# REPUTATION WITHOUT COMMITMENT IN FINITELY-REPEATED GAMES

JONATHAN WEINSTEIN AND MUHAMET YILDIZ

ABSTRACT. In the reputation literature, players have *commitment types* which represent the possibility that they do not have standard payoffs but instead are constrained to follow a particular plan. In this paper, we show that arbitrary commitment types can emerge from incomplete information about the stage payoffs. In particular, any finitely repeated game with commitment types is strategically equivalent to a standard finitely repeated game with incomplete information about the stage payoffs. Then, classic reputation results can be achieved with uncertainty concerning only the stage payoffs.

JEL Numbers: C72, C73.

# 1. INTRODUCTION

The reputation literature relies on the existence of *commitment types*. These types are not strategic but are certain to follow a particular plan. Since the seminal work of Kreps, Milgrom, Roberts, and Wilson (1982) (henceforth, the Gang of Four), it has been wellestablished that inclusion of commitment types may alter predicted outcomes dramatically, as this may entice the original "rational" types to imitate the commitment types, in order to form a reputation for playing according to the committed plan. Building on this insight, a large literature has emerged, with applications in a wide range of areas.<sup>1</sup>

Of course, commitment types can be modeled by using a payoff function that rewards a player who follows a specific plan. For example, the tit-for-tat types used by the Gang of Four in the analysis of finitely repeated prisoners dilemma could be assigned payoff 1 if they follow

Date: First Version: August 2012; This Version: June 2014.

Weinstein: Washington University in St Louis; j.weinstein@wustl.edu. Yildiz: MIT Economics Department; myildiz@mit.edu. We thank the editor, the referees, and seminar participants at Harvard, Koc, Princeton, and Yale Universities and the SITE, especially Dilip Abreu and Stephen Morris, for useful comments.

<sup>&</sup>lt;sup>1</sup>We refer to the textbook of Mailath and Samuelson (2006) for a review. Throughout the paper, we also refer to the same textbook for the existing repeated games results that we mention without a citation.

### JONATHAN WEINSTEIN AND MUHAMET YILDIZ

tit-for-tat and zero otherwise. However, such payoffs cannot arise in a standard repeated game, i.e. as a discounted sum of stage-game payoffs. The only commitment types that arise directly from modified stage-game payoffs, within a standard repeated-game structure, are those who commit to playing the same action throughout the game. In some models, such types have very significant effects,<sup>2</sup> but unlike tit-for-tat types they have no effect on the repeated prisoners dilemma game.

The form of commitment types is important for the interpretation of reputation results. When commitment types must be restricted by fiat to follow a certain plan, or have payoffs which are not a discounted sum of stage-game payoffs, the literature has sometimes referred to them as "crazy" types. A more generous characterization, more in keeping with the current tone of the literature, would say that commitment types reflect psychological anomalies and motivations that lie outside the game, such as maintaining reputation in the context of a super-game. On the other hand, if commitment types arise solely from heterogeneity in stage-game payoffs and beliefs about these payoffs, then reputation formation can occur with full rationality and without resort to such super-game concerns.

In this paper, we show that for any given plan, a commitment type who is *required* to follow this plan can be mimicked by a utility-maximizing type, which we call a "twin". The twin knows it is common knowledge that they play a repeated game (i.e. he comes from a type space in which only the stage-game payoff functions can vary by type), but his unique rationalizable action is to follow the given plan. Moreover, by embedding a collection of such twins into a single type space, every game with commitment types can be converted to a standard repeated game with incomplete information about the stage-game payoff function, such that the twins have prior probabilities almost identical to the commitment types. Therefore, any model of reputation formation in finitely repeated games, where players form a reputation for certain *beliefs* about the stage-game payoffs. This construction requires that we allow sufficient variations in stage-game payoffs and consider a rich set of information structures.

 $\mathbf{2}$ 

<sup>&</sup>lt;sup>2</sup>For example, the existence of such commitment types is sufficient for the seminal analyses of the repeated entry-deterrence models by Kreps and Wilson (1982) and Milgrom and Roberts (1982) and for the Fudenberg and Levine (1996) result that the informed player's payoff is within his Stackelberg payoffs when the uninformed player is short-lived (best-replying myopically).

### REPUTATION WITHOUT COMMITMENT

Of course, one may also wish to restrict the stage-game payoff functions. For example, in a standard prisoners' dilemma game, one might want to assume that it is common knowledge that cooperation is not dominant. Under such restrictions, twins may not exist for some commitment types. Indeed, we also prove an opposing benchmark, showing that one needs some amount of variations in the stage-game payoffs in order to have any reputational effect. We show that if the stage game is dominance-solvable and the stage game payoffs are restricted to a sufficiently small neighborhood of the original stage-game payoff function, then the unique sequential equilibrium of the repeated game with incomplete information prescribes all players to repeat the stage-game solution throughout the game (as in the subgame-perfect equilibrium of the complete information version), regardless of the length of the game.

Therefore, one needs to allow some substantial amount of variation in stage-game payoffs in order to provide an incomplete-information foundation for the commitment types. While the amount of necessary variation may depend on the details of the game and the commitment types at hand, our main result shows that one can always provide such a foundation as long as there is enough variation in allowable stage-game payoff functions.

One limitation is worth emphasizing here. Our construction makes fundamental use of players who do not know their own payoffs. Some of the literature has focused on models with common knowledge that each player knows his own payoffs; Fudenberg, Kreps, and Levine (1988) call this a model with "personal types." We believe an important future step is to determine the extent to which our results can be recovered in a model with personal types.

## 2. Preview of Results

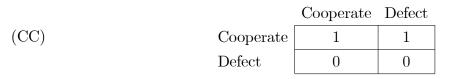
In this section, we preview our main result more carefully on the example analyzed by the Gang of Four: the finitely-repeated prisoner's dilemma game in which player 1 may be committed to tit-for-tat, though this has small ex-ante probability. Consider the repeated game in which the following prisoner's dilemma is repeated  $\bar{t}$  times:

		Cooperate	Defect
(PD)	Cooperate	5, 5	0, 6
	Defect	6,0	1, 1

All the previous moves are publicly observable (perfect monitoring), and the payoff of a player in the repeated game is the sum of his payoffs in the stage game above. A simple application of backward induction in this game yields the play of (Defect, Defect) at every history. Indeed, it is well known that the only Nash equilibrium outcome is playing (Defect, Defect) at every period.

The Gang of Four consider an incomplete information game G in which player 1 may be committed to playing tit-for tat. Player 1 has two types, a "rational" type  $\tau_1^*$ , whose payoffs and available moves are as in the repeated-prisoners' dilemma game above, and a commitment type  $\tau_1^{T4T}$  which can only play tit-for-tat. That is, the latter type must play cooperate in the first round and imitate the last move of player 2 in the subsequent periods. The prior probability of  $\tau_1^{T4T}$  is some small  $\varepsilon > 0$ . Player 2 still has one type  $\tau_2^*$ , which is "rational" as in the original game. The Gang of Four shows that in any sequential equilibrium of the new game, each rational type  $\tau_i^*$  must play Cooperate at all but few periods.

As we mentioned in the introduction, one can replicate the above equilibrium behavior with payoff uncertainty by assigning the payoff function of  $\tau_1^{T4T}$  as 1 at histories at which player 1 plays according to tit-for-tat and 0 at all other histories. Here, the solution concept is sequential equilibrium with the restriction that player 2 assigns probability 1 on  $\tau_1^*$  off the path. Such a payoff function is incompatible with the repeated game payoff structure, and one cannot replicate the commitment to tit-for-tat by simply modifying the stage-game payoff function for  $\tau_1^{T4T}$ . Indeed, such modifications can lead to only two commitment types: the type that plays Cooperate throughout and the type that plays Defect throughout. Commitment to cooperation can be justified by the stage-game payoff function



for example. The inclusion of such simple commitment types cannot affect the behavior of rational types in this game.

However, the austere information structure above is not the only structure we can consider. Our main result (Proposition 1) uses richer type spaces to replicate arbitrary commitment types by payoff types. For any  $\varepsilon' > \varepsilon$ , applying our main result to the game G in the Gang of Four generates a game G' with the following properties.

- **Ex-ante Proximity:** The prior probability of the rational type profile  $(\tau_1^*, \tau_2^*)$  is at least  $1 \varepsilon'$ , and each  $\tau_i^*$  knows that his stage-game payoffs are as in (PD).
- **Repeated-Game Structure:** All types can play all strategies and maximize the sum of stage-game payoffs, which need not be as in (PD).

Strategic Equivalence: G and G' are strategically equivalent in the following sense.

- (1) G' contains types  $\tau_1^*$ ,  $\tau_2^*$ , a twin  $\hat{\tau}_1^{T4T}$  of the tit-for-tat type  $\tau_1^{T4T}$  in G, and a number of other new types (of both players) that we use to encode the beliefs of type  $\hat{\tau}_1^{T4T}$ .
- (2) Though  $\hat{\tau}_1^{T4T}$  is allowed to play any plan of action, tit-for-tat is his unique rationalizable plan.
- (3) The interim beliefs of rational types are equivalent in G and G': rational type  $\tau_1^*$  is certain that he faces the rational type  $\tau_2^*$ , and the rational type  $\tau_2^*$  in turn puts probability  $1 \varepsilon$  on  $\tau_1^*$  and probability  $\varepsilon$  on the twin  $\hat{\tau}_1^{T4T}$  of  $\tau_1^{T4T}$ .

By the strategic equivalence property, the strategic situation the rational types face is the same as in G, except now  $\tau_2^*$  thinks that  $\hat{\tau}_1^{T4T}$  plays tit-for-tat as a result of some rational reasoning under incomplete information rather than as a result of commitment or an unconventional payoff function. Therefore, under the broad set of solution concepts that are invariant to such changes, the solution sets for rational types  $(\tau_1^*, \tau_2^*)$  are identical in Gand G'. The conditional probabilities specified above are achieved by a prior distribution in G' putting probability  $1 - \varepsilon'$  on  $(\tau_1^*, \tau_2^*)$ ,  $\varepsilon (1 - \varepsilon') / (1 - \varepsilon)$  on  $(\hat{\tau}_1^{T4T}, \tau_2^*)$  and the remaining small probability  $(\varepsilon' - \varepsilon) / (1 - \varepsilon)$  on the newly constructed types.

Two points about this construction are worth emphasizing. First, when  $\varepsilon' - \varepsilon$  is small compared to  $\varepsilon$ , the prior probabilities of  $(\tau_1^*, \tau_2^*)$  and  $(\hat{\tau}_1^{T4T}, \tau_2^*)$  are approximately  $1 - \varepsilon$  and  $\varepsilon$ , respectively, with much smaller probability on the new types. Hence, the type spaces of G and G' are nearly identical, and the twin  $\hat{\tau}_1^{T4T}$  assigns much larger probability to the standard type  $\tau_2^*$  than to the new types. Despite this,  $\hat{\tau}_1^{T4T}$  has a unique rationalizable plan because  $\hat{\tau}_1^{T4T}$  believes that his own plan has non-negligible impact on his payoff only if he faces one of the newly constructed types. He finds these types unlikely, but they are likely enough to be his main concern. Second, the unique rationalizable plan emerges under intricate beliefs that require a large number of new types for encoding, especially when the game is long. Nonetheless, we are able to encode such beliefs by putting only a negligible amount of prior probability on the new types.

Our results build on our previous work on non-robustness in repeated games. In Weinstein and Yildiz (2013) we showed that, in any infinitely repeated game, any individually rational and feasible outcome is the *unique* rationalizable outcome of an appropriately chosen perturbation which maintains common knowledge of the repeated-game structure and discounting criterion. A key lemma leading to this result showed that for any plan whatsoever, there is a type who follows this plan as a unique rationalizable action, although he believes in common knowledge of the repeated-game structure. An extension of this lemma to finitely repeated games plays an important role in our construction.

Aside from the obvious differences in motivation and applications, there are two major technical distinctions from our work in Weinstein-Yildiz (2013). First, extending the above lemma from infinitely repeated games to finitely repeated games requires a more difficult construction, as we cannot use future incentives in the last period of a finitely-repeated game. Second, the perturbations allowed here are more constrained. Here, as in the traditional reputation literature, we create a perturbed model which assigns high ex-ante probability to the original model. This is also similar to the perturbations in Kajii-Morris (1997) and other papers on robustness. This *ex ante* notion of perturbation commonly gives very different results from our *interim* framework in Weinstein-Yildiz (2013) and earlier papers, where we allow arbitrary perturbations of interim beliefs in the universal type space.<sup>3</sup> One reason the results here can be achieved with ex-ante perturbations is that our construction centers around perturbing the commitment types, who do not have set beliefs. The main difficulty turns out to be embedding types constructed in the lemma into a common-prior model without affecting the types' rationalizable actions, while keeping the ex-ante probabilities of the new types arbitrarily small.

<sup>&</sup>lt;sup>3</sup>The key difference is that the ex-ante perturbations under a common prior impose additional commonbelief restrictions (Kajii and Morris, 1997), which are crucial in extending the equilibria of the original game to the perturbed one (Monderer and Samet, 1989).

We introduce the basic definitions and formulations in Section 3. In Section 4, we present our construction of a new type space in which the commitment types are replaced by types for which the committed action plan is uniquely rationalizable. In Section 5, we show that, for the original rational types, the constructed game is strategically equivalent to the model with commitment types, under a very broad set of solution concepts. In Section 6, we generalize our result to n-player games in which all players may have commitment types. After presenting our continuity result in Section 7, we offer further remarks on the literature in Section 8 and conclude in Section 9. Some of the more complicated proofs are relegated to the Appendix.

### 3. BASIC DEFINITIONS

We begin with a standard two-player finitely repeated game with perfect monitoring and normal-form stage games; see Section 6 for the *n*-player case. We write  $N = \{1, 2\}$  for the set of players,  $T = \{0, 1, \ldots, \bar{t}\}$  for the set of dates t, and fix a finite set  $A = A_1 \times A_2$ of stage-game action profiles  $a = (a_1, a_2)$ .<sup>4</sup> Note that, since we have perfect monitoring, the non-initial histories in the repeated game are of the form  $h^t = (a^0, \ldots, a^{t-1})$  where  $a^s$ denotes the stage-game strategy profile played at date  $s \in T$ . We write  $h^0$  for the empty initial history, and write H for the set of all non-terminal histories. An outcome path, or terminal history, is a list  $(a^0, \ldots, a^{\bar{t}})$ ; the set of all terminal histories is denoted by Z.

The payoff vector from an outcome path  $(a^0, a^1, \ldots, a^{\bar{t}})$  in a repeated game is simply the sum<sup>5</sup> of the stage game payoffs:

(3.1) 
$$u\left(a^{0}, a^{1}, \dots, a^{\bar{t}}|g\right) = g\left(a^{0}\right) + g\left(a^{1}\right) + \dots + g\left(a^{\bar{t}}\right),$$

where the function  $u = (u_1, u_2)$  denotes the payoffs from the repeated game and the function  $g = (g_1, g_2)$  denotes the payoffs from the stage game. While the particular stage payoffs are not necessarily known, this formula will be common knowledge throughout the games we study here. That is, it is common knowledge that the stage payoff function g is fixed throughout the game and that the players simply maximize the sum of these payoffs.

<sup>&</sup>lt;sup>4</sup>Following the convention in game theory, we write -i for the player  $j \neq i$  and drop the subscript to denote profiles, e.g.,  $x = (x_1, x_2) \in X = X_1 \times X_2$  and  $X_{-1} = X_2$ .

<sup>&</sup>lt;sup>5</sup>Discounting would not affect our results; setting the discount rate to 1 simplifies our derivations.

We write  $\mathcal{G} = [0, 1]^A$  for the set of all possible stage-game payoff functions  $g_i : A \to [0, 1]$ . Here, we put a uniform bound on the stage game payoffs so that small variations of the probability distributions on stage payoffs lead to small variations in expected payoffs, as in the reputation literature. This restriction strengthens our results.

We fix a complete-information repeated game in which it is common knowledge that the stage-game payoffs are a fixed  $(g_1^*, g_2^*)$ . The payoff function in the repeated game is  $u(\cdot|g^*)$ , given by the formula in (3.1). This could, for example, be the repeated prisoner's dilemma game, with  $g^*$  defined as in (PD).

In the complete-information game, a strategy of a player *i* is a mapping  $s_i : H \to A_i$ , which maps each non-terminal history to a strategy in the stage game. Because we analyze incomplete information games, however, we will avoid the word strategy for this mapping and call it instead an *action plan*, reserving the word *strategy* for mappings from types to action plans. (We refer to the strategies in the stage game as *moves*.) The set of all action plans is denoted by  $S_i$ . The outcome path induced by a profile  $(s_1, s_2)$  is denoted by  $z(s_1, s_2)$ . We also allow (behavioral) mixed strategies and write  $\Sigma_i$  for the set of mixed action plans  $\sigma_i : H \to \Delta(A_i)$  for player *i*.

We consider two kinds of elaboration, corresponding to two distinct ways in which the common-knowledge assumption in the complete information game may be relaxed. The first notion of elaboration uses *commitment types*, as is standard in the reputation literature.

**Definition 1.** An  $\varepsilon$ -elaboration with (one-sided) commitment types  $(C, \pi)$  is a Bayesian game such that

- the sets of types for players 1 and 2 are  $\{\tau_1^*\} \cup C$  and  $\{\tau_2^*\}$ , respectively, where  $C \subset S_1$ ;
- Player 2's belief  $\pi$  about player 1's type satisfies  $\pi(\tau_1^*) = 1 \varepsilon$ ;
- the set of plans available to  $\tau_i^*$  is as in the repeated game above, while the only available plan for type  $c \in C$  is c,
- the payoffs are as in the complete information game.

Here, each action plan  $c \in C$  corresponds to a type of Player 1 who can only play c. The incomplete information is only about whether Player 1 can play all action plans or has committed to a particular action plan. The type  $\tau_1^*$  that can play all plans is called the rational type while the types  $c \in C$ , who can play only according to one plan of action, are called *commitment* types. Observe that, since  $C \subset S_1$ , we confine ourselves to pure commitment types.<sup>6</sup>

The second notion of elaboration allows richer type spaces and two-sided incomplete information, but does not allow any payoff function outside of the additive structure in (3.1). Towards stating this notion formally, we define a type space as a list  $(\mathcal{G}, \mathcal{T}, \pi)$  where  $\mathcal{G} \subset \mathcal{G}^*$ is a finite set of payoff function profiles  $g, \mathcal{T} = \mathcal{T}_1 \times \mathcal{T}_2$  is the set of type profiles  $\tau = (\tau_1, \tau_2)$ , and  $\pi \in \Delta(\mathcal{G} \times \mathcal{T})$  is the common prior.<sup>7</sup> We define a *Bayesian repeated game* (without commitment types) as a list  $(N, A, (\mathcal{G}, \mathcal{T}, \pi))$ . We should emphasize that this notation suppresses many important common-knowledge assumptions, such as the fact that the game is repeated, all previous actions are publicly observable (i.e. perfect monitoring), and the payoffs in the repeated game are given by the formula (3.1).

**Definition 2.** An  $\varepsilon$ -elaboration without commitment types of a complete-information game  $g^*$  is a Bayesian repeated game  $(N, A, (\mathcal{G}, \mathcal{T}, \pi))$  with distinguished types  $\tau_1^*, \tau_2^*$  where

(1)  $(g^*, \tau^*) \in \mathcal{G} \times \mathcal{T}$ , (2)  $\pi (g^*, \tau^*) = 1 - \varepsilon$ , and (3)  $\pi (g_i^* | \tau_i^*) = 1$ .

The first and second conditions state that the original complete information game is embedded in the elaboration and has a high ex-ante probability of  $1 - \varepsilon$ . The last condition states that the rational types  $(\tau_1^*, \tau_2^*)$  know their payoffs, and their payoffs are as in the original complete information game. The novelty in this definition is that the elaboration is required to be a Bayesian *repeated* game, i.e., the structure given by the formula (3.1) is common knowledge. In that sense, all the types in an elaboration without commitment types are rational, although we reserve the term rational for types  $(\tau_1^*, \tau_2^*)$  as in the elaborations with commitment types.

<sup>&</sup>lt;sup>6</sup>This is without loss of generality because a belief in a commitment type that plays a mixed strategy in the repeated game is equivalent to a belief in a mixture of pure commitment types (see Section 9 below).

<sup>&</sup>lt;sup>7</sup>Here,  $\Delta(X)$  denotes the set of all probability measures on the finite set X.

Both elaborations above fall under the category of  $\varepsilon$ -elaborations as defined by Kajii and Morris (1997). An  $\varepsilon$ -elaboration without commitment types is a Kajii-Morris elaboration with the additional restriction that the formula (3.1) is common knowledge. While  $\varepsilon$ -elaborations with commitment types were presented above in terms of uncertainty about the set of available strategies, they could also be represented as Kajii-Morris elaborations with a specific simple type space in which the formula (3.1) fails.

Finally, we review a couple of standard concepts in game theory. First, a strategy of a player *i* in a Bayesian repeated game  $(N, A, (\mathcal{G}, \mathcal{T}, \pi))$  is a mapping  $\sigma_i : \mathcal{T}_i \to \Sigma_i$ . Second, interim correlated rationalizability (henceforth ICR) is the outcome of iterated elimination of action plans for types that are never a weak best response, as defined by Dekel, Fudenberg, and Morris (2007). We write  $S_i^{\infty}[\tau_i|G]$  for the set of all interim correlated rationalizable action plans for type  $\tau_i \in \mathcal{T}_i$  in game  $G = (N, A, (\mathcal{G}, \mathcal{T}, \pi))$ . We will give a more detailed definition of ICR later in the construction. We just note here that ICR is the weakest known rationalizability concept for Bayesian games, and all the action plans that are played by a type with positive probability in any equilibrium are ICR for that type.

Third, we say that action plans  $s_i$  and  $s'_i$  are equivalent if  $z(s_i, s_{-i}) = z(s'_i, s_{-i})$  for all action plans  $s_{-i} \in S_{-i}$ , i.e., they lead to the same outcome no matter what strategy the other player plays. Note that  $s_i$  and  $s'_i$  are equivalent iff  $s_i(h^t) = s'_i(h^t)$  for every history  $h^t$  in which *i* played according to  $s_i$  throughout; they may differ only in their prescriptions for histories that they preclude. Hence, in reduced form, action plans can be represented as mappings from the history of other players' play to own stage game actions. We write  $\bar{S}_i$  for the set of reduced-form action plans  $\bar{s}_i$ ; these map each  $(a^l_{-i})_{0 \leq l < t}$  to some action  $a_i \in A_i$ in the stage game. Finally, we introduce the following notation for sets of equivalent action plans: Given any two sets X, Y of action plans, we write  $X \simeq Y$  if for every  $x \in X$  there exists  $y \in Y$  that is equivalent to x, and for every  $y \in Y$  there exists  $x \in X$  that is equivalent to y. In particular,  $X \simeq \{x_i\}$  means that X consists only of strategies that are equivalent to  $x_i$ .

### 4. IRRELEVANCE OF COMMITMENT TYPES

In this section, we state and outline the proof of our main result: any  $\varepsilon$ -elaboration with commitment types can be transformed, for any  $\varepsilon' > \varepsilon$ , into an  $\varepsilon'$ -elaboration without

commitment types, where each commitment type c is replaced with a payoff type  $\tau_1^c$  for which c is uniquely rationalizable. These payoff types follow c not because they are committed or have payoffs that are inconsistent with playing a repeated game but because their reasoning under their information leads them to do so. Moreover, from the point of view of the rational types these are the only types with positive probability, mirroring the elaboration with commitment types. From the point of view of rational types who believe in the ICR concept, the two elaborations are identical. Hence, under ICR (as well as a broader set of solution concepts), the set of solutions for each rational type is identical in the two elaborations.

**Proposition 1.** For any  $\varepsilon, \varepsilon' \in (0,1)$  with  $\varepsilon' > \varepsilon$  and for any  $\varepsilon$ -elaboration G with commitment types  $(C,\pi)$  there exists an  $\varepsilon'$ -elaboration  $G' = (N, A, (\mathcal{G}, \mathcal{T}, \pi'))$  without commitment types in which the commitment types are replaced by types with unique rationalizable action plans, meaning:

- (1)  $\pi'(g^*, \tau_2^* | \tau_1^*) = 1$  and  $\pi'(g^*, \tau_1^* | \tau_2^*) = \pi(g^*, \tau_1^* | \tau_2^*) = 1 \varepsilon$ , and
- (2) for every  $c \in C$  there exists  $\tau_1^c \in \mathcal{T}_1$  such that

$$S_i^{\infty}[\tau_1^c] \simeq \{c\},$$

and  $\pi'(\tau_1^c|\tau_2^*) = \pi(c|\tau_2^*) = \pi(c)$ .

Here, the first condition states that the interim beliefs of rational types regarding their own payoffs and "rationality" of their opponents are identical in the two elaborations. The second condition states that each commitment type c is replaced by a type  $\tau_1^c$  for which following c is uniquely rationalizable in reduced form (i.e.  $S_i^{\infty}[\tau_1^c] \simeq \{c\}$ ), and that the rational type of player 2 in G' assigns the same probability to the type  $\tau_1^c$  as the rational type in G assigns to the commitment type c (i.e.  $\pi'(\tau_1^c | \tau_2^*) \equiv \sum_{g_1} \pi'((g_1, g_2^*), \tau_1^c | \tau_2^*) = \pi(c | \tau_2^*) = \pi(c)$ ).

The equivalence of G' with G is established despite the following constraints:

(1) The repeated-game payoff structure is maintained throughout G'. That is, it is common knowledge throughout that the payoff in the repeated game is the sum of the payoffs in the stage game, and that the stage game is fixed throughout the game. Type  $\tau_1^c$  knows all this and yet follows c as its unique rationalizable plan.

(2) The ex-ante distribution  $\pi'$  in G' can be arbitrarily close to the distribution  $\pi$  in G, in that  $\varepsilon'$  can be arbitrarily close to  $\varepsilon$ .

Proof of Proposition 1. The first step in our construction is the following lemma, which establishes that any given action plan  $s_i$  is the only rationalizable action for a type  $\tau_i^{s_i}$  from some common prior model. (The proof of Lemma 1 is the lengthiest step of the proposition and is relegated to the appendix; we provide a detailed intuition for it later in this section.)

**Lemma 1.** For any  $s_i \in S_i$ , there exists a Bayesian repeated game  $G^{s_i} = (N, A, (\mathcal{G}^{s_i}, \mathcal{T}^{s_i}, \pi^{s_i}))$ with a type  $\tau_i^{s_i} \in \mathcal{T}_i^{s_i}$  such that

- (1)  $\pi^{s_i}(g,\tau) > 0$  for every  $(g,\tau) \in \mathcal{G}^{s_i} \times \mathcal{T}^{s_i}$  and
- (2)  $S_i^{\infty}[\tau_i^{s_i}|G^{s_i}] \simeq \{s_i\}.$

By relabeling if necessary, we take all of the types in the type spaces  $\mathcal{T}^{s_i}$  above to be distinct from each other and from  $\tau^*$ , fixing also a unique type  $\tau_i^{s_i}$  for each  $s_i$ . We construct  $G' = (N, A, (\mathcal{G}', \mathcal{T}', \pi'))$  by setting

$$\begin{split} \mathcal{G}' &= \{g^*, (\mathbf{0}, g_2^*)\} \cup \bigcup_{c \in C} \mathcal{G}^c \\ \mathcal{T}'_i &= \{\tau_i^*\} \cup \bigcup_{c \in C} \mathcal{T}_i^c \qquad (\forall i \in N) \\ \pi'\left(g, \tau\right) &= \begin{cases} 1 - \varepsilon' & \text{if } (g, \tau) = (g^*, \tau^*) \,, \\ \frac{1 - \varepsilon'}{1 - \varepsilon} \pi\left(c\right) & \text{if } (g, \tau) = \left((\mathbf{0}, g_2^*) \,, (\tau_1^c, \tau_2^*)\right) \,, \\ \frac{\varepsilon' - \varepsilon}{(1 - \varepsilon)|C|} \pi^c\left(g, \tau\right) & \text{if } (g, \tau) \in \mathcal{G}^c \times \mathcal{T}^c, \\ 0 & \text{otherwise,} \end{cases} \end{split}$$

where

$$\mathbf{0}\left(a\right) = 0 \qquad \left(\forall a \in A\right).$$

We now observe that G' satisfies the properties in the proposition. Indeed, rational type  $\tau_1^*$  of Player 1 assigns probability 1 on  $(g^*, \tau_2^*)$ . Likewise, we have

$$\pi'\left(\mathcal{G}'\times\{\tau_2^*\}\right) = 1 - \varepsilon' + \frac{1-\varepsilon'}{1-\varepsilon}\sum_{c\in C}\pi(c) = 1 - \varepsilon' + \frac{1-\varepsilon'}{1-\varepsilon}\varepsilon = \left(1-\varepsilon'\right)/\left(1-\varepsilon\right)$$

and therefore, in the interim,  $\tau_2^*$  assigns probability  $1 - \varepsilon$  to  $(g^*, \tau_1^*)$  and probability  $\pi(c)$  to  $\tau_1^c$  for each c. On the other hand, since the beliefs of type  $\tau_1^c$  altered substantially when  $(\mathcal{G}^c, \mathcal{T}^c, \pi^c)$  was incorporated in G', it is not clear that  $\tau_1^c$  follows c as the unique ICR action. The next lemma states that this is indeed the case.

**Lemma 2.** For any  $c \in C$ ,  $i \in N$ , and any  $\tau_i \in T_i^c$ ,  $S_i^{\infty}[\tau_i|G'] = S_i^{\infty}[\tau_i|G^c]$ ; in particular,  $S_1^{\infty}[\tau_1^c|G'] = c$ .

This lemma completes the proof of the proposition; its proof is in the appendix.  $\Box$ 

Our proof has two main steps. The first, found in Lemma 1, is to construct a type space in which a given action plan is uniquely rationalizable for a type. We constructed such a type space in Weinstein and Yildiz (2013) for infinite-horizon repeated games, but without requiring that the constructed type space have a common prior, a property that is essential for our proposition here. In this paper, using the ideas in that construction, we first construct such a type space for finite-horizon games without common prior and then convert it to a common-prior type space, using the ideas and the results developed by Lipman (2003) and Weinstein and Yildiz (2007).

The main economic ideas involved in these constructions come from social learning and reward/punishment mechanisms in repeated games. Our first construction involves types who know their stage-game payoff is a function of their own action alone, but do not initially know their optimal action. They will learn their optimal action from the actions of others. Only some plans are consistent with such beliefs; for instance, no such player could play move  $a^0$  in period 0 and the same move  $a^1 \neq a^0$  in all continuations. For infinite-horizon games, we extended the result to all action plans (including plans that contradict the condition for individual learning), using a reward and punishment mechanism, in Weinstein and Yildiz (2013). It is harder to come up with effective reward and punishment mechanisms for finite-horizon games. After all, one cannot provide any future incentive in the last period. Hence, here, we use a more nuanced construction that combines social learning with a reward and punishment mechanism to extend the result to all action plans in finitely repeated games. In our construction, the player's stage payoffs are additively separable in his action and others' actions. In all periods before the last, his incentives are dominated by the desire to

be rewarded by the other players, while in the last period he has learned his own optimal action and acts accordingly.

The second main step is to incorporate the above type spaces in one common prior model, replacing each commitment type with one of these type spaces. One must do this in such a way that (i) the original complete-information game still has high prior probability  $(1 - \varepsilon')$ , (ii) the interim beliefs of the rational types are as in the original elaboration with commitment types, and (iii) the types' rationalizable behavior in the constructed type space remain the same after incorporating them into common prior model. The conditions (ii) and (iii) oppose each other, making the construction more difficult. To see this, note that (i) and (ii) require that the common prior  $\pi'$  puts a high probability on  $\tau_1^c$ , requiring that probability to be  $\frac{1-\varepsilon'}{1-\varepsilon}\pi(c)$  as in our proof. When  $\varepsilon$  and  $\varepsilon'$  are close, this probability is approximately  $\pi(c)$ . When  $\varepsilon$  and  $\varepsilon'$  are close, this also requires that  $\pi'$  puts a very small probability on  $\mathcal{T}^c$ , the original type profiles in the constructed type space in the first step. That probability can be at most  $(\varepsilon' - \varepsilon) / (1 - \varepsilon)$ , which is negligible with respect to  $\frac{1-\varepsilon'}{1-\varepsilon}\pi(c)$  when  $\varepsilon$  and  $\varepsilon'$  are close. These constraints make the belief of type  $\tau_1^c$  in game G' substantially different from the belief of the type  $\tau_1^c$  in game  $\mathcal{G}^c$ . In our construction, type  $\tau_1^c$  in game G' assigns probability

(4.1) 
$$p^{c} = \frac{|C| (1 - \varepsilon') \pi (c)}{|C| (1 - \varepsilon') \pi (c) + (\varepsilon' - \varepsilon) \pi^{c} (\tau)}$$

on type  $\tau_2^*$ . Note that, for fixed  $\pi(c)$ , when  $\varepsilon' - \varepsilon$  approaches 0,  $p^c$  approaches 1.<sup>8</sup> In contrast,  $\tau_1^c$  in game  $G^c$  assigns zero probability on  $\tau_2^*$ . Consequently, the belief hierarchies of the types in G' can be quite different from the belief hierarchies of the types in  $G^c$  with the same label, which could lead to distinct set of ICR actions. We circumvent this problem with the following trick. We set the beliefs such that, whenever player 2 has type  $\tau_2^*$ , the payoff of type  $\tau_1^c$  is 0 for every move in the stage game, making him indifferent among all outcomes. Since  $p^c < 1$ , his best responses are identical to his best responses conditional on the type of player 2 being other than  $\tau_2^*$ , thereby replicating the best responses of his twin in  $G^c$ . Since this was the only difference between the two type spaces, the rationalizable actions turn out to be identical in games  $G^c$  and G', as shown formally by Lemma 2.

<sup>&</sup>lt;sup>8</sup>Note also that the technique we use in transforming the model without common-prior to the one with common prior also renders  $\pi^c(\tau_1^c)$  small, bringing  $p^c$  near 1 even when  $\varepsilon$  and  $\varepsilon'$  are far apart.

Roughly speaking, from the point of view of rational types, Proposition 1 replaces commitment types by types who follow the same plans as their unique rationalizable plan. Hence, under any rationalizable solution concept, the rational types face the same strategic uncertainty in both games leading the same set of possible behavior. We will next establish such strategic equivalence formally.

# 5. STRATEGIC EQUIVALENCE

In this section, we show that the elaborations G with commitment types and G' without commitment types described in Proposition 1 are "strategically equivalent" for rational types. By this we mean that, for a broad set of solution concepts, the set of solutions for each rational type are identical in games G and G'. Therefore, the same set of behavior can be supported by reputational models regardless of whether one allows commitment types. In other words, the same set of behavior is supported whether one allows payoff functions that are inconsistent with the repeated-game structure or imposes this structure throughout.

Our result here applies to any solution concept that is invariant to replacing commitment types with types that have unique rationalizable action plans (in reduced form). In general Bayesian games, this invariance condition is somewhat stronger than elimination of nonrationalizable strategies, because the new game contains some new types, encoding the beliefs of the types with unique rationalizable plans. We first establish our result for a general class of such invariant solution concepts. We also establish the same strategic equivalence for sequential equilibrium; this requires an additional off-path belief restriction commonly imposed in the reputation literature.

5.1. Strategic Equivalence under Invariant Solutions. The following definitions are standard: A solution concept  $\Sigma$  maps every Bayesian game G to a set  $\Sigma(G)$  of mixed strategies in game G. For any type spaces  $\mathcal{T}$  and  $\mathcal{T}'$  with  $\mathcal{T} \subset \mathcal{T}'$  and any strategy profile  $\sigma$  on  $\mathcal{T}'$ ,  $\sigma_{\mathcal{T}}$  denotes the restriction of  $\sigma$  to  $\mathcal{T}$ . In the following definitions, we also use the convention that two probability distributions that have common support and agree on this support are identical, ignoring any difference in domains. **Definition 3.** A solution concept  $\Sigma$  is said to be *invariant to elimination of non-rationalizable* strategies if

$$\Sigma\left(G\right) = \Sigma\left(G'\right)$$

for any two games G and G' with identical type spaces satisfying (i) if an action plan  $s_i$  is available for a type  $\tau_i$  in game G then  $s_i$  is available for  $\tau_i$  in G' and (ii) if  $s_i$  is not available for  $\tau_i$  in G then  $s_i \notin S_i^{\infty}[\tau_i|G']$ .

**Definition 4.** A solution concept  $\Sigma$  is said to be *invariant to trivial enrichments of the type* spaces if

$$\Sigma(G) = \{\sigma_{\mathcal{T}} | \sigma \in \Sigma(G')\}$$

for any two games G and G' with type spaces  $\mathcal{T}$  and  $\mathcal{T}'$  such that (i)  $\mathcal{T} \subset \mathcal{T}'$ , (ii) every type in  $\mathcal{T}$  has identical set of available action plans in games G and G', and (iii) any type in  $\mathcal{T}$ with multiple action plans has identical interim beliefs in games G and G'.

Note that the transformation in the first definition allows only elimination of non-rationalizable actions and the transformation in the second definition allows only inclusion of new types such that the types who put positive probability to the new types are trivial in that they can play only according to one plan. Proposition 1 implies that under any solution concept that is invariant to the above transformations, elaborations with or without commitment types have the same strategic implications for rational types. Due to its importance, we state this corollary as a proposition:

**Proposition 2.** Let  $\Sigma$  be a solution concept that is invariant to elimination of non-rationalizable strategies and to trivial enrichment of the type spaces. Then, for any  $\varepsilon, \varepsilon' \in (0, 1)$  with  $\varepsilon' > \varepsilon$  and for any  $\varepsilon$ -elaboration G with commitment types, there exists an  $\varepsilon'$ -elaboration  $G' = (N, A, (\mathcal{G}, \mathcal{T}, \pi'))$  without commitment types such that

$$\left\{\sigma\left(\tau^{*}\right)|\sigma\in\Sigma\left(G\right)\right\}=\left\{\sigma\left(\tau^{*}\right)|\sigma\in\Sigma\left(G'\right)\right\},$$

i.e., the set of solutions for rational types are identical in games G and G'.

*Proof.* Note that, in Proposition 1, the elaboration G' can be obtained from G by (1) introducing new types such that only committed types believe in the new types, and (2) allowing commitment types to play any action plan in the repeated game. The first step is a trivial enrichment as in Definition 3 and the second undoes an elimination covered by Definition 4, so the conclusion follows.  $\hfill \Box$ 

5.2. Strategic Equivalence under Sequential Equilibrium. We will next establish the same strategic equivalence under sequential equilibrium, which is defined as follows. Given any Bayesian repeated game with a type space  $(\mathcal{G}, \mathcal{T}, \pi)$ , a belief structure is a list  $\mu = (\mu_{i,\tau_{i},h})_{i \in N, \tau_{i} \in \mathcal{T}_{i}, h \in H}$  of type-specific beliefs  $\mu_{i,\tau_{i},h} \in \Delta (\mathcal{G} \times \mathcal{T}_{-i})$  regarding the underlying payoffs and the other player's types, beliefs that vary with the history of play.<sup>9</sup> An assessment is a pair  $(\tilde{\sigma}, \mu)$  of strategy profile  $\tilde{\sigma} : \mathcal{T} \to \Sigma$  and a belief structure  $\mu$ . An assessment  $(\tilde{\sigma}, \mu)$  is said to be sequentially rational if  $\tilde{\sigma}_i (\cdot | \tau_i)$  is a sequential best response to  $\mu_{i,\tau_{i},h}$  and  $\tilde{\sigma}_{-i}$ , i.e., the restriction of  $\tilde{\sigma}_i (\cdot | \tau_i)$  to the continuation game after every history h is a best response to  $\tilde{\sigma}_{-i}$  and the beliefs  $\mu_{i,\tau_{i},h}$  in the continuation game. An assessment  $(\tilde{\sigma}, \mu)$  is said to be consistent if there exists a sequence  $(\tilde{\sigma}^n, \mu^n) \to (\tilde{\sigma}, \mu)$  such that  $\tilde{\sigma}^n$  assigns positive probability to each available move at every history and  $\mu^n$  is derived from Bayes' rule and  $\tilde{\sigma}^n$ . An assessment  $(\tilde{\sigma}, \mu)$  is said to be a sequential equilibrium if it is sequentially rational and consistent.

In an  $\varepsilon$ -elaboration without commitment types, sequential equilibria are defined as above. In an  $\varepsilon$ -elaboration with commitment types, the definition of course depends on how one formalizes the commitment types. In particular, the definition above implies that Player 2 puts probability 1 on the rational type of Player 1 if the history is not consistent with any commitment type—even when the history is also inconsistent with the strategy of the rational type. This is because the commitment types have only one action, so that only the rational types may tremble. This is an additional assumption when the commitment types are represented by payoff perturbations (violating the additive repeated game structure). In general, the possible off-the-path beliefs can vary depending on the way the commitment types are formulated, but the above assumption is usually maintained. We will keep this additional assumption in our definition for sequential equilibrium without commitment types:

Assumption 1. For every history  $h = (a^0, \ldots, a^{t-1}),$ 

$$\mu_{2,\tau_2^*,h}\left(g^*,\tau_1^*\right) = 1$$

<sup>&</sup>lt;sup>9</sup>A more general definition of a belief structure would also specify the beliefs regarding past actions, but those beliefs are trivial because of perfect monitoring.

whenever h has zero probability under every type  $\tau_1 \neq \tau_1^*$ .

In our analysis we will focus on the behavior of the rational types under sequential equilibrium, which is formally defined as follows.

**Definition 5.** For any elaboration G (with or without commitment types), we write

 $SE^*(G) = \{\sigma(\cdot | \tau^*) | (\sigma, \mu) \text{ is a sequential equilibrium of } G \text{ that satisfies Assumption 1} \}$ for the set of sequential equilibrium action plans for the rational types in G.

We are now ready to state the strategic equivalence result for sequential equilibrium.

**Proposition 3.** For any  $\varepsilon, \varepsilon' \in (0,1)$  with  $\varepsilon' > \varepsilon$  and for any  $\varepsilon$ -elaboration G with commitment types  $(C,\pi)$  there exists an  $\varepsilon'$ -elaboration  $G' = (N, A, (\mathcal{G}, \mathcal{T}, \pi'))$  without commitment types such that

$$SE^{*}\left(G\right) = SE^{*}\left(G'\right),$$

i.e., under Assumption 1, the set of sequential equilibrium action plans for the rational types is same in games G and G'.

*Proof.* We take G' as in Proposition 1.We will show that both conditions  $\sigma(\cdot|\tau^*) \in SE^*(G)$ and  $\sigma(\cdot|\tau^*) \in SE^*(G')$  are characterized by the following conditions, (SR1) and (SR2). First,  $(\sigma, \mu)$  is a sequential equilibrium of G if and only if the following three conditions are satisfied. The consistency condition for  $\tau_2^*$  is

(C) 
$$\mu_{2,\tau_{2}^{*},h}(c) = \mu_{h}^{\sigma\left(\cdot|\tau_{1}^{*}\right)}(c) \equiv \begin{cases} \frac{\pi(c)}{\Pr\left(h|\sigma\left(\cdot|\tau_{1}^{*}\right)\right)(1-\varepsilon)+\sum\limits_{c'\in C^{h}}\pi(c')} & \text{if } c\in C^{h} \\ 0 & \text{otherwise} \end{cases} \quad (\forall h, c)$$

where  $C^h$  is the set of commitment plans  $c \in C$  that is consistent with history h. Of course,  $\mu_h^{\sigma(\cdot|\tau_1^*)}(\tau_1^*) = 1 - \sum_{c \in C} \mu_h^{\sigma(\cdot|\tau_1^*)}(c)$ . The consistency condition for player 1 is trivial, as player 2 has only one type. Note that  $\mu_h^{\sigma(\cdot|\tau_1^*)}$  is a function of  $\sigma(\cdot|\tau_1^*)$ , and hence the following sequential rationality conditions are solely on  $\sigma(\cdot|\tau^*)$ . The sequential rationality conditions are

(SR1):  $\sigma(\cdot|\tau_1^*)$  is a sequential best response to  $\sigma(\cdot|\tau_2^*)$  under  $g_1^*$ , and

(SR2): at each history h,  $\sigma(\cdot|\tau_2^*)$  is conditional best response to the mixed strategy

$$\tilde{\sigma} \equiv \mu_h^{\sigma\left(\cdot|\tau_1^*\right)}\left(\tau_1^*\right)\sigma\left(\cdot|\tau_1^*\right) + \sum_{c \in C} \mu_h^{\sigma\left(\cdot|\tau_1^*\right)}\left(c\right)c$$

under  $g_2^*$ .

Since all the other types are committed to a single plan, there are no other conditions. This shows that  $\sigma(\cdot|\tau^*) \in SE^*(G)$  if and only if (SR1) and (SR2) are satisfied.

To show that  $\sigma(\cdot|\tau^*) \in SE^*(G')$  implies the conditions (SR1) and (SR2), consider any sequential equilibrium  $(\sigma, \mu')$  of G' that satisfies Assumption 1. Firstly, since type  $\tau_1^*$  puts probability one on  $(g^*, \tau_2^*)$ , the sequential rationality condition for that type is (SR1). Secondly, since c is the unique rationalizable action plan of  $\tau_1^c$  in G' (by Lemma 2) on all histories h consistent with c,

(5.1) 
$$\sigma\left(c\left(h\right)|h,\tau_{1}^{c}\right) = 1 \qquad \left(\forall c \in C^{h},\forall h\right).$$

Hence, by Assumption 1 and consistency,

(5.2) 
$$\mu'_{2,\tau_{2}^{*},h}(\tau_{1}^{c}) = \mu_{h}^{\sigma(\cdot|\tau_{1}^{*})}(c) \qquad (\forall h,c) ,$$

which of course also implies that  $\mu'_{2,\tau_2^*,h}(\tau_1^*) = \mu_h^{\sigma(\cdot|\tau_1^*)}(\tau_1^*)$ . By (5.1) and (5.2), under the belief of type  $\tau_2^*$ , player 1 plays according to  $\tilde{\sigma}$  above, and the sequential rationality condition for type  $\tau_2^*$  is (SR2).

To show that (SR1) and (SR2) are sufficient for  $\sigma(\cdot|\tau^*) \in SE^*(G')$ , take any  $\sigma(\cdot|\tau^*)$  that satisfies (SR1) and (SR2). We will construct a sequential equilibrium  $(\sigma, \mu')$  of G' that satisfies Assumption 1. Set  $\mu'_{1,\tau_1^*,h}(g^*, \tau_2^*) = 1$  and  $\mu'_{2,\tau_2^*,h} = \mu_h^{\sigma(\cdot|\tau_1^*)}$ . For each  $c \in C$ , consider a sequential equilibrium  $(\sigma^c, \mu^c)$  of the game in which the action plan of type  $\tau_2^*$  is fixed as  $\sigma(\cdot|\tau_2^*)$ —as moves of nature, and the type space is  $\mathcal{T}^c$  with the interim beliefs in G'. Set  $\sigma(\cdot|\tau_i) = \sigma^c(\cdot|\tau_i)$  and  $\mu'_{i,\tau_i,h} \equiv \mu_{i,\tau_i,h}^c$  for every  $\tau_i \in \mathcal{T}_i^c$  and  $c \in C$ . We now show that  $(\sigma, \mu')$  is a sequential equilibrium of G' and satisfies Assumption 1. Since Lemma 2 applies to the case  $g_2^* = 0$ , in which case  $\sigma(\cdot|\tau_2^*)$  is rationalizable for type  $\tau_2^*$ ,

(5.3) 
$$\sigma^{c}\left(c\left(h\right)|h,\tau_{1}^{c}\right) = 1 \qquad \left(\forall c \in C^{h},\forall h\right).$$

Hence,  $\mu'_{2,\tau_2^*,h}$  is consistent and satisfies Assumption 1. The sequential rationality conditions for rational types are (SR1) and (SR2) by construction and (5.3). The sequential rationality

and consistency for types in  $\mathcal{T}^c$  immediately follows from the construction and the fact that  $(\sigma^c, \mu^c)$  is a sequential equilibrium in the auxiliary game.

The strategic equivalence under sequential equilibrium is somewhat subtle, requiring the lengthy proof above. This is because of the issues relating to the off-the-path beliefs, which play a central role in sequential equilibrium while not being relevant for ICR. If a type  $\tau_1^c$ , who plans to follow c, deviates from c, then his subsequent behavior may be different from c as ICR cannot restrict the behavior at the contingencies that are precluded by one's own strategy. In that case, off the path beliefs of player 2 at the histories that are not consistent with any type could be different. Moreover, consistency may result in unforeseen restrictions on those beliefs as it is applied for types in  $\mathcal{T}^c$  and  $\tau_2^*$  simultaneously. Assumption 1 ensures that Player 2 assigns zero probability to  $\tau_1^c$  whenever Player 1 deviates from c, resulting in beliefs that are identical to those with commitment types, as we show in the proof. Of course, at the histories that are consistent with commitment types, the rational types in the games G and G' face the same uncertainty regarding all relevant aspects, such as whether the other player is rational and which  $c \in C$  he is playing if he is not rational. This leads to the same set of solutions for rational types in both games.

**Remark 1.** The strategic equivalence above implies that the testable predictions with or without commitment types are nearly indistinguishable. Imagine that an empirical or experimental researcher observes outcomes of games that essentially look like a fixed repeated game, as in  $g^*$ , but she does not know the players' beliefs about possible commitments or payoff variations. Using the data, she can obtain an empirical distribution on outcome paths—with some noise. The above strategic equivalence implies that the equilibrium distributions for elaborations with or without commitment types can be arbitrarily close, making it impossible to rule out one model without ruling out the other given the sampling noise (see our online appendix for a formal result along these lines).

# 6. General Case

In this section, we will present the result for the *n*-player case, allowing commitment types for all players. The definitions for the *n*-player case mirror the case of n = 2, and we will not repeat them here. Since we will allow commitment types for all players, an

 $\varepsilon$ -elaboration with commitment types is now defined as a Bayesian game, with common prior  $\pi$ , such that the set of types for each player i is  $\{\tau_i^*\} \cup C_i$  where  $C_i \subset S_i$  can be empty, type  $\tau_i$  can play any action plan while a type  $c_i \in C_i$  can play only  $c_i$ , and the probability  $\pi$  ( $\tau^*$ ) of the rational type profile is  $1 - \varepsilon$ . Note that when  $\varepsilon < 1$ , there some  $C_i$ is non-empty. Note also that the distribution of commitment type is not restricted; they can be correlated for example. Such a Bayesian game can be denoted by  $(C_1, \ldots, C_n, \pi)$  where  $\pi \in \Delta ((\{\tau_1\} \cup C_1) \times \cdots \times (\{\tau_n\} \cup C_n))$  is the prior on the type profiles. Finally, we write  $\Sigma^*$ for the set of solution concepts that are (1) invariant to the elimination of non-rationalizable plans, (2) invariant to trivial enrichments of the type spaces, and (3) include all solutions generated by the sequential equilibria that satisfy Assumption 1. The result is generalized to this case as follows.

**Proposition 4.** For any  $\varepsilon, \varepsilon' \in (0, 1)$  with  $\varepsilon' > \varepsilon$  and for any  $\varepsilon$ -elaboration G with commitment types  $(C_1, \ldots, C_n, \pi)$  there exists a strategically-equivalent  $\varepsilon'$ -elaboration  $G' = (N, A, (\mathcal{G}, \mathcal{T}, \pi'))$ without commitment types in which the commitment types are replaced by types with unique rationalizable action plans:

- (1) for every  $i \in N$ ,  $\pi'(g^*, \tau_{-i}^* | \tau_i^*) = \pi(\tau_{-i}^* | \tau_i^*);$
- (2) for every  $i \in N$  and  $c_i \in C_i$ , there exists  $\tau_i^{c_i} \in \mathcal{T}_i$  such that all ICR action plans of  $\tau_i^{c_i}$  are equivalent to  $c_i$ , and  $\pi' \left( \tau_i^{c_i} | \tau_j^* \right) = \pi \left( c_i | \tau_j^* \right)$  for every  $j \neq i$ ;
- (3) for every  $\Sigma \in \Sigma^*$ ,

$$\{\sigma(\tau^*) | \sigma \in \Sigma(G)\} = \{\sigma(\tau^*) | \sigma \in \Sigma(G')\}.$$

The first two conditions all together state that each commitment type is replaced by a type that follows the committed action profile as his uniquely rationalizable plan, and the interim beliefs of the rational types remain intact under rationalizability. The last condition states that the two games are strategically equivalent for rational types under any invariant solution concept, including sequential equilibria that puts probability one on rational types off the path. An outline of the proof for this result can be found in the appendix.

# 7. Necessity of Commitment under CK of Approximate Payoffs

In the previous sections, while we imposed the constraint that it is always common knowledge that the payoffs are the sum of identical stage-game payoffs, we allowed those payoffs to lie anywhere in the interval [0, 1]. In this section, by contrast, we make the stricter requirement that it is common knowledge that payoffs lie within  $\varepsilon$  of those in the completeinformation game. Under this stricter requirement, we show that commitment types are not dispensable in reputation models. When the stage game is dominance solvable, there is a unique sequential Nash equilibrium outcome, in which the unique rationalizable strategy profile of the stage game is played throughout. Here,  $\varepsilon$  is uniform over all type spaces and the number of repetitions. For example, in the repeated prisoners' dilemma, one cannot have any cooperation without commitment types when it is common knowledge that the payoffs are approximately those in the prisoner's dilemma.

Define the distance between two stage-game payoff functions via the sup norm:

$$d(g',g) = \max_{a} |g'(a) - g(a)|$$

**Proposition 5.** Fix a complete information stage game  $g^*$  which has unique rationalizable profile  $a^*$ . Then, there exists  $\varepsilon > 0$  such that for any  $\varepsilon' > 0$  and any  $\overline{t}$ , every  $\varepsilon'$ elaboration  $(N, A, (\mathcal{G}, \mathcal{T}, \pi))$  without commitment types, satisfying the additional requirement that  $d(g, g^*) < \varepsilon$  for all  $g \in \mathcal{G}$ , has a unique sequential equilibrium in which  $a^*$  is played by all types at all histories.

*Proof.* The elimination process for the finite stage game  $g^*$  is finite. Each time an action is eliminated (again by finiteness) it must be that for some  $\delta > 0$  it is never within  $\delta$  of being a best reply. Choose  $\varepsilon > 0$  so that  $2\varepsilon$  is smaller than the minimum of these  $\delta$ .

Now suppose there is a sequential equilibrium strategy profile  $s^*$  which contradicts the result. Consider one of the latest histories at which any violation of the profile  $a^*$  occurs, and of the violations at this history, consider an action  $a'_i$  which is eliminated first in the elimination process for  $g^*$ , say at stage k. When player i takes this action, he must believe that (a) the profile  $a^*$  is played at all future dates regardless of his action and (b) no action eliminated at stage k - 1 or earlier is played at the current history. But then by (b), the fact that  $a'_i$  is eliminated at stage k, and the choice of  $\varepsilon$ , his action is suboptimal in the stage game; and by (a) his action cannot affect future play. This contradicts the concept of sequential equilibrium.

For example, in a repeated prisoners' dilemma game, if it is common knowledge that payoffs are close to the prisoners' dilemma, then in any sequential equilibrium the players defect throughout the game regardless of the number of repetitions. At some level this is a reflection of general continuity properties of Bayesian Nash equilibrium payoffs with respect to the perturbations of payoffs. Indeed, it is well known that, for any given  $\bar{t}$ , as  $\varepsilon \to 0$ , the Bayesian equilibrium payoffs in  $\varepsilon$ -elaborations of repeated prisoners' dilemma with commitment types approach the payoffs from defection throughout the game. This is in line with the continuity results for Nash equilibrium payoffs with respect to the prior distributions. Hence, for a given  $\bar{t}$ , our result here differs from the existing continuity results only in terms of the perturbations it considers, making the stage payoffs approach to the original game instead of making the probability of types with unrelated payoffs to go to zero. Our result has a major strength however:  $\varepsilon$  is uniform with respect to the number of repetitions. In contrast, for any  $\varepsilon$  probability of a tit-for-tat type, cooperation prevails whenever the number of repetitions are sufficiently large, as famously established by the Gang of Four.

### 8. Remarks

Continuity and Robustness of Equivalence. Since interim correlated rationalizability is upperhemicontinuous (Dekel, Fudenberg, and Morris, 2007), each type  $\tau_i^{c_i}$  with unique rationalizable action  $c_i$  has the same unique rationalizable action on a open neighborhood of parameters and beliefs. Hence, in the elaboration constructed in Proposition 1, we can perturb parameters such as the stage-game payoff functions and beliefs for the newly constructed types as long as the beliefs of rational types are fixed. So, relative to the set of elaborations with the same set of types, where rational players have fixed beliefs, we obtain an open set of  $\varepsilon'$ -elaborations G' without commitment types that are strategically equivalent to G. In particular, type  $\tau_i^{c_i}$  need not be exactly indifferent between his actions conditional on meeting a rational type; this was only a simplifying aspect in our construction.

On the other hand, our result is silent about continuity with respect to variation of the beliefs of the rational types. Such continuity is directly tied to the continuity properties of the solution concept in the original game G, by our strategic equivalence result. Of course, since  $\varepsilon'$  must be larger than  $\varepsilon$  (albeit being arbitrarily close), our result and the reputation

result that it is applied to are relevant only when the solution concept on the original model G is continuous with respect to small variations in  $\varepsilon$  when the commitment types and their relative probabilities with respect to each other are fixed. This is indeed the case for most existing models.<sup>10</sup>

Sensitivity to the Set of Commitment Types. Despite the continuity in the previous paragraph, the equilibrium predictions of reputational models are highly sensitive to the set of commitment types one considers: by varying the set of commitment types one can obtain a rich set of behavior as equilibrium outcomes in long but finitely-repeated games. Indeed, Fudenberg and Maskin (1986) obtain a Folk Theorem in this way. Once again, such sensitivity to the set of commitment types will be inherited by our newly constructed reputation models without commitment types, due to strategic equivalence.

**Short-Lived Players**. The above sensitivity is muted when the uninformed player is *short-lived* (i.e. she myopically best-responds to her belief about the other player's move at every history). In that case, in any Bayesian Nash equilibrium, the payoff of the rational player with commitment types is near his Stackelberg payoff, provided that he is sufficiently patient and has a type that always plays his Stackelberg move (Fudenberg and Levine (1989)). Since our players are all *long-lived*, such an independence result does not hold in the reputation models we consider here. For example, in the repeated prisoners' dilemma game, the Stackelberg type always plays Defect, and the presence of such a type would not have any qualitative impact on the equilibrium behavior. When there is a tit-for-tat type, we would still have cooperation in all but a few rounds. Here, the payoff of rational type exceeds his Stackelberg payoffs, but his payoff could be lower than his Stackelberg payoffs in other games.<sup>11</sup>

We must emphasize that our main result for two-player games would still be true if we assumed that Player 2 is short-lived instead. In that case, for Player 2, we could still generate any action plan that is consistent with her stage payoff being a function of her own action

<sup>&</sup>lt;sup>10</sup>For example, sequential equilibrium is upperhemicontinuous with respect to such scaling of the probabilities of commitment types. Bayesian Nash equilibrium behavior of the rational types is also upperhemicontinuous with respect to all variations of priors (with possibly varying commitment types and relative probabilities) because such variations can be represented as an ex-ante payoff perturbation.

<sup>&</sup>lt;sup>11</sup>This fact has been demonstrated in infinitely repeated games, but we suspect that it can also be shown in long but finitely-repeated games.

only (as in Lemma 3 in the Appendix). Since this is all we need for our Lemma 4 in the Appendix, all of our results would go through as is.

Infinitely Repeated Games. Here, we focus on finitely repeated games. It is actually easier to construct types that are committed to a particular plan of action up to an arbitrary finite horizon as the unique rationalizable plan in infinite-horizon games if one does not insist on the common-prior assumption. Indeed we provided such a result in Weinstein and Yildiz (2013) in another context as we discussed before. It also seems feasible to extend our construction within common-prior assumption to infinitely-repeated games, using finitehorizon truncations. Hence, it seems feasible to obtain a similar result for infinite-horizon games allowing only arbitrarily long but finite-horizon commitments. We do not pursue such results here mainly because the most major results in infinite-horizon reputation literature, such as the above result of Fudenberg and Levine (1989), are based on types that commit to playing a fixed move, and such types can easily be justified within the repeated-game framework.<sup>12</sup>

**Commitment to Mixed Strategies.** In some reputation models, the commitment types are allowed to play a mixed action plan. For the natural case that only the realized moves are observable, such mixed commitment types are incorporated in our paper as follows. A mixed commitment type  $\sigma_i$  induces a probability distribution  $\mu^{\sigma_i}$  on pure action plans of the player in reduced form. From the point of view of the rational type  $\tau_j^*$  of the other player, the commitment type  $\sigma_i$  can be replaced by pure commitment types in the support of  $\mu^{\sigma_i}$ , by putting probability  $\pi(\sigma_i) \mu^{\sigma_i}(s_i)$  on each  $s_i$  in the support of  $\mu^{\sigma_i}$ , where  $\pi(\sigma_i)$  is the probability of  $\sigma_i$  in the original elaboration G and  $\mu^{\sigma_i}(s_i)$  is the probability of  $s_i$  under  $\mu^{\sigma_i}$ . Application of Proposition 4 to the resulting elaboration with pure commitment types yields an elaboration G' without commitment types that is strategically equivalent for the rational types.

**Other games**. The applications in reputation formation are not confined to the repeatedgames framework. Indeed, an important strand of literature explores the role of reputation in bargaining considering types that commit to dynamic plans (see for example Abreu and Gul (2000), Abreu and Pearce (2007), and Wolitzky (2012)). Of course, understanding the

<sup>&</sup>lt;sup>12</sup>When the commitment type plays a mixed move, the resulting pure action plans involve commitment to dynamically varying action plans. An extension of our results could be useful in that case.

#### JONATHAN WEINSTEIN AND MUHAMET YILDIZ

scope of reputation within the structural assumptions of those models is also very important. Here, as a first step, we established a strategic equivalence result for finitely repeated games.

How much variation in stage-game payoffs and type spaces do we need to support a commitment? Our main result establishes that arbitrary commitment types can result from ICR without any restriction on the stage-game payoffs and the type spaces. Moreover, our construction uses only a couple of simple stage-game payoff functions.<sup>13</sup> Feasibility of such stage-game payoff functions is sufficient for supporting arbitrary commitment types by introducing uncertainty on stage game payoffs. On the other hand, Proposition 5 shows that in the limit as the maximal variation in stage-game payoffs is taken to 0, commitment types cannot be generated.

**Personal Types**. Our construction (in the first part of the proof of Lemma 1) makes fundamental use of players who do not know their own payoffs. Some of the literature has focused on models with common knowledge that each player knows his own payoffs; Fudenberg, Kreps, and Levine (1988) call this a model with "personal types." We do not know precisely to what extent our results can be recovered in a model with personal types. There are multiple difficulties. The first is the construction of a type with unique rationalizable plan, as in Lemma 1. This is considerably more difficult when using personal types, and while it is possible to generate commitment types for some non-trivial plans, we do not know if it is possible for all plans. The second difficulty arises when putting the types into a common-prior type space. The technique we used for Lemma 2 relied on the commitment types of Player 1 believing that their payoff is always identically zero when Player 2 is a normal type. With commitment types, this technique cannot be used, as payoffs cannot be correlated with the opponent's type. It is an open question whether some other technique would be successful. Note that this second difficulty only arises if we assume personal types and a common prior.

### 9. CONCLUSION

The reputation literature, one of the main accomplishments of game theory, relies on the existence of commitment types. It is important for the interpretation of the results in this

<sup>&</sup>lt;sup>13</sup>More precisely, it uses the family  $g_i^{a_i^*, a_{-i}^*, \lambda} = \lambda \mathbf{1}_{a_i^*} + (1 - \lambda) \mathbf{1}_{a_{-i}^*}, a_i^* \in A_i, a_{-i}^* \in A_{-i}^*, \lambda \in [0, 1]$ , where  $\mathbf{1}_{a_i^*}$  is the characteristic function of  $a_i^*$ , taking the value of 1 when  $a_i^*$  is played and 0 otherwise.

### REPUTATION WITHOUT COMMITMENT

literature whether one can obtain the same results within a rationalistic framework in which all types can follow the plans that are available to rational types and all types' payoffs satisfy the structural payoff assumptions of the underlying model. If one can obtain the same results within such a framework, we can interpret the results as coming from incomplete information about payoffs. Otherwise, the result must be interpreted as stemming from the factors that are outside of the model, such as irrationality, psychological anomalies, and super-game concerns. In this paper, within the context of finitely repeated games, we have established that one can obtain all results within a rationalistic framework, allowing an interpretation based on incomplete information. This is the case when all stage-game payoffs are allowed. On the other hand, for games with dominance-solvable stage games, we show that reputation cannot have an impact when the stage game payoffs are sufficiently restricted. Hence, the scope of reputation within a rationalistic framework depends on the severity of the additional structural assumptions imposed when there are such assumptions.

### APPENDIX A. OMITTED PROOFS

A.1. **Preliminary Definitions.** In the appendix, we will also consider type spaces without a common prior. Such a type space is a list  $(\mathcal{G}, \mathcal{T}, \pi(\cdot|\cdot))$  where  $\pi(\cdot|\tau_i) \in \Delta(\mathcal{G} \times \mathcal{T}_{-i})$  is the probability distribution of  $\tau_i$ . Here, there need not be a single  $\pi \in \Delta(\mathcal{G} \times \mathcal{T})$  that leads to these interim beliefs by Bayes' rule. Fix any  $G = (N, A, (\mathcal{G}, \mathcal{T}, \pi(\cdot|\cdot)))$ . For each  $i \in N$  and for each belief  $\beta \in \Delta(\mathcal{G} \times S_{-i})$ , we write  $BR_i(\beta)$  for the set of actions  $s_i \in S_i$  that maximize the expected value of  $u_i(z(s_i, s_{-i})|g)$  under the probability distribution  $\beta$ .

Interim correlated rationalizability (ICR) is computed by the following elimination procedure: For each *i* and  $\tau_i$ , set  $S_i^0[\tau_i|G] = S_i$ , and define sets  $S_i^k[\tau_i|G]$  for k > 0 iteratively, by setting  $s_i \in S_i^k[\tau_i|G]$  if  $s_i \in BR_i \left( \max_{\mathcal{G} \times S_{-i}} \beta \right)$  for some  $\beta \in \Delta \left( \mathcal{G} \times \mathcal{T}_{-i} \times S_{-i} \right)$  such that  $\max_{\mathcal{G} \times \mathcal{T}_{-i}} \beta = \pi \left( \cdot |\tau_i| \right)$  and  $\beta \left( s_{-i} \in S_{-i}^{k-1}[\tau_{-i}|G] \right) = 1$ . That is,  $s_i$  is a best response to a belief of  $\tau_i$  that puts positive probability only to the actions that survive the elimination in round k - 1. We write  $S^k[\tau|G] = S_1^k[\tau_1|G] \times S_2^k[\tau_2|G]$ . Then,

$$S_i^{\infty}\left[\tau_i|G\right] = \bigcap_{k=0}^{\infty} S_i^k\left[\tau_i|G\right].$$

The following class of action plans will play an important role in our construction:

**Definition 6.** A plan  $s_i$  is said to be *sure-thing compliant* if there is no partial history h and move  $a_i \in A_i$  such that  $s_i(h, (s_i(h), a_{-i})) = a_i$  for every  $a_{-i}$  but  $s_i(h) \neq a_i$ .

In other words, a plan is sure-thing compliant if whenever the player plays  $a_i$  in all possible continuations next period, he also plays  $a_i$  this period. In the context of a single player with stable preferences who acquires information each period, this would be a consequence of the sure-thing principle of Savage.

A.2. **Proof of Lemma 1.** Our proof has three main steps. First, we will prove it for sure-thing compliant action plans, without requiring the type space to have common prior or the full-support property (property 1 in the statement of the lemma). We then extend this result to all action plans, without requiring the properties on type space once again. Finally, we convert the latter type space to a type space with common prior and full support assumptions without altering the rationalizable actions, proving the lemma. The first step is the following lemma; Weinstein and Yildiz (2013) proved this lemma for infinite-horizon games, and the proof carries over to finitely repeated games with minor modifications. (The proof can be found in the online appendix.)

**Lemma 3** (Weinstein-Yildiz 2013). For any sure-thing compliant action plan  $s_i$ , there exists a game  $\tilde{G} = \left(N, A, \left(\tilde{\mathcal{G}}, \tilde{\mathcal{T}}, \tilde{\pi}(\cdot|\cdot)\right)\right)$  with a type  $\tau_i^{s_i}$  such that  $S_i^{\infty}\left[\tau_i^{s_i}|\tilde{G}\right] \simeq \{s_i\}$ . (The type space does not necessarily have a common prior.)

The next lemma builds on this result to generalize to all action plans.

**Lemma 4.** For any action plan  $s_i$ , there exists a game  $\tilde{G} = \left(N, A, \left(\tilde{\mathcal{G}}, \tilde{\mathcal{T}}, \tilde{\pi}\left(\cdot|\cdot\right)\right)\right)$  with a type  $\tau_i^{s_i}$  such that  $S_i^{\infty}\left[\tau_i^{s_i}|\tilde{G}\right] \simeq \{s_i\}$ . (The type space does not necessarily have a common prior.)

*Proof.* Fix some  $a_{-i}^* \in A_{-i}$ , and define a function  $v_{-i}: A_{-i} \to [0,1]$  by

$$\nu_{-i}(a_{-i}) = \begin{cases} 1 & \text{if } a_{-i} = a_{-i}^*, \\ 0 & \text{otherwise.} \end{cases}$$

For every  $\hat{a}_i \in A_i$ , define a function  $v_i^{\hat{a}_i} : A_i \to [0, 1]$  by

$$v_i^{\hat{a}_i}\left(a_i\right) = \begin{cases} 1 & \text{if } a_i = \hat{a}_i, \\ 0 & \text{otherwise.} \end{cases}$$

Finally, consider the class of stage-game payoff functions  $g_i^{\hat{a}_i}: A \to [0,1]$  for player *i* where

(A.1) 
$$g_i^{\hat{a}_i}(a_i, a_{-i}) = \lambda v_i^{\hat{a}_i}(a_i) + (1 - \lambda) v_{-i}(a_{-i})$$

for some  $\lambda \in (0, 1/(2\bar{t}+1))$ . Here,  $a_{-i}^*$  is a known action profile the other players can take to reward player *i*, while  $\hat{a}_i$  is an action that player *i* can take to increase his payoff. In the type space we construct, type  $\tau_i^{s_i}$  knows his payoffs are of the form defined in (A.1) but does not know the specific  $\hat{a}_i$ . As specified by (A.1), Player *i*'s stage payoffs are additively separable in his action and others' actions. Our choice of  $\lambda$  is small enough so that in all periods before the last, player *i*'s driving incentive is the desire to be rewarded by the other players, while in the last period he has learned his own optimal action  $\hat{a}_i$  and acts accordingly.

Next, for each  $\rho: H \times A_{-i} \to A_i$ , let  $S_{-i}^{\rho}$  be the set of action profiles  $s_{-i}^{\rho}$  satisfying

$$s_{-i}^{\rho}\left(h^{t},\left(a_{i},a_{-i}\right)\right)=a_{-i}^{*}\iff a_{i}=\rho\left(h^{t},a_{-i}\right)$$

for any  $t < \overline{t}$ , any history  $h^t$  and any  $(a_i, a_{-i}) \in A_i$ . Also, let R be the set of functions  $\rho$  satisfying

$$\rho\left(h^{\bar{t}-1}, a_{-i}\right) = s_i\left(h^{\bar{t}-1}\right)$$

for all  $a_{-i}$  and all those  $h^{\bar{t}-1}$  such that player *i* has played according to  $s_i$  throughout. Finally, let  $\hat{S}_{-i} = \bigcup_{\rho \in R} S_{-i}^{\rho}$ .

To sum up, when following a plan in  $S_{-i}^{\rho}$ , at any history  $(h, (a_i, a_{-i}))$ , player -i rewards i by playing  $a_{-i}^*$  if  $a_i = \rho(h, a_{-i})$ . The only restriction on  $\rho$  occurs at date  $\bar{t} - 1$  and in the contingency that i has followed  $s_i$  up to  $\bar{t} - 1$ : he will be rewarded at  $\bar{t}$  if he continues to follow  $s_i$  at  $\bar{t} - 1$ . The set R is symmetric in all other ways. In particular, if player i assigns uniform probability on  $\hat{S}_{-i}$ , he considers it equally likely that each of his moves are rewarded, except possibly at the final stage. Note that the actions in  $\hat{S}_{-i}$  are all sure-thing compliant, because player  $j \neq i$  reacts differently to the rewarded move of player i from all other moves. Thus, by Lemma 3, for each  $s_{-i} \in \hat{S}_{-i}$ , there exists a game  $G^{s_{-i}} = (N, A, (\mathcal{G}^{s_{-i}}, \mathcal{T}^{s_{-i}}, \pi^{s_{-i}}(\cdot | \cdot )))$  with a type  $\tau_{-i}^{s_{-i}}$  such that  $S_{-i}^{\infty} [\tau_{-i}^{s_{-i}} | G^{s_{-i}}] \simeq \{s_{-i}\}.$ 

Let  $\bar{g}_{-i}$  be an arbitrary payoff function for the players other than *i*. Define the game  $\tilde{G} = \left(N, A, \left(\tilde{\mathcal{G}}, \tilde{\mathcal{T}}, \tilde{\pi}\left(\cdot | \cdot \right)\right)\right)$  by

$$\tilde{\mathcal{G}} = \left\{ (g_i^{\hat{a}_i}, \bar{g}_{-i}) | \hat{a}_i \in A_i \right\} \cup \bigcup_{s_{-i} \in \hat{S}_{-i}} \mathcal{G}^{s_{-i}};$$

$$\tilde{\mathcal{T}}_i = \{\tau_i^{s_i}\} \cup \bigcup_{s_{-i} \in \hat{S}_{-i}} \mathcal{T}_i^{s_{-i}}; \quad \tilde{\mathcal{T}}_{-i} = \bigcup_{s_{-i} \in \hat{S}_{-i}} \mathcal{T}_{-i}^{s_{-i}};$$

$$\tilde{\pi} \left( \cdot | \tau_j \right) = \pi^{s_{-i}} \left( \cdot | \tau_j \right) \qquad \left( \forall \tau_j \in \mathcal{T}_j^{s_{-i}}, j \in N, s_{-i} \in \hat{S}_{-i} \right);$$
(A.2)
$$\tilde{\pi} \left( \left( g^{\hat{a}_i \left( (\tau_{-i}^{s_{-i}}), \bar{g}_{-i} \right), \tau_{-i}^{s_{-i}} | \tau_i^{s_i} \right) = 1 / \left| \hat{S}_{-i} \right| \qquad \left( \forall s_{-i} \in \hat{S}_{-i} \right),$$

where we let

$$\hat{a}_{i}(\tau_{-i}^{s_{-i}}) = s_{i} \left( z \left( s_{i}, s_{-i} \right) \right)^{\bar{t}}$$

be the prescribed action at the history reached at the beginning of the last period under the strategy profile  $(s_i, s_{-i})$ . The critical feature of the definition is the belief of the newly introduced type  $\tau_i^{s_i}$ in (A.2). He assigns equal probabilities on types  $\tau_{-i}^{s_{-i}}$  and believes that there is a perfect correlation between the types  $\tau_{-i}^{s_{-i}}$  and the way his own action affects his payoff. If he follows  $s_i$  throughout  $h^{\bar{t}}$ and observes the moves of the other player, he learns what action  $\hat{a}_i(\tau_{-i}^{s_{-i}})$  is best for him, which happens to be the action  $s_i (z (s_i, s_{-i}))^{\bar{t}}$  that he would have played at that history according to  $s_i$ . Note that each of the types other than  $\tau_i^{s_i}$  has a unique rationalizable action in reduced form. Hence, the updated belief of  $\tau_i^{s_i}$  regarding the payoff functions and outcomes is uniquely determined at any history, given that he believes the other players follow rationalizable strategies.

In the rest of the proof, we will show that  $s_i$  is uniquely rationalizable for  $\tau_i^{s_i}$  in reduced form, i.e.,  $S_i^{\infty} \left[ \tau_i^{s_i} | \tilde{G} \right] \simeq \{s_i\}$ . Some additional notation: For any history  $h^t$ , write  $P_t^* \left( h^t \right)$  for the probability that  $a_{-i}^*$  is played at date t conditional on h according to the rationalizable belief of  $\tau_i^{s_i}$ . As noted above, by symmetry,

(A.3) 
$$P_t^*(h^t) = \begin{cases} 1 & \text{if } t = \bar{t} \text{ and } i \text{ follows } s_i \text{ throughout } h^t; \\ 0 & \text{if } t = \bar{t} \text{ and } i \text{ follows } s_i \text{ up to } \bar{t} - 1 \text{ but deviates at } \bar{t} - 1 \text{ in } h^t; \\ 1/|A_{-i}| & \text{otherwise.} \end{cases}$$

Write  $U_i(s'_i|h)$  for the expected payoff of *i* from playing  $s'_i$  under the rationalizable belief of type  $\tau_i^{s_i}$  conditional on history *h*. Write also  $U_i(h^t)$  for the realized expected payoff of  $\tau_i^{s_i}$  up to date *t* at history  $h^t$ .

We now show that  $U_i(s_i|h^t) > U_i(s'_i|h^t)$  for every history  $h^t$  and action plan  $s'_i$  such that i follows  $s_i$  throughout  $h^t$  and  $s'_i(h^t) \neq s_i(h^t)$ . So long as he follows  $s_i$ , every such history  $h^t$  is reached with positive probability under the rationalizable belief of  $\tau_i^{s_i}$ . This therefore will show that the expected payoff from  $s_i$  is strictly higher than any  $s'_i$  that is not equivalent to  $s_i$ . Therefore,  $S_i^{\infty} \left[ \tau_i^{s_i} | \tilde{G} \right] \simeq \{s_i\}.$ 

First consider the case  $t = \bar{t}$ . Conditional on  $h^t$ ,  $\tau_i^{s_i}$  assigns probability 1 on  $g^{s_i(h^t)}$ . If he follows  $s_i$ , playing  $s_i(h^t)$  at  $h^t$ , then his own action contributes  $\lambda$  to his payoff; otherwise, his own action contributes zero to his payoff. Moreover, since he has followed  $s_i$  throughout  $h^t$ , he will be rewarded for sure by the other player at  $\bar{t}$ , contributing  $1 - \lambda$  to his payoff regardless of his own move at  $h^t$ . Hence,

$$U_i\left(s_i|h^t\right) = U_i\left(h^t\right) + 1,$$

and

$$U_i\left(s_i'|h^t\right) = U_i\left(h^t\right) + 1 - \lambda,$$

yielding

$$U_i\left(s_i|h^t\right) - U_i\left(s'_i|h^t\right) = \lambda > 0.$$

Now consider the case  $t < \overline{t}$ . His payoff from following  $s_i$  is

$$U_{i}\left(s_{i}|h^{t}\right) = U_{i}\left(h^{t}\right) + \lambda \sum_{t'=t}^{\bar{t}} E\left[v_{i}^{\hat{a}_{i}(t_{-i})}\left(s_{i}\left(h^{t'}\right)\right)|h^{t}, s_{i}\right] + (1-\lambda)\sum_{t'=t}^{\bar{t}} E\left[P_{t'}^{*}\left(h^{t'}\right)|h^{t}, s_{i}\right]$$
$$\geq U_{i}\left(h^{t}\right) + (1-\lambda)\frac{\bar{t}-t}{|A_{-i}|} + 1.$$

To see the lower bound, note that, so long as he follows  $s_i$ , he gets 1 at date  $\bar{t}$  (as in the previous case) and at least  $(1 - \lambda) / |A_{-i}|$  at each  $t' < \bar{t}$ . (At any  $t' < \bar{t}$ ,  $v_i^{\hat{a}_i(t_{-i})} \ge 0$  and  $P_{t'}^* \ge 1/|A_{-i}|$  when he follows  $s_i$ .) On the other hand, his payoff from following  $s'_i$  is

$$U_{i}(s_{i}'|h^{t}) = U_{i}(h^{t}) + \lambda \sum_{t'=t}^{\bar{t}} E\left[v_{i}^{\hat{a}_{i}(t_{-i})}\left(s_{i}'(h^{t'})\right)|h^{t}, s_{i}'\right] + (1-\lambda) \sum_{t'=t}^{\bar{t}} E\left[P_{t'}^{*}(h^{t'})|h^{t}, s_{i}'\right] \\ \leq U_{i}(h^{t}) + \lambda(\bar{t} - t + 1) + (1-\lambda)\frac{\bar{t} - t + 1}{|A_{-i}|}.$$

The upper bound comes from the fact that  $v_i^{\hat{a}_i} \leq 1$  throughout and  $P_{t'}^* \leq 1/|A_{-i}|$  after a deviation from  $s_i$  (by A.3). Combining the two inequalities, we obtain

$$U_i\left(s_i|h^t\right) - U_i\left(s_i'|h^t\right) \ge (1-\lambda)\left(1 - \frac{1}{|A_{-i}|}\right) - \lambda\left(\bar{t} - t\right) > 0,$$

where the strict inequality follows from  $\lambda < 1/(2\bar{t}+1)$  and  $|A_{-i}| \ge 2$ .

Proof of Lemma 1. By Lemma 4, there exists a game  $\tilde{G} = \left(N, A, \left(\tilde{\mathcal{G}}, \tilde{\mathcal{T}}, \tilde{\pi}\left(\cdot|\cdot\right)\right)\right)$  with a type  $\tilde{\tau}_i$ such that  $S_i^{\infty}\left[\tilde{\tau}_i|\tilde{G}\right] \simeq \{s_i\}$ . This falls short of the conditions of Lemma 1 in that  $\left(\tilde{\mathcal{G}}, \tilde{\mathcal{T}}, \tilde{\pi}\left(\cdot|\cdot\right)\right)$ does not necessarily admit a common prior and the prior could not have a full support (Condition 1) even if it existed. Here, we remedy this problem by converting  $\tilde{G}$  to a common prior game  $G^s = (N, A, (\mathcal{G}^{s_i}, \mathcal{T}^{s_i}, \pi^{s_i}))$  with the desired properties. First, for every  $\lambda \in (0, 1)$ , define  $G^{\lambda} = \left(N, A, \left(\tilde{\mathcal{G}}, \tilde{\mathcal{T}}, \pi^{\lambda}\left(\cdot|\cdot\right)\right)\right)$  by setting

$$\pi^{\lambda}\left(g,\tau_{-j}|\tau_{j}\right) = \frac{\lambda}{\left|\tilde{\mathcal{G}}\times\tilde{\mathcal{T}}_{-j}\right|} + (1-\lambda)\,\tilde{\pi}\left(g,\tau_{-j}|\tau_{j}\right)$$

at each  $(g, \tau_j, \tau_{-j}) \in \tilde{\mathcal{G}} \times \tilde{\mathcal{T}}$ . Now, as  $\lambda \to 0$ ,  $\pi^{\lambda}(g, \tau_{-j}|\tau_j) \to \tilde{\pi}(g, \tau_{-j}|\tau_j)$  everywhere. Together with a continuity result for belief hierarchies by Mertens and Zamir (1985), this implies that the

belief hierarchy of type  $\tilde{\tau}_i$  in game  $G^{\lambda}$  converges to the belief hierarchy of  $\tilde{\tau}_i$  in game  $\tilde{G}$ . Thus, by upperhemicontinuity of ICR (Dekel, Fudenberg, and Morris, 2006), there exists  $\bar{\lambda} > 0$  such that

$$S_i^{\infty}\left[\tilde{\tau}_i|G^{\bar{\lambda}}\right] \subseteq S_i^{\infty}\left[\tilde{\tau}_i|\tilde{G}\right] \simeq \{s_i\}$$

Since  $S_i^{\infty}\left[\tilde{\tau}_i|G^{\bar{\lambda}}\right]$  is non-empty, this implies that

(A.4) 
$$S_i^{\infty} \left[ \tilde{\tau}_i | G^{\bar{\lambda}} \right] \simeq \{ s_i \} \,.$$

Moreover, since  $\tilde{\mathcal{G}} \times \tilde{\mathcal{T}} \times S$  is finite, there exists some finite k such that

(A.5) 
$$S_i^{\infty} \left[ \tilde{\tau}_i | G^{\bar{\lambda}} \right] = S_i^k \left[ \tilde{\tau}_i | G^{\bar{\lambda}} \right].$$

Now, since  $\pi^{\bar{\lambda}}(g, \tau_{-j}|\tau_j) > 0$  everywhere, by the main result of Lipman (2003), there exists a common-prior game  $G^{s_i} = (N, A, (\mathcal{G}^{s_i}, \mathcal{T}^{s_i}, \pi^{s_i}))$  such that the common prior  $\pi^{s_i}$  is positive everywhere and there exists a type  $\tau_i^{s_i} \in \mathcal{T}_i^{s_i}$  whose first k orders of beliefs are identical to that of type  $\tilde{\tau}_i$  in game  $G^{\bar{\lambda}}$ . Dekel, Fudenberg, and Morris (2007) show that  $S^k$  is a function of the first k orders of beliefs, yielding

(A.6) 
$$S_i^k \left[ \tau_i^s | G^s \right] = S_i^k \left[ \tilde{\tau}_i | G^{\bar{\lambda}} \right]$$

Combining (A.4), (A.5) and (A.6), we obtain

$$S_i^{\infty}\left[\tau_i^s | G^s\right] \subseteq S_i^k\left[\tau_i^s | G^s\right] = S_i^k\left[\tilde{\tau}_i | G^{\bar{\lambda}}\right] \simeq \{s_i\}.$$

Since  $S_i^{\infty}[\tau_i^s|G^s] \neq \emptyset$ , this further implies that

$$S_i^{\infty}\left[\tau_i^s | G^s\right] \simeq \{s_i\}$$

as desired.

A.3. **Proof of Lemma 2.** Using induction on k, we will show that  $S_i^k[\tau_i|G'] = S_i^k[\tau_i|G^c]$  for every  $k, \tau_i \in \mathcal{T}_i^c$ , and  $i \in N$ . This is true for k = 0 by definition. Towards an induction, assume that

(A.7) 
$$S_{-i}^{k-1} \left[ \tau_{-i} | G' \right] = S_{-i}^{k-1} \left[ \tau_{-i} | G^c \right] \qquad \left( \forall \tau_{-i} \in T_{-i}^c \right).$$

Take any  $\tau_i \in \mathcal{T}_i^c$  and write  $B(\tau_i|G)$  for the set of all beliefs  $\beta$  of type  $\tau_i$  after round k-1 in game G for any  $G \in \{G', G^c\}$ , where  $\operatorname{marg}_{\mathcal{G} \times \mathcal{T}_{-i}}\beta = \pi(\cdot|\tau_i)$  and  $\beta\left(s_{-i} \in S_{-i}^{k-1}[\tau_{-i}|G]\right) = 1$ . First consider the case  $\tau_i \neq \tau_1^c$ . In that case, by definition,  $\pi'(\cdot|\tau_i) = \pi^c(\cdot|\tau_i)$ . Together with the inductive hypothesis (A.7), this implies that  $B(\tau_i|G') = B(\tau_i|G^c)$ . Therefore,  $s_i \in S_i^k[\tau_i|G']$  if and only if  $s_i \in BR_i\left(\operatorname{marg}_{\mathcal{G}' \times S_{-i}}\beta\right)$  for some  $\beta \in B(\tau_i|G') = B(\tau_i|G^c)$ , and this is the case if and only if  $s_i \in S_i^k[\tau_i|G^c]$ , showing that  $S_i^k[\tau_i|G'] = S_i^k[\tau_i|G^c]$ .

Now consider the case  $\tau_i = \tau_1^c$ . Then,

(A.8) 
$$\pi'(\cdot|\tau_i) = p^c \delta_{((\mathbf{0},g_2^*),\tau_2^*)} + (1-p^c) \pi^c(\cdot|\tau_i)$$

where the probability  $p^c \in (0,1)$  is defined in (4.1), and  $\delta_x$  is the Dirac measure on x, putting probability 1 on  $\{x\}$ . Hence, by the inductive hypothesis (A.7),  $\beta \in B(\tau_i|G')$  if and only if

(A.9) 
$$\beta = p^c \beta \left( \cdot | \tau_2^* \right) + \left( 1 - p^c \right) \beta \left( \cdot | \tau_2^c \right)$$

for some conditional beliefs  $\beta\left(\cdot|\tau_{2}^{*}\right) \in \Delta\left(\left\{\left(\left(\mathbf{0}, g_{2}^{*}\right), \tau_{2}^{*}\right)\right\} \times S_{2}^{k-1}\left[\tau_{2}^{*}|G'\right]\right)$  and

$$\beta\left(\cdot|\mathcal{T}_{2}^{c}\right)\in B\left(\tau_{i}|G^{c}\right)$$

Now, take any  $s_i \in S_i^k[\tau_i|G']$ . Then,  $s_i \in BR_i\left(\max_{\mathcal{G}' \times S_{-i}}\beta\right)$  for some  $\beta \in B(\tau_i|G')$ . By (A.9), for any  $s'_i$ ,

$$p^{c} \cdot 0 + (1 - p^{c}) E_{\beta(\cdot|\mathcal{T}_{2}^{c})} [u_{i}(s_{i}, s_{-i}|g)] = E_{\beta} [u_{i}(s_{i}, s_{-i}|g)]$$

$$\geq E_{\beta} [u_{i}(s'_{i}, s_{-i}|g)]$$

$$= p^{c} \cdot 0 + (1 - p^{c}) E_{\beta(\cdot|\mathcal{T}_{2}^{c})} [u_{i}(s'_{i}, s_{-i}|g)]$$

where  $\beta(\cdot | \mathcal{T}_2^c) \in B(\tau_i | G^c)$ . (Here, the inequality follows from  $s_i$  being a best response, and the equalities follow from (A.9).) Since  $p^c < 1$ , this further implies that

$$E_{\beta\left(\cdot|\mathcal{T}_{2}^{c}\right)}\left[u_{i}\left(s_{i},s_{-i}|g\right)\right] \geq E_{\beta\left(\cdot|\mathcal{T}_{2}^{c}\right)}\left[u_{i}\left(s_{i}',s_{-i}|g\right)\right],$$

showing that  $s_i \in BR_i\left(\max_{\mathcal{G}' \times S_{-i}} \beta\left(\cdot | \mathcal{T}_2^c\right)\right)$ . Therefore,  $s_i \in S_i^k\left[\tau_i | G^c\right]$ .

Conversely, take any  $s_i \in S_i^k[\tau_i|G^c]$ . By definition,  $s_i \in BR_i\left(\max_{\mathcal{G}' \times S_{-i}}\beta\left(\cdot|\mathcal{T}_2^c\right)\right)$  for some  $\beta\left(\cdot|\mathcal{T}_2^c\right) \in B\left(\tau_i|G^c\right)$ . Pick any  $\beta\left(\cdot|\tau_2^*\right) \in \Delta\left(\left\{\left((\mathbf{0}, g_2^*), \tau_2^*\right)\right\} \times S_2^{k-1}[\tau_2^*|G']\right)$ , and define  $\beta \in B\left(\tau_i|G^c\right)$  by (A.9). Now, for any  $s'_i$ ,

$$\begin{split} E_{\beta} \left[ u_i \left( s_i, s_{-i} | g \right) \right] &= p^c \cdot 0 + (1 - p^c) E_{\beta\left( \cdot | \mathcal{T}_2^c \right)} \left[ u_i \left( s_i, s_{-i} | g \right) \right] \\ &\geq p^c \cdot 0 + (1 - p^c) E_{\beta\left( \cdot | \mathcal{T}_2^c \right)} \left[ u_i \left( s_i', s_{-i} | g \right) \right] \\ &= E_{\beta} \left[ u_i \left( s_i', s_{-i} | g \right) \right], \end{split}$$

where the inequality follows from  $s_i$  being a best response, and the equalities follow from (A.9). That is,  $s_i \in BR_i\left(\max_{\mathcal{G}' \times S_{-i}}\beta\right)$ . Therefore,  $s_i \in S_i^k[\tau_i|\mathcal{G}']$ . A.4. **Proof of Proposition 4.** Here we outline the proof of the general proposition. Note first that Lemma 1 applies to general case as well: for each  $i \in N$  and  $c_i \in C_i$  there exists a exists a Bayesian repeated game  $G^{c_i} = (N, A, (\mathcal{G}^{c_i}, \mathcal{T}^{c_i}, \pi^{c_i}))$  with  $\pi^{c_i}$  positive everywhere and with a type  $\tau_i^{c_i} \in \mathcal{T}_i^{c_i}$  for which all ICR actions are equivalent to  $c_i$ . Again all those types can be taken unique and distinct from each other. Write  $\mathcal{T}_i = \{\tau_i^*\} \cup C_i$ , and define mapping  $\tilde{\tau}_i$  on  $\mathcal{T}_i$  by

$$\tilde{\tau}_i \left( \tau_i \right) = \begin{cases} \tau_i^* & \text{if } \tau_i = \tau_i^* \\ \tau_i^{c_i} & \text{otherwise}, \end{cases}$$

and mapping  $\gamma_i$  from  $\tilde{\mathcal{T}}_i = \tilde{\tau}_i(\mathcal{T}_i)$  to stage-game payoff functions by

$$\gamma_i(\tau_i) = \begin{cases} g_i^* & \text{if } \tau_i = \tau_i^* \\ \mathbf{0} & \text{otherwise.} \end{cases}$$

Write also  $\gamma(\tau) = (\gamma_1(\tau_1), \ldots, \gamma_n(\tau_n))$  for  $\tau \in \tilde{\mathcal{T}} = \tilde{\mathcal{T}}_1 \times \cdots \times \tilde{\mathcal{T}}_n$ . We construct  $G' = (N, A, (\mathcal{G}', \mathcal{T}', \pi'))$  by setting

$$\begin{split} \mathcal{G}' &= \prod_{i \in N} \left\{ g_i^*, \mathbf{0} \right\} \cup \bigcup_{i \in N, c_i \in C_i} \mathcal{G}^{c_i} \\ \mathcal{T}'_j &= \left\{ \tau_j^* \right\} \cup \bigcup_{i \in N, c_i \in C_i} \mathcal{T}_j^{c_i} \qquad (\forall j \in N) \\ \pi'\left(g, \tau\right) &= \begin{cases} 1 - \varepsilon' & \text{if } (g, \tau) = (g^*, \tau^*) \,, \\ \frac{1 - \varepsilon'}{1 - \varepsilon} \pi\left(\tau'\right) & \text{if } \tau = \tilde{\tau}\left(\tau'\right) \text{ and } g = \gamma\left(\tau\right) \text{ for some } \tau' \in \mathcal{T} \setminus \{\tau^*\} \,, \\ \frac{\varepsilon' - \varepsilon}{(1 - \varepsilon)\left(|C_1| + \dots + |C_n|\right)} \pi^{c_i}\left(g, \tau\right) & \text{if } (g, \tau) \in \mathcal{G}^{c_i} \times \mathcal{T}^{c_i} \text{ for some } c_i \in C_i \\ 0 & \text{otherwise.} \end{cases}$$

Observe that G' satisfies the properties in the proposition. For any rational type  $\tau_i^*$ ,

$$\pi'\left( au_{j}^{*}
ight)=rac{1-arepsilon'}{1-arepsilon}\pi\left( au_{j}^{*}
ight),$$

and hence  $\pi'\left(\tilde{\tau}_{-j}\left(\tau_{-i}\right)|\tau_{j}^{*}\right) = \pi\left(\tau_{-i}|\tau_{j}^{*}\right)$  for every  $\tau_{-j} \in \mathcal{T}_{-j}$ . On the other hand, for type  $\tau_{i}^{c_{i}}$ , his payoffs vary only when the other types are in  $\mathcal{T}_{-i}^{c_{i}}$ . Hence, as in Lemma 2,  $S_{i}^{\infty}\left[\tau_{i}^{c_{i}}|G'\right] = c_{i}$ .

#### References

- [1] Abreu, Dilip and Faruk Gul (2000): "Bargaining and Reputation," Econometrica 68, 85-117.
- [2] Abreu, Dilip and David Pearce (2007): "Bargaining, Reputation, and Equilibrium Selection in Repeated Games with Contracts," *Econometrica*, 75: 653–710.
- [3] Dekel, Eddie, Drew Fudenberg, and Stephen Morris (2007): "Interim Correlated Rationalizability," *Theoretical Economics*, 2, 15-40.

- [4] Fudenberg, Drew, David Kreps, and David Levine (1988): "On the Robustness of Equilibrium Refinements," Journal of Economic Theory 44, 354-380.
- [5] Fudenberg, Drew and David Levine (1989): "Reputation and Equilibrium Selection in Games with a Patient Player," *Econometrica*, 57, 759-778.
- [6] Fudenberg, Drew and Eric Maskin (1986): "The Folk Theorem in Repeated Games with Discounting or Incomplete Information," *Econometrica*, 54, 533-557.
- [7] Kajii, Atsushi and Stephen Morris (1997): "The Robustness of Equilibria to Incomplete Information," *Econometrica*, 65, 1283-1309.
- [8] Kreps, David and Robert Wilson (1982): "Reputation and imperfect information," Journal of Economic Theory, 27, 253-279.
- Kreps, David, Paul Milgrom, John Roberts and Robert Wilson (1982): "Rational Cooperation in the Finitely-Repeated Prisoners' Dilemma," *Journal of Economic Theory*, 27, 245-252.
- [10] Lipman, Bart (2003): "Finite Order Implications of Common Prior," Econometrica, 71, 1255-1267.
- [11] Mailath, George and Larry Samuelson (2006): Repeated Games and Reputations, Oxford University Press.
- [12] Monderer, Dov, and Dov Samet (1989): "Approximating Common Knowledge with Common Belief," Games and Economic Behavior, 1, 170–190.
- [13] Milgrom, Paul and John Roberts (1982): "Predation, reputation, and entry deterrence," Journal of Economic Theory, 27, 280-312.
- [14] Weinstein, Jonathan and Muhamet Yildiz (2007): "A Structure Theorem for Rationalizability with Application to Robust Predictions of Refinements," *Econometrica*, 75, 365-400.
- [15] Weinstein, Jonathan and Muhamet Yildiz (2013): "Robust predictions in infinite-horizon games—An unrefinable folk theorem," *Review of Economic Studies*, 80, 365-394.
- [16] Wolitzky, Alexander (2012): "Reputational Bargaining with Minimal Knowledge of Rationality," Econometrica, 80, 2047-2088.