



# Subgame-Perfect Implementation Under Information Perturbations

## Citation

Aghion, Phillippe, Drew Fudenberg, Richard Holden, Takashi Kunimoto, and Olivier Tercieux. 2012. Subgame-perfect implementation under information perturbations. *The Quarterly Journal of Economics* 127, no. 4: 1843-1881.

## Published Version

doi:10.1093/qje/qjs026

## Permanent link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:11224965>

## Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Open Access Policy Articles, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#OAP>

## Share Your Story

The Harvard community has made this article openly available.  
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

# Subgame Perfect Implementation Under Information Perturbations\*

Philippe Aghion, Drew Fudenberg, Richard Holden,  
Takashi Kunimoto and Olivier Tercieux<sup>†</sup>

February 23, 2012

## Abstract

We consider the robustness of extensive form mechanisms to deviations from common knowledge about the state of nature, which we refer to as *information perturbations*. First, we show that even under arbitrarily small information perturbations the Moore-Repullo mechanism does not yield (even approximately) truthful revelation and that in addition the mechanism has sequential equilibria with undesirable outcomes. More generally, we prove that any extensive form mechanism is fragile in the sense that if a non-Maskin monotonic social objective can be implemented with this mechanism, then there are arbitrarily small information perturbations under which an undesirable sequential equilibrium also exists. Finally, we argue that outside options can help improve efficiency in asymmetric information environments, and that these options can be thought of as reflecting ownership of an asset.

---

\*This paper builds on two preliminary contributions, respectively by Aghion, Fudenberg and Holden (2009) and Kunimoto and Tercieux (2009).

<sup>†</sup>Aghion: Harvard University, Department of Economics and CIFAR; email: paghion@fas.harvard.edu. Fudenberg: Harvard University, Department of Economics; email: dfudenberg@harvard.edu. Holden: University of New South Wales, Australian School of Business; email: richard.holden@unsw.edu.au. Kunimoto: Department of Economics, Hitotsubashi University, Tokyo, Japan; email: takashikunimoto9@gmail.com. Tercieux: Paris School of Economics, Paris, France; email: tercieux@pse.ens.fr. We thank Oliver Hart, Johannes Horner, John Moore and Andy Skrzypacz for detailed comments on earlier drafts. We are also grateful to Ken Binmore, Yeon-Koo Che, Mathias Dewatripont, Bob Gibbons, Ed Green, Matt Jackson, Philippe Jehiel, Hitoshi Matsushima, Eiichi Miyagawa, Eric Maskin, Roger Myerson, Antonio Penta, Andrew Postlewaite, Jean Tirole, Jorgen Weibull, Ivan Werning, Tom Wilkening, Muhamet Yildiz, seminar participants at Chicago Booth, Harvard, the Paris School of Economics, Stockholm University, the Stockholm School of Economics, Simon Fraser University, Boston University, Bocconi University, the Max Planck Institute in Bonn, and the Canadian Institute for Advanced Research, and the referees and editor of this journal for very useful comments and suggestions. Thanks also to Ashley Cheng for careful proofreading. Financial support from CIFAR (Aghion), from NSF grants SES 0648616, 0954162 (Fudenberg), and from FQRSC, SSHRC of Canada, and the Seimeikai Foundation (Kunimoto) is gratefully acknowledged.

# 1 Introduction

The literature on “complete-information” implementation supposes that players know the payoff-relevant state of the world, and asks which mappings from states to outcomes, i.e., which *social choice rules*, can be implemented by mechanisms that respect the players’ incentives. Although only *Maskin monotonic* social rules are “Nash implementable” (Maskin, 1999), a larger class of social choice rules can be implemented in extensive form games provided that a more restrictive equilibrium notion is used.<sup>1</sup>

This paper considers the robustness of subgame-perfect implementation to arbitrarily small amounts of incomplete information about the state of nature  $\theta$ , which we refer to as “information perturbations”.<sup>2</sup> It is known that refinements of Nash equilibrium are not robust to general small perturbations of the payoff structure (Fudenberg, Kreps and Levine (1988), henceforth FKL), but our results do not follow from theirs as we consider a more restrictive class of perturbations: We fix the map from states to payoffs and perturb the prior distribution over the states of the world and signal structure, so in particular the messages in the mechanism remain cheap talk and do not enter directly into the payoff functions.

Our starting point is the Moore and Repullo (1988, henceforth, MR) result which roughly says that for any social choice rule, one can design a mechanism that yields unique implementation in subgame-perfect equilibria (i.e., for all states of nature, the set of all subgame-perfect equilibria of the induced game yields the desired outcome). In particular, in environments with money, Moore and Repullo propose a simple mechanism (which we call an MR mechanism) inducing truth-telling as the unique subgame-perfect equilibrium. As in MR, our focus is on *exact* implementation, where “exact implementation” means that we require the set of equilibrium outcomes to *exactly* coincide with those picked by the rule.<sup>3</sup>

The requirement of exact implementation can be decomposed into the following two parts: (1) there always exists an equilibrium whose outcome coincides with the given rule; (2) there are no equilibria whose outcomes differ from those of the rule.

Our first result shows that MR mechanisms can only robustly satisfy the first requirement of exact implementation if the rule that is implemented is Maskin monotonic. That is, whenever an

---

<sup>1</sup>Recall that a social choice rule or function  $f$  is *Maskin monotonic* if for any pair of states  $\theta$  and  $\theta'$  such that  $a = f(\theta)$ , and  $a$  never goes down in the preference ranking of any agent when moving from state  $\theta$  to state  $\theta'$ , then necessarily  $a = f(\theta')$ .

<sup>2</sup>It follows from Theorem 14.5 of Fudenberg and Tirole (1991a, p.567) that under our small informational perturbations, for each profile of signals that has strictly positive probability under complete information, some state of nature is *common  $p$ -belief* (Monderer and Samet (1989)) with  $p$  arbitrarily close to one. That is, everybody believes this is the true state with probability at least  $p$ , everybody believes with probability at least  $p$  that everybody believes this is the true state with probability at least  $p$  etc...

<sup>3</sup>Much of the implementation literature studies exact implementation. Virtual implementation (Matsushima (1988) and Abreu and Sen (1991)) uses non-deterministic mechanisms, and only requires that social choice rules be implemented with high probability. As pointed out by Jackson (2001), unlike exact implementation, virtual implementation is not robust to introducing a small amount of nonlinearity in preferences over lotteries. In addition, virtual implementation provides incentives for renegotiation on the equilibrium path: As Abreu and Matsushima (1992) acknowledge, virtual implementation supposes that the social planner can commit ex ante to outcomes that will be known at the time of implementation to be highly inefficient.

MR mechanism implements a non-Maskin monotonic social choice rule, the truth-telling equilibrium ceases to be an equilibrium in some nearby environment. More specifically, we show that an MR mechanism which implements a social choice rule  $f$  under common knowledge (or complete information<sup>4</sup>) about the state of nature does not yield even approximately truthful revelation under arbitrarily small information perturbations, if this  $f$  is not Maskin monotonic.<sup>5</sup>

We then move beyond MR mechanisms to consider *any* extensive-form mechanism. Our second result is concerned with the non-robustness of the second requirement of exact implementation: namely, whenever any mechanism implements a non-Maskin monotonic social choice rule, there exists an undesirable equilibrium in some nearby environment. More specifically, restricting attention to environments with a finite state space, and to mechanisms with finite strategy spaces,<sup>6</sup> then given any mechanism that “subgame-perfect” implements a non-Maskin monotonic social choice rule  $f$  under common knowledge (i.e., whose subgame-perfect equilibrium outcomes in any state  $\theta$  is precisely equal to  $f(\theta)$ ), we can find a sequence of information perturbations (i.e., of deviations from complete information about the state of nature) and a corresponding sequence of sequential equilibria for the mechanism under the corresponding information perturbations, whose outcomes do not converge to  $f(\theta)$  for at least one state  $\theta$ . In other words, there always exist arbitrarily small information perturbations under which an “undesirable” sequential<sup>7</sup> (and hence Perfect Bayesian) equilibrium exists.

Three insights underlie our analysis. The first is that even a small amount of uncertainty about the state at the interim stage, when players have observed their signals but not yet played the game, can loom large ex post once the extensive form game has started and players can partly reveal their private signals through their strategy choice at each node of the game. The second insight is that arbitrarily small information perturbations can turn the outcome of a non-sequential Nash equilibrium of the game with common knowledge of  $\theta$  into the outcome of a sequential equilibrium of the perturbed game. In particular, we know that any extensive-form mechanism that “subgame-perfect” implements a non-Maskin monotonic social choice rule under common knowledge has at least one Nash equilibrium which is not a subgame-perfect equilibrium; we prove that this undesirable Nash equilibrium can be turned into an undesirable sequential equilibrium by only introducing small information perturbations. The third insight is that there is a role for asset ownership to mitigate the investment and trade inefficiencies that arise when the contracting parties have private information ex post about the state of nature  $\theta$ .

Our results are not a straightforward application of those on the robustness of refinements of Nash equilibrium because we consider a smaller class of perturbations. While FKL consider several nested classes of perturbations, even the most restrictive form they analyze allows a player’s payoff in the perturbed game to vary with the realized actions in an arbitrary way. In the mechanism

---

<sup>4</sup>Throughout the paper, we use “complete information” and “common knowledge” interchangeably.

<sup>5</sup>As we shall stress in Section 2.5 below, Maskin monotonicity is precisely the property that the social choice rules usually considered in contract theory do not satisfy.

<sup>6</sup>The Appendix extends the result to the case of countable message spaces.

<sup>7</sup>We remind the reader of the formal definition in Section 4.2.1.

design setting, this implies that some (low-probability) “crazy types” might have a systematic preference for truth-telling. Since the messages and outcome functions of the mechanism are not primitives, but rather endogenous objects to be chosen by the social planner, it may seem natural to restrict the perturbations to be independent of the messages and depend only on the allocation that is implemented.

Our paper contributes most directly to the mechanism design literature, starting with Maskin’s (1999) Nash implementation result and Moore-Repullo’s (1988) subgame-perfect implementation analysis, by showing the non-robustness of subgame-perfect implementation to information perturbations.<sup>8</sup> Our paper is also related to Chung and Ely’s (2003) study of the robustness of undominated Nash implementation. Chung and Ely show that if a social choice rule is not Maskin monotonic but can be implemented in undominated Nash equilibrium<sup>9</sup> under complete information, then there are information perturbations under which an undesirable undominated Nash equilibrium appears. In contrast, we consider extensive-form mechanisms and show that only Maskin monotonic social choice rules can be implemented in the closure of the sequential equilibrium correspondence. In general, the existence of a bad sequential equilibrium in the perturbed game neither implies nor is implied by the existence of a bad undominated Bayesian Nash equilibrium, as undominated Nash equilibria need not be sequential equilibria, and sequential equilibria can use dominated strategies.<sup>10</sup> Hence, although our paper has a similar spirit to Chung and Ely (2003), our argument is quite distinct from theirs.

Our paper also relates to the literature on the hold-up problem. Grossman and Hart (1986) argue that in contracting situations where states of nature are observable but not verifiable, asset ownership (or vertical integration) could help limit the extent to which one party can be held up by the other party, which in turn should encourage ex ante investment by the former. However, vertical integration as a solution to the hold-up problem has been questioned in papers which use or extend the subgame-perfect implementation approach of Moore and Repullo (1988).<sup>11</sup> In particular, Maskin and Tirole (1999a), henceforth MT, show that the non-verifiability of states of nature can be overcome by using a 3-stage subgame-perfect implementation mechanism that induces truth-telling by all parties as the unique equilibrium outcome, and does so in pure strategies. We contribute to this debate in two ways. First we show that the introduction of even small information perturbations greatly reduces the power of subgame-perfect implementation. This suggests that the introduction of incomplete information can significantly change the insights obtained by MT.

---

<sup>8</sup>Other related mechanism design papers include Cremer and McLean (1988), Johnson, Pratt and Zeckhauser (1990), and Fudenberg, Levine and Maskin (1991). These papers show how one can take advantage of the correlation between agents’ signals in designing incentives to approximate the Nash equilibrium under complete information. These papers consider static implementation games with commitment, and look at fairly general information structures, as opposed to our focus on the robustness of subgame-perfect implementation to small perturbations from complete information.

<sup>9</sup>An undominated Nash equilibrium is a Nash equilibrium where no player ever uses a weakly dominated action.

<sup>10</sup>Trembling-hand perfect equilibria cannot use dominated strategies, and sequential and trembling-hand perfect equilibria coincide for generic assignments of payoffs to terminal nodes (Kreps and Wilson [1982]), but the generic payoffs restriction rules out our assumption that messages are cheap talk.

<sup>11</sup>For example, see Aghion-Dewatripont-Rey (1994) and Maskin-Tirole (1999a, 1999b).

Secondly, we show that when there is asymmetric information ex post about the good's valuation, an outside option for the seller permits a more efficient outcome. We argue that this option can be seen as corresponding to ownership of an asset.

The paper is organized as follows. Section 2 uses a simple buyer-seller example to introduce the MR mechanism, to show why truthful implementation using this mechanism is not robust to small information perturbations, and why such perturbations generate an undesirable sequential equilibrium. Section 3 extends our analysis to general MR mechanisms with  $n$  states of nature and transferable utility, and shows that for a given social choice rule  $f$ , truth-telling equilibria are only robust to small information perturbations if this  $f$  is *strategy-proof* (which in turn implies Maskin monotonicity under weak assumptions on preferences).<sup>12</sup> In Section 4, we ask whether *any* extensive form mechanism is robust to small information perturbations. There we prove that for any social choice rule that is not Maskin-monotonic one can find small information perturbations under which an undesirable sequential equilibrium exists. Section 5 considers the case of full informational asymmetry ex post, and shows that asset ownership, by providing outside options, can lead to approximately efficient ex ante investments, whereas contracts or mechanisms with no outside option, cannot. Finally, Section 6 concludes with a few remarks and also suggestions for future research.

## 2 A Hart-Moore example of the Moore-Repullo mechanism

### 2.1 Basic setup

Consider the following simple example from Hart and Moore (2003), which captures the logic of Moore and Repullo (1988)'s subgame-perfect implementation mechanism.

There are two parties, a  $B$ (uyer) and a  $S$ (eller) of a single unit of an indivisible good. If trade occurs then  $B$ 's payoff is

$$V_B = \theta - p,$$

where  $p$  is the price and  $\theta$  is the good's quality.  $S$ 's payoff is

$$V_S = p,$$

thus we normalize the cost of producing the good to zero.

The good can be of either high or low quality. If it is high quality then  $B$  values it at  $\theta_H = 14$ , and if it is low quality then  $B$  values it at  $\theta_L = 10$ . We seek to implement the social choice function whereby the good is always traded ex post, and where the buyer always pays the true  $\theta$  to the seller.

---

<sup>12</sup>If  $f$  is strategy-proof, it is always a weakly dominant strategy for each agent to tell the truth in the direct mechanism associated with  $f$ . See also Definition 1 for a precise definition of strategy-proofness.

## 2.2 Common knowledge

Suppose first that the quality  $\theta$  is observable and common knowledge to both parties. Even though  $\theta$  is not verifiable by a court, so no initial contract between the two parties can be made credibly contingent upon  $\theta$ , truthful revelation of  $\theta$  by the buyer  $B$  and the implementation of the above social choice function can be achieved through the following Moore-Repullo (MR) mechanism:

1.  $B$  announces either a “high” or “low” quality. If  $B$  announces “high” then  $B$  pays  $S$  a price equal to 14 and the game stops.
2. If  $B$  announces “low” and  $S$  does not “challenge”  $B$ ’s announcement, then  $B$  pays a price equal to 10 and the game stops.
3. If  $S$  challenges  $B$ ’s announcement then:
  - (a)  $B$  pays a fine  $F = 9$  to  $T$  (a third party)
  - (b)  $B$  is offered the good for 6
  - (c) If  $B$  accepts the good then  $S$  receives  $F$  from  $T$  (and also a payment of 6 from  $B$ ) and the game stops.
  - (d) If  $B$  rejects at 3b then  $S$  pays  $F$  to  $T$
  - (e)  $B$  and  $S$  each get the item with probability  $1/2$ .

When the true value of the good is common knowledge between  $B$  and  $S$ , this mechanism yields truth-telling as the unique subgame-perfect (and also sequential) equilibrium. To see this, consider first the case  $\theta = \theta_H$ . If  $B$  announces “high” then  $B$  pays 14 and we stop. If, however,  $B$  announces “low” then  $S$  will challenge because at stage 3a,  $B$  pays 9 to  $T$  and, this cost being sunk,  $B$  will still accept the good for price of 6 at stage 3b (since by rejecting he will end up at stage 3e and get  $14/2=7$ , but since the good is worth 14 he gets  $14-6=8$  by accepting). Anticipating this,  $S$  knows that if she challenges  $B$ , she will receive  $9 + 6 = 15$ , which is greater than 10 that she would receive if she did not challenge. Moving back to stage 1, if  $B$  lies and announces “low” when the true state is high, he gets  $14 - 9 - 6 = -1$ , whereas he gets  $14 - 14 = 0$  if he tells the truth, so truth-telling is the unique equilibrium here. Truth-telling is also the unique equilibrium when  $\theta = \theta_L$ : In that case  $S$  will not challenge  $B$  when  $B$  (truthfully) announces “low”, because now  $B$  will refuse the good at price 6 (accepting the good at 6 would yield surplus  $10 - 6 = 4$  to  $B$  whereas by refusing the good and relying on the lottery which assigns the item randomly instead  $B$  can secure a surplus equal to  $10/2 = 5$ ). Anticipating this,  $S$  will not challenge  $B$  because doing so would give her a net surplus equal to  $10/2 - 9 = -4$  which is less than the payment of 10 she receives if she does not challenge  $B$ ’s announcement.

This mechanism (and more generally, the Moore-Repullo mechanisms we describe in Section 3) has two nice and important properties. First, it yields unique implementation in subgame-perfect equilibrium, i.e., for any state of nature, there is a unique subgame-perfect equilibrium which yields

the right outcome. Second, in each state, the unique subgame-perfect equilibrium is appealing from a behavioral point of view since it involves telling the truth. In what follows, we will show that both of these properties fail once we introduce small information perturbations.

## 2.3 The failure of truth-telling with perturbed beliefs about value

### 2.3.1 Pure strategy equilibria

As in the example above, we continue to suppose that the good has possible values  $\theta \in \{\theta_H, \theta_L\}$  with  $\theta_H = 14$  (the high state) and  $\theta_L = 10$  (the low state). However, we now suppose that the players have imperfect information about  $\theta$ . Specifically, we suppose they have a common prior  $\mu$ , with  $\mu(\theta_H) = 1 - \alpha$ ,  $\mu(\theta_L) = \alpha$  for some  $\alpha \in (0, 1)$ , and that each player receives a draw from a signal structure with two possible signals  $s^h$  or  $s^\ell$ , where  $s^h$  is a high signal that is associated with  $\theta_H$ , and  $s^\ell$  is a low signal associated with  $\theta_L$ . We use the notation  $s_B = s_B^h$  (resp.  $s_B = s_B^\ell$ ) to refer to the event in which  $B$  receives the high signal  $s^h$  (resp. the low signal  $s^\ell$ ) and similarly we use the notation  $s_S = s_S^h$  (resp.  $s_S = s_S^\ell$ ) to refer to the event in which  $S$  receives the high signal  $s^h$  (resp. the low signal  $s^\ell$ ). The following table shows the joint probability distribution  $\nu^\varepsilon$  over  $\theta$ , the buyer's signal  $s_B$ , and the seller's signal  $s_S$ :

$\nu^\varepsilon$	$s_B^h, s_S^h$	$s_B^h, s_S^\ell$	$s_B^\ell, s_S^h$	$s_B^\ell, s_S^\ell$
(*) $\theta_H$	$(1 - \alpha)(1 - \varepsilon - \varepsilon^2)$	$(1 - \alpha)\varepsilon$	$(1 - \alpha)\varepsilon^2/2$	$(1 - \alpha)\varepsilon^2/2$
$\theta_L$	$\alpha\varepsilon^2/2$	$\alpha\varepsilon^2/2$	$\alpha\varepsilon$	$\alpha(1 - \varepsilon - \varepsilon^2)$

Note that for all  $\varepsilon$ , the marginal probability distribution of  $\nu^\varepsilon$  on  $\theta$  coincides with  $\mu$ , and that as  $\varepsilon$  converges to zero,  $\nu^\varepsilon$  assigns probability converging to 1 to the signals being correct. Note also that the buyer's signal becomes infinitely more accurate than the seller's signal as  $\varepsilon \rightarrow 0$ . This special feature implies that when deciding whether or not to challenge the buyer if  $S$  and  $B$  were informed of both signals, and the signals disagree, they will conclude that with high probability the state corresponds to  $B$ 's signal.

We will now show that there is no equilibrium in pure strategies in which the buyer always reports truthfully. To simplify the exposition of this example, we keep the payments under the perturbed mechanism the same as in the MR mechanism under common knowledge of the previous subsection, and assume that  $B$  must participate in the mechanism. This is equivalent to assuming that  $B$ 's participation constraint is slack, which in turn can be arranged by a constant ex ante payment and so does not influence the incentives for truth-telling. By way of contradiction, suppose there is a pure strategy equilibrium in which  $B$  reports truthfully, and consider  $B$ 's play when  $s_B = s_B^h$ . Then  $B$  believes that, regardless of what signal player  $S$  gets, the expected value of the good is greater than 10. So  $B$  would like to announce "low" if he expects that  $S$  will not challenge the announcement. If  $B$  does announce "low," then in a fully revealing equilibrium,  $S$  will infer that  $B$  must have received the low signal, i.e.,  $s_B = s_B^\ell$ . But under signal structure (\*),  $S$  thinks that  $B$ 's signal is much more likely to be correct, so  $S$  now believes that there is a large probability that  $\theta = \theta_L$ ; therefore  $S$  will not challenge.



But then, at stage 1, anticipating that  $S$  will not challenge,  $B$  will prefer to announce “low” when he receives the high signal  $s_B^h$ . Therefore, there does not exist a fully revealing equilibrium in pure strategies and consequently, the above social choice function can no longer be implemented through the above MR mechanism in pure strategies.

### 2.3.2 Allowing for mixed strategies

The result that there are no truthful equilibria in pure strategies leaves open the possibility that there are mixed strategy equilibria in which the probability of truthful announcement goes to one as  $\varepsilon$  goes to zero. This is close to the way that the pure-strategy Stackelberg equilibrium can be approximated by a mixed equilibrium of a “noisy commitment game” (van Damme and Hurkens (1997)). We show below that this is not the case under the signal structure (\*).

Let  $\sigma_B^h$  denote the probability that  $B$  announces “low” after receiving the high signal  $s_B^h$ , and let  $\sigma_B^\ell$  be the probability  $B$  announces “high” after receiving the low signal  $s_B^\ell$ , as in the following table:

	High	Low
$s_B^h$	$1 - \sigma_B^h$	$\sigma_B^h$
$s_B^\ell$	$\sigma_B^\ell$	$1 - \sigma_B^\ell$

The corresponding mixing probabilities for player  $S$  are

	Challenge	Don't Challenge
$s_S^h$	$1 - \sigma_S^h$	$\sigma_S^h$
$s_S^\ell$	$\sigma_S^\ell$	$1 - \sigma_S^\ell$

Then for mixed strategy equilibria of the mechanism to converge to the equilibrium under complete information where the buyer announces the valuation truthfully, we should have  $\sigma_B^{\varepsilon,h}, \sigma_B^{\varepsilon,\ell}, \sigma_S^{\varepsilon,h}$  and  $\sigma_S^{\varepsilon,\ell}$  all converge to 0 as  $\varepsilon \rightarrow 0$ . However, this is not the case, as shown by the following:

**Proposition 1.** *Under the information perturbations corresponding to (\*), there is no sequence of equilibrium strategies  $\sigma_B^\varepsilon, \sigma_S^\varepsilon$  such that  $\sigma_B^{\varepsilon,h}, \sigma_B^{\varepsilon,\ell}, \sigma_S^{\varepsilon,h}$  and  $\sigma_S^{\varepsilon,\ell}$  all converge to 0 as  $\varepsilon \rightarrow 0$ .*

*Proof of Proposition 1.* Suppose to the contrary that there is a sequence of equilibrium strategies  $\sigma_B^\varepsilon, \sigma_S^\varepsilon$  such that  $\sigma_B^{\varepsilon,h}, \sigma_B^{\varepsilon,\ell}, \sigma_S^{\varepsilon,h}$  and  $\sigma_S^{\varepsilon,\ell}$  all converge to 0 as  $\varepsilon \rightarrow 0$ . In Stage 1, the expected payoff of player  $B$  who received the low signal  $s_B^\ell$  and plays “High ( $H$ )” tends to  $-4$  while the expected payoff of player  $B$  who received the low signal  $s_B^\ell$  and plays “Low ( $L$ )” tends to 0 (here, player  $B$  makes use of the signal distribution (\*) together with the expectation that the seller’s strategies  $\sigma_S^{\varepsilon,h}$  and  $\sigma_S^{\varepsilon,\ell}$  converge to 0 as  $\varepsilon \rightarrow 0$ ,  $B$  believes with high probability that  $S$  does not “Challenge”). Now, in Stage 1, the expected payoff of player  $B$  who received the high signal  $s_B^h$  and plays “High ( $H$ )” tends to 0 while the expected payoff of player  $B$  who received the high signal  $s_B^h$  and plays “Low ( $L$ )” in the limit is below  $\max\{14 - 6 - 9, 7 - 9\} = -1$  (recall that  $B$  believes with high probability that  $S$  chooses “Challenge”). So for  $\varepsilon$  small, there is no  $\sigma$  that makes player  $B$  indifferent between  $H$  and

$L$ , so player  $B$  plays in pure strategies in Stage 1. And as in argument above about pure-strategy equilibrium, the fact that  $B$ 's signal is much more accurate than  $S$ 's implies that such a strategy profile is not an equilibrium. ■

This shows that one appealing property of the unique equilibrium in the MR mechanism under common knowledge (namely, a good equilibrium is a truthful one) can disappear once we introduce small information perturbations. In the next subsection we show the non-robustness of another appealing property of the MR mechanism under common knowledge, namely that it uniquely implements any desired social choice function.

## 2.4 Existence of persistently bad sequential equilibria

So far we have shown that truth-telling is not a robust equilibrium outcome of the MR mechanism when allowing for information perturbations. But in fact one can go further and exhibit arbitrarily small information perturbations for which the above MR mechanism also has a “bad equilibrium” where the buyer reports “Low” regardless of his signal, which in turn leads to a sequential equilibrium outcome which remains bounded away from the sequential (or subgame-perfect) equilibrium outcome under complete information.

Consider the same MR mechanism as before, with the same common prior  $\mu(\theta_H) = 1 - \alpha$  and  $\mu(\theta_L) = \alpha$ , but with the following perturbation  $\nu^\varepsilon$  of signals about  $\theta$ :

$\nu^\varepsilon$	$s_B^h, s_S^h$	$s_B^h, s_S^\ell$	$s_B^\ell, s_S^h$	$s_B^\ell, s_S^\ell$
$\theta_H$	$(1 - \alpha)(1 - \varepsilon^2)$	$(1 - \alpha)\varepsilon^2/3$	$(1 - \alpha)\varepsilon^2/3$	$(1 - \alpha)\varepsilon^2/3$
$\theta_L$	$\alpha\varepsilon^2$	$\alpha\varepsilon/2$	$\alpha\varepsilon/2$	$\alpha(1 - \varepsilon - \varepsilon^2)$

With this signal structure, both agents believe with high probability that if they receive different signals, the signal corresponding to the low state is correct.

In what follows, we shall construct a sequential equilibrium of the perturbed game with prior  $\nu^\varepsilon$  whose outcome differs substantially from that with complete information.

Consider the following strategy profile of the game with prior  $\nu^\varepsilon$ .  $B$  announces “Low” regardless of his signal. If  $B$  has announced “Low,”  $S$  does not challenge regardless of her signal. Off the equilibrium path, i.e., if  $B$  announced “Low” and  $S$  subsequently challenged, then  $B$  always rejects  $S$ 's offer. These are our candidate strategies for sequential equilibrium. To complete the description of the candidate for sequential equilibrium, we also have to assign beliefs over states and signals for each signal of each player and any history of play. Before playing the game but after receiving their private signals, agents' beliefs are given by  $\nu^\varepsilon$  conditioned on their private signals. Similarly, if  $S$  has the opportunity to move (which in turn requires that  $B$  would have played “Low”), we assume that her posterior beliefs are based on  $\nu^\varepsilon$  together with her private signal. Finally, out of equilibrium, if  $B$  is offered the good for price of 6 (which requires that  $S$  will have challenged), we assume that  $B$  always believes with probability one that the state is  $\theta_L$  and that  $S$  has received the low signal  $s_S^\ell$ .

So what we want to show is that for  $\varepsilon > 0$  sufficiently small, the above strategy profile is *sequentially rational* given the beliefs we just described and that, conversely, these beliefs are *consistent* given the above strategy profile. Here we shall check sequential rationality (the basic intuition for the belief consistency part of the proof is given in footnote 13 below). To establish sequential rationality, we solve the game backwards. At Stage 3, regardless of his signal,  $B$  believes with probability one that the state is  $\theta_L$ . Accepting  $S$ 's offer at price of 6 generates  $10 - 9 - 6 = -5$  and rejecting it generates  $5 - 9 = -4$ . Thus, it is optimal for  $B$  to reject the offer. Moving back to Stage 2, if  $S$  chooses "Challenge,"  $S$  anticipates that with probability one, her offer at price of 6 will be rejected by  $B$  in the next stage, thus  $S$  anticipates that, as  $\varepsilon$  becomes small, the payoff is approximately equal to  $7 - 9 = -2$  if her signal is high (equal to  $s_S^h$ ) and to  $5 - 9 = -4$  if the signal is low (equal to  $s_S^\ell$ ). On the contrary, if  $S$  chooses "Not Challenge,"  $S$  guarantees a payoff of 10. Thus, regardless of her signal, it is optimal for  $S$  not to challenge. Moving back to Stage 1,  $B$  "knows" that  $S$  does not challenge regardless of her signal. Now, suppose that  $B$  receives the high signal  $s_B^h$ . Then, as  $\varepsilon$  becomes small,  $B$  believes with high probability that the true state is  $\theta_H$  so that his expected payoff approximately results in  $14 - 10 = 4$ . This is larger than 0, which  $B$  obtains when announcing "High." Therefore, it is optimal for  $B$  to announce "Low." Obviously, this reasoning also shows that when  $B$  has received the low signal  $s_B^\ell$ , it is optimal for her to announce "Low."<sup>13</sup>

As we will see in the next section, the fact that the MR mechanism cannot induce even approximate truth-telling under information perturbations is closely related to the fact that the social choice function we tried to implement is not *Maskin monotonic*. But before we turn to a more general analysis of the non-robustness of subgame-perfect implementation using MR mechanisms, we review Maskin's necessity result on Nash implementation, and explain why the social choice function we try to implement in this example is not Maskin monotonic.

## 2.5 This example does not satisfy Maskin-monotonicity

### 2.5.1 Maskin's Necessity Result on Nash implementation

Recall that a social choice function  $f$  on state space  $\Theta$  is *Maskin monotonic* if for all pair of states of nature (preference profiles)  $\theta'$  and  $\theta''$  if  $a = f(\theta')$  and

$$\left\{ (i, b) \mid u_i(a; \theta') \geq u_i(b; \theta') \right\} \subseteq \left\{ (i, b) \mid u_i(a; \theta'') \geq u_i(b; \theta'') \right\}$$

---

<sup>13</sup>To establish belief consistency, we need to find a sequence of totally mixed strategies that converges toward the pure strategies described above and so that beliefs obtained by Bayes' rule along this sequence also converge toward the beliefs describe above. It is easy to see that under any sequence of totally mixed strategies converging toward the pure strategies described above, the induced sequence of beliefs about  $\theta$  will converge toward  $\nu^\varepsilon$  conditioned on private signals along the equilibrium path of the pure-strategy equilibrium. When  $B$  is offered the good at price of 6,  $S$  has deviated from the equilibrium path due to the "trembles." Beliefs about  $\theta$  are then determined by the relative probability that  $S$  has trembled after the different signals. For instance, if one chooses a sequence of totally mixed strategies under which it becomes infinitely more likely that  $S$  has trembled after receiving  $s_S^\ell$  rather than when receiving  $s_S^h$ , then  $B$  will assign probability close to one to  $S$  receiving signal  $s_S^\ell$ .

(i.e., no individual ranks  $a$  lower when moving from  $\theta'$  to  $\theta''$ ), then  $a = f(\theta'')$ . Here  $u_i(a; \theta)$  denotes player  $i$ 's utility from outcome  $a$  in state  $\theta$ . A *social choice function* (SCF)  $f$  is said to be *Nash implementable* if there exists a mechanism  $\Gamma = (M, g)$  where  $m = (m_1, \dots, m_n) \in M = M_1 \times \dots \times M_n$  denotes a strategy profile and  $g : M \rightarrow A$  is the outcome function (which maps strategies into outcomes), and if for any  $\theta$  the Nash equilibrium outcome of that mechanism in state  $\theta$  is precisely  $f(\theta)$ . Then, Maskin (1999) shows that if  $f$  is Nash implementable, it must be Maskin monotonic.

Let us summarize the proof, which we shall refer to again below. By way of contradiction, if  $f$  were not Maskin monotonic, then there would exist  $\theta'$  and  $\theta''$  such that for any player  $i$  and any alternative  $b$

$$u_i(f(\theta'); \theta') \geq u_i(b; \theta') \implies u_i(f(\theta'); \theta'') \geq u_i(b; \theta'') \quad (\text{I})$$

and nevertheless  $f(\theta') \neq f(\theta'')$ . But at the same time if  $f$  is Nash-implementable there exists a mechanism  $\Gamma = (M, g)$  such that  $f(\theta') = g(m_{\theta'}^*)$  for some Nash equilibrium  $m_{\theta'}^*$  of the game  $\Gamma(\theta')$ . By definition of Nash equilibrium, we must have

$$u_i(f(\theta'); \theta') = u_i(g(m_{\theta'}^*); \theta') \geq u_i(g(m_i, m_{-i, \theta'}^*); \theta'), \forall m_i.$$

But then, from (I) we must also have

$$u_i(f(\theta'); \theta'') = u_i(g(m_{\theta'}^*); \theta'') \geq u_i(g(m_i, m_{-i, \theta'}^*); \theta''), \forall m_i,$$

so that  $f(\theta')$  is also a Nash equilibrium outcome in state  $\theta''$ . But then if the mechanism implements  $f$ , we must have  $f(\theta') = f(\theta'')$ ; a contradiction.

### 2.5.2 The social choice function in our example is not Maskin monotonic

It is easy to show that the social choice function in our Hart-Moore example is not Maskin monotonic. The set of social outcomes (or alternatives)  $A$  is defined as:<sup>14</sup>

$$A = \{(q, y_B, y_S) \in [0, 1] \times \mathbb{R}^2 \text{ such that } y_B + y_S \leq 0\},$$

where  $q$  is the probability that the good is traded from  $S$  to  $B$ ;  $y_B, y_S$  are the transfers of  $B$  and  $S$  respectively; and the utility functions of the seller and the buyer are respectively:

$$u_S(q, y_B, y_S; \theta) = y_S$$

and

$$u_B(q, y_B, y_S; \theta) = \theta q + y_B.$$

The two states of the world are  $\theta_H$  and  $\theta_L$ , which correspond respectively to the good being

---

<sup>14</sup>The sum  $y_S + y_B$  can be negative to allow for penalties paid to a third party.

of high, and of low quality. We have just seen that if an SCF  $f$  under which trade occurs with probability one is Maskin monotonic, then we must have:

$$f(\theta_H) = f(\theta_L).$$

The social choice function we seek to implement requires that

$$\begin{aligned} f(\theta_L) &= (1, -10, 10), \\ f(\theta_H) &= (1, -14, 14). \end{aligned}$$

Clearly  $f(\theta_L) \neq f(\theta_H)$ , but the buyer ranks outcome  $(1, -10, 10)$  at least as high under  $\theta_L$  as under  $\theta_H$ , while the seller has the same preferences in the two states. Thus,  $f$  is not Maskin monotonic, so Maskin's result implies that this  $f$  is not Nash implementable. It is implementable by a Moore-Repullo (MR) mechanism under common knowledge, but it is not implementable by this mechanism under information perturbations.

Our analysis in the next two sections is motivated by the following questions: (1) Is the nonexistence of truth-telling equilibria in arbitrarily small information perturbations of the above MR mechanism linked to the SCF  $f$  being non Maskin monotonic? (2) Is the existence of a sequence of bad sequential equilibria in arbitrarily small information perturbations of the above MR mechanism, directly linked to  $f$  being non-Maskin monotonic?

In Section 3, we consider a more general version of the MR mechanism and link the failure of MR mechanisms to implement truth-telling in equilibrium under information perturbations to the lack of Maskin monotonicity of the corresponding SCF. Then in Section 4, we consider any sequential mechanism that implements a non-Maskin monotonic SCF (more generally, *social choice correspondences* (SCC), formally defined in Section 4.2.1) under common knowledge, and show that for an arbitrarily small information perturbation of the game there exists a bad sequential equilibrium whose outcome remains bounded away from the good equilibrium outcome under common knowledge, even when the size of the perturbation tends to zero.

### 3 More general Moore-Repullo mechanisms

Moore and Repullo (1988) consider a more general class of extensive form mechanisms, which we shall refer to as "MR mechanisms". Under complete information, Moore and Repullo (1988) consider environments where utilities are transferable and show that truth telling is a unique subgame perfect equilibrium in the MR mechanisms. Since this is the most hospitable environment for subgame-perfect implementation, and because most contracting settings are in economies with money, we focus on it.

### 3.1 Setup

Let there be two players 1 and 2, whose preferences over a social decision  $d \in D$  are given by  $(\theta_1, \theta_2) \in \Theta_1 \times \Theta_2 = \Theta$  where  $\Theta_i = \{\theta_i^1, \dots, \theta_i^n\}$  for each  $i = 1, 2$ .<sup>15</sup> The players have utility functions

$$u_1((d, t_1, t_2); \theta_1) = U_1(d; \theta_1) - t_1$$

and

$$u_2((d, t_1, t_2); \theta_2) = U_2(d; \theta_2) + t_2,$$

where  $d$  is a collective decision,  $t_1$  and  $t_2$  are monetary transfers.<sup>16</sup> Preference characteristics  $(\theta_1, \theta_2)$  are common knowledge between the two parties, but not verifiable by a third party.

Let  $f = (D, T_1, T_2)$  be a social choice function where for each  $(\theta_1, \theta_2) \in \Theta_1 \times \Theta_2$  the social decision is  $d = D(\theta_1, \theta_2)$  and the transfers are  $(t_1, t_2) = (T_1(\theta_1, \theta_2), T_2(\theta_1, \theta_2))$ .

Moore and Repullo (1988) propose the following class of mechanisms. These mechanisms involve two phases, where phase  $i$  is designed so as to elicit truthful revelation of  $\theta_i$ . Each phase in turn consists of three stages. The game begins with phase 1, in which player 1 announces  $\theta_1$  and then carries on with phase 2 in which player 2 announces  $\theta_2$ . Phase 1 proceeds as follows:

1. Player 1 announces a preference  $\theta_1$ , and we proceed to stage 2.
2. If player 2 announces  $\phi_1$  and  $\phi_1 = \theta_1$ , then phase 1 ends and we proceed to phase 2. If player 2's announcement  $\phi_1$  does not agree, (i.e.,  $\phi_1 \neq \theta_1$ ) then player 2 "challenges" and we proceed to stage 3.
3. Player 1 chooses between

$$\{x; t_x + \Delta\}$$

and

$$\{y; t_y + \Delta\},$$

where  $x = x(\theta_1, \phi_1)$  and  $y = y(\theta_1, \phi_1)$  depend on both  $\theta_1$  and  $\phi_1$  and  $\Delta$  is a positive number suitably chosen (see below) and  $(x, y, t_x, t_y)$  are such that

$$U_1(x; \theta_1) - t_x > U_1(y; \theta_1) - t_y$$

and

$$U_1(x; \phi_1) - t_x < U_1(y; \phi_1) - t_y.$$

If player 1 chooses  $\{x; t_x + \Delta\}$ , which proves player 2 wrong in his challenge (in the Hart-Moore example above, this corresponds to the buyer refusing the offer at price 6), then player

<sup>15</sup>Moore and Repullo (1988) allow for an infinite state space but impose bounds on the utility functions.

<sup>16</sup>Because we do not assume that the prior on  $\Theta$  is a product measure, the product structure of  $\Theta = \Theta_1 \times \Theta_2$  is not crucial to our results. To see this, note that given any finite set of states of nature  $\Theta$  and utility functions  $u_i : \Theta \times A \rightarrow \mathbb{R}$  for each player  $i$ , we can identify  $\Theta_i$  with the collection of  $\{u_i(\cdot, \theta)\}_{\theta \in \Theta}$ . Now, define  $\tilde{u}_i : \Theta_1 \times \Theta_2 \times A \rightarrow \mathbb{R}$  as follows: for  $\theta_i = u_i(\cdot, \theta)$  we set  $\tilde{u}_i(\cdot, \theta_i) := u_i(\cdot, \theta)$ . This setting is equivalent to the former one.

1 pays  $t_1 = t_x + \Delta$  and player 2 receives  $t_2 = t_x - \Delta$  and a third party receives  $2\Delta$ . However, if player 1 chooses  $\{y; t_y + \Delta\}$ , which confirms player 2's challenge (in the above Hart-Moore example, this corresponds to the buyer taking up the offer at price 6), then player 1 pays  $t_1 = t_y + \Delta$  and player 2 receives  $t_2 = t_y + \Delta$ . The game ends here.

Phase 2 is the same as phase 1 with the roles of players 1 and 2 reversed (i.e., with player 2 announcing  $\theta_2$  in the first stage of that second phase). We use the notation stage 1.2, for example, to refer to phase 1, stage 2.

The Moore-Repullo argument applies as follows when the state of nature  $\theta$  is common knowledge: If player 1 lies at stage 1.1, then player 2 will challenge, and at stage 1.3 player 1 will find it optimal to choose  $\{y; t_y + \Delta\}$ . If  $\Delta$  is sufficiently large, then at stage 1, anticipating player 2's subsequent challenge, player 1 will find it optimal to announce the truth and thereby implement the social choice function  $f$ . Moreover, player 2 will be happy with receiving  $t_y + \Delta$ . If player 1 tells the truth at stage 1.1 then player 2 will not challenge because she knows that player 1 will choose  $\{x; t_x + \Delta\}$  at stage 1.3 which will cause player 2 to pay the fine of  $\Delta$ .

### 3.2 Perturbing the information structure

We now show that this result does not hold for small perturbations of the information structure of the following form: Each agent  $i = 1, 2$  receives a signal  $s_i^{k,l}$  where  $k$  and  $l$  are both integers in  $\{1, \dots, n\}$ ; the set of signals of player  $i$  is denoted  $S_i$ . We assume that the prior joint probability distribution  $\nu^\varepsilon$  over the product of signal pairs and state of nature is such that, for each  $(k, l)$ :

$$\nu^\varepsilon(s_1^{k,l}, s_2^{k,l}, \theta_1^k, \theta_2^l) = \mu(\theta_1^k, \theta_2^l)[1 - \varepsilon - \varepsilon^2]$$

(\*\*\*)

$$\nu^\varepsilon(s_1^{k,l_1}, s_2^{k_2,l}, \theta_1^k, \theta_2^l) = \mu(\theta_1^k, \theta_2^l) \frac{\varepsilon}{n^2 - 1} \text{ for } (k_2, l_1) \neq (k, l)$$

$$\nu^\varepsilon(s_1^{k_1,l_1}, s_2^{k_2,l_2}, \theta_1^k, \theta_2^l) = \mu(\theta_1^k, \theta_2^l) \frac{\varepsilon^2}{n^4 - n^2} \text{ for } k_1 \neq k \text{ or } l_2 \neq l$$

where  $\mu$  is a *complete information* prior over states of nature and signal pairs (i.e., a prior satisfying  $\mu(s_1^{k_1,l_1}, s_2^{k_2,l_2}, \theta_1^k, \theta_2^l) = 0$  whenever  $(k_i, l_i) \neq (k, l)$  for some player  $i$ ). In the above expressions, we abuse notation and write:  $\mu(\theta_1^k, \theta_2^l)$  for the marg $_{\Theta}(\mu)(\theta_1^k, \theta_2^l)$ . This corresponds to an information perturbation such that each player  $i$ 's signal is much more informative about his own preferences than about those of the other player. Note that in an intuitive sense the prior  $\nu^\varepsilon$  is close to  $\mu$  when  $\varepsilon$  is small; this is also true in a formal sense.<sup>17</sup>

<sup>17</sup>For concreteness we specify the supremum-norm topology when discussing the convergence of the priors. That is, let  $\mathcal{P}$  denote the set of priors over  $\Theta \times S$  with the following metric  $d: \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_+$ : for any  $\mu, \mu' \in \mathcal{P}$ ,

$$d(\mu, \mu') = \max_{(\theta, s) \in \Theta \times S} |\mu(\theta, s) - \mu'(\theta, s)|.$$

So, when we say  $\nu^k \rightarrow \mu$ , we mean that  $d(\nu^k, \mu) \rightarrow 0$  as  $k \rightarrow \infty$ .

We begin by considering pure strategy equilibria. For this purpose, we make use of the concept of strategy-proofness:

**Definition 1.** *An SCF  $f$  is **strategy-proof** if for each player  $i$  and each  $\theta_i$ ,*

$$u_i(f(\theta_i, \theta_{-i}), \theta_i) \geq u_i(f(\theta'_i, \theta_{-i}), \theta_i) \text{ for all } \theta'_i \text{ and } \theta_{-i}.$$

In other words, a social choice function  $f$  is strategy-proof if telling the truth is a weakly dominant strategy through a direct mechanism associated with  $f$  whereby the players are asked to announce their preference parameter. Strategy-proofness implies a weak version of Maskin monotonicity, namely, that for any  $\theta, \theta'$  such that

$$\forall i \in N \text{ and } \forall b \in A \setminus \{f(\theta)\} : u_i(f(\theta); \theta_i) \geq u_i(b; \theta_i) \Rightarrow u_i(f(\theta); \theta'_i) > u_i(b; \theta'_i),$$

we have  $f(\theta) = f(\theta')$ .<sup>18</sup> As a corollary, strategy-proofness also implies the usual Maskin monotonicity condition when preferences over outcomes in  $f(\Theta)$  are strict, where  $f(\Theta)$  denotes the range of  $f$ .

**Theorem 1.** *Suppose that a non strategy-proof SCF  $f$  is implementable by a MR mechanism under complete information. Fix any complete information prior  $\mu$ . There exists a sequence of priors  $\{\nu^\varepsilon\}_{\varepsilon>0}$  that converges to the complete information prior  $\mu$  such that there is no pure equilibrium strategies under which player 1 tells the truth in phase 1 and player 2 tells the truth in phase 2.*

*Proof of Theorem 1.* Under the signal structure (\*\*), if player 2 sees that player 1's announcement about  $\theta_1$  is different from her signal, and she believes player 1 is reporting "truthfully," she disregards her own information on  $\Theta_1$  and follows player 1's announcement (and symmetrically for player 1 vis-a-vis player 2 regarding signals over  $\Theta_2$ ).

Now, suppose that  $f$  is not strategy-proof. Then there is a player, say player 1, and states  $\theta_1^h, \theta_1^k, \theta_2^l$  such that

$$u_1(f(\theta_1^h, \theta_2^l); \theta_1^h) < u_1(f(\theta_1^k, \theta_2^l); \theta_1^h).$$

We claim that there is no pure strategy equilibrium in which player 1 reports truthfully in phase 1 and player 2 reports truthfully in phase 2. By way of contradiction, suppose there is such an equilibrium, and suppose that player 1 gets signal  $s_1^{h,l}$  and player 2 gets signal  $s_2^{h,l}$ . Player 1 would like to announce " $\theta_1^k$ " if she expects that subsequent to such an announcement, player 2 agrees with " $\theta_1^k$ " as well and then tells the truth in phase 2 so that the outcome is  $f(\theta_1^k, \theta_2^l)$ . But this is precisely

---

<sup>18</sup>If  $f(\theta) \neq f(\theta')$ , it must be that there is some player  $i$  and some  $\hat{\theta}_{-i}$  such that  $f(\theta_i, \theta_{-i}) = f(\theta_i, \hat{\theta}_{-i}) \neq f(\theta'_i, \hat{\theta}_{-i})$ , and so in particular  $\theta_i \neq \theta'_i$ . Hence, strategy-proofness of  $f$  implies that for this player  $i$ ,  $u_i(f(\theta_i, \theta_{-i}); \theta_i) = u_i(f(\theta_i, \hat{\theta}_{-i}); \theta_i) \geq u_i(f(\theta'_i, \hat{\theta}_{-i}); \theta_i)$  and  $u_i(f(\theta_i, \theta_{-i}); \theta_i) = u_i(f(\theta_i, \hat{\theta}_{-i}), \theta'_i) \leq u_i(f(\theta'_i, \hat{\theta}_{-i}); \theta'_i)$ , and setting  $b = f(\theta'_i, \hat{\theta}_{-i})$  yields the weak monotonicity condition. Finally, note that if preferences over outcomes in  $f(\Theta)$  are strict, then  $u_i(f(\theta_i, \theta_{-i}), \theta'_i) = u_i(f(\theta_i, \hat{\theta}_{-i}), \theta'_i) < u_i(f(\theta'_i, \hat{\theta}_{-i}), \theta'_i)$  and therefore the above argument yields the usual Maskin monotonicity condition. Our weak monotonicity is closely related to conditions proposed by Dasgupta, Hammond, Maskin (1979). In that paper, strategy-proof social choice functions are characterized via the concept of "independent person-by-person monotonicity" which is stronger than our condition of weak Maskin monotonicity.



what will happen: In a fully revealing equilibrium, player 2 will infer that player 1 must have seen a  $s_1^{k,\tilde{l}}$ -type signal, therefore player 2 will believe with high probability that the state must be  $(\theta_1^k, \theta_2^l)$ . Consequently, player 2 will not challenge player 1's announcement. But then, anticipating this, player 1 will announce " $\theta_1^k$ " and thereby receive  $f(\theta_1^k, \theta_2^l)$  instead of  $f(\theta_1^h, \theta_2^l)$ . This in turn shows that there does not exist a truthfully revealing equilibrium in pure strategies. ■

Theorem 1 links the non-robustness of the MR mechanism to the failure of Maskin monotonicity of the social choice function to be implemented. For instance, in the Hart-Moore example in Section 2, the social choice function is not Maskin monotonic and preferences over  $f(\Theta)$  are strict, so the social choice function in that example is not strategy-proof.

Note that the above result does not preclude the existence of mixed strategy equilibria where truth-telling by one or two players in each phase is robust to small information perturbations. Moreover, the above result provides a necessary condition for the robustness of truth-telling by player  $i$  in phase  $i$ , without requiring truth-telling by player  $j$  as well.

Next, we turn attention to mixed-strategy equilibrium. If we require that both players tell (at least, approximately) the truth in each of the two phases, then *no* social choice function  $f = (D, T_1, T_2)$  can be implemented by the general MR mechanism in such a way that truth-telling by both players in each phase, is a sequential equilibrium outcome which is robust to information perturbations.

More formally, in the Appendix we prove the following:

**Theorem 2.** *Suppose that an SCF  $f$  is implementable by a MR mechanism under complete information. Fix any complete information prior  $\mu$ . There exists a sequence of priors  $\{\nu^\varepsilon\}_{\varepsilon>0}$  that converges to the complete information prior  $\mu$  such that there is no sequence of sequential equilibrium strategy profiles that converges to truth-telling.*

Here is an intuition for why requiring approximate truth-telling by both players in each phase precludes robust implementation by the MR mechanism. Suppose that both players receive a signal which is highly correlated with the true state. Player 1 plays first in phase 1, so if player 1 announces a signal that is highly correlated with some state  $\hat{\theta}$ , then player 2 (playing second in phase 1) will believe that player 1 has told the truth (because by assumption player 1's announcement is close to truthful). But the mechanism is built in such a way that player 2 never wants to challenge player 1 if she thinks that player 1 is telling the truth (otherwise at stage 3 player 2 will be punished), so player 2, if she is not challenging, will also announce  $\hat{\theta}$  and so will not follow her private signal and thus she is not reporting truthfully.

Let us make two remarks at this stage. First, the non-robustness of truth-telling as a sequential equilibrium outcome of the MR mechanism is of interest because truth-telling is cognitively simple, and also because the non-existence of a truthful sequential equilibrium implies the non-existence of a desirable pure equilibrium, and implementation theory has mainly focused on pure-strategy equilibria. Second, neither of the non-robustness results of this section rule out the possibility that some SCF  $f$  can be implemented as the limit of mixed-strategy (non-truthful) sequential equilibrium

outcomes.<sup>19</sup> However, in the next section, we will show that if  $f$  is not Maskin monotonic but yet can be implemented by the MR or by any other extensive form mechanism under common knowledge, then there always exist arbitrarily small information perturbations under which there also exist sequential equilibria with undesirable outcomes.

## 4 Any mechanism

In this section, we go beyond MR mechanisms and consider the set of all extensive form mechanisms. Suppose a non-Maskin monotonic SCF is implemented by a (not necessarily MR) mechanism under complete information. Then, we show that there always exists a “bad” sequential equilibrium in arbitrarily small information perturbations of that mechanism. We begin by presenting the argument in a nutshell, using the Hart-Moore example to illustrate our point. Finally, we proceed to state and establish a more general result.

### 4.1 Overview of the main result

In this subsection we state the main result and provide the reader with an intuition for the proof. The main idea is that introducing just a small amount of incomplete information markedly enlarges the set of (sequential) beliefs that are consistent with Bayesian rationality. As a result, one can turn an arbitrary Nash equilibrium of an extensive-form mechanism that implements a non-Maskin monotonic SCF  $f$  under common knowledge into a sequential equilibrium of the perturbed game.

More specifically, suppose there are  $n$  players, and each player  $i$  has a state dependent utility function  $u_i(a; \theta)$  over outcomes (or alternatives)  $a \in A$ . In the perturbations we consider, players do not observe the state of nature  $\theta$  directly, but are informed about it through private signals. An extensive form mechanism  $\Gamma$  together with a state  $\theta \in \Theta$  defines an extensive form game  $\Gamma(\theta)$ ; Let  $SPE(\Gamma(\theta))$  denote the set of subgame-perfect equilibria of the game  $\Gamma(\theta)$ . A SCF  $f$  is said to be subgame-perfect implementable if there exists a mechanism  $\Gamma = (M, g)$  such that for each state  $\theta$ , every subgame-perfect equilibrium outcome coincides with  $f(\theta)$ . Here is an informal statement of the main result:

**Main Result:** Assume finite state space and finite strategy spaces.<sup>20</sup> Assume, further, that a mechanism  $\Gamma$  subgame-perfect implements a non-Maskin monotonic SCF  $f$  under complete information. Then there exists a sequence of information perturbations parametrized by some  $\varepsilon$  and a corresponding sequence of sequential equilibria of the games induced by  $\Gamma$  under this sequence of perturbations, whose outcomes do not converge to  $f(\theta)$  in some state  $\theta$  as  $\varepsilon \rightarrow 0$ .

In particular, under the usual additional conditions where Maskin monotonicity is sufficient for Nash implementation, this result implies the following: whenever an SCF cannot be implemented

<sup>19</sup>For conditions under which the unique subgame perfect equilibrium outcome of a perfect information game remains an equilibrium outcome in perturbed games, see Takahashi and Tercieux (2011).

<sup>20</sup>In the Appendix we extend the result to the case of countable strategy sets.

using static mechanisms (with Nash equilibrium as the solution concept), there is no hope of implementing it using sequential mechanisms if we want such mechanisms to be robust to information perturbations.

**Intuition for the proof:** Suppose that the SCF  $f$  is not Maskin monotonic. Then, there exist  $\theta'$  and  $\theta''$  such that for any player  $i \in N$  and any alternative  $b \in A$

$$u_i(f(\theta'); \theta') \geq u_i(b; \theta') \implies u_i(f(\theta'); \theta'') \geq u_i(b; \theta'') \quad (\text{I})$$

and nevertheless  $f(\theta') \neq f(\theta'')$ . At the same time, since the extensive form mechanism  $\Gamma$  implements  $f$ , there exists a subgame-perfect equilibrium (SPE)  $m_{\theta'}$  in state  $\theta'$  such that  $g(m_{\theta'}) = f(\theta')$ . But then using the same argument as in the proof of Maskin's theorem summarized in Section 2 above,  $m_{\theta'}$  is also a Nash equilibrium in state  $\theta''$ , and necessarily a “bad” Nash equilibrium since  $f(\theta') \neq f(\theta'')$ .

The remaining part of the proof follows from the fact that one can use information perturbations to “rationalize” this bad Nash equilibrium and turn it into a sequential equilibrium of the perturbed games, in the same way as the construction in Section 2 above showed the non-robustness of the particular MR mechanism considered in that section.

As a concrete example, consider again the MR mechanism studied in Section 2. Under common knowledge of the state, it is a Nash equilibrium for  $B$  to announce  $\theta_L$  at Stage 1 and for  $S$  to never challenge at Stage 2. However, this is a bad Nash equilibrium and it is “not” a sequential equilibrium. In particular, if Stage 3 were to be reached under common knowledge, then  $B$  would just infer that  $S$  deviated from the equilibrium, but never update his beliefs about the true valuation  $\theta$  or about  $S$ 's perception of  $\theta$ .

However, perturbing the signals about  $\theta$  changes the picture radically. Now, if stage 3 is reached, then  $B$  updates his beliefs about which signal  $S$  might have seen. In particular, if  $B$ 's updating puts enough weight on  $S$  having received the low signal  $s_S^l$ , then  $B$  will not take the offer at price 6; then, anticipating this at stage 2,  $S$  will indeed not challenge in equilibrium. Note that by perturbing the signal structure we have enlarged the set of consistent beliefs: under common knowledge it could not be a consistent belief that  $S$  saw the low state  $\theta_L$  if  $B$  “knew” that the state was  $\theta_H$ , but this can become consistent under the perturbation. This is the key to how the perturbation turns a bad (non-sequential) Nash equilibrium of the game with complete information into a sequential equilibrium in the perturbed game.

## 4.2 A more formal statement of the main result

Now, we move from intuition and examples to the formal statement of the result, and refer the reader to the Appendix for the formal proof. In the first reading, the reader can skip the rest of Section 4 here and go directly to Section 5 without losing much of the main idea of the paper.

### 4.2.1 The environment

In what follows, we consider a more general environment, with a finite set  $N = \{1, \dots, n\}$  of players, with  $n \geq 2$ , and a set  $A$  of social alternatives, or outcomes. From now on, we no longer assume that agents have quasi-linear preferences with transferable money, as was needed for MR mechanisms. Each player  $i$  has a state dependent utility function  $u_i : A \times \Theta \rightarrow \mathbb{R}$ , where  $\Theta$  is a finite set of states of nature.<sup>21</sup> Players do not observe the state directly, but are informed of the state via signals. Player  $i$ 's signal set is  $S_i$  which, for simplicity, we identify with  $\Theta$ . A signal profile is an element  $s = (s_1, \dots, s_n) \in S \equiv \times_{i \in N} S_i$ . When the realized signal profile is  $s$ , each player  $i$  observes only his own signal  $s_i$ . We let  $\mu$  denote the prior probability over  $\Theta \times S$ . We write  $\mu(\cdot | s_i)$  for the probability measure over  $\Theta \times S$  conditional on  $s_i$ . Let  $s^\theta$  be the signal profile in which each player's signal is  $\theta$ . *Complete information* refers to the environments in which  $\mu(\theta, s) = 0$  whenever  $s \neq s^\theta$  ( $\mu$  will be then referred to as a complete information prior). Under complete information, the state, and hence the full profile of preferences, is always common knowledge among players.

We will assume for each  $i$  and  $\theta$ , the marginal distribution on  $i$ 's signals places strictly positive weight on each of  $i$ 's signals in every state, i.e.,  $\mu(s_i^\theta) \equiv [\text{marg}_{S_i} \mu](s_i^\theta) > 0$ , so that Bayes' rule is well-defined. Note that in case  $\mu$  is a complete information prior, this implies in particular that for each  $(\theta, s^\theta) \in \Theta \times S : \mu(\theta, s^\theta) > 0$ .

A *social choice correspondence* (SCC) is a set-valued mapping  $\mathcal{F} : \Theta \rightrightarrows A$ . We have focused on social choice functions (SCFs) in the previous sections. In this section, we generalize our arguments to encompass social choice correspondences (SCCs).

Since we consider more general extensive form mechanisms than MR mechanisms, we need to introduce some notation. Most of the notation used here is consistent with Moore and Repullo (1988). The reader is referred to that paper for the definition and notation of extensive form mechanisms. We restrict attention to mechanisms that are multi-stage games with observed actions, meaning at each history  $h$ , all players know the entire history of the play, and if more than one player moves at  $h$ , they do so simultaneously.<sup>22</sup> We also assume that the mechanism has a finite number of stages. The class of mechanisms we consider in the present paper is exactly the same as the one Moore and Repullo (1988) allowed. A *mechanism* is then an extensive game form  $\Gamma = (\mathcal{H}, M, \mathcal{Z}, g)$  where: (1)  $\mathcal{H}$  is the set of all histories; (2)  $M = M_1 \times \dots \times M_n$ ,  $M_i = \times_{h \in \mathcal{H}} M_i(h)$  for all  $i$  where  $M_i(h)$  denotes the set of available messages for  $i$  at history  $h$ ; (3)  $\mathcal{Z}$  describes the history that immediately follows history  $h$  given that the strategy profile  $m$  has been played; and (4)  $g$  is the outcome function that maps the set of terminal histories (denoted  $H_T$ ) into the set of outcomes ( $A$ ).

The following notation will be useful: An element of  $M(h) = M_1(h) \times \dots \times M_n(h)$ , say  $m(h) = (m_1(h), \dots, m_n(h))$  is a message profile at  $h$  while  $m_i(h)$  is  $i$ 's message at  $h$ . If  $\#M_i(h) > 1$  and  $\#M_j(h) > 1$  then players  $i$  and  $j$  move simultaneously after history  $h$ , whereas if  $\#M_i(h) > 1$  and

<sup>21</sup>One can always interpret a partition over  $\Theta$  as corresponding to a particular player  $i$ 's set of types  $\Theta_i$ . Thus the set up considered in the previous sections is indeed a special case of that analyzed in this section.

<sup>22</sup>This includes games of perfect information (sequential and observed moves) as a special case.

$\#M_j(h) = 1$  for all  $j \neq i$  then player  $i$  is the only one to move. Histories and messages are tied together by the property that  $M(h) = \{m : (h, m) \in \mathcal{H}\}$ . An element of  $M_i$  is a pure strategy; and an element of  $M$  is a pure strategy profile.

There is an initial history  $\emptyset \in \mathcal{H}$ , and  $h_t =: (\emptyset, m^1, m^2, \dots, m^{t-1})$  is the history at the end of period  $t$ , where for each  $k$ ,  $m^k \in M(h_k)$ . If for  $t' \geq t + 1$ ,  $h_{t'} = (h_t, m^t, \dots, m^{t'-1})$ , then  $h_{t'}$  follows history  $h_t$ . As  $\Gamma$  contains finitely many stages, there is a set of terminal histories<sup>23</sup>  $H_T \subset \mathcal{H}$  such that  $H_T = \{h \in \mathcal{H} : \text{there is no } h' \text{ following } h\}$ . Given any strategy profile  $m$  and any history  $h$ , there is a unique terminal history denoted by  $h_T[m, h]$ . Formally, let  $\mathcal{Z} : M \times \mathcal{H} \rightarrow \mathcal{H}$  be the mapping where

$$\mathcal{Z}[m, h] = \begin{cases} (h, m(h)) & \text{if } h \notin H_T \\ h & \text{otherwise} \end{cases}$$

is the history that immediately follows  $h$  whenever possible given that strategy profile  $m$  has been played; and so  $h_T[m, h] = \lim_{k \rightarrow \infty} \mathcal{Z}^k[m, h]$  where  $\mathcal{Z}^k[m, h] = \mathcal{Z}[m, \mathcal{Z}^{k-1}[m, h]]$ . Finally, the outcome function  $g : H_T \rightarrow A$  specifies an outcome for each terminal history. We will also denote  $g(m; h)$  the outcome that obtains when players use strategy profile  $m$  starting from history  $h$ , i.e.,  $g(m; h) = g(h_T[m, h])$ . In what follows, we only consider finite mechanisms:

**Assumption A1.**  $M_i(h)$  is finite for each  $i$  and  $h$ .

**Remark 1.** This assumption is useful when using sequential equilibrium and avoids technical complications due to the use of countably infinite (or uncountable) spaces. In the Appendix, we provide additional assumptions on the class of mechanisms so that our result can be extended to countable message spaces. This extension is important because the literature often uses integer games (i.e., games where one dimension of the message space is the set of positive integers) as part of implementing mechanisms.<sup>24</sup>

A mechanism  $\Gamma$  together with a state  $\theta \in \Theta$  defines an extensive game  $\Gamma(\theta)$ . A (pure strategy) Nash equilibrium for the complete information game  $\Gamma(\theta)$  is an element  $m^* \in M$  such that, for each player  $i$ ,  $u_i(g(m^*; \emptyset); \theta) \geq u_i(g((m_i, m_{-i}^*); \emptyset); \theta)$  for all  $m_i \in M_i$ . A (pure strategy) subgame-perfect equilibrium for the game  $\Gamma(\theta)$  is an element  $m^* \in M$  such that, for each player  $i$ ,  $u_i(g(m^*; h); \theta) \geq u_i(g((m_i, m_{-i}^*); h); \theta)$  for all  $m_i \in M_i$  and all  $h \in \mathcal{H} \setminus H_T$ . Recall that  $SPE(\Gamma(\theta))$  denotes the set of subgame-perfect equilibria of the game  $\Gamma(\theta)$  and  $NE(\Gamma(\theta))$  denotes the set of Nash equilibria of the game  $\Gamma(\theta)$ . We say that a mechanism implements an SCC  $\mathcal{F}$  in subgame-perfect equilibrium, or simply SPE-implements  $\mathcal{F}$ , if for each  $(\theta, s^\theta) \in \Theta \times S$ , we have  $g(SPE(\Gamma(\theta)); \emptyset) = \mathcal{F}(\theta)$ .

Given a prior  $\mu$ , the mechanism determines a Bayesian game  $\Gamma(\mu)$  in which each player's type is his signal, and after observing his signal, player  $i$  selects a (pure) strategy from the set  $M_i$ . In what follows, whenever players face uncertainty about the state and other player's signals, they possess

<sup>23</sup>Note that  $M(h) = \{m : (h, m) \in \mathcal{H}\} = \emptyset$  for any  $h \in H_T$ .

<sup>24</sup>Our results do not critically depend on the countability assumption. We believe that our results would hold for arbitrary mechanisms if we were to use perfect Bayesian equilibrium (Fudenberg and Tirole (1991b)) instead of sequential equilibrium as the solution concept.

a probabilistic belief over this uncertainty and with respect to this belief, they aim to maximize expected utility.<sup>25</sup> A strategy profile  $\sigma = (\sigma_1, \dots, \sigma_n)$  lists a strategy for each player  $i$  where  $\sigma_i : S_i \rightarrow M_i$  and  $\sigma_i(h_t, s_i)$  is a message in  $M_i(h_t)$  given history  $h_t$  and signal  $s_i$ . Alternatively, we will sometimes let  $\sigma_i$  be a (mixed) behavior strategy, i.e., a function that maps the set of possible histories and signals into the set of probability distributions over messages:  $\sigma_i(\cdot | h_t, s_i) \in \Delta(M_i(h_t))$  is the probability distribution over  $M_i(h_t)$  given history  $h_t$  and signal  $s_i$ .

With this notation in place we can re-state the definition of *sequential equilibrium* as specialized to these multi-stage games of observed actions. A sequential equilibrium is a profile of assessment (or beliefs)  $\phi$  and strategies  $\sigma$  which satisfy both *consistency* and *sequential rationality*. Here consistency is the requirement that there exists a sequence of totally mixed strategy profiles  $\sigma^n$  converging to  $\sigma$  such that the beliefs  $\phi^n$  computed from  $\sigma^n$  using Bayes' rule converge to  $\phi$ ; Sequential rationality means that for each period  $t$  and history  $h^{t-1}$  up to  $t-1$ , the continuation strategies are optimal for each player  $i$  given the opponents' strategies and his belief  $\phi_i$ . A more formal definition of sequential equilibrium can be found in the Appendix (Section 7.2).

#### 4.2.2 The existence of a bad sequential equilibrium with almost-perfect information

Although we already introduced the definition of Maskin monotonicity for social choice functions in Section 2, we need to extend it to social choice correspondences. A social choice correspondence  $\mathcal{F}$  on a payoff relevant state space  $\Theta$  is *Maskin monotonic* if for all pair of states of nature  $\theta'$  and  $\theta''$  if  $a \in \mathcal{F}(\theta')$  and

$$\{(i, b) \mid u_i(a; \theta') \geq u_i(b; \theta')\} \subseteq \{(i, b) \mid u_i(a; \theta'') \geq u_i(b; \theta'')\} \quad (*)$$

(i.e., no individual ranks  $a$  lower when moving from  $\theta'$  to  $\theta''$ ) then  $a \in \mathcal{F}(\theta'')$ . Henceforth, we assume that  $A$  is a Hausdorff topological space. We are now in a position to provide a more formal statement of our main theorem.

**Theorem 3.** *Assume A1. Suppose that a mechanism  $\Gamma$  SPE implements a non-Maskin monotonic SCC  $\mathcal{F}$ . Fix any complete information prior  $\mu$ . There exists a sequence of priors  $\{\nu^\varepsilon\}_{\varepsilon>0}$  that converges to a complete information prior  $\mu$  and a corresponding sequence of sequential equilibrium assessments and strategy profiles  $\{(\phi^\varepsilon, \sigma^\varepsilon)\}_{\varepsilon>0}$  such that as  $\varepsilon$  tends to 0,  $g(\sigma^\varepsilon(s^\theta); \emptyset) \rightarrow a \notin \mathcal{F}(\theta)$  for some  $\theta \in \Theta$  and some outcome  $a \in A$ .*

*Proof.* See Appendix. ■

**Remark 2.** The above theorem shows the following: under the usual conditions ensuring that Maskin monotonicity is sufficient for Nash implementation, whenever a SCC cannot be implemented using static mechanisms, this SCC cannot be implemented using an extensive-form mechanism that is robust to the introduction of a small amount of incomplete information.

---

<sup>25</sup>All the results extend to more general representations for preferences under uncertainty. The interested reader is referred to Kunimoto and Tercieux (2009) for details.

**Remark 3.** While non-Maskin monotonic SCFs cannot be robustly implemented, things are quite different for Maskin monotonic SCFs. Here we restrict our focus to SCF's rather than SCCs. In the Appendix we extend the argument to the case of SCCs.

What appears as a natural candidate for “robust implementation” of a SCF amounts to constructing a Nash implementable mechanism with the following two properties: (1) there exists at least one strict Nash equilibrium; and (2) the map from information structures to Nash equilibria has a closed graph, so adding a small amount of incomplete information only slightly increases the set of Nash equilibria. In the Appendix, we formalize the above two properties and propose a definition of robust Nash implementation.

To see this, note that the first property ensures that the strict Nash equilibrium continues to be a strict (Bayesian) Nash equilibrium for any nearby environment and hence that there is always a good equilibrium for any nearby environment. The second property in turn ensures that all Nash equilibria will continue to have outcomes that are close to the desired outcome for any nearby environment.

Regarding the first property, the existence of a strict Nash equilibrium in a mechanism that implements a SCF can easily be ensured under a slight strengthening of Maskin monotonicity, namely strong Maskin monotonicity. In the Appendix, we show that this is also the case for SCCs.

As to the second property, in many situations, Nash implementation of Maskin monotonic SCFs can be achieved using finite mechanisms (see Saijo (1988)). Routine arguments then imply that the second property is satisfied.<sup>26</sup>

For the case of infinite mechanisms, the argument is relegated to the Appendix, which provides sufficient conditions under which one can ensure that properties (1) and (2) above are satisfied. There we take care of SCCs as well as SCFs. Interestingly, these sufficient conditions are satisfied by any Maskin monotonic SCF in quasi-linear environments with money.

## 5 Outside options and the hold-up problem

Thus far, we have shown that the mechanisms used by proponents of the “implementation critique” of the property right theory of the firm (e.g., Maskin and Tirole, 1999a) are themselves not robust to small deviations from perfect information and common knowledge. That leaves open the question of what role outside options (e.g., as induced by asset ownership as in Grossman and Hart (1986)) can play in alleviating the hold-up problem in situations that depart more significantly from complete or just symmetric information.

As a first step in this direction, we consider an environment with an ex ante investment stage and where ex post bargaining takes place under one-sided asymmetric information. We present an example where the presence of an outside option allows mechanisms that approximate ex ante efficiency. Moreover, we argue that static or sequential mechanisms without an outside option

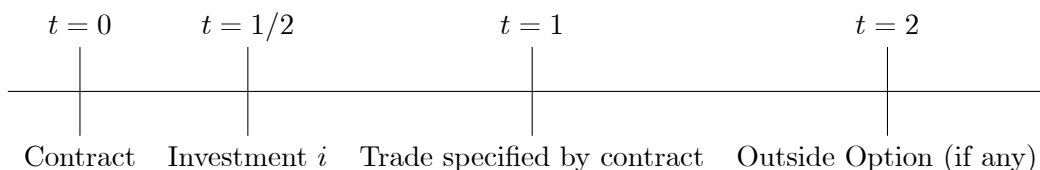
---

<sup>26</sup>This property comes from the following two facts. First, a small change in the prior probability corresponds to a small change in ex ante payoffs. Second, the pure Nash equilibrium correspondence is upper hemi continuous in the space of payoffs.

cannot do as well, which, in turn, we see as a justification for the role of ownership allocation in contracting under incomplete information.

## 5.1 The set-up

Suppose there is a buyer ( $B$ ) and a seller ( $S$ ) of a single unit of an indivisible object with utility  $\tilde{v}$  to the buyer, where  $\tilde{v} \in \{\underline{v}, \bar{v}\}$  and  $\bar{v} > \underline{v} > 0$ . The utility of the seller for the object is assumed to be always zero. Time is discrete, with a contracting period 0 where the good is offered to the buyer at a prespecified price, an investment period 1/2 whereby the seller can increase the buyer's valuation for the good; and a trading period 1. Investment is unobservable as in Grossman-Hart (1986). Moreover, we shall allow for the possibility that an outside option can be exerted in period 2 by one party if trade does not occur in period 1, and focus attention on the case where the outside option yields utility  $\underline{v}$  to whoever has the good at that point. A natural interpretation is that  $\bar{v}$  is the value the buyer and the seller can generate in their relationship and  $\underline{v}$  is the default value that can be generated outside of the relationship. The timing of the events is as follows:



The seller may make an investment in period 1/2 that increases the probability that the good is high quality, as in Che and Hausch (1999). Specifically, suppose that at cost  $c(i)$  the seller achieves  $v = \bar{v}$  with probability  $i$ , where  $c(\cdot)$  is continuous, twice differentiable, and satisfies  $c'(i) > 0$ ,  $c''(i) > 0$ ,  $c(0) = 0$ ,  $c'(1) = +\infty$ , and  $c'(0) < \bar{v} - \underline{v}$ . The buyer will know the value of the good at the beginning of period 1, while the seller will not, so there is one-sided asymmetric information.

## 5.2 Outside options as ownership

One can relate the outside option to the idea of ownership by taking the owner of the good to be the party with the right to exercise the outside option. Thus, under seller ownership, if the seller makes an offer to the buyer but the buyer refuses the offer, then the seller can always choose to always exert his outside option and gets  $\underline{v}$ .

This interpretation as ownership is consistent with other works in the property rights literature, starting with Grossman and Hart (1986), where ownership of the assets of a firm allows the owner to make alternative use of these assets in case of disagreement in the ex post bargaining with the other party(ies). This in turn enhances the owner's ex post bargaining power and therefore it increases the fraction of the ex post production surplus the owner can secure in this bargaining, which, in turn, enhances the owner's investment incentives. In our setting too, ownership of the good will allow the seller to extract a higher price from a high-valuation buyer, and anticipating this, the seller will invest a higher  $i$  in the relationship. What we show below is that no mechanism



(contract) without an outside option can do as well as a contract with outside option to the seller in inducing efficient investment by the seller in period 1/2.

### 5.3 Ex Ante efficiency and outside options

Under our assumptions, the ex-ante efficient outcome is to trade whenever the good is high quality, consume the outside option when the good is low quality,<sup>27</sup> and set investment equal to  $i^*$ , where  $i^* \in (0, 1)$  is the solution to the following first-order condition:

$$\bar{v} - \underline{v} = c'(i^*).$$

The resulting total surplus is then

$$W^* = i^* \bar{v} + (1 - i^*) \underline{v} - c(i^*).$$

We will now show how a mechanism with an outside option can come arbitrarily close to this payoff.

In this setting, a mechanism takes as input the buyer's announced value for the good, and specifies a trade probability  $q$ , transfers  $y_S$  and  $y_B$  to the seller and buyer respectively, a probability  $z_S$  that the seller gets to keep the good if there is no trade, a probability  $z_B$  that the buyer gets the good in that case, and therefore the probability  $1 - z_B - z_S \geq 0$  that the good is destroyed when it is not traded (the mechanism does not specify an investment level, nor condition other outcomes on it, as investment is not observable). Thus the mechanism maps the buyer's announcement  $\tilde{v} \in \{\underline{v}, \bar{v}\}$  into  $A$  where  $A = \{(q, y_B, y_S, z_B, z_S) \in [0, 1] \times \mathbb{R}_+^4 \mid y_S + y_B \leq 0, z_B + z_S \leq 1\}$ . In what follows, we consider the case  $z_S \equiv 1$  (so that the seller gets the outside option whenever there is no trade, regardless of the buyer's announcement), and therefore the mechanism boils down to a mapping  $f(\tilde{v})$  such that  $f(\underline{v}) = (q, \underline{y}_B, \underline{y}_S)$  (when the buyer announces  $\underline{v}$ ) and  $f(\bar{v}) = (\bar{q}, \bar{y}_B, \bar{y}_S)$  (when the buyer announces  $\bar{v}$ ).

Given that  $z_S \equiv 1$ , for  $\varepsilon > 0$  small enough, the mechanism that implements  $(1, -(\bar{v} - \varepsilon), \bar{v} - \varepsilon)$  when the buyer announces valuation  $\bar{v}$ , and  $(0, 0, 0)$  when the buyer announces  $\underline{v}$  satisfies incentive compatibility (it is a strictly dominant strategy for the buyer to report her valuation  $v$  truthfully), individual rationality, and ex post efficiency, i.e., trade occurs if and only if there are social gains from trade.

Now suppose that the buyer and the seller agree on this mechanism with the outside option  $\underline{v}$  allocated to the seller at the contracting stage. Then, moving back to time  $t = 1/2$ , the seller chooses the level of investment to maximize

$$i(\bar{v} - \varepsilon) + (1 - i)\underline{v} - c(i).$$

Given our assumptions, the optimal investment level  $i^*$  (for  $\varepsilon > 0$  small enough) is determined by

---

<sup>27</sup>From the viewpoint of social welfare it does not matter which party gets to use the outside option.

the first-order condition:

$$\bar{v} - \varepsilon - \underline{v} = c'(i^*).$$

From the concavity of the problem, this is approximately the same as the first-best investment when  $\varepsilon$  is small. Thus, a simple contract with seller's ownership can exactly implement an outcome whose total surplus is arbitrarily close to the first best level; this is what we will mean by "approximate ex ante efficiency."

#### 5.4 Ex ante efficiency cannot be approximated without outside options

As in the complete information case, a crucial question is: what exactly can be achieved with contracts/mechanisms that do not use outside options, so that  $z_S = z_B = 0$ . Below, we show that under buyer's private information, any "outside-option-free" contract between the buyer and the seller leads to an outcome which remains bounded away from ex ante efficiency.

First, note that if an SCF  $f$  that maps the true buyer's valuation  $\tilde{v}$  onto a triplet  $f(\tilde{v}) = (\tilde{q}, \tilde{y}_B, \tilde{y}_S)$ , and yields zero continuation utility to both parties if trade does not occur, is to be implemented by some (static or sequential<sup>28</sup>) mechanism in Bayesian Nash equilibrium, it must be at least weakly incentive compatible for the buyer to report truthfully. It is simple to show that  $f$  is incentive compatible if and only if

$$v(\bar{q} - \underline{q}) \leq \underline{y}_B - \bar{y}_B \leq \bar{v}(\bar{q} - \underline{q}). \quad (1)$$

Below we prove that one cannot find SCFs with  $z_S = z_B = 0$  that are incentive compatible and approximately ex ante efficient. To show this, suppose to the contrary that for any  $\varepsilon > 0$  there is an incentive compatible mechanism  $f^\varepsilon$  whose ex ante total surplus is at least  $W^* - \varepsilon$ . Then, the associated probabilities  $i^\varepsilon$  of high quality must converge to  $i^*$ , the probabilities of trade  $\underline{q}^\varepsilon$  and  $\bar{q}^\varepsilon$  must both converge to 1, and the difference in transfers (i.e. money "burnt")  $|\underline{y}_S^\varepsilon - \underline{y}_B^\varepsilon|$  and  $|\bar{y}_S^\varepsilon - \bar{y}_B^\varepsilon|$  must both converge to 0. The incentive compatibility condition (1) then implies that  $|\underline{y}_B^\varepsilon - \bar{y}_B^\varepsilon| \rightarrow 0$ , and this, plus the fact that both  $|\underline{y}_S^\varepsilon - \underline{y}_B^\varepsilon|$  and  $|\bar{y}_S^\varepsilon - \bar{y}_B^\varepsilon| \rightarrow 0$ , implies that  $|\bar{y}_S^\varepsilon - \underline{y}_S^\varepsilon| \rightarrow 0$  as well.

Moving back to time  $t = 1/2$ , the seller will choose investment  $i$  to maximize

$$i\bar{y}_S^\varepsilon + (1 - i)\underline{y}_S^\varepsilon - c(i) = \underline{y}_S^\varepsilon + i(\bar{y}_S^\varepsilon - \underline{y}_S^\varepsilon) - c(i).$$

Because  $|\bar{y}_S^\varepsilon - \underline{y}_S^\varepsilon| \rightarrow 0$  and  $c' > 0$ , the solution  $i^\varepsilon$  to this problem converges to 0, so investment falls far short of the first-best level, which is not consistent with the assumption that the ex ante total surplus converges to  $W^*$ . We conclude that ex ante surplus must be bounded away from efficiency.

---

<sup>28</sup> Approximate ex ante efficiency cannot be achieved by virtual implementation either, since incentive compatibility is also necessary for virtual implementation to work. But precisely we show below that without outside options, one cannot find SCFs that are both, approximately ex ante efficient and incentive compatible.

This establishes our claim that no approximately ex ante efficient SCF can be implemented by a mechanism that does not include an outside option (or some other change to the economic environment), which, in turn, we see as a justification for the role of ownership allocation in contracting under incomplete information.

## 5.5 Summary

Analyzing the hold-up problem in a setting with ex post asymmetric information, as we have done in this section, yields an interesting new insight: outside options such as those induced by asset ownership can help relax incentive compatibility constraints and thereby improve ex ante efficiency compared to what can be achieved through “ownership-free” contracts/mechanisms.

## 6 Concluding remarks

We conclude by making a few additional remarks. First, the bad sequential equilibria in Section 4 survives a standard equilibrium selection criterion. Cho (1987) defines *forward induction equilibrium*, which is an extension of the Cho and Kreps (1987) *intuitive criterion* in signaling games to more general games. The key restriction in this equilibrium concept is that the belief system assigns probability 0 to nodes in some information set  $h$  if this node can be reached only by “bad” deviations, provided that other nodes in  $h$  can be reached by non-bad deviations. Here, “bad deviations” are deviations with the following property: suppose that at any information set where the deviating player can reach by deviating, players are playing best-responses against some arbitrary belief that is consistent with that information set being reached. Then the deviation makes the deviating player strictly worse off compared to his equilibrium payoff. In the Hart-Moore example developed in Section 2, we can show that “Challenge” is never a bad deviation for the seller. To see this, note that when deviating to “Challenge,” the seller may think that an information set under which  $B$  believes that the state  $\theta_H$  may occur with positive probability. Thus we can always pick an appropriate belief (for instance, one that would assign probability 1 to  $\theta_H$ ) under which it is a best reply for  $B$  to accept  $S$ ’s offer if  $S$  challenges. But we know that in such a case “Challenge” by the seller makes her strictly better off compared to the equilibrium, proving that “Challenge” cannot be a bad deviation.

Our second remark is that the non-robustness of subgame-perfect implementation does not mean that implementation is hopeless, but, rather, suggests that we should further explore the implications of Nash implementation. It is well-known that in many important contexts, Nash implementation (or Maskin monotonicity) is quite demanding. For instance, a well-known result by Muller and Satterthwaite (1977) states that any onto and ex post efficient social choice function defined on the domain of all strict preferences is dictatorial when there are at least three outcomes. Maskin (1999) shows that with only two players, this result extends to social choice correspondences. However, it has also been shown that, under some mild domain restrictions, for any social choice function  $f$ , there is a stochastic social choice function that puts probability close to one on the same

outcomes as  $f$  and that is Maskin monotonic (See Abreu and Sen (1991) and Matsushima (1988) for the details of this approach).<sup>29</sup> Indeed, we saw that the social choice function  $f$  we sought to implement in this Hart-Moore example was not Maskin monotonic since  $f(\theta_L) = (1, -10, 10) \neq f(\theta_H) = (1, -14, 14)$ , and therefore not Nash implementable. However, the  $\varepsilon$ -approximation of that SCF defined by

$$f^\varepsilon(\theta_L) = (1 - \varepsilon, -10, 10) \neq f(\theta_H) = (1 - \varepsilon, -14, 14),$$

is Maskin monotonic since for example,  $B$  strictly prefers  $(1 - \varepsilon, -10, 10)$  to  $(1, -10 - 11\varepsilon, 10)$  when  $\theta = \theta_L = 10$  but the reverse is true when  $\theta = \theta_H = 14$ . Hence, even if  $f$  is not Maskin monotonic and therefore not Nash implementable, we can find an  $\varepsilon$ -close stochastic SCF that is Maskin monotonic and therefore Nash implementable for instance in the Moore and Repullo's setting.<sup>30</sup> However, the stochastic nature of this mechanism is problematic in terms of *renegotiation-proofness*. For example, if we consider the above social choice function  $f^\varepsilon$ : with probability  $\varepsilon$ , the planner must induce a bad outcome under which trade does not occur.<sup>31</sup> Given that there are gains from trade, agents will definitely have incentives to renegotiate. If this possibility is explicitly taken into account by the contracting parties, then the social choice function is not going to be Nash implementable anymore. Thus, stochasticity (or randomness) can help to robustly implement nearby efficient social choice functions but it also raises serious renegotiation-proofness issues.

Finally, we feel that laboratory experiments can be useful in assessing the importance of the effect of information perturbations on the likelihood that truth-telling will still occur in equilibrium. Preliminary work by Aghion, Fehr, Holden and Wilkenning (2009) suggests that the effect is potentially large.<sup>32</sup>

---

<sup>29</sup>Here preferences are defined on lotteries over outcomes and agents are assumed to be expected utility maximizers, so typically the restrictions to domains of strict preferences in Muller and Satterthwaite (1977) or in Maskin (1999) are not going to be satisfied.

<sup>30</sup>Note that in the Moore-Repullo setting (i.e., with quasi-linear utilities and arbitrary large transfers), for any social choice function  $f$ , we have the existence of a bad outcome (i.e., an outcome which, in each state of nature, is strictly worse for all players than any outcome in the range of the social choice function). In addition, because for each agent, there is no most preferred outcome,  $f$  also satisfies no-veto-power. Thus by Moore and Repullo (1990, Corollary 3, p.1094)  $f$  is Nash implementable if and only if  $f$  is Maskin monotonic. The stochastic approximation of  $f$  can therefore be implemented with a canonical Maskin mechanism, although since the mechanism uses integer games it is less appealing than the simple MR mechanism.

<sup>31</sup>Renegotiation is less problematic in the case of "exact" Nash implementation since renegotiation then only occurs *out of equilibrium*.

<sup>32</sup>Aghion, Fehr, Holden and Wilkenning (2009) conduct a laboratory experiment testing the robustness of a Moore-Repullo mechanism to information perturbations. The experiment is meant to mimic the Hart-Moore example spelled out in Section 2. Subjects are randomly allocated to the buyer and seller roles, and play the mechanism ten times in a row. In one treatment there is complete information, in the other the subjects each receive a conditionally independent private signal which is 90% accurate-generated by the subjects drawing different colored balls from an urn. In the complete information treatment the proportion of buyers who announce low despite having a high signal declines from around 40% to 10% over the ten rounds. By contrast, in the incomplete information treatment buyers continue to lie more than 40% of the time. In periods 6-10 the average number of lies in the complete information treatment is 24%, whereas it is 42% in the incomplete information treatment.

## References

- [1] Abreu, D, and H. Matsushima (1992), “Virtual Implementation in Iteratively Undominated Strategies: Complete Information”, *Econometrica* 60, 993-1008
- [2] Abreu, D. and A. Sen (1991), “Virtual Implementation in Nash Equilibrium,” *Econometrica*, 59, 997-1021.
- [3] Aghion, P., M. Dewatripont and P. Rey. (1994), “Renegotiation Design with Unverifiable Information,” *Econometrica* 62, 257-282.
- [4] Aghion, P., E. Fehr, R. Holden and T. Wilkening. (2009), “Subgame Perfect Implementation: A Laboratory Experiment,” unpublished working paper.
- [5] Aghion P., D. Fudenberg, and R. Holden (2009) “Subgame Perfect Implementation With Almost Perfect Information,” NBER working paper w15167.
- [6] Aliprantis, C and K. Border (1999): *Infinite Dimensional Analysis*, 2nd. ed., Springer Verlag, Berlin.
- [7] Che, Y. and D. Hausch (1999), “Cooperative Investments and the Value of Contracting,” *American Economic Review* 89, 125-147.
- [8] Cho, I.K. (1987), “A Refinement of Sequential Equilibrium,” *Econometrica* 55, 1367-1389
- [9] Cho, I.K. and D. Kreps (1987), “Signaling Games and Stable Equilibria”, *Quarterly Journal of Economics*, 102, 179-221.
- [10] Chung, K.S and J. Ely (2003), “Implementation with Near-Complete Information,” *Econometrica* 71, 857-871.
- [11] Cremer, J. and R.P. McLean (1988), “Full Extraction of the Surplus in Bayesian and Dominant Strategy Auctions”, *Econometrica* 56, 1247-1257.
- [12] van Damme, E. and S. Hurkens (1997), “Games with Imperfectly Observable Commitment,” *Games and Economic Behavior* 21, 282-308.
- [13] Dasgupta, P., P. Hammond and E. Maskin (1979), “The Implementation of Social Choice Rules: Some General Results on Incentive Compatibility,” *Review of Economic Studies*, 46, 185-216.
- [14] Dekel, E. and D. Fudenberg (1990), “Rational Play Under Payoff Uncertainty,” *Journal of Economic Theory* 52, 243-267.
- [15] Fudenberg, D., D. Kreps, and D.K. Levine (1988), “On the Robustness of Equilibrium Refinements,” *Journal of Economic Theory* 44, 354-380.

- [16] Fudenberg, D., D.K. Levine, and E. Maskin (1991), "Balanced-Budget Mechanisms for Adverse Selection Problems," unpublished working paper.
- [17] Fudenberg D. and J. Tirole (1991a), *Game Theory*, MIT Press
- [18] Fudenberg, D. and J. Tirole (1991b), "Perfect Bayesian and Sequential Equilibrium," *Journal of Economic Theory* 53 (1991), 236-260.
- [19] Grossman, S, and O. Hart (1986), "The Costs and Benefits of Ownership: A Theory of Vertical and Lateral Integration," *Journal of Political Economy* 94, 691-719.
- [20] Hart, O. and J. Moore (2003), "Some (Crude) Foundations for Incomplete Contracts," unpublished working paper.
- [21] Hendon, E., H.J. Jacobsen and B. Sloth (1996), "The One-Shot Deviation Principle for Sequential Rationality," *Games and Economic Behavior* 12, 274-282.
- [22] Jackson, M (2001), "A Crash Course in Implementation Theory," *Social Choice and Welfare*, 18, 655-708.
- [23] Johnson, S, J.W. Johnson and R.J. Zeckhauser (1990), "Efficiency Despite Mutually Payoff-Relevant Private Information: The Finite Case," *Econometrica* 58, 873-900.
- [24] Kreps, D.M. and R. Wilson (1982), "Sequential Equilibria," *Econometrica* 50, 863-894
- [25] Kunimoto, T (2010), "How Robust is Undominated Nash Implementation?," mimeo
- [26] Kunimoto, T. and O. Tercieux (2009) "Implementation with Near-Complete Information: The Case of Subgame Perfection," mimeo
- [27] Maskin, E (1999), "Nash Equilibrium and Welfare Optimality," *Review of Economic Studies* 66, 23-38.
- [28] Maskin, E. and J. Moore (1999), "Implementation and Renegotiation," *Review of Economic Studies*, 66, 39-56.
- [29] Maskin, E. and J. Tirole (1999), "Unforeseen Contingencies and Incomplete Contracts," *Review of Economic Studies* 66, 83-114.
- [30] Maskin, E. and J. Tirole (1999b), "Two Remarks on the Property-Rights Literature," *Review of Economic Studies* 66, 139-149.
- [31] Matsushima, H (1988), "A New Approach to the Implementation Problem," *Journal of Economic Theory*, 45, 128-144.
- [32] Monderer, D. and D. Samet (1989), "Approximating Common Knowledge with Common Beliefs," *Games and Economic Behavior* 1, 170-190.

- [33] Moore, J. and R. Repullo (1988), "Subgame Perfect Implementation," *Econometrica* 56, 1191-1220.
- [34] Moore, J. and R. Repullo (1990), "Nash Implementation: A Full Characterization," *Econometrica*, 58, 1083-1099.
- [35] Muller, E. and M. A. Satterthwaite (1977), "The Equivalence of Strong Positive association and Strategy-Proofness," *Journal of Economic Theory*, 14, 412-418.
- [36] Oury, M. and O. Tercieux (2009), "Continuous Implementation," working paper, available at <http://www.pse.ens.fr/tercieux/ContUTP.pdf>.
- [37] Saijo, T. (1988), "Strategy Space Reduction in Maskin's Theorem: Sufficient Conditions for Nash Implementation," *Econometrica*, 56, 693-700.
- [38] Selten, R (1965), "Spieltheoretische Behandlung Eines Oligopolmodells mit nachfragetrageit," *Zeitschrift fur die gesamte Staatswissenschaft*, 121, 301-24.
- [39] Takahashi, S. and Tercieux, O. (2011), "Robust Equilibria in Sequential Games under Almost Common Certainty of Payoffs," mimeo.

## 7 Appendix

### 7.1 Proof of Theorem 2

We fix a SCF  $f$  which is implemented under complete information by the MR mechanism  $\Gamma^{MR}$ . We let  $\mu$  be a complete information prior and show that for the sequence of priors  $\nu^\varepsilon$  (indexed by  $\varepsilon > 0$ ) as specified in (\*\*\*) of Section 3.2, there is no sequence of equilibrium strategy profiles converging to truth-telling. Let  $\Gamma^{MR}(\nu^\varepsilon)$  be an incomplete information game associated with the MR mechanism and a prior  $\nu^\varepsilon$ . By way of contradiction, assume that for each  $\varepsilon > 0$ , there exists a profile of mixed equilibrium strategies of the game  $\Gamma^{MR}(\nu^\varepsilon)$  such that as  $\varepsilon$  goes to 0, the probability that both players report their signals truthfully converges to 1. Fix such a sequence of mixed equilibrium strategy profiles. We then use the following notation to describe equilibrium play in the games  $\Gamma^{MR}(\nu^\varepsilon)$ :

- $\sigma_{k,l}^{\varepsilon,j}$  denotes the probability that player 1 with signal  $s_1^{k,l}$  announces  $\theta_1^j$  at Stage 1 of Phase 1;
- $\lambda_{k,l}^{\varepsilon,j}[\hat{\theta}_1]$  denotes the probability that player 2 with signal  $s_2^{k,l}$  announces  $\theta_1^j$  at Stage 2 of Phase 1 given that at Stage 1 of Phase 1, player 1 has announced  $\hat{\theta}_1$
- $\rho_{k,l}^{\varepsilon,j}$  denotes the probability that player 2 with signal  $s_2^{k,l}$  announces  $\theta_2^j$  at Stage 1 of Phase 2; and
- $\tau_{k,l}^{\varepsilon,j}[\hat{\theta}_2]$  denotes the probability that player 1 with signal  $s_1^{k,l}$  announces  $\theta_2^j$  at Stage 2 of Phase 2 given that at Stage 1 of Phase 2, player 2 has announced  $\hat{\theta}_2$ .

Using the above notation, our hypothesis to derive a contradiction is summarized as follows: for all  $k \neq j$ , all  $l$  and all announcements  $\hat{\theta}_1$ ,  $\sigma_{k,l}^{\varepsilon,j}$  and  $\lambda_{k,l}^{\varepsilon,j}[\hat{\theta}_1]$  converge to 0 as  $\varepsilon \rightarrow 0$ ; and for all  $l \neq j$ , all  $k$  and all announcement  $\hat{\theta}_2$ ,  $\rho_{k,l}^{\varepsilon,j}$  and  $\tau_{k,l}^{\varepsilon,j}[\hat{\theta}_2]$  converge to 0 as  $\varepsilon \rightarrow 0$ .

We will use the following claim about the properties of the MR mechanism under complete information:

**Claim 1.** *For truth-telling to be the unique subgame-perfect equilibrium of the MR mechanism under complete information, it must be that for each  $\theta = (\theta_1, \theta_2)$  and each  $\phi_1$ ,*

$$u_1(f(\theta_1, \theta_2); \theta_1) > U_1(y(\phi_1, \theta_1); \theta_1) - t_{y(\phi_1, \theta_1)} - \Delta, \quad (2)$$

and

$$u_2(f(\theta_1, \theta_2); \theta_2) > U_2(x(\theta_1, \phi_1); \theta_2) + t_{x(\theta_1, \phi_1)} - \Delta. \quad (3)$$

*Proof of Claim 1.* Suppose first that the inequality (2) goes the other way, that is for some  $\theta = (\theta_1, \theta_2)$  and some  $\phi_1$ , we have

$$u_1(f(\theta_1, \theta_2); \theta_1) < U_1(y(\phi_1, \theta_1); \theta_1) - t_{y(\phi_1, \theta_1)} - \Delta.$$



Then, under complete information where the true state is  $\theta$ , we claim that truthtelling is not a subgame-perfect equilibrium: player 1 has an incentive to deviate by claiming some  $\phi_1 \neq \theta_1$  (and player 2 challenges player 1's report at Stage 2 under truthtelling) in order to reach Stage 3 where he would pick  $\{y(\phi_1, \theta_1), t_{y(\phi_1, \theta_1)} + \Delta\}$ . This contradicts the hypothesis that truthtelling is a subgame perfect equilibrium of the MR mechanism under complete information.

Now, suppose instead that for some  $\theta = (\theta_1, \theta_2)$ , and some  $\phi_1 \neq \theta_1$ , we have

$$u_1(f(\theta_1, \theta_2); \theta_1) = U_1(y(\phi_1, \theta_1); \theta_1) - t_{y(\phi_1, \theta_1)} - \Delta.$$

In this case, we claim that there is a subgame-perfect equilibrium at  $\theta = (\theta_1, \theta_2)$  where player 1 does not report truthfully. To see this, we propose the following strategy profile  $\sigma^*$ : At Stage 1 of Phase 1, player 1 reports  $\phi_1 \neq \theta_1$ ; player 2 reports the true state  $\theta_1$  at Stage 2 irrespective of player 1's announcement; and at Stage 3, player 1 always plays his optimal action. Note here that player 1's optimal play at Stage 3 depends on what he reported at Stage 1. In Phase 2, both players always report truthfully and player 2 plays his optimal action at stage 3. Here again, player 2's optimal action at Stage 3 depends on what he reported at Stage 1. Given the continuation strategy profile from Stage 2 induced by  $\sigma^*$ , player 1 is indifferent between reporting  $\phi_1$  and  $\theta_1$  at Stage 1, and so (if truthtelling is a subgame perfect equilibrium) this  $\sigma^*$  is indeed a subgame-perfect equilibrium at  $\theta = (\theta_1, \theta_2)$ . This contradicts the uniqueness of truthtelling as a subgame perfect equilibrium of the MR mechanism under complete information.

Similarly, we must have that for each  $\theta = (\theta_1, \theta_2)$  and each  $\phi_1$ ,

$$u_2(f(\theta_1, \theta_2); \theta_2) > U_2(x(\theta_1, \phi_1); \theta_2) + t_{x(\theta_1, \phi_1)} - \Delta.$$

By way of contradiction, we argue why this must be the case. Suppose first that for some  $\theta = (\theta_1, \theta_2)$  and some  $\phi_1$ , we have

$$u_2(f(\theta_1, \theta_2); \theta_2) < U_2(x(\theta_1, \phi_1); \theta_2) + t_{x(\theta_1, \phi_1)} - \Delta.$$

Then, under complete information where the true state is  $\theta = (\theta_1, \theta_2)$ , we claim that truthtelling is not an equilibrium: player 2 has an incentive to deviate by claiming some  $\phi_1 \neq \theta_1$  in order to reach stage 3 where player 1 would pick  $\{x(\theta_1, \phi_1), t_{x(\theta_1, \phi_1)} + \Delta\}$ . This contradicts the hypothesis that truthtelling is a subgame perfect equilibrium of the MR mechanism under complete information.

Now, suppose instead that for some  $\theta = (\theta_1, \theta_2)$ , and some  $\phi_1 \neq \theta_1$ , we have

$$u_2(f(\theta_1, \theta_2); \theta_2) = U_2(x(\theta_1, \phi_1); \theta_2) + t_{x(\theta_1, \phi_1)} - \Delta.$$

In this case, we claim that there is a subgame-perfect equilibrium at  $\theta = (\theta_1, \theta_2)$  where player 2 does not report truthfully. To see this, we construct the following strategy profile  $\sigma^{**}$ : At Stage 1 of Phase 1, player 1 always reports  $\theta_1$  truthfully; player 2 reports a false state  $\phi_1$  if player 1 has claimed  $\theta_1$  and otherwise challenges with  $\theta_1$ ; and at Stage 3, player 1 always plays his optimal

action. Note here that player 1's optimal play at Stage 3 depends on what he reported at Stage 1. In Phase 2, both players always report truthfully and player 2 plays his optimal action at stage 3. Here again, player 2's optimal action at Stage 3 depends on what he reported at Stage 1. Since player 1 would choose  $\{x(\theta_1, \phi_1), t_{x(\theta_1, \phi_1)} + \Delta\}$  at Stage 3, player 2 is indifferent between reporting  $\theta_1$  and  $\phi_1$  at Stage 2 after player 1 reported  $\theta_1$ . This shows that (if truthtelling is a subgame perfect equilibrium)  $\sigma^{**}$  is a subgame perfect equilibrium where player 2 does not report truthfully. However, this contradicts the uniqueness of truthtelling as a subgame perfect equilibrium of the MR mechanism under complete information. This completes the proof of the claim. ■

Now, let us fix the prior  $\nu^\varepsilon$  (as defined in (\*\*\*) of section 3.2.). Consider the case where player 1 receives  $s_1^{k,l}$ . Clearly,  $\nu^\varepsilon(\theta_1^k, \theta_2^l, s_2^{k,l} | s_1^{k,l}) \rightarrow 1$  as  $\varepsilon \rightarrow 0$ . Hence, at stage 1, by continuity of expected payoffs with respect to beliefs, the expected equilibrium payoff of player 1 for announcing  $\theta_1^k$  converges (as  $\varepsilon$  vanishes) to

$$u_1(f(\theta_1^k, \theta_2^l); \theta_1^k),$$

while if he lies by claiming  $\phi_1 \neq \theta_1^k$  at Stage 1, his expected equilibrium payoff converges to something (weakly) smaller than

$$U_1(y(\phi_1, \theta_1^k); \theta_1^k) - t_{y(\phi_1, \theta_1^k)} - \Delta.$$

By Equation (2) and choosing  $\varepsilon > 0$  small enough, there is no way that the equilibrium strategies  $\{\sigma_{k,l}^{\varepsilon,j}, \lambda_{k,l}^{\varepsilon,j}[\hat{\theta}_1], \rho_{k,l}^{\varepsilon,j}, \tau_{k,l}^{\varepsilon,j}[\hat{\theta}_2]\}_{k,l,j,\hat{\theta}_1,\hat{\theta}_2}$  can make player 1's best response indifferent at Stage 1. Hence, for  $\varepsilon > 0$  small enough, player 1 with signal  $s_1^{k,l}$  plays pure strategies at Stage 1 of Phase 1. This reasoning holds for an arbitrary choice of  $s_1^{k,l}$  so that player 1 plays in pure strategies irrespective of his signal.

Note now that player 1 with signal  $s_1^{k,l}$  could deviate and claim that  $\theta_1^{k'}$  is the true state where  $k' \neq k$ . In this case, because player 1 plays in pure strategies (and hence, the equilibrium is fully revealing), in the first phase, after observing  $\theta_1^{k'}$ , player 2 believes with probability one that player 1 has received a signal of the form  $s_1^{k',l'}$  for some  $l'$ . We claim that player 2 with signal  $s_2^{k,l}$  will not challenge: indeed, by construction of  $\nu^\varepsilon$ , player 2 with signal  $s_2^{k,l}$  believes with high probability that  $\theta = (\theta_1, \theta_2)$  where  $\theta_2 = \theta_2^l$  is the true state. If player 2 challenges with  $\theta_1^k$ , by construction of  $\nu^\varepsilon$ , he expects player 1 to choose  $\{x(\theta_1^{k'}, \theta_1^k), t_{x(\theta_1^{k'}, \theta_1^k)} + \Delta\}$  at Stage 3. On the other hand, if he does not challenge, his expected payoff would tend to  $u_2(f(\theta_1^{k'}, \theta_2^l); \theta_2^l)$  as  $\varepsilon$  vanishes. Hence, by Equation (3), for  $\varepsilon > 0$  small, player 2 will be better off by not challenging. Thus, we get that  $\lambda_{k,l}^{\varepsilon,k}[\theta_1^{k'}] = 0$ , which is a contradiction. This completes the proof of Theorem 2.

## 7.2 Proof of Theorem 3

We first introduce some notation. Given a prior  $\mu$  over  $\Theta \times S$ , we write  $\mu(\theta)$  for  $[\text{marg}_\Theta \mu](\theta)$ , and given  $s_{-i} \in S_{-i}$ , we will write  $\mu(s_{-i})$  as  $[\text{marg}_{S_{-i}} \mu](s_{-i})$ . Finally, given an arbitrary countable space  $X$ ,  $\delta_x$  will denote the probability measure that puts probability 1 on  $\{x\} \subset X$ .

Let  $\mu$  be any complete information prior, and assume that a mechanism  $\Gamma$  SPE-implements a

non-Maskin monotonic SCC  $\mathcal{F}$ . By hypothesis  $\mathcal{F}$  is not Maskin monotonic, so there are  $\theta', \theta''$  and  $a \in \mathcal{F}(\theta')$  satisfying (\*) in the definition of Maskin monotonicity while  $a \notin \mathcal{F}(\theta'')$ . We now fix this particular  $\theta', \theta''$  and  $a$  throughout.

Since the mechanism  $\Gamma$  SPE-implements  $\mathcal{F}$ , there exists a pure strategy subgame-perfect equilibrium  $m_{\theta'}^*$  in  $\Gamma(\theta')$  such that  $g(m_{\theta'}^*) = a$ . Fix one such equilibrium. Clearly,  $m_{\theta'}^*$  is a Nash equilibrium of  $\Gamma(\theta')$ . From (\*) in the definition of Maskin monotonicity, it follows that  $m_{\theta'}^*$  is also a Nash equilibrium of  $\Gamma(\theta'')$ . Recall that  $\mathcal{H}$  denotes the set of all possible histories. For each  $t \geq 0$ , let  $h_t^*$  be the history induced by  $m_{\theta'}^*$  up to date  $t$  and let  $\mathcal{H}^*$  denote the set of all such histories for any  $t$ . In addition, for each player  $i$ , let  $\mathcal{H}_{-i}^*$  be the set of histories  $h$  along which every player  $j \neq i$  has chosen the message  $m_{\theta',j}^*(h)$ ; formally,  $\mathcal{H}_{-i}^* \equiv \{h \in \mathcal{H} : h = (\emptyset, m^1, m^2, \dots, m^{t-1}) \text{ for some } t \text{ and } m_j^{t'} = m_{j,\theta'}^{*,t'} \text{ for all } t' \leq t-1 \text{ and all } j \neq i\}$ . Note that  $h_t^* \in \mathcal{H}_{-i}^*$  for each  $t \geq 1$ .

Consider the following family of information structures  $\nu^\varepsilon$ . For each player  $i$ , let  $\tau_i$  represent the profile of signals  $s = (s_1, \dots, s_n)$  defined by  $s_i = s_i^{\theta''}$  and  $s_j = s_j^{\theta'}$  for all  $j \neq i$ . For all  $i$ ,  $\nu^\varepsilon$  is given by<sup>33</sup>

$$\begin{aligned} \nu^\varepsilon(\theta', \tau_i) &= \frac{\varepsilon}{n} \mu(\theta', s^{\theta'}); \\ \nu^\varepsilon(\theta', s^{\theta'}) &= (1 - \varepsilon) \mu(\theta', s^{\theta'}); \text{ and} \\ \nu^\varepsilon(\tilde{\theta}, s^{\tilde{\theta}}) &= \mu(\tilde{\theta}, s^{\tilde{\theta}}) \quad \forall \tilde{\theta} \neq \theta'. \end{aligned}$$

In this information structure when the state is anything other than  $\theta'$  or  $\theta''$ , the state is common knowledge. Furthermore, when a player observes  $s^{\theta'}$ , he knows that the state is  $\theta'$ . Obviously,  $\nu^\varepsilon \rightarrow \mu$  as  $\varepsilon \rightarrow 0$ . The support of  $\nu^\varepsilon$  is denoted

$$\text{supp}(\nu^\varepsilon) = \{(\tilde{\theta}, s^{\tilde{\theta}}) : \tilde{\theta} \in \Theta\} \cup \{(\theta', \tau_i) : i \in N\}.$$

Before we prove Theorem 3, we introduce some notation and the formal definition of sequential equilibrium. A system of beliefs of agent  $i$  is defined as a function  $\phi_i : S_i \times \mathcal{H} \rightarrow \Delta(\Theta \times S_{-i})$ . Let  $\phi_i[(\theta, s_{-i}) | s_i, h_t]$  denote agent  $i$ 's belief that  $(\theta, s_{-i})$  is realized when agent  $i$ 's signal is  $s_i$  and the observed history is  $h_t$ . We will henceforth abuse notation and sometimes consider  $\phi_i[(\theta, s_{-i}) | s_i, h_t]$  as an element of  $\Delta(\Theta \times S)$ . We also say a vector of beliefs  $\phi = (\phi_1, \dots, \phi_n)$  is *Bayes consistent* with a strategy profile  $\sigma$  if beliefs are updated from one stage to the next using Bayes' rule whenever possible (see Fudenberg and Tirole (1991a) for its precise definition). An assessment is a pair  $(\phi, \sigma)$  consisting of a profile of beliefs and a pure behavior strategy profile. We formally define sequential equilibrium.

**Definition 2.** *A sequential equilibrium is an assessment  $(\phi, \sigma)$  that satisfies condition (S) and (C):*

<sup>33</sup>This sequence of perturbations is similar to that used by Chung and Ely (2003). However, because sequential equilibrium requires verifying sequential rationality conditions that are not imposed by undominated Nash equilibrium, the body of proof is very different from that in Chung and Ely (2003).

**(S) Sequential rationality:** for all  $i \in N$ ,  $s_i \in S_i$ ,  $h_t \in \mathcal{H}$  :

$$\sum_{(\theta, s_{-i}) \in \Theta \times S_{-i}} \phi_i[\theta, s_{-i} | s_i, h_t] \{u_i(g(\sigma(s); h_t); \theta) - u_i(g((\sigma'_i(s_i), \sigma_{-i}(s_{-i})); h_t); \theta)\} \geq 0$$

for each  $\sigma'_i$ .

**(C) Consistency:** there exists a sequence of totally mixed strategy profiles  $(\sigma_1^k, \dots, \sigma_n^k)$  converging to  $(\sigma_1, \dots, \sigma_n)$  with Bayes consistent beliefs  $\phi^k$  converging to  $\phi$ .<sup>34</sup>

Now we come back to the proof and in particular, build a sequential equilibrium  $(\phi^\varepsilon, \sigma^\varepsilon)$  of  $\Gamma(\nu^\varepsilon)$  where  $g(\sigma^\varepsilon(s^{\theta''}); \emptyset) = a$  for each  $\varepsilon > 0$  small enough. This will show that there exist a sequence of priors  $\{\nu^\varepsilon\}_{\varepsilon > 0}$  that converges to  $\mu$  and a corresponding sequence of sequential equilibria  $\{(\phi^\varepsilon, \sigma^\varepsilon)\}_{\varepsilon > 0}$  such that  $g(\sigma^\varepsilon(s^{\theta''}); \emptyset) \rightarrow a \notin \mathcal{F}(\theta'')$  as  $\varepsilon$  goes to 0. This will complete the proof.

In the sequel, we will omit the dependence of  $\sigma^\varepsilon$  with respect to  $\varepsilon$  and simply write  $\sigma$  for  $\sigma^\varepsilon$ . In the following lines, we define a strategy  $\sigma$  and a family of systems of beliefs  $\Phi$  so that  $g(\sigma(s^{\theta''}); \emptyset) = a$ . In addition, we will show that  $(\phi, \sigma)$  is a sequential equilibrium of  $\Gamma(\nu^\varepsilon)$  for some  $\phi \in \Phi$ . We define  $\Phi$  and  $\sigma$  as follows:

**Definition of  $\sigma$ :**

**Σ1.** For any player  $i$  and any  $h_t \in \mathcal{H}^*$  or  $h_t \notin \mathcal{H}_{-i}^*$ ,  $\sigma_i(h_t, s_i^{\theta''}) = m_{i, \theta'}^*(h_t)$ ;<sup>35</sup>

**Σ2.** For any player  $i$ , any  $h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ ,  $\sigma_i(h_t, s_i^{\theta''}) = \bar{m}_i(h_t)$  where  $\bar{m}_i$  satisfies for any  $h_t$ ,

$$\begin{aligned} h_t \in \mathcal{H}^* \text{ or } h_t \notin \mathcal{H}_{-i}^* &\Rightarrow \bar{m}_i(h_t) = m_{i, \theta'}^*(h_t); \\ h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^* &\Rightarrow \bar{m}_i(h_t) \in \arg \max_{\tilde{\theta}} \sum \nu^\varepsilon(\tilde{\theta} | s_i^{\theta''}) u_i(g((m'_i, m_{-i, \theta'}^*); h_t); \tilde{\theta}), \end{aligned}$$

where the max is taken over all pure messages  $m'_i \in M_i$  that differs from  $\bar{m}_i$  only at  $h$ .<sup>36</sup> By A1 there exists such  $\bar{m}_i$ ;

**Σ3.** For any player  $i$  and any  $h_t \in \mathcal{H}$ ,  $\sigma_i(h_t, s_i^{\theta'}) = m_{i, \theta'}^*(h_t)$ ;

**Σ4.** And for any  $h_t \in \mathcal{H}$ ,  $\sigma_i(h_t, s_i^{\tilde{\theta}}) = m_{\tilde{\theta}, i}^*(h_t)$  for  $\tilde{\theta} \neq \theta', \theta''$  where  $m_{\tilde{\theta}, i}^*$  is an arbitrary pure strategy subgame-perfect equilibrium of  $\Gamma(\tilde{\theta})$ . (This is well-defined since  $\mathcal{F}$  is implementable in subgame-perfect equilibrium under complete information.)

<sup>34</sup>Convergence in the definition of consistency is taken uniformly over messages and histories. Given that the set of messages (and so the set of histories) can be countably infinite, two natural convergence notions can be used: *point-wise* convergence or *uniform* convergence. The set of sequential equilibria is smaller when one assumes uniform convergence. Hence, the use of uniform convergence strengthens our main result.

<sup>35</sup>Note that players here send the messages that  $m$  prescribes for state  $\theta'$  when their signal suggests that the state is  $\theta''$ .

<sup>36</sup>Note that the maximization above is over all pure messages  $m'_i \in M_i$  that differs from  $\bar{m}_i$  only at  $h$ . Hence, since player  $i$  may be playing at several stages, it might be the case that this maximization depends on what player  $i$  is playing at further histories, and these further histories may be outside  $\mathcal{H}_{-i}^* \setminus \mathcal{H}^*$  (for instance in case a player  $j$  different of  $i$  does not play according to  $m_{j, \theta'}^*$  at some subsequent history). This is why we also have to define  $\bar{m}_i$  outside  $\mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ .

**Definition of  $\Phi$ :**

$\phi \in \Phi$  if and only  $\phi$  satisfies the following three properties.

**$\Phi 1$ .** Fix any  $i \in N$ , any  $h_t \notin \mathcal{H}_{-i}^*$ ,

$$\phi_i \left[ \cdot | s_i^{\theta''}, h_t \right] = \delta_{(\theta', s_{-i}^{\theta'})}$$

and

$$\text{supp} \left( \phi_i \left[ \cdot | s_i^{\theta'}, h_t \right] \right) \subseteq \text{supp} \left( \nu^\varepsilon \left[ \cdot | s_i^{\theta'} \right] \right)$$

and for all  $l \neq i$  with  $h_t \in \mathcal{H}_{-l}^* \setminus \mathcal{H}_{-i}^*$

(i.e., player  $l$  has deviated from the path prescribed by  $m_{\theta'}^*$ )

$$\phi_i[(\theta', \tau_l) | s_i^{\theta'}, h_t] = 0.$$

**$\Phi 2$ .** For any  $i \in N$ , any  $h_t \in \mathcal{H}_{-i}^*$ , any  $s_i \in \{s_i^{\theta'}, s_i^{\theta''}\}$ ,

$$\phi_i[\cdot | s_i, h_t] = \nu^\varepsilon(\cdot | s_i).$$

**$\Phi 3$ .** For any  $i \in N$ , any  $h_t \in \mathcal{H}$  and any  $\tilde{s}_i \notin \{s_i^{\theta'}, s_i^{\theta''}\}$ ,  $\phi_i \left[ \cdot | s_i^{\tilde{\theta}}, h_t \right] = \delta_{(\tilde{\theta}, s_{-i}^{\tilde{\theta}})}$  where  $\delta_x$  denotes the probability measure that puts probability 1 on  $\{x\}$ .

Note that  $h_T[\sigma(s^{\theta''}), \emptyset] = h_T[m_{\theta'}^*, \emptyset]$  and so,  $\sigma$  generates  $g(\sigma(s^{\theta''}); \emptyset) = g(m_{\theta'}^*; \emptyset) = a$ . Hence, it only remains to show that  $(\phi, \sigma)$  constitutes a sequential equilibrium for some  $\phi \in \Phi$ . In Section 7.2.1, we show that  $(\phi, \sigma)$  satisfies sequential rationality for any  $\phi \in \Phi$ ; and we establish that  $(\phi, \sigma)$  satisfies consistency for some  $\phi \in \Phi$  in Section 7.2.2.

### 7.2.1 Sequential rationality

Fix any  $\phi \in \Phi$ . Sequential rationality of  $(\phi, \sigma)$  will be proved by Claims 2 and 3 below.

**Claim 2.** For any  $i \in N$ ,  $s_i \neq s_i^{\theta''}$ ,  $h_t \in \mathcal{H}$  :

$$\sum_{(\tilde{\theta}, s_{-i})} \phi_i[(\tilde{\theta}, s_{-i}) | s_i, h_t] \left[ u_i(g(\sigma(s); h_t); \tilde{\theta}) - u_i(g(\sigma'_i(s_i), \sigma_{-i}(s_{-i})); h_t); \tilde{\theta}) \right] \geq 0$$

for each  $\sigma'_i$ .

This claim 2 states that for any player  $i$  with any signal  $s_i \neq s_i^{\theta''}$ ,  $\sigma_i$  is a best response to  $\sigma_{-i}$  given his belief  $\phi_i$ . This will be checked by considering three classes of histories: (1) Histories where all players have played according to the equilibrium  $m_{\theta'}^*$  (i.e., in  $\mathcal{H}^*$ ); (2) histories where player  $i$  has not played according to  $m_{\theta'}^*$  but all other players have (i.e., in  $\mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ ); and finally (3) histories where some player other than  $i$  has not played according to  $m_{\theta'}^*$  (i.e., outside  $\mathcal{H}_{-i}^*$ ).

In particular, in the non-trivial case where  $s_i = s_i^{\theta'}$ , we will show that for any of these histories  $h_t$ , whenever player  $i$  follows  $\sigma_i$  against  $\sigma_{-i}$ , player  $i$  believes with probability one that the outcome

will be given by  $g(m_{\theta'}^*; h_t)$ , while if player  $i$  deviates from  $\sigma_i(s_i)$  to some  $m'_i$ , player  $i$  believes with probability one that the outcome will be given by  $g(m'_i, m_{-i, \theta'}^*; h_t)$ . Because  $m_{\theta'}^*$  is a subgame-perfect equilibrium in the complete information game  $\Gamma(\theta')$  and player  $i$  with signal  $s_i^{\theta'}$  believes with probability one that  $\theta'$  is the true state, this will prove the claim.

*Proof of Claim 2.* Fix any player  $i$ . Claim 1 is obvious for  $s_i^{\tilde{\theta}} \neq s_i^{\theta'}$  because by  **$\Phi 3$** ,  $\phi_i[\cdot | s_i^{\tilde{\theta}}, h_t] = \delta_{(\tilde{\theta}, s_{-i}^{\tilde{\theta}})}$  and so state  $\tilde{\theta}$  is common knowledge. By  **$\Sigma 4$** , we can further conclude that  $\sigma(s^{\tilde{\theta}}) = m_{\tilde{\theta}}^*$  is a subgame-perfect equilibrium in the complete information game  $\Gamma(\tilde{\theta})$ . Hence, we focus on the case where  $s_i = s_i^{\theta'}$ . By construction,  $\nu^\varepsilon(\theta' | s_i^{\theta'}) = 1$  and so this player knows the state is  $\theta'$ , and he knows the profile of signals is either  $s^{\theta'}$  or  $\tau_k$  for some  $k \neq i$ . We partition the set of all histories into three classes  $\mathcal{H}^*$ ;  $\mathcal{H}_{-i}^* \setminus \mathcal{H}^*$  and  $\mathcal{H} \setminus \mathcal{H}_{-i}^*$  and consider the following three cases: Case (1)  $h_t \in \mathcal{H}^*$ ; Case (2)  $h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ ; and Case (3)  $h_t \notin \mathcal{H}_{-i}^*$ .

- Case (1):  $h_t \in \mathcal{H}^*$

In this case, each player has played according to  $m_{\theta'}^*$  and if players  $j \neq i$  received signals of either  $s_j^{\theta'}$  or  $s_j^{\theta''}$ , by  **$\Sigma 1$**  and  **$\Sigma 3$** , this will continue to be the case as long as all players conform to  $\sigma$ . So when players are playing strategy  $\sigma$ , and the profile of signals received is  $s^{\theta'}$  or  $\tau_k$ , for  $k \neq i$  any subsequent history also falls into  $\mathcal{H}^*$ . Thus,  $g(\sigma(s^{\theta'}); h_t) = g(\sigma(\tau_k); h_t) = g(m_{\theta'}^*; h_t)$ .

Now suppose player  $i$  deviates to a strategy  $\sigma'_i$  so that  $\sigma'_i(s_i^{\theta'}) = m'_i$ . Clearly, since  $m'_i \neq \sigma_i(s_i^{\theta'})$ , there is a date at which player  $i$  does not play according to  $m_{i, \theta'}^*$ . Thus, by  **$\Sigma 1$**  and  **$\Sigma 3$** , when the profile of signals received is either  $s^{\theta'}$  or  $\tau_k$  for  $k \neq i$ , any subsequent history of  $h_t$  either falls in  $\mathcal{H}^*$  (player  $i$  has played according to  $m_{i, \theta'}^*$  so far) or does not fall in  $\mathcal{H}_{-k}^*$  for each  $k \neq i$  (at some point in this history, player  $i$  has not played according to  $m_{i, \theta'}^*$ ). In each of these cases, again by  **$\Sigma 1$**  and  **$\Sigma 3$** , player  $i$ 's opponents are playing according to  $m_{-i, \theta'}^*$ . So we get <sup>37</sup>

$$g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}^{\theta'}); h_t) = g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(\tau_k); h_t) = g(m'_i, m_{-i, \theta'}^*; h_t).$$

Here again, since  $m_{\theta'}^*$  is a subgame-perfect equilibrium in the complete information game  $\Gamma(\theta')$ , we have

$$u_i(g(m_{\theta'}^*; h_t); \theta') \geq u_i(g(m'_i, m_{-i, \theta'}^*; h_t); \theta').$$

Thus, we get  $u_i(g(\sigma(s^{\theta'}); h_t); \theta') \geq u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}^{\theta'}); h_t); \theta')$  and  $u_i(g(\sigma(\tau_k); h_t); \theta') \geq u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(\tau_k); h_t); \theta')$  for each  $k \neq i$ . Now since by  **$\Phi 2$** ,  $\phi_i[\cdot | s_i^{\theta'}, h_t]$  may assign strictly positive weight only to  $(\theta', s_{-i}^{\theta'})$  and  $(\theta', \tau_k)$  for each  $k \neq i$ , we can conclude

$$\sum_{(\tilde{\theta}, s_{-i})} \phi_i[(\tilde{\theta}, s_{-i}) | s_i^{\theta'}, h_t] \left[ u_i(g(\sigma_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}); h_t); \tilde{\theta}) - u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}); h_t); \tilde{\theta}) \right] \geq 0.$$

- Case (2):  $h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$

---

<sup>37</sup>We abuse notation because we should use  $\sigma_{-i}(\tau_l \setminus s_i^{\theta'})$  instead of  $\sigma_{-i}(\tau_l)$ .

Since  $h_t \in \mathcal{H}_{-i}^*$  and  $h_t \notin \mathcal{H}^*$ , only player  $i$  has not played according to  $m_{i,\theta'}^*$ . Then, it is clear that  $h_t$  does not fall in  $\mathcal{H}_{-k}^*$  for each  $k \neq i$  (recall that  $\mathcal{H}_{-k}^*$  is the set of histories under which every player  $j$  other than  $k$  has played according to  $m_{j,\theta'}^*$ ). It is also clear that any subsequent history does not fall in  $\mathcal{H}_{-k}^*$  for each  $k \neq i$ . By  **$\Sigma 1$**  and  **$\Sigma 3$** , we thus obtain that each player  $k$  other than  $i$  will play according to  $m_{k,\theta'}^*$  at any subsequent history when receiving signal  $s_k^{\theta'}$  or  $s_k^{\theta''}$ . Hence,

$$g(\sigma(s^{\theta'}); h_t) = g(\sigma(\tau_k); h_t) = g(m_{\theta'}^*; h_t).$$

Consider the case where player  $i$  deviates to a strategy  $\sigma'_i$  so that  $\sigma'_i(s_i^{\theta'}) = m'_i$ . Here, since (by a similar argument as above) any history that player  $i$  can achieve by deviating does not fall in  $\mathcal{H}_{-k}^*$  for each  $k \neq i$ , each player  $k$  other than  $i$  will be playing according to  $m_{k,\theta'}^*$  at any subsequent history whether he receive  $s_k^{\theta'}$  or  $s_k^{\theta''}$ , which implies

$$g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}^{\theta'}); h_t) = g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(\tau_k); h_t) = g(m'_i, m_{-i,\theta'}^*; h_t).$$

Since  $m_{\theta'}^*$  is a subgame-perfect equilibrium in the complete information game  $\Gamma(\theta')$ , we already have  $u_i(g(m_{\theta'}^*; h_t); \theta') \geq u_i(g(m'_i, m_{-i,\theta'}^*; h_t); \theta')$ . Thus, we also get

$$\begin{aligned} u_i(g(\sigma(s^{\theta'}); h_t); \theta') &\geq u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}^{\theta'}); h_t); \theta') \quad \text{and} \\ u_i(g(\sigma(\tau_k); h_t); \theta') &\geq u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(\tau_k); h_t); \theta') \quad \text{for each } k \neq i. \end{aligned}$$

Now, since by  **$\Phi 2$**  we know that  $\phi_i[\cdot | s_i^{\theta'}, h_t]$  assigns a strictly positive weight only to  $(\theta', s_{-i}^{\theta'})$  and  $(\theta', \tau_k)$  for each  $k \neq i$ , we can conclude

$$\sum_{(\tilde{\theta}, s_{-i})} \phi_i[(\tilde{\theta}, s_{-i}) | s_i^{\theta'}, h_t] \left[ u_i(g(\sigma(s^{\theta'}), \sigma_{-i}(s_{-i}), h_t); \tilde{\theta}) - u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}), h_t); \tilde{\theta}) \right] \geq 0.$$

- Case (3):  $h_t \notin \mathcal{H}_{-i}^*$

In this case, at least one player  $j \neq i$  has not played according to  $m_{j,\theta'}^*$ .

By  **$\Sigma 3$** , we know that when each player  $j$  receives signal  $s_j^{\theta'}$ , then these players play according to  $m_{j,\theta'}^*$ , so  $\sigma(s^{\theta'}) = m_{\theta'}^*$ . Thus, at history  $h_t$ , the outcome achieved by playing  $\sigma$  when the profile of signals is  $s^{\theta'}$  must be the same as the one when playing  $m_{\theta'}^*$ , i.e.,

$$g(\sigma(s^{\theta'}); h_t) = g(m_{\theta'}^*; h_t).$$

In addition, for each  $l \neq i$  with  $h_t \notin \mathcal{H}_{-l}^*$ , by definition, some player  $j$  other than  $l$  has not played according to  $m_{j,\theta'}^*$  and obviously this will continue to be the case at any subsequent histories. Hence, any subsequent histories does not belong to  $\mathcal{H}_{-l}^*$  either. At any such histories, we know by  **$\Sigma 1$** , that player  $l$  will be playing according to  $m_{l,\theta'}^*$  when he receives  $s_l^{\theta''}$  while when players  $j$  other than  $l$  receive signal  $s_j^{\theta'}$ , by  **$\Sigma 3$**  they will also be playing according

to  $m_{j,\theta'}^*$ . Hence, we get that the outcome achieved from history  $h_t$  when playing  $\sigma$  and when the profile of signals received is  $\tau_l$  is equal to the outcome achieved from history  $h_t$  when playing  $m_{\theta'}^*$ . Otherwise stated, for each  $l \neq i$  with  $h_t \notin \mathcal{H}_{-l}^*$ , we have

$$g(\sigma(\tau_l); h_t) = g(m_{\theta'}^*; h_t).$$

Now, when player  $i$  deviates say to a strategy  $\sigma'_i$  so that  $\sigma'_i(s_i^{\theta'}) = m'_i$ , using the argument above, when the other players receive signal profile  $s_{-i}^{\theta'}$ , we know that the outcome achieved is

$$g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}^{\theta'}); h_t) = g(m'_i, m_{-i,\theta'}^*; h_t).$$

while for each  $l \neq i$  with  $h_t \notin \mathcal{H}_{-l}^*$ , we know that

$$g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(\tau_l); h_t) = g(m'_i, m_{-i,\theta'}^*; h_t).$$

Since  $m_{\theta'}^*$  is a subgame-perfect equilibrium in the complete information game  $\Gamma(\theta')$ , we have  $u_i(g(m_{\theta'}^*; h_t); \theta') \geq u_i(g(m'_i, m_{-i,\theta'}^*; h_t); \theta')$ . Thus, we get

$$u_i(g(\sigma(s^{\theta'}); h_t); \theta') \geq u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}^{\theta'}); h_t); \theta')$$

and for each  $l \neq i$  such that  $h_t \notin \mathcal{H}_{-l}^*$ ,  $u_i(g(\sigma(\tau_l); h_t); \theta') \geq u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(\tau_l); h_t); \theta')$ . Because by  $\Phi\mathbf{1}$ ,  $\phi_i[\cdot | s_i^{\theta'}, h_t]$  may assign strictly positive weight only to  $(\theta', s_{-i}^{\theta'})$  and  $(\theta', \tau_l)$  for each  $l \neq i$  such that  $h_t \notin \mathcal{H}_{-l}^*$ , we can conclude

$$\sum_{(\tilde{\theta}, s_{-i})} \phi_i[(\tilde{\theta}, s_{-i}) | s_i^{\theta'}, h_t] \left[ u_i(g(\sigma(s^{\theta'}), s_{-i}); h_t; \tilde{\theta}) - u_i(g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i})); h_t); \tilde{\theta}) \right] \geq 0.$$

This completes the proof of that claim.

■

**Claim 3.** For any  $i \in N$ ,  $s_i = s_i^{\theta''}$ , and  $h_t \in \mathcal{H}$ :

$$\sum_{(\tilde{\theta}, s_{-i})} \phi_i[(\tilde{\theta}, s_{-i}) | s_i, h_t] \left[ u_i(g(\sigma(s); h_t); \tilde{\theta}) - u_i(g(\sigma'_i(s_i), \sigma_{-i}(s_{-i})); h_t); \tilde{\theta}) \right] \geq 0$$

for each  $\sigma'_i$ .

This claim 3 states that for any player  $i$  with signal  $s_i^{\theta''}$ ,  $\sigma_i$  is a best response to  $\sigma_{-i}$  given his belief  $\phi_i$ . Here again we consider the same partition of histories as in Claim 2. When  $h_t$  is a history where each player has played according to  $m_{\theta'}^*$  (i.e.,  $h_t \in \mathcal{H}^*$ ), player  $i$  assigns positive probability to both  $\theta''$  and  $\theta'$ . However, we will show that here again player  $i$  believes with probability one that the other players will be playing according to  $m_{-i,\theta'}^*$ , whether he deviates or not. Hence, if he does not deviate and  $h_t \in \mathcal{H}^*$ , he gets  $a$  while if he deviates to  $m'_i$  he gets  $g(m'_i, m_{-i,\theta'}^*; h_t)$ . Because



$m_{\theta'}^*$  is a subgame-perfect equilibrium in  $\Gamma(\theta')$ , we know that the deviation is not profitable if  $\theta'$  is the true state, and Maskin monotonicity (condition  $(*)$ ) implies that this is also not profitable if the state is  $\theta''$ . Since these are the only states to which player  $i$  assigns strictly positive probability, this will complete the argument for this class of histories.

The easy case occurs when  $h_t$  is a history where a player other than  $i$  has not played according to  $m_{\theta'}^*$  (i.e.,  $h_t \notin \mathcal{H}_{-i}^*$ ). In such a case, player  $i$  believes with probability one that  $\theta'$  is the true state. In addition we will check that whenever player  $i$  uses  $\sigma_i$  against  $\sigma_{-i}$ , player  $i$  believes with probability one that the outcome will be given by  $g(m_{\theta'}^*; h_t)$ , while if player  $i$  deviates from  $\sigma_i(s_i)$  to  $m'_i$ , player  $i$  believes with probability one that the outcome will be given by  $g(m'_i, m_{-i, \theta'}^*; h_t)$ . Here again, the fact that  $m_{\theta'}^*$  is a subgame-perfect equilibrium in the complete information game will lead to the desired result. Finally, in the last case where player  $i$  has not played according to  $m_{\theta'}^*$  while all other players have (i.e.,  $h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ ), we will also check that player  $i$  assigns probability one to his opponent playing  $m_{-i, \theta'}^*$ . But  $\sigma_i$  has been constructed (see  **$\Sigma 2$** ) so that playing  $\sigma_i$  is better than any one-shot deviation. Then the one-shot deviation principle for sequential equilibrium will complete the proof of Claim 3. Taken together, Claims 2 and 3 establish sequential rationality of  $(\phi, \sigma)$ .

*Proof of Claim 3.* This claim will be proved by studying three different cases depending on the type of history we consider: (1)  $h_t \in \mathcal{H}^*$ ; (2)  $h_t \notin \mathcal{H}_{-i}^*$ ; and (3)  $h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ .

- Case (1):  $h_t \in \mathcal{H}^*$

In this case, each player has played according to  $m_{\theta'}^*$ . Note that, by  **$\Sigma 1$**  and  **$\Sigma 3$** , if each player  $j$  received signals of either  $s_j^{\theta'}$  or  $s_j^{\theta''}$ , this will continue to be the case as long as all players conform to  $\sigma$ . So when players are playing strategy  $\sigma$ , and player  $i$ 's opponents received either signal profile  $s_{-i}^{\theta'}$  or  $s_{-i}^{\theta''}$ , any subsequent history also falls into  $\mathcal{H}^*$ . Thus,

$$g(\sigma(s_i^{\theta''}, s_{-i}^{\theta''}); h_t) = g(\sigma(s_i^{\theta'}, s_{-i}^{\theta'}); h_t) = g(m_{\theta'}^*; h_t).$$

Now suppose that player  $i$  deviates to a strategy  $\sigma'_i$  so that  $\sigma'_i(s_i^{\theta''}) = m'_i$ . Since  $m'_i \neq \sigma_i(s_i^{\theta''})$ , there must exist a date at which player  $i$  does not play according to  $m_{i, \theta'}^*$ . Thus, by  **$\Sigma 1$**  and  **$\Sigma 3$** , when player  $i$ 's opponents receive signal  $s_{-i}^{\theta'}$  or  $s_{-i}^{\theta''}$ , any subsequent history of  $h_t$  either falls in  $\mathcal{H}^*$  (player  $i$  has played according to  $m_{i, \theta'}^*$  so far) or does not fall in  $\mathcal{H}_{-k}^*$  for each  $k \neq i$  (at some point in this history, player  $i$  has not played according to  $m_{i, \theta'}^*$ ). In each of these cases, by  **$\Sigma 1$**  and  **$\Sigma 3$** , player  $i$ 's opponents are playing according to  $m_{-i, \theta'}^*$ . So we get

$$g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}^{\theta''}); h_t) = g(\sigma'_i(s_i^{\theta'}), \sigma_{-i}(s_{-i}^{\theta''}); h_t) = g(m'_i, m_{-i, \theta'}^*; h_t). \quad (4)$$

Here again, since  $m_{\theta'}^*$  is a subgame-perfect equilibrium in the complete information game  $\Gamma(\theta')$ , we have

$$u_i(g(m_{\theta'}^*; h_t); \theta') \geq u_i(g(m'_i, m_{-i, \theta'}^*; h_t); \theta').$$

Thus, we also get

$$u_i(g(\sigma(s_i^{\theta''}, s_{-i}^{\theta'}); h_t); \theta') \geq u_i(g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}^{\theta'}); h_t); \theta'). \quad (5)$$

The above inequality, together with (4), also implies

$$u_i(g(\sigma(s_i^{\theta''}, s_{-i}^{\theta''}); h_t); \theta') \geq u_i(g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}^{\theta''}); h_t); \theta').$$

Since  $g(\sigma(s_i^{\theta''}, s_{-i}^{\theta''}); h_t) = g(m_{\theta'}^*; h_t^*) = a$  and we have assumed that  $\theta'$  and  $\theta''$  are two states satisfying (\*) in the definition of Maskin monotonicity, we get that

$$u_i(g(\sigma(s_i^{\theta''}, s_{-i}^{\theta''}); h_t); \theta'') \geq u_i(g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}^{\theta''}); h_t); \theta''). \quad (6)$$

Now, since by **\Phi2**,  $\phi_i[\cdot | s_i^{\theta''}, h_t]$  assigns a strictly positive weight only to  $(\theta', s_{-i}^{\theta'})$  and  $(\theta'', s_{-i}^{\theta''})$ , we conclude (5) and (6) imply that:

$$\sum_{(\tilde{\theta}, s_{-i})} \phi_i[(\tilde{\theta}, s_{-i}) | s_i^{\theta''}, h_t] \left[ u_i(g(\sigma(s_i^{\theta''}, s_{-i}); h_t); \tilde{\theta}) - u_i(g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}); h_t); \tilde{\theta}) \right] \geq 0.$$

- Case (2):  $h_t \notin \mathcal{H}_{-i}^*$

In this case, at least one player  $j \neq i$  has not played according to  $m_{j, \theta'}^*$ ; This is still the case for any subsequent histories, so that they all fall outside  $\mathcal{H}_{-i}^*$ . By **\Sigma1**, if player  $i$  plays according to  $\sigma_i$ , from  $h_t$ , he will play according to  $m_{i, \theta'}^*$ . Now, by **\Sigma3**, we know that when player  $j$  other than  $i$  receives signal  $s_j^{\theta'}$ , then he plays according to  $m_{j, \theta'}^*$ . Thus, the outcome achieved when the profile of signals is  $(s_i^{\theta''}, s_{-i}^{\theta'})$  must be the same as the outcome achieved when  $m_{\theta'}^*$  is played. That is, we obtain

$$g(\sigma(s_i^{\theta''}, s_{-i}^{\theta'}); h_t) = g(m_{\theta'}^*; h_t).$$

Suppose player  $i$  deviates to a strategy  $\sigma'_i$  so that  $\sigma'_i(s_i^{\theta''}) = m'_i$ . Since, if the other players are receiving signal profile  $s_{-i}^{\theta'}$ , they will all be playing according to  $m_{-i, \theta'}^*$ , we obtain

$$g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}^{\theta'}); h_t) = g(m'_i, m_{-i, \theta'}^*; h_t).$$

Since  $m_{\theta'}^*$  is a subgame-perfect equilibrium in the complete information game  $\Gamma(\theta')$ , we have  $u_i(g(m_{\theta'}^*; h_t); \theta') \geq u_i(g(m'_i, m_{-i, \theta'}^*; h_t); \theta')$ . Thus, we also get

$$u_i(g(\sigma(s_i^{\theta''}, s_{-i}^{\theta'}); h_t); \theta') \geq u_i(g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}^{\theta'}); h_t); \theta').$$

Because by **Φ1**,  $\phi_i[(\theta', s_{-i}^{\theta'}) | s_i^{\theta''}, h_t] = 1$ , so we can conclude

$$\sum_{(\tilde{\theta}, s_{-i})} \phi_i[(\tilde{\theta}, s_{-i}) | s_i^{\theta''}, h_t] \left[ u_i(g(\sigma(s_i^{\theta''}, s_{-i}); h_t); \tilde{\theta}) - u_i(g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}); h_t); \tilde{\theta}) \right] \geq 0.$$

- Case (3):  $h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$

Since  $h_t \in \mathcal{H}_{-i}^*$  and  $h_t \notin \mathcal{H}^*$ , only player  $i$  has not played according to  $m_{i, \theta'}^*$ . Then  $h_t$  does not fall in  $\mathcal{H}_{-k}^*$  for each  $k \neq i$  (recall that  $\mathcal{H}_{-k}^*$  is the set of histories under which every player  $j$  other than  $k$  has played according to  $m_{j, \theta'}^*$ ). It is also clear that any subsequent history does not fall in  $\mathcal{H}_{-k}^*$  for each  $k \neq i$ . By **Σ1** and **Σ3**, whether player  $i$ 's opponents have received  $s_{-i}^{\theta'}$  or  $s_{-i}^{\theta''}$ , they all play according to  $m_{-i, \theta'}^*$ . By **Φ2** we know that  $\phi_i[\cdot | s_i^{\theta''}, h_t] = \nu^\varepsilon(\cdot | s_i^{\theta''})$  assigns a strictly positive weight only to  $(\theta', s_{-i}^{\theta'})$  and  $(\theta'', s_{-i}^{\theta''})$ . In addition, we have that for any  $h \in \mathcal{H}^*$  or  $h \notin \mathcal{H}_{-i}^* : \sigma_i(h, s_i^{\theta''}) = m_{i, \theta'}^*(h, s_i^{\theta''})$ . Since  $h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ , we conclude with **Σ2** that:

$$\sum_{(\tilde{\theta}, s_{-i})} \nu^\varepsilon(\tilde{\theta}, s_{-i} | s_i^{\theta''}) \left[ u_i(g(\sigma(s_i^{\theta''}, s_{-i}); h_t); \tilde{\theta}) - u_i(g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}); h_t); \tilde{\theta}) \right] \geq 0$$

for any  $\sigma'_i$  that differs from  $\sigma_i$  only at  $h_t$ . By this and case (1) and (2), we know that at any history players have no profitable one-shot deviation, by the one-shot deviation principle (see Hendon, Jacobsen, and Sloth (1996)<sup>38</sup>) this yields:

$$\sum_{(\tilde{\theta}, s_{-i})} \nu^\varepsilon(\tilde{\theta}, s_{-i} | s_i^{\theta''}) \left[ u_i(g(\sigma(s_i^{\theta''}, s_{-i}); h_t); \tilde{\theta}) - u_i(g(\sigma'_i(s_i^{\theta''}), \sigma_{-i}(s_{-i}); h_t); \tilde{\theta}) \right] \geq 0$$

for any  $\sigma'_i$ . This completes the proof.

■

## 7.2.2 Consistency

In this section, we show that for some  $\phi \in \Phi$ ,  $(\phi, \sigma)$  satisfies consistency.

To show this part, we first fix  $\sigma$  as defined above and consider the following sequence  $\{(\phi^k, \sigma^k)\}_{k=0}^\infty$  of assessments. Let  $\eta_k > 0$  for each  $k$  and  $\eta_k \rightarrow 0$  as  $k \rightarrow \infty$ . For each player  $i$ ,  $h_t \in \mathcal{H}$ , and signal  $s_i$ , let  $\xi_i(h_t, s_i, \cdot)$  be any strictly positive prior over  $M_i(h_t) \setminus \{\sigma_i(s_i, h_t)\}$  and define  $\sigma_i^k$  as

$$\sigma_i^k(m_i^t | h_t, s_i^{\theta''}) = \begin{cases} 1 - \eta_k^{T \times n} & \text{if } m_i^t = \sigma_i(h_t, s_i^{\theta''}); \\ \eta_k^{T \times n} \times \xi_i(h_t, s_i^{\theta''}, m_i^t) & \text{otherwise} \end{cases}$$

<sup>38</sup>Hendon, Jacobsen, and Sloth (1996) assume that for each  $i$  and  $h$ ,  $M_i(h)$  is finite, which is our A1. It is easy to check that their argument goes through in case  $M_i(h)$  is countably infinite. This fact is implicitly used in Section 7.3.

where  $T$  is the (finite) length of the longest final history, and for any signal  $s_i \neq s_i^{\theta''}$  :

$$\sigma_i^k(m_i^t | h_t, s_i) = \begin{cases} 1 - \eta_k & \text{if } m_i^t = \sigma_i(h_t, s_i); \\ \eta_k \times \xi_i(h_t, s_i, m_i^t) & \text{otherwise} \end{cases}.$$

Let  $\phi^k$  be the unique consistent belief associated with each  $\sigma^k$ . It is easy to check that  $\sigma^k$  converges to  $\sigma$  and also that  $\phi^k$  converges.<sup>39</sup> Let  $\phi \equiv \lim_{k \rightarrow \infty} \phi^k$ . In what follows, we show that  $\phi$  satisfies  **$\Phi 1$** ,  **$\Phi 2$**  and  **$\Phi 3$** . This will show that  $(\phi, \sigma)$  satisfies consistency, and  $\phi \in \Phi$  as claimed.

To do so, we explicitly compute each  $\phi^k$  and study its limit as  $k$  tends to infinity. In general for each  $(\tilde{\theta}, \tilde{s}_{-i}) \in \Theta \times S_{-i}$ , each  $h_t = (m^1, \dots, m^{t-1}) \in \mathcal{H}$ , and each  $\tilde{s}_i \in S_i$ , we have

$$\phi_i^k[(\tilde{\theta}, \tilde{s}_{-i}) | \tilde{s}_i, h_t] = \frac{\nu^\varepsilon(\tilde{\theta}, \tilde{s}_{-i}, \tilde{s}_i) \times \prod_{t'=1}^{t-1} [\sigma^k(m^{t'} | h_{t'}, \tilde{s})]}{\sum_{(\hat{\theta}, s'_{-i})} \nu^\varepsilon(\hat{\theta}, s'_{-i}, \tilde{s}_i) \times \prod_{t'=1}^{t-1} [\sigma^k(m^{t'} | h_{t'}, s'_{-i}, \tilde{s}_i)]}.$$

In the above formula for each  $t' \leq t$ ,  $h_{t'}$  stands for the truncation of  $h_t$  to the first  $t'$  elements, i.e.,  $h_{t'} = (m^1, \dots, m^{t'-1})$ .

**Claim 4.**  $\phi$  satisfies  **$\Phi 1$** .

Claim 4 says that, for any player  $i$  who sees signal  $s_i^{\theta''}$  and has an opportunity to play after some other player has not played according to  $m_{\theta'}^*$  (i.e.,  $h_t \notin \mathcal{H}_{-i}^*$ ), then under  $\phi \equiv \lim_{k \rightarrow \infty} \phi^k$ , player  $i$  believes with probability one that the state is  $\theta'$ , and that the other players have received  $s_{-i}^{\theta'}$ . In order to show that, we observe that if every player other than  $i$  has received a signal  $s_j \in \{s_j^{\theta'}, s_j^{\theta''}\}$ , then at such a history some player  $j$  other than  $i$  has deviated from  $\sigma$ . Then, since under the sequence of totally mixed strategies built above, it is (infinitely) more likely (as  $\eta_k$  tends to 0) that a deviation occurred at  $s_j^{\theta'}$  rather than at  $s_j^{\theta''}$ . In the limit, Bayes' rule will then put probability one on  $s_j^{\theta'}$  and given that the prior  $\nu^\varepsilon$  assigns strictly positive weight only to  $(\theta'', s_{-i}^{\theta''})$  and  $(\theta', s_{-i}^{\theta'})$ , Bayes rule will then put probability arbitrarily close to one on  $(\theta', s_{-i}^{\theta'})$ . In case player  $i$  received the private signal  $s_i^{\theta'}$ , if  $h_t$  is a history under which all players other than  $l$  have played according to  $m_{\theta'}^*$  (i.e.  $h_t \in \mathcal{H}_{-l}^*$ ), then the deviating player is  $l$  and again using a similar argument as above, we show that player  $i$  must assign probability 0 to player  $l$  receiving  $s_l^{\theta''}$  and so to  $\tau_l$ .

Consider player  $i$  at history  $h_t \notin \mathcal{H}_{-i}^*$ . The proof is reduced to checking the following two cases:

*Proof of Claim 4. Case 1:*  $s_i = s_i^{\theta''}$

---

<sup>39</sup>As will become clear from the proof, the sequence  $\{\phi^k\}_k$  does converge. Moreover, convergence in the definition of consistency is taken uniformly over messages and histories. In the case where  $M_i(h)$  is countably infinite (we will discuss this case in Section 7.3 in the Appendix), two natural convergence notions can be used: *point-wise* convergence or *uniform* convergence. The set of sequential equilibria is smaller when one assumes uniform convergence. Hence, the use of uniform convergence strengthens our result.

Recall that  $\nu^\varepsilon(\cdot, s_i^{\theta''})$  assigns a strictly positive weight only to  $(\theta'', s_{-i}^{\theta''})$  and  $(\theta', s_{-i}^{\theta'})$ . Hence,

$$\begin{aligned}
& \phi_i^k[(\theta', s_{-i}^{\theta'})|s_i^{\theta''}, h_t] \\
&= \frac{\nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta''}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta'})}{\nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta''}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta'}) + \nu^\varepsilon(\theta'', s_{-i}^{\theta''}, s_i^{\theta''}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta''})} \\
&= \frac{\nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta''})}{\nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta''}) + \nu^\varepsilon(\theta'', s_{-i}^{\theta''}, s_i^{\theta''}) \times \frac{\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta''})}{\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta'})}}.
\end{aligned}$$

We now show that the ratio below converges to 0 as  $k \rightarrow \infty$ :

$$\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta''}) \Big/ \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta'}) \rightarrow 0 \text{ as } k \rightarrow \infty.$$

This will show that  $\phi_i^k[(\theta', s_{-i}^{\theta'})|s_i^{\theta''}, h_t] \rightarrow 1$  and  $\phi_i^k[(\theta'', s_{-i}^{\theta''})|s_i^{\theta''}, h_t] \rightarrow 0$  as  $k \rightarrow \infty$ .

Note first that in case every player  $j$  other than  $i$  receives signal  $s_j \in \{s_j^{\theta'}, s_j^{\theta''}\}$ , there must exist a player  $\hat{j} \neq i$  and a date  $\hat{t} \leq t-1$  so that  $\hat{j}$  has not played according to  $\sigma_{\hat{j}}$ , i.e.  $\sigma_{\hat{j}}(h_{\hat{t}}, s_{\hat{j}}) \neq m_{\hat{j}}^{\hat{t}}$ . To see this, we proceed by contradiction and assume that  $\sigma_{-i}(h_{t'}, s_{-i}) = m_{-i}^{t'}$  for all  $t' \leq t-1$ . This implies that whenever  $h_{t'-1} \in \mathcal{H}_{-i}^*$ , we must have  $h_{t'} \in \mathcal{H}_{-i}^*$ , because  $h_{t'-1} \in \mathcal{H}_{-i}^*$  implies that either  $h_{t'-1} \in \mathcal{H}^*$  (i.e., no player has deviated) or  $h_{t'-1} \notin \mathcal{H}_{-j}^*$  for all  $j \neq i$  (i.e., player  $i$  has deviated). In either case,  $\sigma_{-i}(h_{t'-1}, s_{-i}) = m_{-i, \theta'}^*(h_{t'-1})$  is obtained by **\Sigma1** and **\Sigma3**. Since we have assumed that  $\sigma_{-i}(h_{t'-1}, s_{-i}) = m_{-i}^{t'-1}$ , we get  $m_{-i}^{t'-1} = m_{-i, \theta'}^*(h_{t'-1})$ , which proves that  $h_{t'} \in \mathcal{H}_{-i}^*$ . Since  $h_1 = \emptyset \in \mathcal{H}^* \subseteq \mathcal{H}_{-i}^*$ , this simple inductive argument shows that  $h_t \in \mathcal{H}_{-i}^*$ , a contradiction.

By construction of  $\sigma^k$ , this implies that for some  $\hat{j} \neq i$  and  $\hat{t} \leq t-1$ :

$$\sigma_{\hat{j}}^k(m_{\hat{j}}^{\hat{t}}|h_{\hat{t}}, s_{\hat{j}}^{\theta''}) = \eta_k^{T \times n} \xi_{\hat{j}}(h_{\hat{t}}, s_{\hat{j}}^{\theta''}, m_{\hat{j}}^{\hat{t}}). \tag{7}$$

Now, we have:

$$\begin{aligned}
& \frac{\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta''})}{\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta'})} \leq \frac{\eta_k^{T \times n} \times \xi_{\hat{j}}(h_{\hat{t}}, s_{\hat{j}}^{\theta''}, m_{\hat{j}}^{\hat{t}}) \times 1}{\prod_{j \neq i} \prod_{t'=1}^{t-1} \eta_k \xi_j(h_{t'}, s_j^{\theta'}, m_j^{t'})} \\
&= \frac{\eta_k^{T \times n}}{\eta_k^{(t-1)(n-1)}} \times \frac{\xi_{\hat{j}}(h_{\hat{t}}, s_{\hat{j}}^{\theta''}, m_{\hat{j}}^{\hat{t}})}{\prod_{j \neq i} \prod_{t'=1}^{t-1} \xi_j(h_{t'}, s_j^{\theta'}, m_j^{t'})} \rightarrow 0 \text{ (as } k \rightarrow \infty).
\end{aligned}$$

Here, the inequality is assured by (7) and the construction of  $\sigma^k$  that, for all  $j$  and  $t' \leq t-1$ ,  $\sigma_j^k(m_j^{t'}|h_{t'}, s_j^{\theta'}) \geq \eta_k \times \xi_j(h_{t'}, s_j^{\theta'}, m_j^{t'})$ .

**Case 2:**  $s_i = s_i^{\theta'}$

Recall that  $\nu^\varepsilon(\cdot, s_i^{\theta'})$  assigns a strictly positive weight only to  $(\theta', s_{-i}^{\theta'})$  and  $(\theta', \tau_l)$  for each  $l \neq i$ .

Hence,

$$\begin{aligned}
& \phi_i^k[(\theta', \tau_l) | s_i^{\theta'}, h_t] \\
&= \frac{\nu^\varepsilon(\theta', \tau_l) \times \prod_{j \neq l, i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'}) \times \prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''})}{\sum_{z \neq i} \nu^\varepsilon(\theta', \tau_z) \times \prod_{j \neq z, i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'}) + \nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta'}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'})} \\
&= \frac{\nu^\varepsilon(\theta', \tau_l)}{\sum_{z \neq i} \nu^\varepsilon(\theta', \tau_z) \times c_z(k) + \nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta'}) \times \prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta'}) / \prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''})}
\end{aligned}$$

for some positive functions  $c_z(k)$ . We now show that if  $h_t \in \mathcal{H}_{-l}^*$ , then the ratio below converges to  $+\infty$  as  $k \rightarrow \infty$ :

$$\prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta'}) / \prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''}) \rightarrow +\infty \text{ as } k \rightarrow \infty.$$

This will show that  $\phi_i^k[(\theta', \tau_l) | s_i^{\theta'}, h_t] \rightarrow 0$  for all  $l$  if  $h_t \in \mathcal{H}_{-l}^*$ ; and hence that  $\phi$  satisfies  $\Phi 1$ . Assume that  $h_t \in H_{-l}^*$  for some  $l$ , as we already claimed, if every player  $j$  other than  $i$  has received a signal  $s_j \in \{s_j^{\theta'}, s_j^{\theta''}\}$ , there is a player  $\hat{j} \neq i$  and a date  $\hat{t} \leq t-1$  so that  $\hat{j}$  has not played according to  $\sigma_{\hat{j}}$ , i.e.,  $\sigma_{\hat{j}}(h_{\hat{t}}, s_{\hat{j}}) \neq m_{\hat{j}}^{\hat{t}}$ . Now, since  $h_t \in \mathcal{H}_{-l}^*$ , we claim that  $\hat{j} = l$ . Indeed,  $h_t \in \mathcal{H}_{-l}^*$  means that any player  $j$  other than  $l$  has played according to  $m_{j, \theta'}^*$ . So if player  $l$  had played according to  $\sigma_l$  (i.e., for all  $t' : \sigma_l(h_{t'}, s_l) = m_l^{t'}$ ), repeated applications of  $\Sigma 1$  and  $\Sigma 3$  would yield to  $h_t = h_t^* \in \mathcal{H}_{-i}^*$  which is false by assumption.

By construction of  $\sigma^k$ , this implies that there exists  $\hat{t} \leq t-1$  such that  $\sigma_l(h_{\hat{t}}, s_l) \neq m_l^{\hat{t}}$  and so:

$$\sigma_l^k(m_l^{\hat{t}} | h_{\hat{t}}, s_l^{\theta''}) = \eta_k^{T \times n} \xi_l(h_{\hat{t}}, s_l^{\theta''}, m_l^{\hat{t}}). \tag{8}$$

Now, we have

$$\frac{\prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta'})}{\prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''})} \geq \frac{\eta_k^{t-1} \prod_{t'=1}^{t-1} \xi_l(h_{t'}, s_l^{\theta'}, m_l^{t'})}{\eta_k^{T \times n} \xi_l(h_{\hat{t}}, s_l^{\theta''}, m_l^{\hat{t}}) \times 1} \rightarrow \infty \text{ (as } k \rightarrow \infty).$$

Where the inequality is assured by (8) and (assuming without loss of generality that  $\eta_k$  is small) we use the fact that by construction, for all  $t' \leq t-1$ ,  $\sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta'}) \geq \eta_k \times \xi_l(h_{t'}, s_l^{\theta'}, m_l^{t'})$ . ■

**Claim 5.**  $\phi$  satisfies  $\Phi 2$ .

Claim 5 says that if a player  $i$  gets signal  $s_i^{\theta'}$  or  $s_i^{\theta''}$  then at a history  $h_t$  under which each of his opponent has played according to  $m_{\theta'}^*$ ,  $\phi$  is the same as his beliefs given only by his private signal.

To prove this, we show that if every player  $j \neq i$  has received a signal  $s_j \in \{s_j^{\theta'}, s_j^{\theta''}\}$  then at histories where all players other than  $i$  have played according to  $m_{\theta'}^*$ , each player other than  $i$  has played according to  $\sigma$  at each previous stage. This ensures that for any  $h_t \in \mathcal{H}_{-i}^*$ , no player other than  $i$  has deviated from the candidate for sequential equilibrium strategy  $\sigma$  and so player  $i$ 's beliefs must be given by his private signal.

*Proof of Claim 5.* Consider player  $i$  at history  $h_t \in \mathcal{H}_{-i}^*$ . Here again, the proof is reduced to checking the following two cases.

**Case 1:**  $s_i = s_i^{\theta''}$

Recall that  $\nu^\varepsilon(\cdot, s_{-i}^{\theta''})$  assigns a strictly positive weight only to  $(\theta'', s_{-i}^{\theta''})$  and  $(\theta', s_{-i}^{\theta'})$ . Hence,

$$\begin{aligned} & \phi_i^k[(\theta'', s_{-i}^{\theta''}) | s_i^{\theta''}, h_t] \\ &= \frac{\nu^\varepsilon(\theta'', s_{-i}^{\theta''}, s_i^{\theta''}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta''})}{\nu^\varepsilon(\theta'', s_{-i}^{\theta''}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta''}) + \nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta''}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'})} \\ &= \frac{\nu^\varepsilon(\theta'', s_{-i}^{\theta''}, s_i^{\theta''})}{\nu^\varepsilon(\theta'', s_{-i}^{\theta''}) + \nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta''}) \times \frac{\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'})}{\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta''})}} \end{aligned}$$

We now show that the ratio below converges to 1 as  $k \rightarrow \infty$ :

$$\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'}) / \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta''}) \rightarrow 1 \quad \text{as } k \rightarrow \infty.$$

This will show that  $\phi_i^k[(\theta'', s_{-i}^{\theta''}) | s_i^{\theta''}, h_t] \rightarrow \nu^\varepsilon((\theta'', s_{-i}^{\theta''}) | s_i^{\theta''})$  and  $\phi_i^k[(\theta', s_{-i}^{\theta'}) | s_i^{\theta''}, h_t] \rightarrow \nu^\varepsilon((\theta', s_{-i}^{\theta'}) | s_i^{\theta''})$ .

Note now that if players  $j \neq i$  receive signal  $s_j \in \{s_j^{\theta'}, s_j^{\theta''}\}$ , then for all  $t' \leq t-1$ ,  $\sigma_j(h_{t'}, s_j) = m_j^{t'}$ . To see this, note that for any  $t' \leq t-1$ :  $h_{t'} \in \mathcal{H}_{-i}^*$ , thus, either every player has played according to  $m_{\theta'}^*$  (i.e.,  $h_{t'} \in \mathcal{H}^*$ ) or player  $i$  has not played according to  $m_{i, \theta'}^*$  (i.e.,  $h_{t'} \notin \mathcal{H}_{-i}^*$  for all  $j \neq i$ ). In each of these cases we know, by  $\Sigma\mathbf{1}$  and  $\Sigma\mathbf{3}$ , that  $\sigma_j$  prescribes to play according to  $m_{j, \theta'}^*$ . Since  $h_{t'} \in \mathcal{H}_{-i}^*$  this implies that  $\sigma_j(h_{t'}, s_j) = m_{j, \theta'}^*(h_{t'}) = m_j^{t'}$ .

By construction of  $\sigma^k$ , this in turn implies that for all  $j \neq i$  and  $t' \leq t-1$ :

$$\sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'}) = 1 - \eta_k \quad \text{and} \quad \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta''}) = 1 - \eta_k^{T \times n}.$$

Thus,

$$\prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'}) \bigg/ \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta''}) \rightarrow 1 \text{ as } k \rightarrow \infty.$$

**Case 2:**  $s_i = s_i^{\theta'}$

Recall that  $\nu^\varepsilon(\cdot, s_i^{\theta'})$  assigns a strictly positive weight only to  $(\theta', s_{-i}^{\theta'})$  and  $(\theta', \tau_l)$  for  $l \neq i$ .

Hence,

$$\begin{aligned} & \phi_i^k[(\theta', s_{-i}^{\theta'}) | s_i^{\theta'}, h_t] \\ &= \frac{\nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta'}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'})}{\nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta'}) \times \prod_{j \neq i} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'}) + \sum_{l \neq i} \nu^\varepsilon(\theta', \tau_l) \times \prod_{j \neq i, l} \prod_{t'=1}^{t-1} \sigma_j^k(m_j^{t'} | h_{t'}, s_j^{\theta'}) \times \prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''})} \\ &= \frac{\nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta'})}{\nu^\varepsilon(\theta', s_{-i}^{\theta'}, s_i^{\theta'}) + \sum_{l \neq i} \nu^\varepsilon(\theta', \tau_l) \times \frac{\prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''})}{\prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta'})}} \end{aligned}$$

We now show that for each  $l \neq i$ , the ratio below converges to 1 as  $k \rightarrow \infty$ :

$$\prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''}) \bigg/ \prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta'}) \rightarrow 1 \text{ as } k \rightarrow \infty.$$

This will show that  $\phi_i^k[(\theta', s_{-i}^{\theta'}) | s_i^{\theta'}, h_t] \rightarrow \nu^\varepsilon((\theta', s_{-i}^{\theta'}) | s_i^{\theta'})$  and similar reasoning shows that for each  $l \neq i$ :  $\phi_i^k[(\theta', \tau_l) | s_i^{\theta'}, h_t] \rightarrow \nu^\varepsilon((\theta', \tau_l) | s_i^{\theta'})$ , and hence,  $\phi$  satisfies  $\Phi 2$ . ■

Now, by similar reasoning as in the case above, we get that for all  $l \neq i$  and  $t' \leq t-1$ :

$$\sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta'}) = 1 - \eta_k \text{ and } \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''}) = 1 - \eta_k^{T \times n}.$$

Thus,

$$\prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta''}) \bigg/ \prod_{t'=1}^{t-1} \sigma_l^k(m_l^{t'} | h_{t'}, s_l^{\theta'}) \rightarrow 1 \text{ as } k \rightarrow \infty.$$

Finally, observing that for  $s_i^{\tilde{\theta}} \notin \{s_i^{\theta'}, s_i^{\theta''}\}$ ,  $\nu^\varepsilon(\cdot, s_i^{\tilde{\theta}})$  assigns a weight one to  $(\tilde{\theta}, s_{-i}^{\tilde{\theta}})$ , we have established the following claim, which completes the proof of Theorem 3.

**Claim 6.**  $\phi$  satisfies  $\Phi 3$ .

### 7.3 Theorem 3 extends to countable messages

Here we extend Theorem 3 to mechanisms that have countably infinite message spaces. This extension is important because some of the literature on implementation theory uses “integer games” where each player has to announce an integer and becomes the dictator when his integer is the



largest one, as in Maskin (1999) and in Moore and Repullo (1988).

**Assumption A2.**  $M_i(h)$  is countable for each  $i$  and  $h$ .

The next assumption says that against any profile of strategy in the complete information game, in the neighborhood of complete information, each player  $i$  has a non-empty set of best responses. This condition is vacuously satisfied under A1, so Theorems 3 and 4 show that if a mechanism can implement a non-Maskin monotonic social choice correspondence (SCC) both under complete information and under small information perturbations, then under this mechanism players must not have well-defined best responses. In addition, we show in the supplemental materials that when the state space is finite (this is our case), Moore and Repullo's general mechanism has well-defined best-responses (under weak assumptions) and so our argument also applies there.

**Assumption A3.** *The sequential mechanism  $\Gamma$  has well-defined best replies: for any player  $i$ , any  $\theta \in \Theta$ , any  $m_{-i} \in M_{-i}$ , there exists  $\bar{\xi}(i, \theta, m_{-i}) > 0$  such that for any  $\beta \in \Delta(\Theta)$  with  $\beta(\theta) \geq 1 - \bar{\xi}(i, \theta, m_{-i})$ , for any  $m_i \in M_i$  we have for all  $h \in \mathcal{H}$ :*

$$\arg \max_{\tilde{\theta}} \sum \beta(\tilde{\theta}) u_i(g((m'_i, m_{-i}); h); \tilde{\theta}) \neq \emptyset$$

where the max is taken over all pure messages  $m'_i \in M_i$  that differ from  $m_i$  only at  $h$ .

**Remark 1.** *If the mechanism is not finite but the set of outcomes is, A3 is also vacuously satisfied. We also note that A3 is not needed for sequential mechanisms in which each player moves only once.*<sup>40</sup>

**Theorem 4.** *Assume A2 and A3. Suppose that a mechanism  $\Gamma$  SPE-implements a non-Maskin monotonic SCC  $\mathcal{F}$ . Fix any complete information prior  $\mu$ . There exist a sequence of priors  $\{\nu^\varepsilon\}_{\varepsilon>0}$  that converges to  $\mu$  and a corresponding sequence of sequential equilibria  $\{(\phi^\varepsilon, \sigma^\varepsilon)\}_{\varepsilon>0}$  such that as  $\varepsilon$  tends to 0,  $g(\sigma^\varepsilon(s^\theta); \emptyset) \rightarrow a \notin \mathcal{F}(\theta)$  for some  $\theta \in \Theta$  and some  $a \in A$ .*

*Proof.* The proof is essentially the same as the proof of Theorem 3 where we only consider finite mechanisms. So, we claim that there are essentially only two changes we need to extend the proof of Theorem 3 to the case of countably infinite message spaces. First, in the beginning of the proof of Theorem 3, we have to choose  $\varepsilon > 0$  small enough to apply A3. Second, we will show that A3 guarantees that  $\Sigma 2$  (which is introduced in the proof of Theorem 3) is well defined. This will be proved in the next subsection. ■

### 7.3.1 A3 guarantees that $\Sigma 2$ is well-defined

Fix  $\varepsilon > 0$  small enough so that  $\nu^\varepsilon(\theta' | s_i^{\theta'}) \geq 1 - \bar{\xi}(i, \theta, m_{-i, \theta}^*)$ . We shall claim that A3 guarantees that one can construct  $\bar{m}_i$  needed for  $\Sigma 2$ . First, for any  $h_t \in \mathcal{H}^*$  or  $h_t \notin \mathcal{H}_{-i}^*$ , we set  $\bar{m}_i(h_t) =$

<sup>40</sup>One can directly check this in the definition of strategy  $\sigma$  ( $\Sigma 2$ ) used in the proof of Theorem 3. More specifically, it can be checked there that for each player, A3 is only used at histories where this player has to choose a message and at which he has previously deviated from the equilibrium. By construction, there is no such a history.

$m_{i,\theta}^*(h_t)$ . Second, we define  $\bar{m}_i$  by induction on the set of histories in  $\mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ . Take any history  $h_t \in \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$  so that there is no subsequent history that falls into  $\mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ . Since we already defined  $\bar{m}_i(h_t) = m_{i,\theta}^*(h_t)$  for any  $h_t \notin \mathcal{H}_{-i}^* \setminus \mathcal{H}^*$ ,  $\bar{m}_i$  has been defined for any subsequent histories. By A3 we obtain

$$\arg \max_{\tilde{\theta}} \sum \nu^\varepsilon(\tilde{\theta} | s_i^{\theta'}) u_i(g((m'_i, m_{-i,\theta}^*)); h_t); \tilde{\theta} \neq \emptyset$$

where the max is taken over all pure messages  $m'_i \in M_i$  that differ from  $\bar{m}_i$  only at  $h_t$  and are identical at any subsequent histories (what happens before  $h_t$  is obviously irrelevant).

Now set

$$\bar{m}_i(h_t) \in \arg \max_{\tilde{\theta}} \sum \nu^\varepsilon(\tilde{\theta} | s_i^{\theta'}) u_i(g((m'_i, m_{-i})); h_t); \tilde{\theta}.$$

This establishes that one can inductively construct  $\bar{m}_i$  so that  $\bar{m}_i$  satisfies the properties needed for  $\Sigma 2$ .

### 7.3.2 A3 is satisfied in the Moore-Repullo canonical mechanism

We will review some of the main results of Moore and Repullo (1988) here.

**Definition 3** (Moore and Repullo (1988)). *A social choice correspondence  $\mathcal{F}$  satisfies Condition C if, for every pair of profiles  $\theta, \phi \in \Theta$  with  $a \in \mathcal{F}(\theta) \setminus \mathcal{F}(\phi)$ , there exists a finite sequence*

$$\sigma(\theta, \phi; a) \equiv \{a_0 = a, a_1, \dots, a_k, \dots, a_l, a_{l+1}\} \subset A,$$

with  $l = l(\theta, \phi; a) \geq 1$ , such that:

1. for each  $k = 0, \dots, l-1$ , there is some particular player  $j(k) = j(k|\theta, \phi; a)$ , for whom

$$u_{j(k)}(a_k; \theta) \geq u_{j(k)}(a_{k+1}; \theta);$$

2. there is some player  $j(l) = j(l|\theta, \phi; a)$  for whom

$$u_{j(l)}(a_l; \theta) \geq u_{j(l)}(a_{l+1}; \theta) \text{ and } u_{j(l)}(a_{l+1}; \phi) > u_{j(l)}(a_l; \phi).$$

Further,  $l(\theta, \phi; a)$  is uniformly bounded by some  $\bar{l} < \infty$ .

Assuming Condition C holds, let  $\mathcal{Q}(\mathcal{F})$  be a class of subsets  $Q$  of  $A$ . A typical  $Q$  is defined as follows:

For each pair of profiles  $\theta$  and  $\phi$  in  $\Theta$ , and for each  $a \in \mathcal{F}(\theta) \setminus \mathcal{F}(\phi)$ , select one sequence  $\sigma(\theta, \phi; a)$  satisfying (1) and (2) in Condition C. Then let  $Q$  be the union of the elements in these sequences.

$\mathcal{Q}(\mathcal{F})$  comprises the  $Q$ 's constructed from all possible selections.

**Definition 4.** A social choice correspondence  $\mathcal{F}$  satisfies Condition  $C^+$  if it satisfies Condition  $C$  and the following condition as well: there exists a particular  $Q^+ \in \mathcal{Q}(\mathcal{F})$ , and a particular set  $B \subseteq A$  containing  $Q^+$ , such that the following is true for each  $\theta \in \Theta$ :

- Each player  $i$  has nonempty maximal set  $B_i^*(\theta) \subseteq B$  under  $\theta$ , i.e.,  $B_i^*(\theta) = \arg \max_{a \in B} u_i(a; \theta)$ .
- $B_i^*(\theta) \cap B_j^*(\theta) = \emptyset$  for each  $\theta \in \Theta$  and each  $i, j \in N$  with  $i \neq j$
- $B_i^*(\theta) \cap Q^+ = \emptyset$  for each  $i$  and each  $\theta$ .

Let the selected sequences  $\sigma(\theta, \phi; a) \in Q^+$  be labelled  $\sigma^+(\theta, \phi; a)$ . Define the Moore-Repullo canonical mechanism  $\Gamma^{MR} = (M, g)$  as follows.

**Stage 0:** each player  $i$  announces some triplet  $m_{i,0} = (\theta^i, a^i, n_0^i)$ , where  $\theta^i \in \Theta$ ,  $a^i \in \mathcal{F}(\theta^i)$ , and  $n_0^i$  is a nonnegative integer. There are three possibilities to consider:

1. all  $n$  players agree on a common profile  $\theta$  and outcome  $a \in \mathcal{F}(\theta)$ , then outcome  $a$  is chosen. STOP
2. If only  $n - 1$  players agree on a common profile  $\theta$  and outcome  $a \in \mathcal{F}(\theta)$ , and if the remaining player  $i$  announces a profile  $\phi$ , and
  - (a) if  $a \in \mathcal{F}(\phi)$ , then outcome  $a$  is implemented; STOP
  - (b) if  $a \notin \mathcal{F}(\phi)$  but  $i$  is not the agent  $j(0)$  prescribed in  $\sigma^+(\theta, \phi; a)$ , then outcome  $a$  is implemented; STOP
  - (c) if  $a \notin \mathcal{F}(\phi)$  and  $i = j(0)$ , then go to Stage 1.
3. If neither (1) nor (2) apply, then the player with the highest integer  $n_0^i$  is allowed to choose an outcome from  $B$ . Ties are broken by selecting from the players who announced the highest number according to who has the smallest  $i$ . STOP

**Stage  $k = 1, \dots, l$ :** each player  $i$  can either raise a “flag,” or announce a nonnegative integer  $n_k^i \in \mathbb{N}$ , i.e.,  $m_{i,k} \in M_{i,k} \in \{\text{flag}\} \cup \mathbb{N}$ . Again there are three possibilities to consider:

1. If  $n - 1$  or more flags are raised, then the agent  $j(k - 1)$  prescribed in  $\sigma^+(\theta, \phi; a)$  is allowed to choose an outcome from  $B$ . STOP
2. If  $n - 1$  or more players announce zero, and
  - (a) if the player  $j(k)$  prescribed in  $\sigma^+(\theta, \phi; a)$  is one of those who announce zero, then implement outcome  $a_k$  from sequence  $\sigma^+(\theta, \phi; a)$ ; STOP
  - (b) if  $j(k)$  does not announce zero, then
    - i. if  $k < l$ , go to Stage  $k + 1$ ;

- ii. if  $k = l$ , implement outcome  $a_{l+1}$  from sequence  $\sigma^+(\theta, \phi; a)$ . STOP
- (c) If neither (1) nor (2) apply, then the player who announces the highest integer  $n_k^i$  is allowed to choose an outcome from  $B$ . STOP

**Theorem 5** (Moore and Repullo (1988)). *If a social choice correspondence  $\mathcal{F}$  satisfies Condition  $C^+$ , and  $n \geq 3$ , then  $\mathcal{F}$  can be implemented in subgame-perfect equilibrium.*

Moore and Repullo (1988) show the above theorem by using the mechanism described above. We note that this mechanism satisfies A3 if the set of outcomes  $A$  is finite or when each player's preferences over  $A$  are strict and utilities are bounded. Furthermore, the above mechanism satisfies A3 whenever (i) the set  $B$  given in Condition  $C^+$  is a compact set of outcomes; (ii)  $u_i : A \times \Theta \rightarrow \mathbb{R}$  is continuous in  $a$ .<sup>41,42</sup> It is worth noting that many researchers assume (i) and (ii) after appealing to Moore and Repullo's (1988) result. This is the case for instance in Moore and Repullo (1988)'s examples of risk-sharing (Section 6.1) or the production contract example (Section 6.2). More importantly, it is also the case in Maskin and Tirole (1999a)'s proof of the irrelevance Theorem. Hence our non-robustness result (Theorem 4) also apply to Maskin and Tirole's irrelevance Theorem.

#### 7.4 Sufficiency for Robust Implementation: The Case of Social Choice Correspondences (SCCs)

In Remark 3 of Section 4, we argue that Maskin monotonic social choice functions are robustly implementable. Here we extend this argument to the case of social choice correspondences.

We need to strengthen Maskin monotonicity to the following:

**Definition 5.** *An SCC  $\mathcal{F}$  satisfies **strong Maskin Monotonicity** if for every SCF  $f$  selected from  $\mathcal{F}$  and every pair of states  $\theta'$  and  $\theta''$  such that*

$$\{(i, b) \mid u_i(f(\theta'); \theta') > u_i(b; \theta')\} \subseteq \{(i, b) \mid u_i(f(\theta'); \theta'') \geq u_i(b; \theta'')\}$$

*then  $f(\theta') \in \mathcal{F}(\theta'')$ .*

Strong Maskin monotonicity is equivalent to Maskin monotonicity in many economic environments.<sup>43</sup> For example, consider environments in which there is a private good that is both desirable and continuously transferable. Another example is an environment in which agents have strict preferences. The next definition is the no-veto-power condition, which is widely used in the literature.

<sup>41</sup>Then, for any  $\beta \in \Delta(\Theta)$ ,  $\arg \max_{a \in B} \sum_{\tilde{\theta}} \beta(\tilde{\theta}) u_i(a; \tilde{\theta}) \neq \emptyset$ . We note that a one-shot deviation of player  $i$  at stage  $k$  in  $\Gamma^{MR}$  allows player  $i$  possibly to fall into an integer game at stage  $k$  where he can get any outcome in  $B$ ; if he cannot fall into this integer game, he can only induce a finite number of outcomes, say  $B_k$ , by deviating. In any case, he has a most preferred deviation, i.e.,  $\arg \max_{a \in B} \sum_{\tilde{\theta}} \beta(\tilde{\theta}) u_i(a; \tilde{\theta}) \neq \emptyset$ ,  $\arg \max_{a \in B \cup B_k} \sum_{\tilde{\theta}} \beta(\tilde{\theta}) u_i(a; \tilde{\theta}) \neq \emptyset$ , and  $\arg \max_{a \in B_k} \sum_{\tilde{\theta}} \beta(\tilde{\theta}) u_i(a; \tilde{\theta}) \neq \emptyset$ . Then A3 is satisfied whenever (i) and (ii) hold.

<sup>42</sup>Note that A2 need not be satisfied for these mechanisms since  $B$  need not be countable. A2 was introduced only to define sequential equilibrium in a simple manner. If one uses perfect Bayesian equilibrium instead, we believe that A2 is not required.

<sup>43</sup>What we mean by "strong" is that we replace the first weak inequality of (\*) in the definition of Maskin monotonicity with a strict one. This notion also appears in Chung and Ely (2003).

**Definition 6.** An SCC  $\mathcal{F}$  satisfies **no-veto-power** if whenever there is an alternative  $c \in A$  such that for at least  $n - 1$  players  $i$ ,  $u_i(c; \theta) \geq u_i(b; \theta)$  for every  $b \in A$ , we have  $c \in \mathcal{F}(\theta)$ .

We need one extra condition together with strong Maskin monotonicity and no-veto power. This is the no-worst-alternative condition as defined by Cabrales and Serrano (2011):

**Definition 7.** An SCC  $\mathcal{F}$  satisfies the **no-worst-alternative** (NWA) condition if for each agent  $i \in N$ ,  $\theta \in \Theta$  and  $f$  selected from  $\mathcal{F}$ , there exists  $z(i, \theta, f) \in A$  such that  $u_i(f(\theta); \theta) > u_i(z(i, \theta, f); \theta)$ .

Let  $\mathcal{P}$  denote the set of priors over  $\Theta \times S$  with the following metric  $d : \mathcal{P} \times \mathcal{P} \rightarrow \mathbb{R}_+$ : for any  $\mu, \mu' \in \mathcal{P}$ ,

$$d(\mu, \mu') = \max_{(\theta, s) \in \Theta \times S} |\mu(\theta, s) - \mu'(\theta, s)|.$$

So, when we say  $\nu^k \rightarrow \mu$ , we mean that  $d(\nu^k, \mu) \rightarrow 0$  as  $k \rightarrow \infty$ . When  $\Theta \times S$  is a finite state space, Theorem 14.5 of Fudenberg and Tirole (1991a) shows that when  $\nu^k \rightarrow \mu$  as  $k \rightarrow \infty$ , there exists  $\{p^k\}_{k=1}^\infty$  such that (1)  $p^k \rightarrow 1$  as  $k \rightarrow \infty$ ; (2)  $\nu^k(\{(\theta, s^\theta)\}_{\theta \in \Theta}) \geq p^k$  for each  $k$ ; and (3) for each  $k$ , it is common  $p^k$ -belief at any profile of signals  $s^\theta$  that  $\theta$  has realized.<sup>44</sup>

We propose the following definition of robust implementation:

**Definition 8.** An SCC  $\mathcal{F}$  is **robustly implementable** under the complete information prior  $\mu$  if there exists a mechanism  $\Gamma = (M, g)$  satisfying the following two properties: for any SCF  $f$  selected from  $\mathcal{F}$  and any sequence of priors  $\{\nu^\varepsilon\}_{\varepsilon > 0}$  converging to  $\mu$ , (1) there is a sequence of sequential equilibria  $\{\sigma^\varepsilon\}_{\varepsilon > 0}$  in  $\{\Gamma(\nu^\varepsilon)\}_{\varepsilon > 0}$  satisfying  $\lim_{\varepsilon \rightarrow 0} g(\sigma^\varepsilon(s^\theta); \emptyset) = f(\theta)$  for every  $\theta \in \Theta$ ; and (2) for any sequence of sequential equilibria  $\{\sigma^\varepsilon\}_{\varepsilon > 0}$  in  $\{\Gamma(\nu^\varepsilon)\}_{\varepsilon > 0}$ , we have  $\lim_{\varepsilon \rightarrow 0} g(\sigma^\varepsilon(s^\theta); \emptyset) \in \mathcal{F}(\theta)$  for every  $\theta \in \Theta$ .

**Remark 4.:** The first requirement of robust implementation says that for any SCF  $f$  selected from a given SCC  $\mathcal{F}$  and any environment near  $\mu$ , there is an equilibrium whose outcome is close to that given by  $f$  whenever a signal profile  $s$  has strictly positive probability under  $\mu$  (i.e.,  $s = s^\theta$  for some  $\theta$ ). The second requirement says that for any environment near  $\mu$ , whenever a signal profile  $s$  has strictly positive probability under  $\mu$ , equilibrium outcomes are close to that of  $\mathcal{F}$ . Both requirements are robust analogs of the two standard requirements of implementation.<sup>45</sup> Roughly speaking, the first requirement embodies a version of lower hemi-continuity of the equilibrium correspondence and the second embodies a version of upper hemi-continuity.<sup>46</sup> As is clear from the proof of Theorem 3, to show that Maskin monotonicity is necessary for SCCs to be robustly implemented, we only used the second property of robust implementation and do not exploit the full strength of robust implementation. Finally, the subsequent argument provides sufficient conditions under which a

<sup>44</sup>See Monderer and Samet (1989) and/or Fudenberg and Tirole (1991a) for the precise definition of common  $p$ -belief.

<sup>45</sup>See, for instance, Maskin (1999) for the definition of Nash implementation.

<sup>46</sup>Property (2) in our definition says that the correspondence from priors to equilibrium outcomes has a closed graph. In general, this is not equivalent to upper hemi-continuity. However, the closed graph property of the equilibrium outcomes correspondence implies upper hemi-continuity if the range of the correspondence is compact (see Aliprantis and Border (1999)).

static mechanism yields robust implementation. Hence, the result would hold if we were to replace sequential equilibrium by Nash equilibrium in the above statement.

We are now ready to state the result:

**Theorem 6.** *Suppose there are at least three players, i.e.,  $|N| = n \geq 3$ . If an SCC  $\mathcal{F}$  satisfies strong Maskin monotonicity, no-veto-power and the NWA condition, then  $\mathcal{F}$  is robustly implementable.*

*Proof.* We construct an implementing mechanism  $\Gamma = (M, g)$ .<sup>47</sup> For each  $i \in N$ , let  $M_i = (\Theta \times \mathcal{F}) \cup (\mathbb{Z}_+ \times A)$  where  $\mathbb{Z}_+$  is the set of nonnegative integers. That is, each agent is asked to report *either* a state and a social choice function or an integer and an alternative. Let  $m^{\theta, f}$  denote the message profile  $((\theta, f), (\theta, f), \dots, (\theta, f))$ , and  $m^{\theta, f} \setminus m_i$  the profile obtained from  $m^{\theta, f}$  by substituting  $m_i$  for agent  $i$ . We set  $g(m^{\theta, f}) = f(\theta)$ . If  $m_i = (\theta', f')$ , and if there exists an alternative  $c \in A$  such that  $u_i(c; \theta') > u_i(f(\theta); \theta')$  but  $u_i(f(\theta); \theta) > u_i(c; \theta)$ , then we set  $g(m^{\theta, f} \setminus m_i) = c$ . (If there is more than one such  $c$ , select one arbitrarily). For all other cases, we set  $g(m^{\theta, f} \setminus m_i) = z(i, \theta, f(\theta))$  as defined for the NWA condition.

Consider any other profile of messages  $m$ . If each  $m_i$  consists of a state and a social choice function, then choose  $g(m)$  to be an arbitrary element of  $\mathcal{F}(\Theta)$ . If at least one agent has announced an integer and an alternative, set  $g(m)$  to be the alternative named by the agent whose named integer is the greatest (breaking ties by choosing the lowest index among those who announced the greatest integer).

The rest of the proof can be completed by the following three steps: in Step 1, we show that for any SCF  $f$  selected from  $\mathcal{F}$ , there exists a good equilibrium whose outcome coincides with that of  $f$  for any nearby environment. In Step 2, we show that any Nash equilibrium outcome is socially desirable. In Step 3, we show that this continues to be the case in nearby environments.

For any complete information prior  $\mu$ , let  $U(\mu)$  denote a neighborhood around  $\mu$  with respect to metric  $d$ .

**Step 1:** Let  $\mu$  be a complete information prior. For each SCF  $f$  selected from  $\mathcal{F}$ , there exists a neighborhood  $U(\mu)$  for which there exists a strict Bayesian Nash equilibrium  $\sigma$  of the game  $\Gamma(\nu)$  for each  $\nu \in U(\mu)$  such that  $g(\sigma(s^\theta)) = f(\theta)$  for every  $\theta \in \Theta$ .

For each SCF  $f$  selected from  $\mathcal{F}$  and  $\theta \in \Theta$ , consider the truthful strategy of agent  $i$  as  $m_i^{\theta, f} = (\theta, f)$ . This yields  $g(m^{\theta, f}) = f(\theta)$ . By construction, if in state  $\theta$ , agent  $i$  sends message  $m_i \neq m_i^{\theta, f}$ ,

$$u_i(g(m^{\theta, f}); \theta) > u_i(g(m^{\theta, f} \setminus m_i); \theta).$$

Hence,  $m^{\theta, f}$  is a *strict* Nash equilibrium of the game  $\Gamma(\theta)$ . Define  $\sigma_i(s_i^\theta) = (\theta, f)$  for each  $s_i^\theta \in S_i$  as agent  $i$ 's strategy of the game  $\Gamma(\mu)$ . Then  $\sigma$  is a strict Nash equilibrium of the game  $\Gamma(\mu)$ . Define

$$A[\sigma_{-i}] = \left\{ a \in A \mid \exists s_{-i} \in S_{-i}, \exists \sigma'_i \text{ such that } g(\sigma'_i(s_i), \sigma_{-i}(s_{-i})) = a \right\}$$

as the set of possible outcomes that can be induced by agent  $i$ 's strategy  $\sigma'_i$  against  $\sigma_{-i}$ . By

<sup>47</sup>The proof here is a modification of that of Theorem 2 of Chung and Ely (2003).

construction of  $\Gamma$  and the finiteness of  $S$ ,  $A[\sigma_{-i}]$  is finite. It is important to note that each agent can only induce a finite number of outcomes, while the set of strategies may be infinite. By the continuity of expected utility and the finiteness of  $S$ ,  $N$ , and  $A[\sigma_{-i}]$ , there is a neighborhood  $U(\mu)$  such that  $\sigma$  continues to be a strict Bayesian Nash equilibrium of the game  $\Gamma(\nu)$  for every  $\nu \in U(\mu)$ .

**Step 2:** Let  $\mu$  be a complete information prior and  $\sigma$  be a Nash equilibrium of the game  $\Gamma(\mu)$ . Then,  $g(\sigma(s^\theta)) \in \mathcal{F}(\theta)$  for every  $\theta \in \Theta$ .

Suppose  $\sigma$  is a Nash equilibrium of  $\Gamma(\mu)$ . Assume further that in  $\sigma(s^\theta)$ , each player announces the same state and SCF  $(\theta', f')$ . Then,  $g(\sigma(s^\theta)) = f'(\theta')$ . In this case, we claim that  $f'(\theta') \in \mathcal{F}(\theta)$ . If this is not the case, by strong Maskin monotonicity, there exist a player  $i$  and an alternative  $a$  such that  $u_i(a; \theta) > u_i(f'(\theta'); \theta)$  but  $u_i(f'(\theta'); \theta') > u_i(a; \theta')$ . By construction of  $\Gamma$ , we can conclude that  $g(\sigma(s^\theta) \setminus (\theta, f')) = a$ . Thus,  $\sigma(s^\theta)$  would not be a Nash equilibrium of  $\Gamma(\theta)$ . For any other profile  $\sigma(s^\theta)$ , there must be at least  $n - 1$  agents who can deviate from  $\sigma(\theta)$  and bring about a profile in which there are at least 3 distinct messages. Thus, by construction of  $\Gamma$ , each of these agents could dictatorially choose his most preferred alternative from  $A$  in state  $\theta$ . But since  $\sigma(s^\theta)$  is a Nash equilibrium of  $\Gamma(\theta)$ , it must be that for each of these players  $i$ ,  $u_i(g(\sigma(s^\theta)); \theta) \geq u_i(a; \theta)$  for every  $a \in A$ . Since  $\mathcal{F}$  satisfies no-veto-power,  $g(\sigma(s^\theta)) \in \mathcal{F}(\theta)$ .

**Step 3:** Let  $\mu$  be a complete information. Suppose that  $\sigma$  is a strategy profile such that  $g(\sigma(s^\theta)) \notin \mathcal{F}(\theta)$  for some  $\theta \in \Theta$ . It is enough for our purpose to show that there must exist a neighborhood  $U(\mu)$  such that  $\sigma$  is not a Bayesian Nash equilibrium of the game  $\Gamma(\nu)$  for every  $\nu \in U(\mu)$ .

Suppose  $\sigma$  is given such that  $g(\sigma(s^\theta)) \notin \mathcal{F}(\theta)$  for some  $\theta \in \Theta$ . This implies that  $\sigma$  is not a Nash equilibrium of  $\Gamma(\theta)$ . Hence, there exists an agent  $i$  and a strategy  $\sigma'_i$  such that

$$u_i(g((\sigma'_i, \sigma_{-i})(s^\theta)); \theta) > u_i(g(\sigma(s^\theta)); \theta).$$

By the continuity of expected utility and the finiteness of  $N$ ,  $S$ , and  $A[\sigma_{-i}]$ , there exists a neighborhood  $U(\mu)$  such that for any  $\nu \in U(\mu)$ ,

$$\sum_{\tilde{\theta} \in \Theta} \sum_{s_{-i} \in S_{-i}} \nu(\tilde{\theta}, s_{-i} | s_i^\theta) \left[ u_i(g(\sigma'_i(s_i^\theta), \sigma_{-i}(s_{-i})); \tilde{\theta}) - u_i(g(\sigma_i(s_i^\theta), \sigma_{-i}(s_{-i})); \tilde{\theta}) \right] > 0.$$

This implies that  $\sigma$  is not a Bayesian Nash equilibrium of  $\Gamma(\nu)$  for every  $\nu \in U(\mu)$ . ■