# Mechanism Design without Commitment; General Solution and Application to Bilateral Bargaining* †

Hannu Vartiainen‡

HECER and University of Helsinki

May 5, 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.

In the literature on commitment in mechanism design, the usual approach is to appeal to the *incrutability 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.[1]

Renegotiation-proofness is not, however, an entirely satisfactorily 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. 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 inscrutability principle, our research strategy is to reduce, after each history, the continuation equilibria to a single direct

---

[1]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).

incentive compatible mechanism.[2]  In order to do this, we separate the two tasks of 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.[3]  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. We appeal to the $\varepsilon-$optimality criterion that impose an $\varepsilon-$cost for deviations. Such costs can be interpreted as a consequence of redesign of the mechanism. 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 an $\varepsilon-$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 *in*dependent 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

---

[2]Incrutability principle: any equilibrium of the mechanism selection can be represented as a direct single stage mechanism that is truthfully played and obediently implemented.

[3]Assuming public signals restrics away private communication. This is a simplification. See Skreta (2006, 2010) for analyses of mechanism design without commitment but with private communication.

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

Our second result studies mechanisms that can be implemented wihtout 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 affirmitive: 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 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 (2995) 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 critical question in the context of bargaining. The famous Coase Theorem asserts that, in the absence of commitment, the uniformed 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 Holmstöm 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 $\{1, ..., n\}$ 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 $\Theta_{-i} = \times_{j \neq i} \Theta_i$ with a typical element $\theta_{-i} = (\theta_j)_{j \neq i}$.[4] 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 $\text{supp}(p)$.[5] Agent $i$'s vNM utility functions $u$ is of form $u_i : X \times \Theta_i \to \mathbb{R}$.

---

[4]That is, $p_i(\theta_i) = \sum_{\theta_{-i}} p(\theta_i, \theta_{-i})$.

[5]$\text{supp}(p) = \{\theta : p(\theta) > 0\}$.

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* $f : \Theta \to \Delta X$ that specifies a probability distribution over outcomes for each type profile. Denote by

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

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 $f : \Theta \to \Delta X$, agent $i$'s expected payoff is

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

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

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

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. We separate these tasks. A mechanism is then a composite function

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

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

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

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 $\hat{\theta} = (\hat{\theta}_1, ..., \hat{\theta}_n)$, 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 $r$.[6]

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

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

---

[6] 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.

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 $r$.

Denote in particular by $1_s$ a constant information processing device that generates signal $s \in S$ with probability one under all type profiles. That is,

$$1_s = r \text{ such that } r(s : \theta) = 1, \text{ for all } \theta \in \Theta.$$

Assuming that the agents report their types truthfully, the signal $s$ generated by a mechanism $r$ induces a posterior distribution $p(r, s) \in \Delta\Theta$ such that, whenever $s \in r(\text{supp}(p))$,

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

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

Many composite mechanisms induce the same outcome function. A particular example of a composite mechanims 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.

## 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 $f = g \circ r : \Theta \to \Delta X$. 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

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

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

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

$$(1)$$

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

$$IC(p) = \{r \in \Phi : g \circ r \text{ is incentive compatible under } p\}$$

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 all not equivalent in terms of the planner's incentives. Different mechansims reveal different amount of information to the planner and hence may affect her strategic possibilties 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

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

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 $\sigma$ is defined ion $H$ and has to satisfy

$$\sigma(h) \in IC(p_{|h}), \quad \text{for all } h \in H. \tag{2}$$

One the one hand, $\sigma(h) \in IC(p_{|h})$ means that the planner continues the game by designing a new mechanism given the belief $p$ and history $h$. One the other hand, $\sigma(h, (r, s)) = 1_{g(s)}$ means that the planner does not change the suggested outcome; she chooses the constant mechanism that implements $g(s)$. 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 $\sigma$ would be followed ex post:

$$C^\sigma(h) = \left\{ r \in IC(p_{|h}) : \sigma(h, (r, s)) = 1_s, \ \text{for all} \ \ s \in r(\Theta) \right\}.$$

Choice set $C^\sigma(h)$ is 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 (Strategic Consistency)** *Choice rule $\sigma$ is* stregically consistent *if $\sigma(h) \in C^\sigma(h)$, for all $h \in H$.*

Given a mechanism $r$ and a choice rule $\sigma$, denote by $r \circ \sigma$ the mechanism that results from first truthfully playing $r$ and then following $\sigma$. That is

$$(r \circ \sigma)(s : \theta) = \sum_{s' \in S} r(s' : \theta)\sigma(h, (r, s'))(s : \theta), \text{ for all } s \in S, \text{ for all } \theta \in \Theta.$$

This compound mechanism need not be incentive compatible. However, *if* it is, then the planner could first use $r$ and then follow $\sigma$ in order to implement $r \circ \sigma$. Since it should matter whether she implemens $r \circ \sigma$ directly or via first employing $r$ and then following $\sigma$, the mechanism $r \circ \sigma$ should belong to the choice set of the planner at $h$.

**Definition 2 (Structural Consistency)** *Choice rule $\sigma$ is structurally consistent if $r \circ \sigma \in IC(p_{|h})$ implies $r \circ \sigma \in C^{\sigma}(h)$, for all $r \in IC(p_{|h})$, for all $h \in H$.*

Our final condition reflects local optimality.

**Definition 3 ($\varepsilon$−Optimality)** *For $\varepsilon > 0$, the choice rule $\sigma$ is $\varepsilon$−optimal if $v(\sigma(h), p_{|h}) \geq v(r, p_{|h}) - \varepsilon$, for all $r \in C^{\sigma}(h)$, for all $h \in H$.*

In words, for any chosen mechanism of the planner there should be no other mechanism that is more profitable for her than what she can commit to, given the $\varepsilon$−cost of changing the strategy. $\varepsilon$−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 $\sigma$. Conversely, without $\varepsilon$−optimality, $\sigma$ could not be convinsingly committed to since the planner is able to make a reliable one-time deviation at some history.

The $\varepsilon$−cost has natural interpretation in the cost of changing the planned mechanism. Imagine that planner declares which mechanism she implements after all contingencies. If there is an $\varepsilon$−cost of reneging from this agreement, the mechanism that she can reliably commit to are characterized by the $\varepsilon$−optimality condition.

The notion of ■-equilibria is important in the general theory of stochastic games. For example, there are simple examples of stochastic games that do not Nash nor subgame perfect Nash equilibrium equilibrium but do possess an ■-equilibrium for any ■ strictly bigger than 0. Consequently, existence results in many natural classes of stochastic games require $\varepsilon$−threshold for deviations (see e.g. Flesch et al 2010).

**Interpretation: A redesing game**   We now extend the single stage mechanism selection game by allowing the planner to redesign the mechanism repeatedly, once an outcome has been realized. We argue that a Bellman optimal and consistent mechanism selection strategy $\sigma$ can be interpreted as reduced form expression of a weak Perfect Bayesian Equilibrium (PBE) of a natural mechanism design game, that reflects the planner's inability to commit to the rules.

Consider the following multistage game:

At each stage $t = 0, 1, 2, ...$, the planner announces a mechanism $r^t$, where $r^t : \Theta \to \Delta S$.[7] Given $r^t$, the agents choose a message profile $m \in \Theta$ which produces a lottery over messages $r^t(\cdot : m) \in \Delta S$. 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 play to period $t + 1$, in which case the game repeats itself.[8]

This is a proper extensive form game with incomplete information and imperfect monitoring. The game has many (weak) PBEa, where (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

---

[7] Allowing messages spaces larger than $\times_i \Theta_i$ would not affect the results.

[8] 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.

beliefs in each period after observing the signal that this period's information processing device generates. Off-equilibrium beliefs are chosen to justify the optimization behavior in (ii).

In particular there is always the babbling equilibrium where messages have no meaning, information does not accumulate, and the planner's decisions are made under the initial belief. However, such an equilibrium fails to be responsive to the planner's ability induce truthtelling behavior, whenever such behavior is consistent with incetnives - the basic doctrine in the Myersonian mechanism design literature. The aim of the current model to capture the consistency conditions that specifiy when can a mechanism deommmitted to, even when the agents are responsive to the echanism and the planner cannot commit not to exploit theirninforation.

Note that the assumption of consistency is consistent with incentive compatibility. Under the hypothesis that the designer follows the choice rule $\sigma$ :

- if $r \in C^\sigma(h)$, then, since the mechanism $r$ will *not* be redesigned ex post and $r$ is incentive compatible, truthful reporting in $r$ can be sustained in equilibrium.

- if $r \notin C^\sigma(h)$, then, since the mechanism $r$ will be redesigned ex post and the choice rule is structurally consistent, truthful reporting in $r$ cannot be sustained in equilibrium.

Note that restricting attention to mechanism design stratgy $\sigma$ that, at each stage, implements a mechanism truthfully (or implements obediently the recommendation of the last stage mechanism) is without loss of generality. If, at history $h$, a public perfect equilibrium that eventually induces a type $\theta$ conditioned lottery over final outcomes and posterior beliefs, then there is an incentive compatible one-stage mechanism that induces the same type dependent lotteries over outcomes and posterior beliefs that generates at least as high payoff to the planner. This follows by the revelation principle and the fact that the set $S$ of signals is infinite: finitely many consecutive mechanism have the physical ability to ganerate exactly the same posterior beliefs as a single stage mechanism. 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.

## 4  Existence

We now state the main result of the paper.

**Theorem 1** *For any $\varepsilon > 0$, there is a strategically consistent, structurally consistent, and $\varepsilon-$optimal mechanism design strategy $\sigma$.*

The remainder of this section is devoted to proving the result. First, let

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

and denote
$$\bar{v}(p) = v(x(p), p)$$

Use the notation $P(r, p)$ the set of posterior beliefs, generated with positive probability by the mechanism $r$ from the initial beliefs $p$ :
$$P(r, p) = \cup_{s \in r(\text{supp}(p))} \{p_{|(r,s)}\}.$$

We will prove the existence via a series of subresults. The proof will rely on a fixed point argument. Wirst we develop an iterative procedure that rules out beliefs under which we can be sure that the planner cannot implement a constant mechanism.

Denote by $\rho$ a *response plan* $\rho : S \to R$ such that $\rho(s) \in IC(p_{|(r,s)})$ for all $s \in r(\Theta)$. Then denote by $r \circ \rho$ the compound mechanism such that
$$(r \circ \rho)((s, s') : \theta) = r(s : \theta)\rho(s)(s' : \theta), \text{ for all } (s, s') \in S^2, \text{ for all } \theta \in \Theta.$$

In particular, for any $s \in S$, denote by $\rho_s$ the *simple* response plan with a property that $\rho_s(s') = 1_{s'}$ if $s' \neq s$.

Let $B$ be a subset of beliefs, i.e. $B \subseteq \Delta$. Given belief $p$, we say that a response plan $\rho$ $B-blocks$ mechanism $r \in IC(p)$ if

- $v(\rho(s), p_{|(r,s)}) - \varepsilon \geq \bar{v}(p_{|(r,s)})$, for all $s$

- $P((r \circ \rho), p) \subseteq B$,

- $(r \circ \rho) \notin IC(p)$ or $v((r \circ \rho), p) - \varepsilon < \bar{v}(p)$.

Further, we say that a mechanism $r$ $B-covers$ belief $p$ if

- if $v(r, p) - \varepsilon \geq \bar{v}(p)$ and $r \in IC(p)$,

- $P(r, p) \subseteq B$,

- there is no *simple* response plan $\rho_s$ that $B-blocks$ $r$, given $p$.

That is, when a mechanism $r$ is blocked by a response plan $\rho_s$, choosing first $r$ and then continuing with $\rho_s$ is incentive compatible only if the eventual expected payoff for the planner is less profitable than implemetning a constant mechanism under $p$. Thus choosing $r$ cannot be justified under the hopothesis that $\rho_s$ is employed ex post. Conversely, when a mechanism $r$ covers belief $p$, there is no such way to punish the planner. Thus the planner cannot commit to implemetning a constant mechanism in the beginning.

Denote by $UC(B)$ the *uncovered* beliefs under $B$, i.e. ones that are not covered under $B$ by any $r$. Let $UC(\Delta) = UC^0$ and $UC^{t+1} = UC(UC^t)$ for all $t$.

By construction $UC^t \subseteq UC^{t+1}$ for all $t$. Denote
$$UC^* = \cap_t UC^t.$$

Then $UC^*$, the *ultimate uncovered set,* has the property that no belief in $p$ is $UC^*-$covered by any mechanism $r \in IC(p)$.

**Lemma 1** *The set $UC^*$ is nonempty.*

**Proof.** It suffices to show that there is $p$ that is not covered in any iteration of $UC^t$. Take a degenerate distribution $p_\theta$ such that $p_\theta(\theta) = 1$ and $p_\theta(\theta') = 1$ for all $\theta' \neq \theta$. We show that $p_\theta$ is not covered in any iteration. To this end, for any $r$,

$$v(r, p_\theta) = v(r, \theta) = \sum_s r(s:\theta)v(g(s), \theta) \leq v(x(p_\theta), \theta) = \bar{v}(p_\theta).$$

Thus the items in definition covering cannot be met by any $r$. ∎

We first state that the covering operation iterates.

**Lemma 2** *Let $r$ $B-cover$ $p$ and let $\rho_s(s)$ $B-cover$ $p_{|(r,s)}$. Then $(r \circ \rho_s)$ $B-covers$ $p$.*

**Proof.** Since $r$ $B-$covers belief $p$, any response plan $\rho_s$ such that $v(\rho_s(s), p_{|(r,s)}) - \varepsilon \geq \bar{v}(p_{|(r,s)})$ and $P((r \circ \rho), p) \subseteq B$ also satisfies $(r \circ \rho_s) \in IC(p)$ and $v((r \circ \rho), p) - \varepsilon \geq \bar{v}(p)$. ∎

We now argue that any $p$ not in $UC^*$ is not only $UC^t-$covered by some $r^t$ for all $t$ but also $UC^*-$covered by some $r$. This property guarantees that once a belief outside $UC^*$ is reached, there a reliable way to invoke a mechanism that induces beliefs back in $UC^*$.

**Lemma 3** *If $p \notin UC^*$, then there is $r$ that $UC^*-covers$ $p$.*

**Proof.** Let $p$ be $UC^t-$covered, i.e. $p \notin UC^{t+1}$. Assume that there is no $r$ that $UC^t-$covers $p$, for some $t$. We prove a contradiction. By repeatedly applying Lemma 2, we can construct a sequence $\{\rho_{s^k}^k\}$ of response plans such that $(r \circ \rho_{s^0}^0 \circ ... \rho_{s^k}^k)$ $UC^t-$covers $p$, for all $k = 0, 1, ...$ .

Denote $(r \circ \rho_{s^0}^0 \circ ... \rho_{s^k}^k) = \bar{\rho}^k$. We now claim that there is a mechanism $\bar{\rho}(p)$ such that $\bar{\rho}^k(p) \to \bar{\rho}(p)$. Then also $\bar{\rho}(p)$ $UC^t-$covers $p$. To this end, denote by $\bar{S}_k \subseteq S$ the set of absorbing signals under $\bar{\rho}^k(p)$, i.e.

$$\bar{S}_k = \{s : p_{|(\bar{\rho}^k, s)} \notin UC^{t+1}\}.$$

Let the probability $\beta^k$ be defined by

$$\beta^k = \sum_\theta \sum_{s \in \bar{S}_k} p(\theta) \bar{\rho}^k(p)(s:\theta).$$

Since, by the construction of $\bar{\rho}^k(p)$,

$$P(\bar{\rho}^k(p), p) \cap UC^{t+1} \subseteq P(\bar{\rho}^{k+1}(p), p) \cap UC^{t+1}, \text{ for all } k = 0, 1, ... , \tag{3}$$

also

$$\beta^k \geq \beta^{k+1}, \quad \text{ for all } k = 0, 1, ... . \tag{4}$$

13

It now suffices that $\beta^k \to_k 0$.

The payoff generated by $\rho_t^k(p)$ satisfies

$$
\begin{aligned}
v(\rho_t^k(p), p) &= \sum_\theta \sum_s p(\theta)\, v(g(s), \theta) \rho_t^k(s:\theta) \\
&= \sum_\theta \sum_s p(\theta)\, v(g(s), \theta) \rho_t^{k-1}(s:\theta) + \sum_\theta \sum_s p(\theta)\, \rho^{k-1}(p)(s:\theta) \left[ v(\rho(p)(p_{|(r \circ \rho^{k-1}, s)}), p_{|(r \circ \rho^{k-1}, s)}) \right. \\
&\geq \sum_\theta \sum_s p(\theta)\, v(g(s), \theta) \rho_t^{k-1}(s:\theta) + \sum_\theta \sum_{s \in \bar{S}_k} p(\theta)\, \rho^{k-1}(p)(s:\theta)\varepsilon \\
&= \sum_\theta \sum_s p(\theta)\, v(g(s), \theta) \rho_t^{k-1}(p)(s:\theta) + \beta^k \varepsilon,
\end{aligned}
$$

where the inequality follows from the defition of covering. Hence, by induction on $k$,

$$
\begin{aligned}
v(\rho_t^k(p), p) &\geq \sum_\theta \sum_{s \in S^*} p(\theta)\, v(g(s), \theta) \rho(p)(s:\theta) + \left( k - \sum_{\ell=0}^k \beta^\ell \right) \varepsilon \\
&\geq v(\rho(p), p) + k\beta^k \varepsilon,
\end{aligned}
$$

where the second inequality follows from (4). Since $v(\rho_t^k(p), p)$ and $v(\rho(p), p)$ lie in a bounded space, and since $\varepsilon > 0$, the second term of the final expression must converge to a finite limit as $k$ tends to infinity. Hence $\beta^k \to 0$, as desired.

Since (**??**) holds, the mechasnism $\rho_t^\infty(p)$ is well defined, $UC^t$−covers $p$, and has the property

$$
P(\rho_t^\infty, p) \subseteq UC^{t+1}.
$$

By induction on $t$,

$$
\rho_t^\infty \circ \dots \circ \rho_\infty^\infty
$$

$UC^*$−covers any $p \notin UC^*$. ∎

We now construct a mechanism design strategy that meets the two desiderata by using the notions of $UC^*$ and $UD(UC^*)$. We first identify a way to punish the planner when the beliefs are in the set $UC^*$.

**Lemma 4** *Let $p \in UC^*$ and $r \in IC(p)$. Then there is $\rho_s$ that $UC^*$−blocks $r$.*

**Proof.** By the construction of $UC^*$, there is no $r$ that $UC^*$−covers $p$. Thus $r$ is blocked by a response plan. ∎

We now partition the set $H$ of histories into distinct "phases" $H_{p,r}$. To this end, we first define two operators. For any pair $(p, r)$ such that $p \in UC^*$, $v(r, p) - \varepsilon \geq \bar{v}(p)$, $r \in IC(p)$, and $P(r, p) \subseteq UC^*$, identify a simple response plan $\rho_s^{p,r}$ that $UC^*$−blocks $r$ given $p$, such that there is no other simple response plan $\rho_s$ that $UC^*$−blocks $r$ with the property that

$$
v(\rho_s, p) - \varepsilon \geq v(\rho_s^{p,r}, p).
$$

14

By Lemma 4, and since the planner's payoffs are bounded, such a $\rho_s^{p,r}$ exists.

Similarly, for each $p \notin UC^*$, identify a mechanism $r^p$ that $UC^*-$covers $p$, such that there is no other mechanism $r$ that $UC^*-$covers $p$ with the property that

$$v(r,p) - \varepsilon \geq v(r^p, p).$$

By Lemma 3, and since the planner's payoffs are bounded, such a $r^p$ exists.

Finally, for any pair $(p,r)$ such that $p \in UC^*$, construct a response plan $\rho^{p,r}$ that $UC^*-$blocks $r$ given $p$ such that there is no other response plan $\rho$ that $UC^*-$blocks $r$ with the property that

$$v(\rho(s), p_{|(r,s)}) - \varepsilon \geq v(\rho^{p,r}(s), p_{|(r,s)}).$$

To see why such a response plan exists, note that at each $p_{|(r,s)} \notin UC^*$ there is $r^{p_{|(r,s)}}$ that $UC^*-$covers $p_{|(r,s)}$.

Let $\emptyset \in H_{p^0, 1_{s^0}}$ for some $s^0$ and, for any $h \in H_{p,r}$, let

$$(h, (r', s)) \in \begin{cases} H_{p_{|h,(r',s)}, 1_s}, & \text{if } P(r', p_{|h}) \subseteq UC^* \text{ and } (r \circ r') \notin IC(p) \text{ or } v(r \circ r', p) - \varepsilon \leq \bar{v}(p), \\ H_{p, r \circ r'}, & \text{otherwise.} \end{cases}$$

(5)

Recursively,

$$h \in H_{p,r} \text{ implies that there is } s \text{ such that } p_{|(r,s)} = p_{|h}$$

Construct a *mechanism design strategy* $\sigma$ that is measurable with respect to the partition $\{H_{p,r}\}_{p,r}$ such that, for any $h \in H_{p,r}$,

$$\sigma(h, (r', s')) = \begin{cases} 1_s, & \text{if } P(r', p_{|h}) \subseteq UC^* \text{ and } (r \circ r') \notin IC(p) \text{ or } v(r \circ r', p) - \varepsilon \leq \bar{v}(p), \\ \rho_s^{p,r}(s), & \text{otherwise.} \\ r^{p_{|h,(r',s')}}, & \text{if } p_{|h,(r',s')} \notin UC^* \end{cases}$$

(6)

That is, the mechanism design strategy depends on the phase, the current belief of the planner, and the status quo outcome. By Lemma **??**, $G^*$ 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 (5) and (6) is structurally consistent, strategically consistent, and $\varepsilon-$optimal.*

**Proof.** First we describe the relevant choice sets for each history history $h' = (h, (r', s))$. Let $h \in H_{p,r}$.

1. If $P(r', p_{|h}) \subseteq UC^*$ and $(r \circ r') \notin IC(p)$ or $v(r', p) - \varepsilon \leq \bar{v}(p)$, then

$$C^\sigma(h') = \left\{ r'' \in IC(p_{|h'}) : P(r'', p_{|h'}) \subseteq UC^* \text{ and } v(1_s \circ r'', p_{|h'}) - \varepsilon \leq \bar{v}(p_{|h'}) \right\}.$$

2. Otherwise,

$$C^\sigma(h') = \left\{ r'' \in IC(p_{|h'}) : \begin{array}{l} P(r'', p_{|h'}) \subseteq UC^*, \text{ and} \\ (r \circ r' \circ r'') \notin IC(p) \text{ or } v(r \circ r' \circ r'', p) - \varepsilon \leq \bar{v}(p) \end{array} \right\}.$$

15

We then check that $\sigma$ is strategically consistent:

1. If $P(r', p_{|h}) \subseteq UC^*$ and $v(r', p) - \varepsilon \leq \bar{v}(p)$, then, since $p_{|h'} \in P(r', p_{|h})$, $P(1_s, p_{|h'}) = \{p_{|h'}\} \subseteq UC^*$ and $v(1_s \circ 1_s, p_{|h'}) - \varepsilon \leq \bar{v}(p_{|h'})$. Hence $1_s \in C^\sigma(h')$.

2. If not $P(r', p_{|h}) \subseteq UC^*$ and $v(r', p) - \varepsilon \leq \bar{v}(p)$, then, since $\rho_{p, r \circ r'}$ $UC^*-$blocks $r$ given $p$, $P(\rho_{p, r \circ r'}(s), p) \subseteq UC^*$ and $(r \circ r' \circ \rho_{p, r \circ r'}(s)) \notin IC(p)$ or $v(r \circ r' \circ \rho_{p, r \circ r'}(s), p) - \varepsilon \leq \bar{v}(p)$. Hence $\rho_{p, r \circ r'}(s) \in C^\sigma(h')$.

To prove that $\sigma$ is $\varepsilon-$optimal, we proceed case by case.

1. If $P(r', p_{|h}) \subseteq UC^*$ and $v(r', p) - \varepsilon \leq \bar{v}(p)$, then $v(r'', p_{|h'}) - \varepsilon \leq \bar{v}(p_{|h'})$ for all $r'' \in C^\sigma(h')$. Thus $1_s \in \arg\max_{s'} v(1_{s'}, p_{|h'}) = \bar{v}(p_{|h'})$.

2. If not $P(r', p_{|h}) \subseteq UC^*$ and $v(r', p) - \varepsilon \leq \bar{v}(p)$, then, by the construction of $\rho_{p, r \circ r'}(s)$, there is no $\rho'$ that $UC^*-$blocks $r$ such that

$$\rho' = \begin{cases} r' \text{ s.t. } v(r', p_{|h, (r \circ r', s')}) - \varepsilon > v(\rho_{p, r \circ r'}(s'), p_{|h, (r \circ r', s')}), & \text{if } s' = s, \\ \rho_{p, r}(s'), & \text{if } s' \neq s. \end{cases}$$

Thus there is no $r' \in C^\sigma(h')$ such that $v(r', p_{|h, (r \circ r', s)}) - \varepsilon > v(\rho_{p, r \circ r'}(s), p_{|h, (r \circ r', s)})$.

Finally, we conclude that $\sigma$ is structurally consistent. Consider $r'' \notin C^\sigma(h')$. Then

$$C^\sigma(h', (r'', s')) = \left\{ r''' \in IC(p_{|h', (r'', s')}) : \begin{array}{l} P(r''', p_{|h', (r'', s')}) \subseteq UC^*, \text{ and} \\ (r'' \circ r''') \notin IC(p_{|h'}) \text{ or } v(r'' \circ r''', p_{|h'}) - \varepsilon \leq \bar{v}(p_{|h'}) \end{array} \right\}.$$

Thus

$$\begin{aligned} & \left\{ (r'' \circ r''') \in IC(p_{|h'}) : r''' \in C^\sigma(h', (r'', s')) \right\} \\ \subseteq \ & \left\{ (r'' \circ r''') \in IC(p_{|h'}) : P((r'' \circ r''), p_{|h'}) \subseteq UC^* \text{ and } v(r'' \circ r'', p_{|h'}) - \varepsilon \leq \bar{v}(p_{|h'}) \right\} \\ \subseteq \ & C^\sigma(h'), \end{aligned}$$

implying structural consistency. ∎

## 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 strenghtened.

**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

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

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[9]

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

Assume that the set of feasible outcomes

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

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.[10] 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

$$\sum_{\theta_{-i}} p(\theta) \left[ \sum_{s} \tilde{u}_i(g(s), \theta_i) r(s : \theta) - \sum_{s} \tilde{u}_i(g(s), \theta_i) r(s : \theta_{-i}, \theta'_i) \right] \geq 0, \quad (7)$$

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

---

[9] *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).

[10] Veto-incentive compatibility is due to Forges (1998), and is closely related to IC* of Matthews and Postlewaite (1989).

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).[11]

# 5 Application: Bilateral Bargaining

Since Myerson and Satterthewaite (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 $\Theta_1 = \Theta_2 = \left\{0, \frac{1}{K}, ..., \frac{K-1}{K}, 1\right\}$, 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 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 $(a, m) : \Theta \to \{0,1\} \times \mathbb{R}$. 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 distrbution $p$, 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

$$
\sum_{\theta_2} p\left(\theta_2 : \theta_1\right) [a(\theta'_1, \theta_2)\theta_1 - m(\theta'_1, \theta_2)],
$$

$$
\sum_{\theta_1} p\left(\theta_1 : \theta_2\right) [m(\theta_1, \theta'_2) - a(\theta_1, \theta'_2)\theta_2].
$$

A direct mechanism is incentive compatible if

$$
\sum_{\theta_2} p\left(\theta_2 : \theta_1\right) [a(\theta)\theta_1 - m(\theta)] \geq \sum_{\theta_2} p(\theta_2 : \theta_1)[a(\theta'_1, \theta_2)\theta_1 - m(\theta'_1, \theta_2)], \text{ for all } \theta_1, \theta'_1 \in \Theta_1,
$$

$$
\sum_{\theta_1} p\left(\theta_1 : \theta)\right)[m(\theta) - a(\theta)\theta_2] \geq \sum_{\theta_1} p(\theta_1 : \theta_2)[m(\theta_1, \theta'_2) - a(\theta_1, \theta'_2)\theta_2], \text{ for all } \theta_2, \theta'_2 \in \Theta_2,
$$

---

[11] Choose $\theta_i = \theta'_i$ in (7). We only need EXP-IR and IC in the remainder of the paper.

and it is (interim) individually rational if

$$\sum_{\theta_2} p\left(\theta_2 : \theta_1\right)\left[a(\theta)\theta_1 - m(\theta)\right] \geq 0, \text{ for all } \theta_2 \in \Theta_2,$$

$$\sum_{\theta_1} p\left(\theta_1 : \theta_2\right)\left[m(\theta) - a(\theta)\theta_2\right] \geq 0, \quad \text{for all } \theta_1 \in \Theta_1$$

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 absolutly 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 starategy 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 indepedently 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 distrbuted 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 $\gamma \in (0, 1]$,*

$$if \quad c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) \geq 0, \quad then \quad a^\gamma(\theta_1, \theta_2) = 1 \tag{8}$$

$$if \quad c_1(\theta_1, \gamma) - c_2(\theta_2, \gamma) < 0, \quad then \quad a^\gamma(\theta_1, \theta_2) = 0. \tag{9}$$

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 (8)-(9). Then, for any $(\theta_1, \theta_2) \in \Theta_1 \times \Theta_2$,*

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

*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 (9).

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 5** *Let $\theta_1$ and $\theta_2$ be independently distrbuted 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 (8)-(9) such that the transfer rule $m^\gamma$ prescribes zero monetary transfer when trade does not take place, i.e.*

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

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 consistent and Bellman optimal 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

$$S^* = \{\langle \theta_1, \theta_2 \rangle : \theta_1 \geq \theta_2\} \cup \{0\}. \tag{10}$$

Consider the following information processing device $r^* : \Theta_1 \times \Theta_2 \to \Delta S^*$. For any $t$,

let $\kappa(t) = \#\{t' : t \geq t'$ and $c_1(t,\gamma) \leq c_2(t',\gamma)$ or $t \geq t'$ and $c_1(t',\gamma) \leq c_2(t,\gamma)\}$. Then

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

(11)

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 $g^* : S^* \to \{0,1\} \times \mathbb{R}$ such that, for any $s \in S^*$,

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

(12)

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

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

Hence, by Lemma 5,

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

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

**Lemma 6** *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 $\mathrm{supp}(p_{|h}) = \{(\bar{t},\bar{t}),(\bar{t},\underline{t}),(\underline{t},\underline{t})\}$. 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 in on induction on the cardinality of $\mathrm{supp}(p_{|h'})$. First, let $g^* : S^* \to \{0,1\} \times \mathbb{R}$ be defined by

$$
g^*(s) = \begin{cases}
(1,\bar{t}), & \text{if } s = \langle \bar{t},\bar{t}\rangle, \\
(1,(\bar{t}+\underline{t})/2), & \text{if } s = \langle \bar{t},\underline{t}\rangle, \\
(1,\underline{t}), & \text{if } s = \langle \underline{t},\underline{t}\rangle, \\
(0,0), & \text{if } s = 0.
\end{cases}
\tag{13}
$$

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

$$
h' \in \begin{cases}
H^1, & \text{if } h \in H^0 \text{ and } \{(\bar{t},\bar{t}),(\underline{t},\underline{t})\} \subseteq \mathrm{supp}(p_{h'}), \\
H^0 & \text{otherwise.}
\end{cases}
\tag{14}
$$

Construct a choice rule $\sigma^*$ such that, for any $h' = (h,(r,s)) \in H$,

$$
\sigma^*(h') = \begin{cases}
\begin{cases}
\text{stop}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\bar{t}),(\bar{t},\underline{t}),(\underline{t},\underline{t})\}, s=0, \text{ and } h \in H^0, \\
1_0, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\bar{t}),(\bar{t},\underline{t}),(\underline{t},\underline{t})\}, s \neq 0, \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 } \mathrm{supp}(p_{|h}) = \{(\bar{t},\bar{t}),(\bar{t},\underline{t}),(\underline{t},\underline{t})\} \text{ and } h \in H^1,
\end{cases} \\[2em]
\begin{cases}
\text{stop}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\bar{t}),(\underline{t},\underline{t})\}, s=0, \text{ and } h \in H^0, \\
1_0, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\bar{t}),(\underline{t},\underline{t})\}, s \neq 0, \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 } \mathrm{supp}(p_{|h}) = \{(\bar{t},\bar{t}),(\underline{t},\underline{t})\} \text{ and } h \in H^1,
\end{cases} \\[2em]
\begin{cases}
\text{stop}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\underline{t}),(\underline{t},\underline{t})\} \text{ and } s = \langle \underline{t},\underline{t}\rangle, \\
1_{\langle \underline{t},\underline{t}\rangle}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\underline{t}),(\underline{t},\underline{t})\} \text{ and } s \neq \langle \underline{t},\underline{t}\rangle,
\end{cases} \\[1em]
\begin{cases}
\text{stop}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\bar{t}),(\bar{t},\underline{t})\} \text{ and } s = \langle \bar{t},\bar{t}\rangle, \\
1_{\langle \bar{t},\bar{t}\rangle}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\bar{t}),(\bar{t},\underline{t})\} \text{ and } s \neq \langle \bar{t},\bar{t}\rangle,
\end{cases} \\[1em]
\begin{cases}
\text{stop}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\bar{t})\} \text{ and } s = \langle \bar{t},\bar{t}\rangle, \\
1_{\langle \bar{t},\bar{t}\rangle}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\bar{t})\} \text{ and } s \neq \langle \bar{t},\bar{t}\rangle,
\end{cases} \\[1em]
\begin{cases}
1_{\langle \bar{t},\underline{t}\rangle}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\underline{t})\} \text{ and } s \neq \langle \bar{t},\underline{t}\rangle, \\
\text{stop}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\bar{t},\underline{t})\} \text{ and } s = \langle \bar{t},\underline{t}\rangle,
\end{cases} \\[1em]
\begin{cases}
\text{stop}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\underline{t},\underline{t})\} \text{ and } s = \langle \underline{t},\underline{t}\rangle, \\
1_{\langle \underline{t},\underline{t}\rangle}, & \text{if } \mathrm{supp}(p_{|h'}) = \{(\underline{t},\underline{t})\} \text{ and } s \neq \langle \underline{t},\underline{t}\rangle.
\end{cases}
\end{cases}
\tag{15}
$$

**Lemma 7** *Let* $supp(p_{|\emptyset}) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$ *with* $\bar{t} > \underline{t}$. *There is a consistent, Bellman optimal, and admissible 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. $supp(p_{|h}) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$ and $h \in H^0$. Then $C^{\sigma^*}(h) \backslash \{1_0\} = \emptyset$.

   To see this, suppose on the contrary that $r \in C^{\sigma^*}(h) \backslash \{1_0\}$. Then:

   (a) $\{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\} \not\subseteq supp(p_{|h,(r,s)})$ for all $s \in r(supp(p_{|h}))$, by the construction of $\sigma^*(h, (r, s))$,

   (b) $g^*(s) \in \{(1, \underline{t}), (1, (\bar{t} + \underline{t})/2), (1, \bar{t})\}$, for all $s \in r(supp(p_{|h}))$, by (a), and the construction of $\sigma^*(h, (r, s))$,

   (c) $g^*(s) = (1, \bar{t})$, for all $s \in r(\bar{t}, t)$ for all $t \in \{\bar{t}, \underline{t}\}$, by (b) and individual rationality,

   (d) $g^*(s) = (1, \underline{t})$, for all $s \in r(t, \underline{t})$ for all $t \in \{\bar{t}, \underline{t}\}$, by (b) and individual rationality,

   (e) $g^*(s) = (1, \bar{t})$, for all $s \in r(\underline{t}, t)$ for all $t \in \{\bar{t}, \underline{t}\}$, by (c) and incentive compatibility,

   (f) $g^*(s) = (1, \underline{t})$, for all $s \in r(t, \bar{t})$ for all $t \in \{\bar{t}, \underline{t}\}$, by (d) and incentive compatibility.

   (g) By (c) and (e), $g^*(s) = (1, \bar{t})$ 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 (d) and (f), $g^*(s) = (1, \bar{t})$ for all $s \in r(t, t')$ for all $(t, t') \in \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$, a contradiction.

2. $supp(p_{|h}) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$ and $h \in H^1$. Then $C^{\sigma^*}(h) = \{r \in IC(p_{|h}) : g^*(s) = \sigma^*(h, (r, s))$ for all $s \in supp(p_{|h})\}$.

3. $supp(p_{|h}) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}$ and $h \in H^0$. Then $C^{\sigma^*}(h) = \{r \in IC(p_{|h}) : g^*(s) = \sigma^*(h, (r, s))$ for all $s \in supp(p_{|h})\}$.

4. $supp(p_{|h}) = \{(\bar{t}, \bar{t}), (\underline{t}, \underline{t})\}$ and $h \in H^1$. Then $C^{\sigma^*}(h) = \{r \in IC(p_{|h}) : g^*(s) = \sigma^*(h, (r, s))$ for all $s \in supp(p_{|h})\}$.

5. $supp(p_{|h}) = \{(\bar{t}, \underline{t}), (\underline{t}, \underline{t})\}$. Then $C^{\sigma^*}(h) = \{1_{\langle \underline{t}, \underline{t} \rangle}\}$, by incentive compability, individual rationality, and the construction of $\sigma^*(h, \cdot)$.

6. $supp(p_{|h}) = \{(\bar{t}, \bar{t}), (\bar{t}, \underline{t})\}$. Then $C^{\sigma^*}(h) = \{1_{\langle \bar{t}, \bar{t} \rangle}\}$, by incentive compability, individual rationality, and the construction of $\sigma^*(h, \cdot)$.

7. $supp(p_{|h}) = \{(\bar{t}, \bar{t})\}$. Then $C^{\sigma^*}(h) = \{1_{\langle \underline{t}, \underline{t} \rangle}\}$, by the construction of $\sigma^*(h, \cdot)$.

8. $supp(p_{|h}) = \{(\bar{t}, \underline{t})\}$. Then $C^{\sigma^*}(h) = \{1_{\langle \bar{t}, \underline{t} \rangle}\}$, by the construction of $\sigma^*(h, \cdot)$.

24

9. $\mathrm{supp}(p_{|h}) = \{(\underline{t},\underline{t})\}$. Then $C^{\sigma^*}(h) = \{1_{\langle \overline{t},\overline{t}\rangle}\}$, by the construction of $\sigma^*(h, \cdot)$.

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

- In all cases, if $\sigma^*(h) = stop$, then consistency is automatically implied. Suppose below that $\sigma^*(h) \neq stop$.

- In Cases 1, 3, 5, 6, 7, 8, and 9, consistency is implied by the fact that $C^{\sigma^*}(h) = \{\sigma^*(h)\}$.

- In Case 2, $\sigma^*(h) = r^*$ such that

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

There are two cases, $s = \langle \overline{t}, \underline{t}\rangle$ and $s = 0$. In the former, $g^*(\langle \overline{t}, \underline{t}\rangle) = (1, (\overline{t} + \underline{t})/2)$ and $\mathrm{supp}(p_{|h,(r,\langle \overline{t},\underline{t}\rangle)}) = \{(\overline{t}, \underline{t})\}$. By construction $\sigma^*(h, (r^*, \langle \overline{t}, \underline{t}\rangle)) = g^*(\langle \overline{t}, \underline{t}\rangle)$. In the latter case, $g^*(0) = (0, 0)$ and $\mathrm{supp}(p_{|h,(r,0)}) = \{(\overline{t}, \overline{t}), (\underline{t}, \underline{t})\}$. 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

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

Since $\mathrm{supp}(p_{|h}) = \{(\overline{t}, \overline{t}), (\underline{t}, \underline{t})\}$, only signals $s = \langle \overline{t}, \overline{t}\rangle$ and $s = \langle \underline{t}, \underline{t}\rangle$ materialize with positive probability. In the former case, $g^*(\langle \overline{t}, \overline{t}\rangle) = (1, \overline{t})$ and $\mathrm{supp}(p_{|h,(r^{**},\langle \overline{t},\overline{t}\rangle)}) = \{(\overline{t}, \overline{t})\}$. By construction, $\sigma^*(h, (r^{**}, \langle \overline{t}, \overline{t}\rangle)) = g^*(\langle \overline{t}, \overline{t}\rangle)$ when $s = \langle \overline{t}, \overline{t}\rangle$. In the latter case, $g^*(\langle \underline{t}, \underline{t}\rangle) = (1, \underline{t})$ and $\mathrm{supp}(p_{|h,(r^{**},\langle \underline{t},\underline{t}\rangle)}) = \{(\underline{t}, \underline{t})\}$. By construction, $\sigma^*(h, (r^{**}, \langle \underline{t}, \underline{t}\rangle)) = g^*(\langle \underline{t}, \underline{t}\rangle)$, implying consistency when $s = \langle \underline{t}, \underline{t}\rangle$. 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 when $s = 0$.

To see that $\sigma^*$ is Bellman optimal when $\sigma^*(h) \neq stop$:

- In Case 1, $C^{\sigma^*}(h) \backslash \{1_0\} = \emptyset$, and hence $\sigma^*(h) = 1_0$ is a Bellman optimal choice in this case.

- In Case 2, $v(g^* \circ r^*, p_{|h}) = (\overline{t} - \underline{t})p(\overline{t}, \underline{t}) \geq v(g^* \circ r, p_{|h})$, for all budget balanced mechanisms $g^* \circ r$. Hence $\sigma^*(h) = r^*$ is a Bellman optimal choice in this case.

- In Case 3, $v(g^* \circ 1_0, p_{|h}) = 0 = v(g^* \circ r, p_{|h})$, for all budget balanced mechanisms $g^* \circ r$. Hence $\sigma^*(h) = 1_0$ 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, $C^{\sigma^*}(h) = \{\sigma^*(h)\}$, and hence $\sigma^*(h)$ is a Bellman optimal choice in these cases.

To see that $\sigma^*$ is Bellman optimal when $\sigma^*(h) = stop$ :

- In Case 1, $\widehat{C}^{\sigma^*}(h)\backslash\{1_0\} \subseteq C^{\sigma^*}(h)\backslash\{1_0\} = \emptyset$, and hence $\sigma^*(h) = stop$ is a Bellman optimal choice in this case.

- In Case 3, $v(g(0), p_{|h}) = 0 = v(g^* \circ r, p_{|h})$, for all budget balanced mechanisms $g^* \circ r$. Hence $\sigma^*(h) = stop$ is a Bellman optimal choice in this case.

- In Cases 5, 6, 7, 8, and 9, $g(s) \in X$ becomes implemented such that $C^{\sigma^*}(h) = \{1_s\}$. Hence $\sigma^*(h) = stop$ is a Bellman optimal choice in these cases.

■

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 consitute a Bellman optimal and consistent rule in the continuation game. By Theorem 1, such a continuation strategy does exists. 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 chice rule, staring from $p_{|\emptyset}$

# A  Appendix

## A.1  Omitted proofs of Section 5

**Proof of Proposition 2.**  Denote

$$a_1(\theta_1) = \sum_{\theta_2} p_2(\theta_2)a(\theta_1, \theta_2),$$

$$a_2(\theta_2) = \sum_{\theta_1} p_1(\theta_1)a(\theta_1, \theta_2),$$

and use the shorthand

$$V_1(\theta_1) = \sum_{\theta_2} p(\theta_2 : \theta_1)[a(\theta_1, \theta_2)\theta_1 - m(\theta_1, \theta_2)],$$

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

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

$$a_1(\theta_1')(\theta_1 - \theta_1') \leq V_1(\theta_1) - V_1(\theta_1') \leq a_1(\theta_1)(\theta_1 - \theta_1').$$
$$a_2(\theta_2')(\theta_2 - \theta_2') \geq V_2(\theta_2') - V_2(\theta_2) \geq a_2(\theta_2)(\theta_2 - \theta_2').$$

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

$$V_1(\theta_1) \geq \sum_{t \leq \theta_1'} a_1(t) + V_1(0),$$

$$V_2(\theta_2) \geq \sum_{t \geq \theta_2'} a_2(t) + V_2(1).$$

Let

$$P_i(\theta_i) = \sum_{t \leq \theta_i} p_i(t),$$

$$A_i(\theta_i) = \sum_{t \leq \theta_i} a_i(t).$$

Then

$$P_1(\theta_1)A_1(\theta_1) = \sum_{t \leq \theta_1} P_1(t)[A_1(t) - A_1(t')] + \sum_{t \leq \theta_1} [P_1(t) - P_1(t')]A_1(t')$$

$$= \sum_{t \leq \theta_1} P_1(t)a_1(t) + \sum_{t \leq \theta_1} p_1(t)A_1(t').$$

Thus

$$\sum_t p_1(t)A_1(t') = P_1(1)A_1(1) - \sum_t P_1(t)a_1(t) \qquad (16)$$

$$= \sum_t a_1(t)(1 - P_1(t)).$$

And similarly for the agent 2 :

$$\sum_t p_2(t)[A_2(1) - A_2(t')] = A_2(1) - \sum_t p_2(t)A_2(t') \qquad (17)$$

$$= \sum_t a_2(t)P_2(t).$$

The planner's problem can be written

$$\max_{a(r(\cdot))} \sum_{\theta_1} \sum_{\theta_2} p_1(\theta_1)p_2(\theta_2)(\theta_1 - \theta_2)a(r(\theta_1, \theta_2))$$

s.t.

$$\sum_{\theta_1}\sum_{\theta_2} p_1(\theta_1)p_2(\theta_2)(\theta_1-\theta_2)a(r(\theta_1,\theta_2)) = \sum_{\theta_1} p_1(\theta_1)V_1(\theta_1) + \sum_{\theta_2} p_2(\theta_2)V_2(\theta_2) \qquad (18)$$

$$A_1(\theta_1) \geq V_1(\theta_1) - V_1(0) \geq A_1(\theta_1'), \text{ for all } \theta_1 \qquad (19)$$

$$A_2(1) - A_2(\theta_2) \geq V_2(\theta_2) - V_2(1) \geq A_2(1) - A_2(\theta_2'), \text{ for all } \theta_2 \qquad (20)$$

$$V_1(\theta_1) \geq 0 \text{ for all } \theta_1, \; V_2(\theta_2) \geq 0 \text{ for all } \theta_2 \qquad (21)$$

where (18) is the ex ante budget balance condition, (19) and (20) are the incentive compatibility constraints, and (21) is the participation constraint.

Since the right hand side inequalities of (19) and (20) imply

$$\sum_{\theta_1} p_1(\theta_1)V_1(\theta_1) \geq \sum_{\theta_1} p_1(\theta_1)A_1(\theta_1') + V_1(0),$$

$$\sum_{\theta_2} p_2(\theta_2)V_2(\theta_2) \geq \sum_{\theta_2} p_2(\theta_2)[A_2(1) - A_2(\theta_2')] + V_2(1),$$

(16), (17), and (18) result in

$$V_1(0) + V_2(1) + \sum_{\theta_1} a_1(\theta_1)(1 - P_1(\theta_1)) + \sum_{\theta_2} a_2(\theta_2)P_2(\theta_2)$$

$$\leq \sum_{\theta_1}\sum_{\theta_2} p_1(\theta_1)p_2(\theta_2)(\theta_1 - \theta_2)a(r(\theta_1,\theta_2)),$$

or, more compactly,

$$\sum_{\theta_1}\sum_{\theta_2} p_1(\theta_1)p_2(\theta_2)\left[\left(\theta_1 - \frac{1 - P_1(\theta_1)}{p_1(\theta_1)}\right) - \left(\theta_2 + \frac{P_2(\theta_2)}{p_2(\theta_2)}\right)\right]a(\theta_1,\theta_2) \geq 0. \qquad (22)$$

Maximizing the objective function with respect to (22), and interpreting $\gamma/(1-\gamma)$ as the Lagrange multiplier, gives the desired programme. Since, at the optimum, (22) holds as equality, the solution to the programme also meets the left hand side inequalities of (19) and (20). Since this implies that $a_1$ is increasing and $a_2$ is decreasing, it also follows that the participation constraint (21) 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 5.** 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 :

$$m_1^\gamma(\theta_1) = a_1^\gamma(\theta_1)\theta_1 - \sum_{t<\theta_1'} a_1^\gamma(t), \quad \text{for all } \theta_1 \in T$$

$$m_2^\gamma(\theta_2) = a_2^\gamma(\theta_2)\theta_2 + \sum_{t>\theta_2'} a_2^\gamma(t), \quad \text{for all } \theta_2 \in T.$$

The ex ante budget balance of the incentive efficient mechanism implies

$$\sum_{\theta_1} p_1(\theta_1) m_1^\gamma(\theta_1) = \sum_{\theta_2} p_2(\theta_2) m_2^\gamma(\theta_2). \tag{23}$$

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

$$p_1(1)m(1,0) + \sum_{t<1} p_1(t) m_1^\gamma(t) = m_2^\gamma(0), \tag{24}$$

$$p_2(0)m(1,0) + \sum_{t>0} p_2(t) m_2^\gamma(t) = m_1^\gamma(1). \tag{25}$$

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 (24) and (25), 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,

$$\sum_{\theta_1} p_1(\theta_1) \sum_{\theta_2} p_2(\theta_2) m(\theta_1, \theta_2) = \sum_{\theta_2} p_2(\theta_2) \sum_{\theta_1} p_1(\theta_1) m(\theta_1, \theta_2). \tag{26}$$

By construction,

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

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

Now letting $m(1,0)$ be defined by (24), 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 (26) and (23),

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

Thus $m(1,0)$ also satisfies (25). ∎

29

# 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 STRAUSZ, 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] FLESCH, J., KUIPERS, J., SCHOENMAKERS G., AND K. VRIEZE (2010), Subgame perfection in positive recursive games with perfect information, *Mathematics of Operations Research* **35,** 193-207

[11] FORGES, F. (1994), Posterior efficiency, *Games and Economic Behavior* **6**, 238-61.

[12] FREIXAS, X., GUESNERIE, R., AND TIROLE, J. (1985), Planning under incomplete information ans the ratchet effect, *Review of Economic Studies* **52,** 173-92

[13] GERARDI, D., HÖRNER, J., AND L. MAESTRI (2011), The role of commitment in bilateral trade, manuscript, Collegio Carlo Alberto

[14] GREEN, J. AND LAFFONT, J.-J.. (1987), Posterior implementability in a two-player decision problem, *Econometrica* **55,** 69-94

[15] GRESIK, T. (1991), Ex ante efficient, ex post individually rational trade, *Journal of Economic Theory* **53,** 131-45.

[16] GRESIK, T. (1996), Incentive efficient equilibria of two-party sealed-bid bargaining games, *Journal of Economic Theory* **68,** 26-48

[17] HOLMSTROM, B. AND MYERSON, R. (1983), Efficient and durable decision rules with incomplete information, *Econometrica* **51**, 1799-819,

[18] KIYOTAKI, N. (2011), A mechanism design approach to financial frictions, manuscript, Princeton University

[19] KRISHNA, V. AND PERRY, M. (1998), Efficient mechanism design, working paper, Penn State University.

[20] LAFFONT, J.-J. AND TIROLE, J. (1990), Adverse selection and renegotiation in procurement , *Review of Economic Studies* **57,** 597-625

[21] LAGUNOFF, R. (1995), Resilient allocation rules for bilateral trade, *Journal of Economic Theory* **66**, 463-87

[22] LAGUNOFF, R. (1992), Fully endogenous mechanism selection on finite outcome sets, *Economic Theory* **2**, 465-80

[23] MCADAMS, D. AND SCHWARZ, M. (2006), Credible sales mechanisms and intermediaries, *American Economic Review* forthcoming

[24] MCAFEE, P. AND VINCENT, D. (1997), Sequentially optimal auctions, *Games and Economic Behavior* **18,** 246-76

[25] MYERSON, R. (1991), *Game Theory, Analysis of Conflict*, Cambridge MA, Harvard University Press

[26] MYERSON, R., (1979), Incentive compatibility and the bargaining problem, *Econometrica* **47**, 61–73.

[27] MYERSON, R. (1982), Optimal coordination mechanisms in generalized principle-agent problems, *Journal of Mathematical Economics* **28,** 67–81.

[28] NEEMAN Z. AND G. PAVLOV (2012), Renegotiation-proof mechanism design, forthcoming in *Journal of Economic Theory.*

[29] SEGAL, I., AND M. WHINSTON (2002), The Mirrlees approach to mechanism design with renegotiation *Econometrica* **70**, 1-45.

[30] SKRETA V. (2011), Optimal auction design under non-commitment, manuscript, NYU Stern School of Business

[31] SKRETA V., (2006), Sequentially optimal mechanisms, *Review of Economic Studies* **73**

[32] VARTIAINEN, H. (2010), Auction design without commitment, forthcoming in *Journal of the European Economic Association*