamma$  prescribes zero monetary transfer when trade does not take place, i.e.*

$$a^\gamma(\theta_1, \theta_2) = 0 \quad \text{implies} \quad m^\gamma(\theta_1, \theta_2) = 0.$$

Our question is whether there is a Bellman optimal and consistent mechanism choice rule that permits implementation of a compound mechanism that is outcome equivalent with the incentive efficient mechanism  $(a^\gamma, m^\gamma)$ . We shall show that this is the case.

We are now ready to state the desired result: the agents can commit to implementing the Myerson-Satterthwaite incentive efficient mechanism in the bilateral bargaining context even in the absence of external commitment devices. This entails that the agents design an information processing device through which their communication takes place in a way that they cannot commit not to continue bargaining after it becomes clear that the inefficient no-trade outcome will become implemented.

**Theorem 2** *Let  $\theta_1$  and  $\theta_2$  be independently distributed with regular distribution functions  $p_1$  and  $p_2$ , respectively. Then there is a Bellman optimal and consistent mechanism selection strategy  $\sigma$  that implements an incentive efficient mechanism under  $(p_1, p_2)$ , when  $p_{|\emptyset} = (p_1, p_2)$ .*

The remainder of this section proves the result. Our key task is to construct an information processing device which provides just the right amount of information for the agents to commit to the inefficient no-trade outcome.

There are many ways to for the information precessing device  $r$  to provide enough information for the mechanism to work properly. Our central task is to design  $r$  in such a way that it blocks further negotiation but still permits implementation the outcomes prescribed by the incentive efficient mechanism  $(a^\gamma, m^\gamma)$ .

Let (a subset of) the signal space be defined by ordered pairs

$$S^* = \{\langle \theta_1, \theta_2 \rangle : \theta_1 \geq \theta_2\} \cup \{0\}. \quad (20)$$

Consider the following information processing device  $r^* : \Theta_1 \times \Theta_2 \rightarrow \Delta S^*$ . For any  $t$ , let  $\kappa(t) = \#\{t' : t \geq t' \text{ and } c_1(t, \gamma) \leq c_2(t', \gamma) \text{ or } t \geq t' \text{ and } c_1(t', \gamma) \leq c_2(t, \gamma)\}$ . Then

$$r^*(\cdot : \theta_1, \theta_2) = \begin{cases} 1_{\langle \theta_1, \theta_2 \rangle}, & \text{if } \theta_1 > \theta_2, \\ \frac{1}{\kappa(t)} \left( \sum_{t': t' \leq t \text{ and } c_1(t, \gamma) \leq c_2(t', \gamma)} 1_{\langle t, t' \rangle} + \sum_{t': t' \geq t \text{ and } c_1(t', \gamma) \leq c_2(t, \gamma)} 1_{\langle t', t \rangle} \right), & \text{if } \theta_1 = \theta_2 = t. \\ 1_0, & \text{if } \theta_1 < \theta_2. \end{cases} \quad (21)$$

That is, a signal  $\langle \theta_1, \theta_2 \rangle$  such that  $c_1(\theta_1, \gamma) \geq c_2(\theta_2, \gamma)$  may be sent only by the type pair  $(\theta_1, \theta_2)$ , and a signal  $\langle \theta_1, \theta_2 \rangle$  such that  $c_1(\theta_1, \gamma) < c_2(\theta_2, \gamma)$  and  $\theta_1 \geq \theta_2$  may be sent by type pairs  $(\theta_1, \theta_2)$ ,  $(\theta_1, \theta_1)$ , or  $(\theta_2, \theta_2)$ . A signal "0" may only be send by a type pair  $(\theta_1, \theta_2)$  such that  $\theta_1 < \theta_2$ .

Further, define an implementation device  $g^* : S^* \rightarrow \{0, 1\} \times \mathbb{R}$  such that, for any  $s \in S^*$ ,

$$g^*(s) = \begin{cases} (1, m^\gamma(\theta_1, \theta_2)), & \text{if } s = \langle \theta_1, \theta_2 \rangle \text{ and } c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) \geq 0, \\ (0, 0), & \text{if } s = \langle \theta_1, \theta_2 \rangle \text{ and } c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) < 0, \\ (0, 0), & \text{if } s = 0. \end{cases} \quad (22)$$

By construction, the compound mechanism  $g^* \circ r^*$  satisfies

$$\begin{aligned} \text{if } c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) &\geq 0, & \text{then } g^*(r^*(\theta_1, \theta_2)) &= (1, m^\gamma(\theta_1, \theta_2)), \\ \text{if } c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) &< 0, & \text{then } g^*(r^*(\theta_1, \theta_2)) &= (0, 0). \end{aligned}$$

Hence, by Lemma 8,

$$g^*(r^*(\cdot)) = (a^\gamma, m^\gamma)(\cdot).$$

By Proposition 2,  $(a^\gamma, m^\gamma)$  is an incentive efficient mechanism when  $p_1$  and  $p_2$  are regular. We conclude:

**Lemma 9** *Let  $p_1$  and  $p_2$  be regular distributions. Then the mechanism  $g^* \circ r^*$  is incentive efficient.*

Our aim is to show that the mechanism  $g^* \circ r^*$  can be committed to under regular distributions  $(p_1, p_2)$ . To show this, construct a mechanism design strategy  $\sigma^*$  that is consistent and Bellman optimal, and implements  $g^* \circ r^*$  when  $(p_1, p_2)$  is taken as the initial belief  $p|_{\emptyset}$ .

Note first that when  $s = \langle \theta_1, \theta_2 \rangle$  such that  $c_1(\theta_1, \gamma) \geq c_2(\theta_2, \gamma)$  or  $s = 0$ , the implemented outcome  $g^*(s)$  is ex post efficient. Since there is no mechanism that surplus dominates such an outcome, the only issue is whether the planner can commit to the inefficient no-trade outcome, i.e. when  $s = \langle \theta_1, \theta_2 \rangle$  such that  $\theta_1 > \theta_2$  and  $c_1(\theta_1, \gamma) < c_2(\theta_2, \gamma)$ . We need to consider the posterior belief that is induced by such a signal.

Note that an information processing device  $r^*$  may send a signal  $s = \langle \bar{t}, \underline{t} \rangle$  such that  $\bar{t} > \underline{t}$  and  $c_1(\bar{t}, \gamma) < c_2(\underline{t}, \gamma)$  under the following ordered pairs of types  $(\theta_1, \theta_2) : (\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})$ . This implies that the signal  $\langle \bar{t}, \underline{t} \rangle$  induces a posterior belief  $p|_{r^*, \langle \bar{t}, \underline{t} \rangle}$  such that

$$\text{supp}(p|_{r^*, \langle \bar{t}, \underline{t} \rangle}) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}.$$

Our task is to construct a mechanism selection rule  $\sigma^*$  such that there is no credible way to continue bargaining under the belief  $p|_{r^*, \langle \bar{t}, \underline{t} \rangle}$  even though a mutually profitable trading opportunity exists with strictly positive probability.

We construct a  $\sigma^*$  that satisfies Bellman optimality and consistency on  $\text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$ . For any  $h'$ , let  $\sigma^*(h')$  depend on the distribution  $p|_{h'} \in \Delta\{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$ . Our construction is on induction on the cardinality of  $\text{supp}(p|_{h'})$ . First, let  $g^* : S^* \rightarrow \{0, 1\} \times \mathbb{R}$  be defined by

$$g^*(s) = \begin{cases} (1, \bar{t}), & \text{if } s = \langle \bar{t}, \bar{t} \rangle, \\ (1, (\bar{t} + \underline{t})/2), & \text{if } s = \langle \bar{t}, \underline{t} \rangle, \\ (1, \underline{t}), & \text{if } s = \langle \underline{t}, \underline{t} \rangle, \\ (0, 0), & \text{if } s = 0. \end{cases} \quad (23)$$

Partition first the set of public histories  $H$  into two sets  $H^0$  and  $H^1$  such that, for any  $h' = (h, (r, s)) \in H$ ,

$$h' \in \begin{cases} H^1, & \text{if } h \in H^0 \text{ and } \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\} \subseteq \text{supp}(p_{h'}), \\ H^0 & \text{otherwise.} \end{cases} \quad (24)$$

Construct a choice rule  $\sigma^*$  such that

$$\sigma^*(h) = \begin{cases} 1_0, & \text{if } \text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\} \text{ and } h \in H^0, \\ r^*(\theta) := \begin{cases} \langle \bar{t}, \underline{t} \rangle, & \text{if } \theta = (\bar{t}, \underline{t}), \\ 0, & \text{if } \theta \neq (\bar{t}, \underline{t}), \end{cases} & \text{if } \text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\} \text{ and } h \in H^1, \\ 1_0, & \text{if } \text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\} \text{ and } h \in H^0, \\ r^{**}(\theta) := \begin{cases} \langle \underline{t}, \underline{t} \rangle, & \text{if } \theta = (\underline{t}, \underline{t}), \\ \langle \bar{t}, \bar{t} \rangle, & \text{if } \theta = (\bar{t}, \bar{t}), \\ 0, & \text{if } \theta \notin \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}, \end{cases} & \text{if } \text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\} \text{ and } h \in H^1, \\ 1_{\langle \underline{t}, \underline{t} \rangle}, & \text{if } \text{supp}(p|_h) = \{(\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}, \\ 1_{\langle \bar{t}, \bar{t} \rangle}, & \text{if } \text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}, \\ 1_{\langle \bar{t}, \bar{t} \rangle}, & \text{if } \text{supp}(p|_h) = \{(\bar{t}, \bar{t})\}, \\ 1_{\langle \bar{t}, \underline{t} \rangle}, & \text{if } \text{supp}(p|_h) = \{(\bar{t}, \underline{t})\}, \\ 1_{\langle \underline{t}, \underline{t} \rangle}, & \text{if } \text{supp}(p|_h) = \{(\underline{t}, \underline{t})\}. \end{cases} \quad (25)$$

**Lemma 10** *Let  $\text{supp}(p|_\emptyset) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$  with  $\bar{t} > \underline{t}$ . There is a Bellman optimal and consistent choice rule  $\sigma^*$  such that  $\sigma^*(h) = (0, 0)$ .*

**Proof.** First we describe the choice set  $C^{\sigma^*}(h)$  for each public history  $h$ . There are 9 distinct cases:

1.  $\text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$  and  $h \in H^0$ . Then  $C^{\sigma^*}(h) = \emptyset$ .

To see this, suppose on the contrary that  $r \in C^{\sigma^*}(h)$ . Then:

- (a)  $\text{supp}(p|_{h,(r,s)}) \subset \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$ , for all  $s \in r(\text{supp}(p|_h))$ , by the construction of  $\sigma^*(h, (r, s))$ ,
- (b)  $\text{supp}(p|_{h,(r,s)}) \neq \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}$ , by the construction of  $\sigma^*(h, (r, s))$ ,
- (c)  $g^*(s) \in \{(1, \underline{t}), (1, (\bar{t} + \underline{t})/2), (1, \bar{t})\}$ , for all  $s \in r(\text{supp}(p|_h))$ , by (a), (b), and the construction of  $\sigma^*(h, (r, s))$ ,
- (d)  $g^*(s) = (1, \bar{t})$ , for all  $s \in r(\bar{t}, t)$  for all  $t \in \{\bar{t}, \underline{t}\}$ , by (c) and individual rationality,
- (e)  $g^*(s) = (1, \underline{t})$ , for all  $s \in r(t, \underline{t})$  for all  $t \in \{\bar{t}, \underline{t}\}$ , by (c) and individual rationality,
- (f)  $g^*(s) = (1, \bar{t})$ , for all  $s \in r(\underline{t}, t)$  for all  $t \in \{\bar{t}, \underline{t}\}$ , by (d) and incentive compatibility,
- (g)  $g^*(s) = (1, \underline{t})$ , for all  $s \in r(t, \bar{t})$  for all  $t \in \{\bar{t}, \underline{t}\}$ , by (e) and incentive compatibility.
- (h) By (d) and (f),  $g^*(s) = (1, \bar{t})$  for all  $s \in r(t, t')$  for all  $(t, t') \in \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$  and, by (e) and (g),  $g^*(s) = (1, \bar{t})$  for all  $s \in r(t, t')$  for all  $(t, t') \in \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$ , a contradiction.

2.  $\text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$  and  $h \in H^1$ . Then  $C^{\sigma^*}(h) = \{r \in IC(p|_h) : g^*(s) = \sigma^*(h, (r, s))\}$ .
3.  $\text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}$  and  $h \in H^0$ . Then  $C^{\sigma^*}(h) \subseteq \{r \in IC(p|_h) : g(r(\bar{t}, \bar{t})) \in \{(0, 0), (1, \bar{t})\}, g(r(\underline{t}, \underline{t})) \in \{(0, 0), (1, \underline{t})\}\}$ , by incentive compability, individual rationality, and the construction of  $\sigma^*(h, \cdot)$ .
4.  $\text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}$  and  $h \in H^1$ . Then  $C^{\sigma^*}(h) \subseteq \{r \in IC(p|_h) : g(r(\bar{t}, \bar{t})) \in \{(0, 0), (1, \bar{t})\}, g(r(\underline{t}, \underline{t})) \in \{(0, 0), (1, \underline{t})\}\}$ , by incentive compability, individual rationality, and the construction of  $\sigma^*(h, \cdot)$ .
5.  $\text{supp}(p|_h) = \{(\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$ . Then  $C^{\sigma^*}(h) = \{r : g^* \circ r = 1_{(1, \underline{t})}\}$ , by incentive compability, individual rationality, and the construction of  $\sigma^*(h, \cdot)$ .
6.  $\text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t})\}$ . Then  $C^{\sigma^*}(h) = \{r : g^* \circ r = 1_{(1, \bar{t})}\}$ , by incentive compability, individual rationality, and the construction of  $\sigma^*(h, \cdot)$ .
7.  $\text{supp}(p|_h) = \{(\bar{t}, \bar{t})\}$ . Then  $C^{\sigma^*}(h) = \{r : g^* \circ r = 1_{(1, \bar{t})}\}$ , by the construction of  $\sigma^*(h, \cdot)$ .
8.  $\text{supp}(p|_h) = \{(\bar{t}, \underline{t})\}$ . Then  $C^{\sigma^*}(h) = \{r : g^* \circ r = 1_{(1, (\bar{t} + \underline{t})/2)}\}$ , by the construction of  $\sigma^*(h, \cdot)$ .
9.  $\text{supp}(p|_h) = \{(\underline{t}, \underline{t})\}$ . Then  $C^{\sigma^*}(h) = \{r : g^* \circ r = 1_{(1, \underline{t})}\}$ , by the construction of  $\sigma^*(h, \cdot)$ .

To see that  $\sigma^*$  is consistent:

- In Cases 1, 3, 5, 6, 7, 8, and 9,  $\sigma^*(h) \in X$ , implying consistency in these cases.
- In Case 2,  $\sigma^*(h) = r^*$  such that

$$r^*(\theta) = \begin{cases} \langle \bar{t}, \underline{t} \rangle, & \text{if } \theta = (\bar{t}, \underline{t}), \\ 0, & \text{if } \theta \neq (\bar{t}, \underline{t}), \end{cases}$$

There are two cases,  $s = \langle \bar{t}, \underline{t} \rangle$  and  $s = 0$ . In the former,  $g^*(\langle \bar{t}, \underline{t} \rangle) = (1, (\bar{t} + \underline{t})/2)$  and  $\text{supp}(p|_{h, (r, \langle \bar{t}, \underline{t} \rangle)}) = \{(\bar{t}, \underline{t})\}$ . By construction  $\sigma^*(h, (r^*, \langle \bar{t}, \underline{t} \rangle)) = g^*(\langle \bar{t}, \underline{t} \rangle)$ . In the latter case,  $g^*(0) = (0, 0)$  and  $\text{supp}(p|_{h, (r, 0)}) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}$ . By construction, since  $(h, (r, 0)) \in H^0$ ,  $\sigma^*(h, (r^*, 0)) = g^*(0)$ , implying consistency in Case 2.

- In Case 4,  $\sigma^*(h) = r^{**}$  such that

$$r^{**}(\theta) = \begin{cases} \langle \underline{t}, \underline{t} \rangle, & \text{if } \theta = (\underline{t}, \underline{t}), \\ \langle \bar{t}, \bar{t} \rangle, & \text{if } \theta = (\bar{t}, \bar{t}), \\ 0, & \text{if } \theta \notin \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}. \end{cases}$$

Since  $\text{supp}(p|_h) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}$ , only signals  $s = \langle \bar{t}, \bar{t} \rangle$  and  $s = \langle \underline{t}, \underline{t} \rangle$  materialize with positive probability. In the former case,  $g^*(\langle \bar{t}, \bar{t} \rangle) = (1, \bar{t})$  and  $\text{supp}(p|_{h, (r^{**}, \langle \bar{t}, \bar{t} \rangle)}) =$

$\{(\bar{t}, \bar{t})\}$ . By construction,  $\sigma^*(h, (r^{**}, \langle \bar{t}, \bar{t} \rangle)) = g^*(\langle \bar{t}, \bar{t} \rangle)$ . In the latter case,  $g^*(\langle \underline{t}, \underline{t} \rangle) = (1, \underline{t})$  and  $\text{supp}(p_{|h, (r^{**}, \langle \underline{t}, \underline{t} \rangle)}) = \{(\underline{t}, \underline{t})\}$ . By construction,  $\sigma^*(h, (r^{**}, \langle \underline{t}, \underline{t} \rangle)) = g^*(\langle \underline{t}, \underline{t} \rangle)$ , implying consistency in Case 5. Associate belief  $p_{|h, (r^{**}, 0)} = p_{|h}$  to the off-equilibrium signal  $s = 0$ . Then, since  $(h, (r^{**}, 0)) \in H^0$ ,  $\sigma^*(h, (r^{**}, 0)) = g^*(0)$ , implying consistency in Case 2.

To see that  $\sigma^*$  is Bellman optimal:

- In Case 1,  $C^{\sigma^*}(h) = \widehat{C}^{\sigma^*}(h) = \emptyset$ , and hence  $\sigma^*(h) = (0, 0) \in X$  is a Bellman optimal choice in this case.
- In Case 2,  $v(g^* \circ r^*, p_{|h}) = (\bar{t} - \underline{t})p(\bar{t}, \underline{t}) \geq v(g^* \circ r, p_{|h})$ , for all budget balanced mechanisms  $g^* \circ r$ . Hence  $\sigma^*(h) = r^*$  is a Bellman optimal choice in this case.
- In Case 3,  $v(1_{(0,0)}, p_{|h}) = 0 = v(g^* \circ r, p_{|h})$ , for all budget balanced mechanisms  $g^* \circ r$ . Hence  $\sigma^*(h) = (0, 0)$  is a Bellman optimal choice in this case.
- In Case 4,  $v(g^* \circ r^{**}, p_{|h}) = 0 = v(g^* \circ r, p_{|h})$ , for all budget balanced mechanisms  $g^* \circ r$ . Hence  $\sigma^*(h) = r^{**}$  is a Bellman optimal choice in this case.
- In Cases 5, 6, 7, 8, and 9,  $v(\sigma^*(h), p_{|h}) = v(g^* \circ r, p_{|h})$  for all budget balanced mechanisms  $g^* \circ r$ , and hence  $\sigma^*(h) \in X$  is a Bellman optimal choice in this case.

■

To complete the description of  $\sigma^*$  that meets the conditions of Theorem 2, let off-equilibrium path  $h$  mechanism selection rule  $\sigma^*(h)$  be anything that would constitute a Bellman optimal and consistent rule in the continuation game. By Theorem 1, such a continuation strategy does exist. Since under  $p_{|\emptyset}$  the mechanism  $\sigma^*(\emptyset)$  is the second best, the planner has no incentive to deviate it in the first stage. Hence the constructed  $\sigma^*$  is Bellman optimal and consistent choice rule, starting from  $p_{|\emptyset}$

## A Appendix

### A.1 Omitted proofs of Section 5

**Proof of Proposition 2.** Denote

$$\begin{aligned} a_1(\theta_1) &= \sum_{\theta_2} p_2(\theta_2) a(\theta_1, \theta_2), \\ a_2(\theta_2) &= \sum_{\theta_1} p_1(\theta_1) a(\theta_1, \theta_2), \end{aligned}$$

and use the shorthand

$$\begin{aligned} V_1(\theta_1) &= \sum_{\theta_2} p(\theta_2 : \theta_1) [a(\theta_1, \theta_2)\theta_1 - m(\theta_1, \theta_2)], \\ V_2(\theta_2) &= \sum_{\theta_1} p(\theta_1 : \theta_2) [m(\theta_1, \theta_2) - a(\theta_1, \theta_2)\theta_2]. \end{aligned}$$

Denoting  $t'$  the immediate predecessor of  $t$ , incentive compatibility of a mechanism implies

$$\begin{aligned} a_1(\theta'_1)(\theta_1 - \theta'_1) &\leq V_1(\theta_1) - V_1(\theta'_1) \leq a_1(\theta_1)(\theta_1 - \theta'_1), \\ a_2(\theta'_2)(\theta_2 - \theta'_2) &\geq V_2(\theta'_2) - V_2(\theta_2) \geq a_2(\theta_2)(\theta_2 - \theta'_2). \end{aligned}$$

Thus  $a_1$  is increasing,  $a_2$  is decreasing, and

$$\begin{aligned} V_1(\theta_1) &\geq \sum_{t \leq \theta'_1} a_1(t) + V_1(0), \\ V_2(\theta_2) &\geq \sum_{t \geq \theta'_2} a_2(t) + V_2(1). \end{aligned}$$

Let

$$\begin{aligned} P_i(\theta_i) &= \sum_{t \leq \theta_i} p_i(t), \\ A_i(\theta_i) &= \sum_{t \leq \theta_i} a_i(t). \end{aligned}$$

Then

$$\begin{aligned} P_1(\theta_1)A_1(\theta_1) &= \sum_{t \leq \theta_1} P_1(t)[A_1(t) - A_1(t')] + \sum_{t \leq \theta_1} [P_1(t) - P_1(t')]A_1(t') \\ &= \sum_{t \leq \theta_1} P_1(t)a_1(t) + \sum_{t \leq \theta_1} p_1(t)A_1(t'). \end{aligned}$$

Thus

$$\begin{aligned} \sum_t p_1(t)A_1(t') &= P_1(1)A_1(1) - \sum_t P_1(t)a_1(t) \\ &= \sum_t a_1(t)(1 - P_1(t)). \end{aligned} \tag{26}$$

And similarly for the agent 2 :

$$\begin{aligned} \sum_t p_2(t)[A_2(1) - A_2(t')] &= A_2(1) - \sum_t p_2(t)A_2(t') \\ &= \sum_t a_2(t)P_2(t). \end{aligned} \tag{27}$$

The planner's problem can be written

$$\begin{aligned} & \max_{a(r(\cdot))} \sum_{\theta_1} \sum_{\theta_2} p_1(\theta_1) p_2(\theta_2) (\theta_1 - \theta_2) a(r(\theta_1, \theta_2)) \\ & \text{s.t.} \end{aligned}$$

$$\sum_{\theta_1} \sum_{\theta_2} p_1(\theta_1) p_2(\theta_2) (\theta_1 - \theta_2) a(r(\theta_1, \theta_2)) = \sum_{\theta_1} p_1(\theta_1) V_1(\theta_1) + \sum_{\theta_2} p_2(\theta_2) V_2(\theta_2) \quad (28)$$

$$A_1(\theta_1) \geq V_1(\theta_1) - V_1(0) \geq A_1(\theta'_1), \text{ for all } \theta_1 \quad (29)$$

$$A_2(1) - A_2(\theta_2) \geq V_2(\theta_2) - V_2(1) \geq A_2(1) - A_2(\theta'_2), \text{ for all } \theta_2 \quad (30)$$

$$V_1(\theta_1) \geq 0 \text{ for all } \theta_1, \quad V_2(\theta_2) \geq 0 \text{ for all } \theta_2 \quad (31)$$

where (28) is the ex ante budget balance condition, (29) and (30) are the incentive compatibility constraints, and (31) is the participation constraint.

Since the right hand side inequalities of (29) and (30) imply

$$\begin{aligned} \sum_{\theta_1} p_1(\theta_1) V_1(\theta_1) & \geq \sum_{\theta_1} p_1(\theta_1) A_1(\theta'_1) + V_1(0), \\ \sum_{\theta_2} p_2(\theta_2) V_2(\theta_2) & \geq \sum_{\theta_2} p_2(\theta_2) [A_2(1) - A_2(\theta'_2)] + V_2(1), \end{aligned}$$

(26), (27), and (28) result in

$$\begin{aligned} & V_1(0) + V_2(1) + \sum_{\theta_1} a_1(\theta_1) (1 - P_1(\theta_1)) + \sum_{\theta_2} a_2(\theta_2) P_2(\theta_2) \\ & \leq \sum_{\theta_1} \sum_{\theta_2} p_1(\theta_1) p_2(\theta_2) (\theta_1 - \theta_2) a(r(\theta_1, \theta_2)), \end{aligned}$$

or, more compactly,

$$\sum_{\theta_1} \sum_{\theta_2} p_1(\theta_1) p_2(\theta_2) \left[ \left( \theta_1 - \frac{1 - P_1(\theta_1)}{p_1(\theta_1)} \right) - \left( \theta_2 + \frac{P_2(\theta_2)}{p_2(\theta_2)} \right) \right] a(\theta_1, \theta_2) \geq 0. \quad (32)$$

Maximizing the objective function with respect to (32), and interpreting  $\gamma/(1-\gamma)$  as the Lagrange multiplier, gives the desired programme. Since, at the optimum, (32) holds as equality, the solution to the programme also meets the left hand side inequalities of (29) and (30). Since this implies that  $a_1$  is increasing and  $a_2$  is decreasing, it also follows that the participation constraint (31) is met whenever  $V_1(0) \geq 0$  and  $V_2(1) \geq 0$  which hold as equality at the optimum. Finally, the optimality of  $a^\gamma$  under regular  $p_1$  and  $p_2$  follows by maximizing the objective function pointwisely. ■

**Proof of Lemma 8.** Our task is to construct an  $m(\cdot)$  that prescribes zero monetary transfer when trade does not take place. That is

$$m(r(\theta_1, \theta_2)) = 0 \quad \text{whenever} \quad c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) < 0.$$

Denote by  $a^\gamma$  the incentive efficient allocation rule under Lagrange multiplier  $\gamma$ . Denote by  $m_1^\gamma$  and  $m_2^\gamma$  the implied expected transfers from 1 and to 2 :

$$\begin{aligned} m_1^\gamma(\theta_1) &= a_1^\gamma(\theta_1)\theta_1 - \sum_{t < \theta_1'} a_1^\gamma(t), \quad \text{for all } \theta_1 \in T \\ m_2^\gamma(\theta_2) &= a_2^\gamma(\theta_2)\theta_2 + \sum_{t > \theta_2'} a_2^\gamma(t), \quad \text{for all } \theta_2 \in T. \end{aligned}$$

The ex ante budget balance of the incentive efficient mechanism implies

$$\sum_{\theta_1} p_1(\theta_1)m_1^\gamma(\theta_1) = \sum_{\theta_2} p_2(\theta_2)m_2^\gamma(\theta_2). \quad (33)$$

Construct  $m(\cdot)$  such that

$$\begin{aligned} m(\theta_1, 0) &= m_1^\gamma(\theta_1), \quad \text{for all } \theta_1 < 1, \\ m(1, \theta_2) &= m_2^\gamma(\theta_2), \quad \text{for all } \theta_2 > 0, \\ m(\theta_1, \theta_2) &= 0, \quad \text{for all } (\theta_1, \theta_2) \text{ such that } \theta_1 < 1 \text{ and } \theta_2 > 0. \end{aligned}$$

To complete the description of  $m$ , let  $m(1, 0)$  satisfy

$$p_1(1)m(1, 0) + \sum_{t < 1} p_1(t)m_1^\gamma(t) = m_2^\gamma(0), \quad (34)$$

$$p_2(0)m(1, 0) + \sum_{t > 0} p_2(t)m_2^\gamma(t) = m_1^\gamma(1). \quad (35)$$

Then  $m(\cdot)$  prescribes zero transfer under no-trade and

$$\begin{aligned} m_1(\theta_1) &= m_1^\gamma(\theta_1), \quad \text{for all } \theta_1 \in T, \\ m_2(\theta_2) &= m_2^\gamma(\theta_2), \quad \text{for all } \theta_2 \in T. \end{aligned}$$

Thus  $m$  is consistent with the incentive efficient allocation  $a^\gamma$  rule. However, since a single variable  $\bar{m}$  is determined by two equations (34) and (35), we need to verify that a desired  $m(1, 0)$  does exist. The remainder of the proof establishes this.

First, fix any  $m(1, 0)$  that completes the description of  $m$ . Since the order of summation does not matter,

$$\sum_{\theta_1} p_1(\theta_1) \sum_{\theta_2} p_2(\theta_2)m(\theta_1, \theta_2) = \sum_{\theta_2} p_2(\theta_2) \sum_{\theta_1} p_1(\theta_1)m(\theta_1, \theta_2). \quad (36)$$

By construction,

$$\begin{aligned} \sum_{\theta_1} p_1(\theta_1) \sum_{\theta_2} p_2(\theta_2)m(\theta_1, \theta_2) &= \sum_{t < 1} p_1(t)m_1^\gamma(t) + p_1(1) \left( p_2(0)m(1, 0) + \sum_{t > 0} p_2(t)m_2^\gamma(t) \right) \\ \sum_{\theta_2} p_2(\theta_2) \sum_{\theta_1} p_1(\theta_1)m(\theta_1, \theta_2) &= p_2(0) \left( p_1(1)m(1, 0) + \sum_{t < 1} p_1(t)m_1^\gamma(t) \right) + \sum_{t > 0} p_2(t)m_2^\gamma(t) \end{aligned}$$

Now letting  $m(1, 0)$  be defined by (34), it follows that

$$\sum_{\theta_2} p_2(\theta_2) \sum_{\theta_1} p_1(\theta_1) m(\theta_1, \theta_2) = \sum_{\theta_2} p_2(\theta_2) m_2^\gamma(\theta_2).$$

By (36) and (33),

$$\sum_{\theta_1} p_1(\theta_1) \sum_{\theta_2} p_2(\theta_2) m(\theta_1, \theta_2) = \sum_{\theta_1} p_1(\theta_1) m_1^\gamma(\theta_1).$$

Thus  $m(1, 0)$  also satisfies (35). ■

## References

- [1] AGHION, P., DEWATRIPONT, M., AND P. REY (1994), Renegotiation design with unverifiable information, *Econometrica* **62**, 257-282.
- [2] AUSUBEL, L. AND DENECKERE, R. (1989), A direct mechanism characterization of sequential bargaining with one-sided incomplete information, *Journal of Economic Theory* **48**, 18-46
- [3] AUSUBEL, L. AND DENECKERE, R. (1993), Efficient sequential bargaining, *Review of Economic Studies* **60**, 435-461.
- [4] BALIGA, S., CORCHÓN, L. AND SJÖSTRÖM, T. (1997). The theory of implementation when the planner is a player, *Journal of Economic Theory* **77**, 15-33.
- [5] BERGEMANN, D AND MORRIS, S. (2005), Robust mechanism design, *Econometrica* **73**, 1771-1813
- [6] BESTER, H. AND STRAUZ, R. (2001), Contracting with imperfect commitment and the revelation principle: the single agent case, *Econometrica* **69**, 1077-98
- [7] COASE. R. (1972), Durability and monopoly, *Journal of Law and Economics* **15**, 143-9.
- [8] DEWATRIPONT, M. (1989), Renegotiation and information revelation over time: the case of optimal labor contracts, *Quarterly Journal of Economics* **104**, 589-619.
- [9] EVANS, R. (2012), Mechanism design with renegotiation and costly messages, *Econometrica*.**80**, 2089-2104
- [10] FORGES, F. (1994), Posterior efficiency, *Games and Economic Behavior* **6**, 238-61.
- [11] FREIXAS, X., GUESNERIE, R., AND TIROLE, J. (1985), Planning under incomplete information and the ratchet effect, *Review of Economic Studies* **52**, 173-92

- [12] GERARDI, D., HÖRNER, J., AND L. MAESTRI (2011), The role of commitment in bilateral trade, manuscript, Collegio Carlo Alberto
- [13] GREEN, J. AND LAFFONT, J.-J.. (1987), Posterior implementability in a two-player decision problem, *Econometrica* **55**, 69-94
- [14] GRESIK, T. (1991), Ex ante efficient, ex post individually rational trade, *Journal of Economic Theory* **53**, 131-45.
- [15] GRESIK, T. (1996), Incentive efficient equilibria of two-party sealed-bid bargaining games, *Journal of Economic Theory* **68**, 26-48
- [16] HOLMSTROM, B. AND MYERSON, R. (1983), Efficient and durable decision rules with incomplete information, *Econometrica* **51**, 1799-819,
- [17] KIYOTAKI, N. (2011), A mechanism design approach to financial frictions, manuscript, Princeton University
- [18] KRISHNA, V. AND PERRY, M. (1998), Efficient mechanism design, working paper, Penn State University.
- [19] LAFFONT, J.-J. AND TIROLE, J. (1990), Adverse selection and renegotiation in procurement , *Review of Economic Studies* **57**, 597-625
- [20] LAGUNOFF, R. (1995), Resilient allocation rules for bilateral trade, *Journal of Economic Theory* **66**, 463-87
- [21] LAGUNOFF, R. (1992), Fully endogenous mechanism selection on finite outcome sets, *Economic Theory* **2**, 465-80
- [22] MCADAMS, D. AND SCHWARZ, M. (2006), Credible sales mechanisms and intermediaries, *American Economic Review* forthcoming
- [23] MCAFEE, P. AND VINCENT, D. (1997), Sequentially optimal auctions, *Games and Economic Behavior* **18**, 246-76
- [24] MYERSON, R. (1991), *Game Theory, Analysis of Conflict*, Cambridge MA, Harvard University Press
- [25] MYERSON, R., (1979), Incentive compatibility and the bargaining problem, *Econometrica* **47**, 61-73.
- [26] MYERSON, R. (1982), Optimal coordination mechanisms in generalized principle-agent problems, *Journal of Mathematical Economics* **28**, 67-81.
- [27] NEEMAN Z. AND G. PAVLOV (2012), Renegotiation-proof mechanism design, forthcoming in *Journal of Economic Theory*.
- [28] SEGAL, I., AND M. WHINSTON (2002), The Mirrlees approach to mechanism design with renegotiation *Econometrica* **70**, 1-45.

- [29] SKRETA V. (2011), Optimal auction design under non-commitment, manuscript, NYU Stern School of Business
- [30] SKRETA V., (2006), Sequentially optimal mechanisms, *Review of Economic Studies* **73**
- [31] VARTIAINEN, H. (2010), Auction design without commitment, forthcoming in *Journal of the European Economic Association*