

# Non-binding agreements and forward induction reasoning\*

**Emiliano Catonini<sup>†</sup>**

November 2015

In dynamic games, players may observe a deviation from a pre-play, possibly incomplete, non-binding agreement before the game is over. The attempt to rationalize the deviation may lead players to revise their beliefs about co-players' behavior in the continuation of the game. This instance of forward induction reasoning is based on interactive beliefs about not just rationality but also the compliance with the agreement itself. I study the effects of such rationalization on the self-enforceability of the agreement. Accordingly, outcomes of the game are deemed to be implementable by some agreement or not. Conclusions depart substantially from what the equilibrium refinement tradition suggests. A non subgame perfect equilibrium may represent a self-enforcing agreement, while a subgame perfect equilibrium may not. The incompleteness of the agreement can be crucial to implement an equilibrium outcome. However, every game possesses an outcome which is compatible with both forward induction and subgame perfection.

Keywords: Agreements, Self-Enforceability, Forward Induction, Strong- $\Delta$ -Rationalizability, Selective Rationalizability, Stability.

**J.E.L. Classification:** C72, C73, D86.

---

\*I am grateful to Pierpaolo Battigalli for introducing me to the mysteries of epistemic game theory. Thank you also to Adam Brandenburger, Alfredo Di Tillio, Amanda Friedenberg, Mattia Landoni, Elena Manzoni, Burkhard Schipper, Dimitrios Tsomocos, Yi-Chun Chen, Xiao Luo, Satoru Takahashi, Amanda Jakobsson, Madhav Aney, Shurojit Chatterji and Atsushi Kajii for precious suggestions.

<sup>†</sup>Higher School of Economics, ICEF, emiliano.catonini@gmail.com

# 1 Introduction

When the players of a dynamic game can communicate before the game starts, they are likely to exploit this opportunity to reach a possibly incomplete agreement<sup>1</sup> about how to play. In most cases, the context allows them to reach only a non-binding agreement, which cannot be enforced by a court of law. The only way a non-binding agreement can affect the behavior of players is through the beliefs it is able to induce in their minds. This paper sheds light on which agreements players can believe in and, among them, which agreements players will comply with. Moreover, in an implementation perspective, the paper investigates which outcomes of the game can be ensured by *some* agreement. The paper will not deal with the pre-play bargaining phase. Yet, the evaluation of their credibility has a clear and strong feedback on which agreements are likely to be reached.

In this paper I take the view that players will believe in the agreement only if this is compatible with reasonable assumptions about rationality,<sup>2</sup> beliefs in rationality and their interaction with the beliefs in the agreement of all orders. Ann will believe in the agreement only if Bob may comply with it in case he is rational, he believes in the agreement, he believes that Ann is rational and believes in the agreement (which may add non-agreed upon restrictions on what Bob expects Ann to do), and so on. Moreover, I take the view that deviations, or more generally past actions, are not interpreted as mistakes but as intentional choices. Suppose that for Bob, in case he is rational and believes in the agreement, some move makes sense only if he plans to play a certain action thereafter. Ann, upon observing such move, will believe that Bob will play that action (and Bob may use the move to signal this). This instance of forward induction reasoning is based not just on the belief in Bob's rationality but also on its interaction with the belief that Bob believes in the agreement. Example 2 in Section 2 provides a case in point. Consider now a move that Bob, if he is rational and believes in the agreement, cannot find profitable whatever he plays thereafter. Example 1 in Section 2 provides a situation of this kind. Then Ann cannot maintain both beliefs that Bob is rational and that he believes in the agreement. Which of the

---

<sup>1</sup>The mathematical representation of agreements in this paper can be given also different interpretations. For instance, the agreement can represent a set of public announcements (from a subset of players)

<sup>2</sup>The notion of rationality employed in this paper imposes expected utility maximization but it does not impose by itself any restriction on beliefs. See Section 3 for details.

two she will (try to)<sup>3</sup> maintain determines the choice of the appropriate rationalizability concept to capture these lines of strategic reasoning. *Strong- $\Delta$ -rationalizability* (Battigalli, [5]; Battigalli and Siniscalchi, [11]) captures the hypothesis that the beliefs in the agreement are given higher *epistemic priority* than the beliefs in rationality; that is, in a contingency which some player would not have allowed under the beliefs in the agreement and rationality of order, respectively,  $n$  and  $n - 1$ ,<sup>4</sup> the co-players abandon the belief in rationality of order  $n$ .<sup>5</sup> When instead the co-players want to retain all orders of belief in rationality that are per se compatible with the observed behavior and rather drop the orders of belief in the agreement that are at odds with them, *selective rationalizability* is the appropriate tool (see Catonini [13]). Since the agreement originates from just cheap talk among players, I suggest that selective rationalizability is the appropriate solution concept in most situations. However, for robustness purposes and theoretical insight, it will be useful to evaluate agreements with both rationalizability concepts. Strong- $\Delta$ -rationalizability and selective rationalizability deliver either an empty set (if believing in the agreement is not compatible with the strategic reasoning hypotheses) or the possible behavioral implications of the agreement (otherwise). Some behavioral implications of a credible agreement may be inconsistent with the agreement itself: a player may or may not comply with the agreement depending on her conjectures where the incomplete agreement and strategic reasoning do not uniquely pin down the moves of the co-players.<sup>6</sup>

For notational simplicity and for the proofs of some results, the focus is restricted to the class of finite games with complete information and observable actions.<sup>7</sup> However, the methodology can be applied to all dynamic games with a countable set of non-terminal histories<sup>8</sup> and perfect recall, hence possibly infinite horizon.<sup>9</sup> Which

---

<sup>3</sup>The observed behavior may contradict also some order of belief in rationality per se.

<sup>4</sup>For  $n = 1$  I mean rationality itself. The first-order-belief in the agreement has bite only when associated with rationality. Likewise, the belief in the agreement of order  $n$  has bite only when associated with the belief in rationality of order  $n - 1$ . See Section 3 for details.

<sup>5</sup>All the beliefs in the agreement of higher orders can be maintained, since they have no behavioral consequence without the belief in rationality of the lower order, so they cannot be falsified by observation. This remark is due to Battigalli and Prestipino [8].

<sup>6</sup>Or when indifferences apply.

<sup>7</sup>Games where every player always knows the current history of the game, i.e. - allowing for truly simultaneous moves - information sets are singletons. For instance, all repeated games with perfect monitoring are games with observable actions.

<sup>8</sup>The limitation to a countable set of non-terminal histories allows to use Conditional Probability Systems (see Section 3), which require a countable set of conditioning events.

<sup>9</sup>Battigalli and Prestipino [8] have the most general analysis of strong- $\Delta$ -rationalizability in finite

agreements are *credible* and will be complied with? Which outcomes of the game can be achieved through some agreement? To answer these questions, the concepts of *self-enforceability* (of agreements) and *implementability* (of outcomes) are defined.

First, a very natural class of agreements is analyzed. A *path agreement* does not attempt to prescribe behavior at contingencies that follow a violation of the agreement itself. Thus, it just corresponds to agreeing on an outcome to be achieved. Example 2 in Section 2 analyzes a path agreement. Independently of the epistemic priority hypothesis, only for a strongly rationalizable,<sup>10</sup> subgame perfect equilibrium (henceforth, SPE) outcome, the corresponding path agreement can be self-enforcing, but (unless it is the only one) not even its credibility is guaranteed.

What about non-subgame perfect equilibrium outcomes? The corresponding path agreement may be credible too. Moreover, a non subgame perfect equilibrium outcome may be implementable with appropriate off-the-path threats (whereas a SPE one may not). Example 1 in Section 2 is a case in point. For two-players game, I provide a full characterization of implementable outcomes: an outcome is implementable under epistemic priority to rationality (resp., to the agreement) if and only if it is induced by a "strict"<sup>11</sup> equilibrium in strongly rationalizable (resp., sequentially rational) strategies. The incompleteness of the agreement can be crucial to implement an equilibrium outcome, also when players are more than two. The last examples of Section 4.1 and Section 4.2 illustrate how.

At this point, one may wonder whether there always exists (the support of) a SPE outcome (distribution) that can be implemented by some agreement. The answer is negative under the assumption that players do not agree on mixed actions. However, a preliminary investigation in this direction leads to a result of broader interest: in every game with observable actions, there are always a SPE and an equilibrium in strongly rationalizable strategies which induce the same distribution over outcomes.<sup>12</sup> That is, backward induction (plus equilibrium reasoning)<sup>13</sup> and forward induction (as

---

games: incomplete information, imperfectly observable actions, chance moves. Selective rationalizability can be extended in the same direction too.

<sup>10</sup>Strong rationalizability is a modification of Extensive Form Rationalizability (Pearce, [27]), which Battigalli and Siniscalchi [10] use to characterize the behavioral implications of Rationality and Common Strong Belief in Rationality.

<sup>11</sup>i.e., without best replies to the equilibrium conjecture which would induce a different outcome.

<sup>12</sup>In games with perfect information and no relevant ties, strong rationalizability yields as unique outcome the SPE one, as proved by Battigalli and Siniscalchi [9] and Heifetz and Perea [23].

<sup>13</sup>Chen and Micali [14] already proved that backward induction without equilibrium reasoning, i.e. purely as an elimination procedure, has an overlap with forward induction in terms of predicted

captured by strong rationalizability)<sup>14</sup> never give disjoint predictions.

While the rationalizability literature provides concepts and tools for the analysis, the equilibrium refinement literature was the first to introduce forward induction considerations into equilibrium reasoning. Stable equilibria by Kohlberg and Mertens [24] (henceforth K&M) is a set-valued solution concept that captures instances of forward induction triggered here by a path agreement. Osborne [26] shows that *equilibrium paths that can be upset by a convincing deviation*<sup>15</sup> are not stable. Govindan and Wilson [20] refine sequential equilibria with a weaker notion of forward induction.<sup>16</sup> However, these works have two common shortcomings. First, they never question subgame perfection as a must-have for a "strategically stable" solution.<sup>17</sup> Second, the strategic reasoning that leads to play (or reject) such equilibria is unclear (or limited). The rationalizability approach adopted in this paper, which is backed by epistemic foundations, allows to eliminate both shortcomings. First, there is no constraint about how precisely and on which kind of equilibrium or disequilibrium behavior players agree. Second, there is transparency about which particular agreements, beliefs and epistemic assumptions trigger different lines of reasoning.

In this sense, this work can also be interpreted as the axiomatic realization of a program akin to K&M (see [24], p. 1020).<sup>18</sup> Epistemic priority to rationality implies,

---

outcomes. In an earlier paper, Battigalli [4] proves the same result in a smaller class of games.

<sup>14</sup>Strong-rationalizability captures "unrestricted" forward induction, based on the beliefs in rationality only. The existence of such Nash allows to easily claim the existence, at least in 2-players games, of an implementable SPE outcome distribution with agreements on mixed actions. In this sense, also "restricted" forward induction based on the agreement and subgame perfection do not give disjoint predictions. In 2-players games, Govindan and Wilson [20] prove the existence of sequential equilibria which capture a weak notion of forward induction based on the beliefs in the equilibrium path, see footnote 16.

<sup>15</sup>Particular sequences of equilibria of the stage game of a finitely repeated game. I characterize them as non credible agreements in the Online Appendix.

<sup>16</sup>In the refined sequential equilibria, players believe in strategies of the opponent which are best replies to a weakly sequential equilibrium with the same outcome distribution, but without restrictions on beliefs. Thus, beliefs in rationality and in the path above the second order are not captured. (I conjecture that second order beliefs are, whereas the authors observe explicitly strong belief in rationality only.)

<sup>17</sup>In some cases, the actual realization of K&M's ideal program departs from subgame perfection. Yet, in the attempt to capture it, valuable equilibria are disregarded: all the non-SPE outcomes implemented in the examples of this paper are unstable, although they also pass the credibility test, so they are compatible with forward induction a la K&M (see this paragraph). K&M regard the inability to imply subgame perfection as a weakness of stability, and "hope that in the future some appropriately modified definition of stability will, in addition, imply connectedness and backwards induction." This paper suggests the opposite direction.

<sup>18</sup>K&M write: "We agree that an ideal way to discuss which equilibria are stable, and to delineate

in generic games, iterated weak dominance.<sup>19</sup> Implementable outcomes are proved (and not assumed) to be admissible and Nash. Full-fledged forward induction reasoning is captured and clarified.<sup>20</sup> Agreements on *pure* actions provide clear motivation and intuitive implementation, whereas stability requires hard-to-guess mixed strategies even for the most intuitive outcomes. Finally, I do not disregard non subgame perfect equilibrium outcomes. Players are not necessarily required to coordinate on expected behavior after a deviation. A credible threat must be optimal against a plan of the deviator which is compatible with what the deviation signals (after all possible steps of reasoning) about her future intentions. In the equilibrium literature, forward induction is based on the following interpretation of a deviation: the deviator believed in the path (i.e. that the co-players would have followed the path), but does not believe in the threat. But then, that the threat is a best reply to a plan of the deviator which is a best reply to the threat itself is of no additional value. Moreover, threats compatible with both forward induction and subgame perfection typically do not exist (see the first example of Section 4.1): if a threat compatible with forward induction is made explicit, then that *another* threat is subgame perfect is of no additional value either. If the threat is not explicit, subgame perfection will be obtained as a result and not as an assumption (see Theorem 2), but very few SPE paths are self-enforcing. Selective rationalizability is based on an agnostic interpretation of deviations: when a rational player displays disbelief in the agreement, the co-players are free to assume that the agreement has been believed to any partial extent compatible with the observed behavior.<sup>21</sup> In addition, a "theoretical" rather than "literal" use of path agreements, in the fashion of the Iterated Intuitive Criterion by Cho and Kreps [15], allows to evaluate the robustness of threats against the interpretation above (see Example 2 in Section 2 for clarifications).<sup>22</sup> The implementation of all the outcomes

---

this common feeling, would be to proceed axiomatically. However, we do not yet feel ready for such an approach; we think the discussion in this section will abundantly illustrate the difficulties involved." Thirty years later, the achievements of epistemic game theory allow to overcome many of these difficulties.

<sup>19</sup>In generic games, iterated weak dominance is equivalent to strong rationalizability (see Battigalli and Siniscalchi [10]), of which selective rationalizability is a refinement.

<sup>20</sup>In K&M, forward induction is defined on the normal form. This makes it hard to understand to what extent forward induction reasoning is captured, especially at information sets that do not immediately follow a unilateral deviation from the equilibrium path, and for many steps of reasoning.

<sup>21</sup>Strong- $\Delta$ -rationalizability assumes instead that the agreement is always fully believed.

<sup>22</sup>Selective rationalizability can be modified to capture a finer epistemic priority ordering, with the beliefs in the compliance with the path in between the beliefs in rationality and the beliefs in the threats. This would automatize such robustness control.

in the examples of this paper is robust in this sense.

To explain intuitively these ideas, Section 2 discusses two examples (formally solved in the Online Appendix). In the first, players can profitably agree on a non subgame perfect equilibrium of the game, also when they assign priority to rationality. In the second, players would like to achieve a desirable SPE outcome, but forward induction reasoning makes the corresponding agreement not credible. Section 3 introduces the theoretical framework and the analytic tools for the formal treatment of Section 4. Section 5 concludes with directions for future research. The proofs of the theorems where observability of actions and finiteness are actually used are postponed to the Appendix.

## 2 Main examples

**Example 1** In a city, two parties can form a coalition for the election of the mayor if they choose the same positioning on a few issues. If they both choose a *Radical* positioning, their coalition will win and split equally a surplus of 10. If they both choose a *Moderate* positioning, their coalition will win and the surplus to split grows to 12. The problem is that Party 2 may be tempted to take a radical positioning even if Party 1 chooses a moderate one. In this case,  $P_2$  would not win at the first round but would proceed to the second round (a ballot) against a third candidate. In the last public speeches,  $P_1$  can ask its voters to *Support*  $P_2$  or the *Alternative* candidate and  $P_2$  can make a political *Offer* to  $P_1$ 's voters, which costs 2, or *Not*.  $P_2$  wins for sure if  $P_1$  supports and with probability 1/2 if not but the offer is made.  $P_1$  earns 2 if the supported candidate wins, unless  $P_2$  wins without making the offer.

$P_1 \backslash P_2$	$M$	$R$
$M$	(6, 6)	..
$R$	(0, 0)	(5, 5)

---
-->

$P_1 \backslash P_2$	$N$	$O$
$A$	(2, 0)	(1, 4)
$S$	(0, 10)	(2, 8)

The subgame has only one equilibrium, where  $A$  and  $N$  are played with probability 1/3. Since the payoff of  $P_2$  in this equilibrium is higher than 6,  $(M, M)$  is not a SPE outcome. Yet, it is the outcome of a Nash equilibrium where  $P_1$  plays  $M$  and  $A$  and  $P_2$  plays  $M$  and  $N$ .

The two parties meet before declaring their positioning to the public and try to reach an agreement. Can they credibly agree on the Nash equilibrium? The answer is yes and it is robust to the two different epistemic priority assumptions. Suppose that  $P_2$  deviates to  $R$ . Is it credible that  $P_1$  replies with  $A$ ? A rational  $P_2$  believing in the agreement would not have played  $R$ . Hence  $P_1$  must either believe that  $P_2$  does not believe in the agreement or that it is irrational. If  $P_1$  drops the belief that  $P_2$  is rational (epistemic priority to the agreement), it can expect any move and  $A$  is a best reply to  $N$ . If  $P_2$  believes that  $P_1$  reasons in this way, it can believe that  $P_1$  would reply to the deviation with  $A$ . The agreement is credible and, once believed, players will play  $(M, M)$ , consistently with the agreement. If  $P_1$  drops the belief that  $P_2$  believes in the agreement (epistemic priority to rationality), it can expect any rational move:  $(R, N)$  is a best reply to all conjectures that put probability 1 on  $S$  and  $A$  is a best reply to  $N$  in the subgame.<sup>23</sup> Again, if  $P_2$  believes that  $P_1$  reasons in this way, it can believe that  $P_1$  would reply to the deviation with  $A$ ; so the agreement is credible and players will play  $(M, M)$ . That is, the complete agreement on the Nash profile is self-enforcing. To keep the games simple, situations where the incompleteness of the agreement is instead necessary to implement a desirable Nash outcome are provided in the last examples of Section 4.1 and 4.2.

Note that the threat is credible also if  $P_2$ , after the deviation, is assumed to be rational and have believed in the path (i.e. that  $P_1$  would have played  $M$ ) although not in the threat. Rationalizing the deviation in this way,  $P_1$  can expect  $P_2$  to play  $N$  and hence best reply with  $A$ . Expecting this rationalization,  $P_2$  can believe in  $A$  and hence not deviate.

The agreement on the SPE path  $(R, R)$  is self-enforcing. In other cases, forward induction reasoning based on the beliefs in the path may rule out the off-the-path beliefs that are necessary to prevent a deviation.<sup>24</sup> This is the case in Example 2.

**Example 2** The duopolists of the cola market,  $A$  and  $B$ , have to decide their marketing strategy before two big sport events, for which the population will gather

---

<sup>23</sup>Also  $(R, O)$  is rational. But  $S$  is best reply to  $O$  and  $N$  is best reply to  $S$ , so we can run in circle and the incentive to play  $A$  remains compatible with the beliefs in rationality of all orders.

<sup>24</sup>Forward induction arguments can also go in favour of the agreement. Suppose that after a unilateral deviation from a path agreement, the deviator can get a higher payoff than under the path only if the co-player plays a certain action. Suppose that this action is best reply only to an action that prevents the deviator to get a higher payoff than under the path. Then this action will not be played and the possibility of deviation is ruled out via forward induction.



in front of the television. There are 10 million buyers: 2 of them are somewhat loyal to brand  $A$ , 2 of them are somewhat loyal to brand  $B$ , the others just follow the advertisements before the event (if both or none brands do it, they split equally). At the current price, each million of buyers brings a profit of 1. Advertising costs 2. There is also another marketing strategy, which consists of a discount in the supermarkets before the event. Every buyer switches to the brand making discounts, but the profit drops to 0.2 per million buyers.

The game is a twice repeated prisoner dilemma with a punishment action, because advertisement ( $D$ ) is a best reply to both no advertisement ( $C$ ) and advertisement, while the aggressive discount campaign ( $P$ ) is best reply only to the other firm doing the same, and the profits of both firms fall anyway.

$A \setminus B$	$C$	$D$	$P$
$C$	5, 5	2, 6	0, 2
$D$	6, 2	3, 3	0, 2
$P$	2, 0	2, 0	1, 1

There is a SPE where the two firms collude in the first stage. Suppose that the two marketing directors, Ann and Bob, agree not to advertise their products before the first event and to do it before the second event. It is understood that the agreement falls through if it was violated for the first event. All this sums up to agreeing on the SPE path.

The agreement does not rule out punishment  $P$  after a deviation; therefore, it seems possible (although not guaranteed) that players fear it and comply with the path. Instead, by forward induction reasoning, punishment is actually ruled out and thus the path agreement is not credible. Bob, if he is rational and believes that Ann will comply with the agreement, will play  $D$  in the first stage only if he does not expect  $P$  in the second. Ann, if she believes this and observes  $D$ , will then (by forward induction) expect Bob to play  $D$  again and play  $D$  herself. Bob, if he believes that Ann would interpret the deviation in this way, will actually play  $D$  in the first stage. Ann, if she believes that Bob expects such interpretation, will predict the deviation and not believe in the agreement.

Two objections may be raised at this point.

First, Ann could interpret the deviation as follows: "Bob believed that  $I$  would

have not complied with the agreement, and best replied by not complying himself." But then, if the beliefs of Ann are dynamically consistent (as I will require), she must believe that Bob does not trust her from the start: the deviation of Bob, assuming he is rational, is not at odds with the belief that Ann complies with the agreement.

Second, anticipating that the path agreement would not work, players could explicitly threaten the punishments beforehand. In the next section I propose a few reasons why this may not happen. However, even if players reach the complete agreement, the non credibility of the path agreement could make the belief in the complete agreement fall for the following reason. Suppose that after a deviation, the co-players pose themselves the following dilemma: the deviator did not believe in the agreement on path (i.e. she expected someone else to deviate) or does not believe in our post-deviation threat? If the co-players rationalize the deviation in the second way, the instances of forward induction reasoning triggered by a path agreement take place.

Battigalli and Siniscalchi ([11] and [12]) show that the Iterated Intuitive Criterion test is equivalent to non-emptiness of strong- $\Delta$ -rationalizability under the restrictions that correspond to the equilibrium outcome distribution. Theorem 1 in Section 4.1 will prove that non-emptiness of selective rationalizability and of strong- $\Delta$ -rationalizability are equivalent for path restrictions. Therefore, the "credibility of the path agreement" test can be seen as a generalization<sup>25</sup> of the Iterated Intuitive Criterion, and provides for it the motivation above. Moreover, only strongly rationalizable outcomes can pass the test (see Proposition 1). The test is applied to "equilibrium paths that can be upset by a convincing deviation" in the Online Appendix, and it gives negative response. When the test gives positive response, it also provides conjectures under which players comply with the path and which are compatible with the rationalization of deviations discussed above. From such conjectures, it *may* be possible to derive an agreement that implements the outcome. This possibility is discussed in Section 4.2.

---

<sup>25</sup>Once extended to outcome distributions and games with incomplete information.

### 3 Agreements, beliefs and strategic reasoning

#### 3.1 Preliminaries

Let  $I$  be the finite set of *players*. For any collection  $(X_i)_{i \in I}$  and  $J \subseteq I$ , I write  $X_J := \times_{j \in J} X_j$ ,  $X := X_I$  and  $X_{-i} := X_{I \setminus \{i\}}$ ; for any profile  $x \in X$ , I write  $x_J := \text{Proj}_{X_J} x$ .

Let  $(\bar{A}_i)_{i \in I}$  be the finite sets of *actions* potentially available to each player and  $\bar{H} \subseteq \cup_{t=1, \dots, T} \bar{A}^t \cup \{\emptyset\}$ , where  $T$  is the finite horizon of the game, the set of *histories*, with the following properties:

- $h^0 := \{\emptyset\} \in \bar{H}$  (root of the game);
- for any  $h = (a^1, \dots, a^t) \in \bar{H}$  and  $l < t$ ,  $h' = (a^1, \dots, a^l) \in \bar{H}$  and I write  $h' \prec h$ ,<sup>26</sup>
- calling  $Z := \{z \in \bar{H} : \forall h \in \bar{H}, z \not\prec h\}$  the set of *terminal histories* and  $H := \bar{H} \setminus Z$  the set of *non-terminal histories*, for every  $i \in I$  there exists a non-empty-valued correspondence  $A_i : H \rightrightarrows \bar{A}_i$  such that for every  $h \in H$ ,  $(h, a) \in \bar{H}$  if and only if  $a \in A_i(h)$ .

Finally, for every  $i \in I$ , let  $u_i : Z \rightarrow \mathbb{R}$  be the payoff function;  $\Gamma = \langle I, \bar{H}, (u_i)_{i \in I} \rangle$  is a *finite game with complete information and observable actions*.

Now I can derive the following objects. A *strategy* is a function  $s_i : H \rightarrow \bar{A}_i$  such that for every  $h \in H$ ,  $s_i(h) \in A_i(h)$ . The set of all strategies is denoted by  $S_i$ . The set of strategies that are compatible with a history  $h = (a^1, \dots, a^t) \in \bar{H}$  is defined as:

$$S_i(h) := \{s_i \in S_i : s_i(h^0) = a_i^1 \wedge \forall l < t, s_i((a^1, \dots, a^l)) = a_i^{l+1}\}.$$

For any  $\bar{S}_i \subset S_i$ ,  $\bar{S}_i(h) := S_i(h) \cap \bar{S}_i$  and  $\bar{S}_i[h] := \{a_i \in \bar{A}_i : \exists s_i \in \bar{S}_i(h), s_i(h) = a_i\}$ .

On the other hand, the set of histories that are compatible with a set of strategy (sub-)profiles  $\bar{S}_J \subseteq S_J$  is defined as  $\bar{H}(\bar{S}_J) := \{h \in \bar{H} : \bar{S}_J(h) \neq \emptyset\}$ . Then, the sets of terminal histories and non-terminal histories that are compatible with  $\bar{S}_J$  are respectively  $\zeta(\bar{S}_J) := \bar{H}(\bar{S}_J) \cap Z$  and  $H(\bar{S}_J) := \bar{H}(\bar{S}_J) \cap H$ .

Given a strategy  $s_i \in S_i$ , the set of realization equivalent strategies is  $[s_i] := \{s'_i \in S_i : \forall h \in H(s_i), s'_i(h) = s_i(h)\}$ . The immediate predecessor of a history  $h \in \bar{H} \setminus \{h^0\}$  is denoted by  $p(h)$ .<sup>27</sup>

<sup>26</sup> $\bar{H}$  endowed with the precedence relation  $\prec$  is a tree with root  $h^0$ .

<sup>27</sup>The history  $h'$  such that  $h = (h', a)$  for some  $a \in \bar{A}$ .

In games with observable actions, every player always knows the current history of the game. That is, *information sets* are singletons and coincide with histories. I will use the term information set instead of history when it is the appropriate one in games where the two do not coincide. Note moreover that to simplify notation every player is required to play an action at every history: when a player is not truly active at a history, her set of feasible actions consists of just one "wait" action.

### 3.2 Agreements, complete agreements and path agreements

I consider pre-play, non-binding *agreements* among players of the following form.

**Definition 1 (Agreement)** *An agreement is a profile of correspondences  $e = (e_i : H \rightrightarrows \bar{A}_i)_{i \in I}$ , such that for every  $i \in I$  and  $h \in H$ ,  $\emptyset \neq e_i(h) \subseteq A_i(h)$ .*

An agreement specifies at every information set the *pure* actions among which players are expected to choose. Note that the agreement can restrict the behavior of players also at information sets that are supposedly precluded by the agreement itself. When this is not the case and the agreement allows to reach only one outcome of the game, I call it a *path agreement*.

**Definition 2 (Path Agreement)** *Fix  $z = (a^1, \dots, a^t) \in Z$ . A path agreement on  $z$  is an agreement  $e = (e_i)_{i \in I}$  such that for every  $i \in I$ ,  $e_i(h^0) = a_i^1$ ; for every  $l < t$ ,  $e_i((a^1, \dots, a^l)) = a_i^{l+1}$ ; for every  $h \not\prec z$ ,  $e_i(h) = A_i(h)$ .*

Path agreements are particularly interesting for many reasons. First, they correspond to a very natural kind of agreement: choosing an outcome as a goal. Second, punishments are not explicitly threatened in many real-life situations: discussing what to do in case the partner defects is an inconvenient way to start a relationship. Third, players may anticipate that punishments would not be believed: the belief in an agreement that has already been violated is very likely to fall. Or, punishments may be believed only if compatible with the rationalization of deviations from the path discussed in the previous section. Both for its self-enforceability as real agreement and for its credibility as test, a path agreement has an additional desirable property: conclusions are robust to the epistemic priority choice (see Theorem 1).

A *complete agreement* assigns one action to each player at each information set. Thus, it can correspond to an entire (subgame perfect) equilibrium of the game.

The wide range of agreements between path and complete will turn out to be useful too. The next subsection provides the tools for the evaluation of any agreement.

### 3.3 Beliefs, rationality and rationalizability

The elementary belief that an agreement may be able to induce is that the co-players will comply with it. This belief can be represented as a restricted set of conjectures about which strategies the co-players are going to play. In this dynamic framework, beliefs are modeled as Conditional Probability Systems (Renyi, [29]; henceforth, CPS). Here I define CPS's directly for the problem at hand.

**Definition 3** *A Conditional Probability System on  $(S_{-i}, (S_{-i}(h))_{h \in H})$  is a mapping  $\mu(\cdot|\cdot) : 2^{S_{-i}} \times \{S_{-i}(h)\}_{h \in H} \rightarrow [0, 1]$  satisfying the following axioms:*

*CPS-1 for every  $C \in (S_{-i}(h))_{h \in H}$ ,  $\mu(\cdot|C)$  is a probability measure on  $S_{-i}$ ;*

*CPS-2 for every  $C \in (S_{-i}(h))_{h \in H}$ ,  $\mu(C|C) = 1$ ;*

*CPS-3 for every  $E \in 2^{S_{-i}}$  and  $C, D \in (S_{-i}(h))_{h \in H}$ , if  $E \subseteq D \subseteq C$ , then  $\mu(E|C) = \mu(E|D)\mu(D|C)$ .*

*The set of all CPS's on  $(S_{-i}, (S_{-i}(h))_{h \in H})$  is denoted by  $\Delta^H(S_{-i})$ .*

For brevity, the conditioning events will be indicated with just the information set, which represents all the information acquired by players through observation.

I say that players believe in the agreement if at every information set they believe in what the agreement prescribes at the subsequent ones. In the Online Appendix, an example shows the important reason why players are not required to believe in strategies that comply with the agreement at information sets which have not been reached and cannot be reached anymore.

**Definition 4 (Belief in the agreement)** *Fix an agreement  $e = (e_i)_{i \in I}$ . For every  $i \in I$ , let  $\Delta_i^e$  be the set of all  $\mu_i = (\mu_i(\cdot|h))_{h \in H} \in \Delta^H(S_{-i})$  such that for every  $h \in H$ :*

$$\text{Supp} \mu_i(\cdot|h) \subseteq \left\{ (s_j)_{j \neq i} \in S_{-i}(h) : \forall j \neq i, \forall \tilde{h} \succeq h, s_j(\tilde{h}) \in e_j(\tilde{h}) \right\}.$$

Note that any two players have the same restrictions about any third player. For a path agreement on  $z \in Z$ ,  $\Delta_i^e$  coincides with the set of  $\mu_i \in \Delta^H(S_{-i})$  such that  $\mu_i(S_{-i}(z)|h^0) = 1$ : path agreements will be treated keeping this specification in mind.

For any agreement  $e = (e_i)_{i \in I}$  and each player  $i \in I$ ,  $\Delta_i^e$  is compact:<sup>28</sup> this allows to claim that selective rationalizability and strong- $\Delta$ -rationalizability with  $(\Delta_i^e)_{i \in I}$  have the epistemic characterizations of [13] and, for instance, [8].

I consider players who reply rationally to their conjectures. By rationality I mean that players, at every information set, choose an action that maximizes expected utility given the conjecture about how co-players will play and the expectation to reply rationally again in the continuation of the game. This is equivalent (see Battigalli [3]) to playing a *sequential best reply* to the CPS.

**Definition 5** Fix  $\mu_i \in \Delta^H(S_{-i})$ . A strategy  $s_i \in S_i$  is a *sequential best reply* to  $\mu_i$  if for every  $h \in H(s_i)$ ,  $s_i$  is a *continuation best reply* to  $\mu_i(\cdot|h)$ , i.e. for every  $\tilde{s}_i \in S_i(h)$ ,

$$\sum_{s_{-i} \in \text{Supp}\mu_i(\cdot|h)} u_i(\zeta(s_i, s_{-i}))\mu_i(s_{-i}|h) \geq \sum_{s_{-i} \in \text{Supp}\mu_i(\cdot|h)} u_i(\zeta(\tilde{s}_i, s_{-i}))\mu_i(s_{-i}|h).$$

I say that a strategy  $s_i$  is *rational* if it is a sequential best reply to some  $\mu_i \in \Delta^H(S_{-i})$ . The set of sequential best replies to  $\mu_i$  is denoted by  $\rho(\mu_i)$ . The set of best replies to a conjecture  $\nu_i \in \Delta(S_{-i})$  in the normal form of the game is denoted by  $r(\nu_i)$ . Note that if  $s_i \in \rho(\mu_i)$ ,  $[s_i] \subseteq \rho(\mu_i)$ , and if  $s_i \in r(\nu_i)$ ,  $[s_i] \subseteq r(\nu_i)$ . Moreover, the following relationship between continuation and sequential best replies holds: if  $s_i$  is a continuation best reply to  $\mu_i(\cdot|h)$  and  $h \in H(\rho(\mu_i))$ , there exists  $\hat{s}_i \in \rho(\mu_i)(h)$  such that for every  $\tilde{h} \succeq h$ , if  $\tilde{h} \in H(\text{Supp}\mu_i(\cdot|h))$ , then  $\hat{s}_i(\tilde{h}) = s_i(\tilde{h})$ . The reason is that a continuation best reply at  $h$  has to be a continuation best reply at every history  $h'$  that can be reached with positive probability from  $h$ , because the conjecture at  $h'$  is the conjecture at  $h$  conditional on  $h'$ .

Here I take the view that players refine further their conjectures through strategic reasoning based on beliefs in rationality and beliefs in the belief in the agreement. If players assign *epistemic priority to (the beliefs in) the agreement* (see the Introduction and [13]), the appropriate solution concept is *strong- $\Delta$ -rationalizability*. I translate its ultimate definition by Battigalli and Prestipino [8] in the framework of this paper.

**Definition 6 (strong- $\Delta$ -rationalizability)** Fix an agreement  $e = (e_i)_{i \in I}$ .

<sup>28</sup>For each  $h \in H$ , the set of probability measures that assign probability 1 to a particular subset of  $S_{-i}(h)$  is compact in  $\Delta(S_{-i})$ , endowed with the topology of weak convergence. Thus, the set of arrays of probability measures that believe in the remainder of the agreement at every information set is compact in  $(\Delta(S_{-i}))^H$ , endowed with the product topology. Since  $\Delta_i^e$  is the intersection between such set and the compact set  $\Delta^H(S_{-i})$ , it is compact itself.

(Step 0) For every  $i \in I$ , let  $S_{i,\Delta^e}^0 = S_i$ .

(Step  $n > 0$ ) For every  $i \in I$  and  $s_i \in S_i$ , let  $s_i \in S_{i,\Delta^e}^n$  if and only if there exists  $\mu_i \in \Delta_i^e$  such that:

D1  $s_i \in \rho(\mu_i)$

D2  $\forall q = 0, \dots, n-1, \forall h \in H, S_{-i,\Delta^e}^q(h) \neq \emptyset \Rightarrow \mu_i(S_{-i,\Delta^e}^q | h) = 1$  (i.e.  $\mu_i$  strongly believes  $S_{-i,\Delta^e}^q$ );

Finally let  $S_{i,\Delta^e}^\infty = \bigcap_{n \geq 0} S_{i,\Delta^e}^n$ . The profiles in  $S_{\Delta^e}^\infty$  are called strongly- $\Delta$ -rationalizable.

Strong rationalizability (Battigalli and Siniscalchi, [10]) can be seen as a special case of strong- $\Delta$ -rationalizability where no restriction applies ( $\Delta_i^e = \Delta^H(S_{-i})$ ) and it will be denoted by dropping the subscript  $\Delta^e$ .

The  $n$ -th step of strong- $\Delta$ -rationalizability captures the first  $n$  assumptions below:

1. players are rational and believe that co-players will comply with the agreement (and believe that everyone else will comply with the agreement, and so on);
2. players believe that 1 holds as long as not contradicted by observation (i.e. at every information set that can be reached if 1 holds);
3. players believe that 1 and 2 hold as long as not contradicted by observation;
4. ...

The sentence in brackets in assumption 1 means that players commonly believe in the agreement at every information set, regardless of the compatibility with beliefs in rationality of any order. This is not a necessary assumption to characterize strong- $\Delta$ -rationalizability: the belief in the agreement of order  $n$  has no behavioral implication once the belief in rationality of order  $n-1$  is dropped. Both epistemic characterizations can be found in [8].

Strong- $\Delta$ -rationalizability does two things.

First, it constitutes a compatibility test for the belief in the agreement with the strategic reasoning hypotheses. If strong- $\Delta$ -rationalizability delivers an empty set, the agreement does not pass the test. This happens when a player at some step allows an information set only with strategies that do not comply with the remainder

of the agreement, hence they are not compatible with the first-order-belief restrictions of the co-players.

Second, if the agreement passes the test, the strongly- $\Delta$ -rationalizable strategy profiles coincide with the behavioral implications of rationality, (common) belief in the agreement, and common strong belief in, jointly, rationality and (common) belief in the agreement.

If players assign *epistemic priority to rationality* (see the Introduction and [13]), the appropriate solution concept is *selective rationalizability* [13].

**Definition 7 (selective rationalizability)** *Fix an agreement  $e = (e_i)_{i \in I}$ . Let  $((S_j^m)_{j \in I})_{m \geq 0}$  be the strong rationalizability procedure.*

(Step 0) *For every  $i \in I$ , let  $S_{i,R\Delta^e}^0 = S_i$ .*

(Step  $n > 0$ ) *For every  $i \in I$  and  $s_i \in S_i$ , let  $s_i \in S_{i,R\Delta^e}^n$  if and only if there exists  $\mu_i \in \Delta_i^e$  such that:*

*S1  $s_i \in \rho(\mu_i)$ ;*

*S2  $\forall q = 0, \dots, n-1, \forall h \in H, S_{-i,R\Delta^e}^q(h) \neq \emptyset \implies \mu_i(S_{-i,R\Delta^e}^q | h) = 1$ ;*

*S3  $\forall q \geq 0, \forall h \in H, S_{-i}^q(h) \neq \emptyset \implies \mu_i(S_{-i}^q | h) = 1$ .*

*Finally let  $S_{i,R\Delta^e}^\infty = \bigcap_{n \geq 0} S_{i,R\Delta^e}^n$ . The profiles in  $S_{R\Delta^e}^\infty$  are called *selectively-rationalizable*.*

The  $n$ -th step of selective rationalizability captures the first  $n$  assumptions below:

1. players believe in the agreement, are rational, believe that co-players are rational (as long as not contradicted by observation), and so on;
2. players believe that 1 holds as long as not contradicted by observation;
3. players believe that 1 and 2 hold as long as not contradicted by observation;
4. ...

The first step already captures common strong belief in rationality [10]. This is necessary to guarantee that at every information set players give epistemic priority to the beliefs in rationality of all order which are compatible with the information



set (see the Discussion section of [13] for details). Hence, selective rationalizability refines strong rationalizability by selecting the strategies that are compatible with the beliefs in the agreement and consequent forward induction reasoning.

Selective rationalizability accomplishes the same two tasks of strong- $\Delta$ -rationalizability. A non-empty set of selectively-rationalizable strategy profiles coincides with the behavioral implications of (i) rationality and common strong belief in rationality, (ii) belief in the agreement, and common strong belief in, jointly, (i) and (ii).

As already argued, it seems more plausible that players will assign epistemic priority to rationality in presence of a non-binding agreement. However, for robustness purposes and theoretical insight, the next section evaluates agreements with both rationalizability tools.

## 4 Self-enforceability and implementability

### 4.1 Credibility and self-enforceability

In order to evaluate a given agreement, two features have to be investigated. First, whether the agreement is credible or not. Second, if the agreement is credible, whether players will comply with it or not.

An agreement is *credible* if it passes the appropriate rationalizability test. For brevity, the generic rationalizability procedure will be indicated with the subscript  $e$ ; it has to be replaced with  $\Delta^e$  or  $R\Delta^e$  to obtain the correct expressions under the chosen epistemic priority assumption.

**Definition 8 (Credibility)** *An agreement  $e = (e_i)_{i \in I}$  is credible if  $S_e^\infty \neq \emptyset$ .*

Credibility does not imply that players will comply with the agreement, but only that they may do so. If players, once they refine their conjectures according to the agreement, always have the strict incentive to comply with it, the agreement is *self-enforcing*.

**Definition 9 (Self-enforceability)** *An agreement  $e = (e_i)_{i \in I}$  is self-enforcing if it is credible and for every  $i \in I$ ,  $s_i \in S_{i,e}^\infty$  and  $h \in H(S_e^\infty) \cap H(s_i)$ ,  $s_i(h) \in e_i(h)$ .*

The definition requires explicitly that *every* rationalizable strategy complies with the agreement at the information sets which can be reached under (itself and) the

rationalizable strategy profiles, so that violation of the agreement will actually occur. Credibility implies that players also believe in the compliance with the agreement at the other information sets, so that the desired behavioral consequences apply.

One might think that a self-enforcing agreement under priority to rationality will be, a fortiori, self-enforcing under priority to the agreement. An example in the appendix of [13] shows that this is not the case. Instead, for a path agreement credibility and self-enforceability are robust to the epistemic priority assumption.

**Theorem 1** *A path agreement is credible/self-enforcing under priority to rationality if and only if it is credible/self-enforcing under priority to the agreement.*

Hence, all the conditions for credibility and self-enforceability of path agreements can be stated regardless of the epistemic priority assumption, and the proofs can rely on either selective rationalizability or strong- $\Delta$ -rationalizability.

The first condition claims that a path can be credible only if it is induced by some strongly rationalizable strategy profile.<sup>29</sup>

**Proposition 1** *Fix  $z \in Z$ . If  $z \notin \zeta(S^\infty)$ , the corresponding path agr. is not credible.*

**Proof.** Selective rationalizability is a refinement of strong rationalizability. ■

Thus, only a strongly rationalizable outcome can pass the credibility test.<sup>30</sup>

This necessary condition for credibility becomes sufficient for self-enforceability when the path is the only strongly rationalizable one.

**Proposition 2** *If  $\zeta(S^\infty)$  is a singleton, the corresponding path agr. is self-enforcing.*

**Proof.** Let  $\zeta(S^\infty) = \{z\}$ . For every  $i \in I$ ,  $S_i^\infty \subseteq S_i(z)$ . Thus, for every  $j \in I$  and  $\mu_j \in \Delta^H(S_{-j})$  that satisfies S3,  $\mu_j(S_{-j}(z)|h^0) = 1$ , so  $\mu_j \in \Delta_j^e$ . Hence, the restrictions are immaterial, so that selective rationalizability coincides with strong rationalizability after convergence.<sup>31</sup> ■

<sup>29</sup>Moreover, results of Battigalli and Siniscalchi [11] imply that only a self-confirming equilibrium outcome (Fudenberg and Levine [18]; called conjectural equilibrium in Battigalli [2]) can pass the credibility test. However, for implementability the outcome will need to be Nash.

<sup>30</sup>By virtue of Remark 7 in the Appendix, this property holds true if the test is extended to distributions over outcomes rather than just one outcome. The extension to incomplete information games would prove that the Iterated Intuitive Criterion selects strongly rationalizable outcomes.

<sup>31</sup>Selective rationalizability without actual restrictions could be empty when there exists  $i \in I$  and  $s_i \in S_i^\infty$  such that for every  $\mu_i$  that strongly believes  $(S_{-i}^n)_{n \geq 0}$ ,  $s_i \notin \rho(\mu_i)$ . This is never the case under finiteness or mild regularity conditions.

Analogously, while SPE paths can be self-enforcing or not, non SPE paths can be credible but not self-enforcing.

**Theorem 2** *Fix  $z \in Z$ . If the corresponding path agreement is self-enforcing, then there exists a SPE<sup>32</sup>  $(\sigma_i)_{i \in I} \in \times_{i \in I} \Delta(S_i)$  such that for every  $i \in I$ ,  $\sigma_i(S_i(z)) = 1$ .*

Differently than for strongly rationalizable paths, when the game has only one SPE path, it needs not be self-enforcing.

$A$	-	-	→	$A/B$	$L$	$C$	$R$
$O$	↓		$I$	$U$	(5, 0)	(0, 5)	(0, 0)
	(4, ·)			$M$	(0, 5)	(5, 0)	(0, 0)
				$D$	(3, 0)	(3, 0)	(3, 3)

In the subgame, there is no equilibrium where both  $U$  and  $M$  are played: if Ann is indifferent between them, she prefers  $D$ . If  $U$  was played and  $M$  not, Bob would not play  $L$ , so  $U$  cannot be a best reply, and analogously for  $M$ . Thus, the unique equilibrium is  $(D, R)$  and it induces Ann to choose the outside option. Yet, Ann may be sufficiently confident of  $L$  or  $C$  to enter the subgame and play  $U$  or  $M$ .<sup>33</sup> This also implies that the subgame perfect threat  $R$  is not credible, being a best reply only to  $D$ . The path is however credible. It can be shown that a path is credible when it is the unique SPE path or, more generally, when it is robust to any choice of off-the-path equilibria:<sup>34</sup> this is the notion of credibility used by Gossner [19] in his work on incomplete codes.

When SPE paths are more than one, none of them may be even credible.<sup>35</sup> Consider the twofold repetition of the following game. The players must perform a task that yields a profit of 3 to each of them at the total effort cost of 2. If at least one

---

<sup>32</sup>SPE in mixed strategies are defined in the Appendix, using additional notation. For every SPE in mixed strategy there is a SPE in behavioral strategies that induces the same outcome distribution, thus the Theorem holds also with SPE in behavioral strategies.

<sup>33</sup>So that Bob in turn may reply with  $L$  or  $C$  and all 4 actions remain rationalizable.

<sup>34</sup>For a path not to be credible, some player at some step of the rationalizability procedure must have the incentive to abandon it for every reaction she may expect to her deviation. The inner recursive step in the proof of Lemma 10 in the Appendix shows that in such case a (probabilistic) SPE path of the subgame would have survived until that step. This shows that the non credible path cannot fall in the class I am considering.

<sup>35</sup>I conjecture that with agreements on mixed actions, a credible SPE outcome distribution exists. This would extend the main result of Govindan and Wilson [20] to full-fledged forward induction reasoning.

player works, the task is performed; if only one player works, she pays the total cost of effort; if they both work, they share the effort cost equally.

$A \setminus B$	<i>Work</i>	<i>FreeRide</i>
<i>W</i>	2, 2	1, 3
<i>FR</i>	3, 1	0, 0

No deterministic SPE path is even credible. If the path prescribes the same Nash in both stages, the unhappy player can signal with a deviation the intention to switch to the preferred equilibrium in the second stage. If the path prescribes to play one Nash in the first stage and the other Nash in the second stage, the player whose preferred equilibrium is played in the first stage can deviate from it to signal the intention to play it in the second stage. Similar paths have been already classified by Osborne [26] for 2-players, finitely-repeated games as *equilibrium paths that can be upset by a convincing deviation*; in the Online Appendix, the definition and the formal reasoning that makes them not credible are provided.

For complete agreements the equivalence between self-enforceability under priority to the agreement and under priority to rationality does not hold. When the epistemic priority falls on the agreement, a "strict" SPE is self-enforcing. I say that an equilibrium  $s \in S$  is *strict* when for every  $i \in I$ ,  $r_i(s_{-i}) \subseteq S_i(\zeta(s))$ . The proof of this result is in the Online Appendix.

When the epistemic priority falls on rationality, instead, not all SPE are credible agreements. A SPE outcome may not be strongly rationalizable, hence also not selectively rationalizable. Moreover, even if the outcome is strongly rationalizable, the SPE strategy profile may be not.<sup>36</sup> Yet, the outcome may be induced also by some Nash equilibrium in strongly rationalizable strategies. Or, like in the first example of Section 2, players may be interested in a non subgame perfect equilibrium. For Nash equilibria, also when the priority falls on the agreement, the complete agreement may not work, because an off-the-path threat may not be optimal against any clear-cut intention of the deviator or of a third player. Here I illustrate an example of the first

---

<sup>36</sup>Even in a perfect information game like the variation of the centipede in Reny [28], the agreement corresponding to the only SPE is not credible, because among the strongly rationalizable strategy profiles, despite inducing the SPE outcome, there is not the SPE one (forward induction and backward induction moves do not coincide)

kind; an example of the second kind is provided at the end of the next subsection.

$Ann$	-	-	→	$A \setminus B$	$L$	$C$	$R$
$O$	↓	$I$		$U$	$(9, 0)$	$(8, 3)$	$(1, 2)$
	$(4, 4)$			$D$	$(8, 3)$	$(9, 0)$	$(0, 2)$

In this game,  $O$  is not a SPE outcome because in any equilibrium of the subgame,  $R$  cannot be played with probability higher than  $1/2$  (otherwise Ann would play  $U$  and Bob  $C$ ). The Nash equilibria  $(O, U, R)$  and  $(O, D, R)$  are not self-enforcing agreements: whatever Ann promises to do in case she deviates to  $I$ ,  $R$  is not a best reply. Yet, if Bob refuses to "listen" to Ann's post-deviation promises, the incomplete agreement on  $O$  and  $R$  is self-enforcing (regardless of the epistemic priority assumption).

Therefore, when players want to implement some equilibrium of the game and are prone to discuss and believe off-the-path threats, it is more appropriate to tackle the problem through the alternative perspective: can the outcome be implemented by *some* agreement? This issue is addressed in the next subsection.

## 4.2 Implementability

Taking the opposite, implementation perspective, which outcomes of the game can be achieved through *some* agreement? I say that a set of outcomes is *implementable* if for some agreement the appropriate rationalizability procedure delivers a non-empty subset of it.

**Definition 10 (Implementability)** *A set of outcomes  $P \subset Z$  is implementable if there exists a credible agreement  $e = (e_i)_{i \in I}$  such that  $\zeta(S_e^\infty) \subseteq P$ .*

In the Online Appendix I discuss the relationship between implementability and two much weaker but somehow related solution concepts: *Extensive Form Best Response Sets* (Battigalli and Friedenberg, [7]) and *Mutually acceptable courses of actions* (Greenberg et al., [21]).

Focusing on a single outcome, one should first check if the corresponding path agreement is self-enforcing.<sup>37</sup> If this is not the case, regardless of the epistemic priority assumption, implementability by loosening restrictions can be excluded.

---

<sup>37</sup>Given the results of Section 4.1, one can hope so only for a strongly rationalizable, SPE outcome.

**Theorem 3** Fix  $z \in Z$  and the corresponding path agreement  $e = (e_i)_{i \in I}$ . If  $z$  is implemented by some agreement  $e' = (e'_i)_{i \in I}$  such that for every  $i \in I$  and  $h \in H$ ,  $e_i(h) \subseteq e'_i(h)$ , then  $e = (e_i)_{i \in I}$  is self-enforcing.

Thus, if players want to reach a given outcome, without willing or trusting the possibility to coordinate in case of deviation, they cannot do any better than agreeing on the path. This is a sort of "revelation principle" for agreements design: leaving some mystery about on-the-path moves cannot be of any help for the goal.

Allowing for off-the-path restrictions, is it possible to implement a non-equilibrium outcome? The answer is negative under both epistemic priority assumptions. For the remainder of this section, the needed additional arguments for the priority to rationality/selective rationalizability case are provided in square brackets.

**Proposition 3** Fix  $z \in Z$ . If  $z$  is implemented by  $e = (e_i)_{i \in I}$ , then there exists a strict Nash equilibrium  $s \in S_e^\infty$  such that  $\zeta(s) = z$ .

**Proof.** Fix  $i \in I$ . By credibility there exists  $\mu_i \in \Delta_i^e$  such that  $\mu_i(S_{-i,e}^\infty | h^0) = 1$ .<sup>38</sup> So, for every  $j \neq i$ , there exists  $s_j \in S_{j,e}^\infty$  such that for every  $h \in H$ ,  $s_j(h) \in e_j(h)$ . By implementability  $s_j \in S_j(z)$ .

Fix  $\mu_i \in \Delta_i^e$  that strongly believes  $(S_{-i,e}^n)_{n=0}^\infty$  [and  $(S_{-i}^n)_{n=0}^\infty$ ] such that  $\mu_i(s_{-i} | h^0) = 1$ . Suppose by contradiction that  $\emptyset \neq r(\mu_i(\cdot | h^0)) \setminus S_i(z) \ni \tilde{s}_i$ . Then there would exist  $\hat{s}_i \in \rho(\mu_i)$  such that for every  $h \in H(\tilde{s}_i) \cap H(s_{-i})$ ,  $\tilde{s}_i(h) = \hat{s}_i(h)$ ,<sup>39</sup> so that  $\hat{s}_i \notin S_i(z)$  too, contradicting implementability. So  $s_i \in r(s_{-i}) = r(\mu_i(\cdot | h^0)) = S_i(z)$ . ■

The condition is not sufficient. In the last example,  $(R_1, L_2)$  is a strict Nash outcome that cannot be implemented by any agreement, because it is sustained by  $R_3$  and  $R_4$  but  $R_3$  is not rational against  $R_4$ . If players are more than two, off-the-path threats of two different players, even if rational, may be incompatible with each other. With two players, still, a threat could be rational but not optimal against any clear-cut intention of the deviator. In the second case, there are clear conditions under which the problem can be solved by leaving the agreement incomplete (see the end of the previous subsection for an example).

<sup>38</sup>This may be false when there exist  $i \in I$  and  $s_i \in S_{i,e}^\infty$  such that for every  $\mu_i$  that strongly believes  $(S_{i,e}^n)_{n \geq 0}$ ,  $s_i \notin \rho(\mu_i)$ . This is never the case under finiteness or mild regularity conditions.

<sup>39</sup>See the relationship between continuation best replies and sequential best replies in Section 3.

**Proposition 4** *Consider a 2-players game and a strict Nash equilibrium  $(s_1, s_2)$ . If  $s_1$  and  $s_2$  are rational [strongly rationalizable],  $z := \zeta(s)$  is implementable.*

**Proof.** For every  $i \in I$  and  $h \in H(s_i)$ , let  $e_i(h) = s_i(h)$ , otherwise, let  $e_i(h) = A_i(h)$ . Since  $s_i$  is rational [ $s_i \in S_i^\infty$ ], there exists  $\mu_i \in \Delta^H(S_{-i})$  [that strongly believes  $(S_{-i}^n)_{n=0}^\infty$ ]<sup>40</sup> such that  $s_i \in \rho(\mu_i)$ . For every  $\bar{\mu}_i \in \Delta_i^e$ ,  $\bar{\mu}_i([s_{-i}] | h^0) = 1$  and for every  $s'_{-i} \in [s_{-i}]$ , since by strictness  $r(s_{-i}) \subseteq S_i(z)$ ,  $r(s'_{-i}) \subseteq S_i(z)$  too. Hence,  $S_{i,e}^1 \subseteq S_i(z)$ , yielding implementability in case of credibility. Then, for every  $i \in I$  and  $m \in \mathbb{N}$  such that  $s \in S_e^{m-1}$ , there exists  $\tilde{\mu}_i \in \Delta_i^e$  that strongly believes  $(S_{-i,e}^n)_{n=0}^{m-1}$  [and  $(S_{-i}^n)_{n=0}^\infty$ ] such that  $\tilde{\mu}_i(s_{-i} | h^0) = 1$  and for every  $h \notin H(S_{-i}(z))$ ,  $\tilde{\mu}_i(\cdot | h) = \mu_i(\cdot | h)$ . For every  $h \prec z$ ,  $\tilde{\mu}_i(s_{-i} | h) = 1$  and  $s_i \in r(s_{-i})$ ; for every  $h \in H(s_i) \setminus H(S_{-i}(z))$ ,  $s_i$  is a continuation best reply to  $\tilde{\mu}_i(\cdot | h)$  by construction; else,  $h \notin H(S_i(z))$ , so  $h \notin H(s_i)$ . Thus,  $s_i \in \rho(\tilde{\mu}_i) \subseteq S_{i,e}^m$ . Inductively,  $s \in S_e^\infty \neq \emptyset$ . ■

Hence, in 2-players games I can fully characterize the set of implementable outcomes with equilibria in *pure* strategies as follows.

**Theorem 4** *Consider a 2-players game. An outcome  $z \in Z$  is implementable under priority to the agreement [rationality] if and only if there exists a rational [strongly rationalizable] strict Nash equilibrium  $s \in S$  such that  $\zeta(s) = z$ .*

**Proof.** The "if" direction coincides with Proposition 4. The "only if" direction coincides with Proposition 3 once observed that if  $s \in S_{\Delta^e}^\infty$ ,  $s$  is rational and if  $s \in S_{R\Delta^e}^\infty \subseteq S^\infty$ ,  $s$  is strongly rationalizable. ■

Then, since a rationalizable strategy is obviously rational, differently than for self-enforceability, implementability under priority to rationality implies implementability under priority to the agreement (at least in two-players games).

**Corollary 5** *In a 2-players game, if an outcome is implementable under priority to rationality, then it is implementable under priority to the agreement.*

In games with more than two players, there is no equilibrium equivalent of implementable outcomes. Here I provide a 3-players game where a non subgame perfect

---

<sup>40</sup>Under priority to rationality, the existence is guaranteed by finiteness or mild regularity conditions. Otherwise, the Proposition can be restated requiring such CPS's to exist.

equilibrium outcome is implementable thanks to an incomplete agreement (and incompleteness is not only on the side of the deviator).

				<i>Ann</i>	--	---	--	→	(4, 4, 4)
				<i>I</i>					
				↓					
<i>Cleo</i>	<i>Bob</i>								
	<i>M1</i>	<i>L</i>	<i>C</i>	<i>R</i>	--	<i>M2</i>	<i>L</i>	<i>C</i>	<i>R</i>
<i>Ann</i>	<i>U</i>	(8, 3, 0)	(9, 0, 0)	(1, 2, 1)		<i>U</i>	(8, 3, 0)	(9, 3, 0)	(1, 2, 0)
	<i>M</i>	(9, 3, 0)	(8, 3, 0)	(0, 2, 0)		<i>D</i>	(9, 0, 0)	(8, 3, 0)	(0, 2, 1)

Ann moves first and chooses between a "fair" outside option and a simultaneous moves game with Bob and Cleo, where Cleo chooses the matrix. In any equilibrium of the subgame, *R* cannot be played with probability higher than 1/2, otherwise Ann would choose *U* and then Bob would switch to *L*. Thus, *O* is not a SPE outcome. Still, Bob can credibly threaten to play *R*. For this to hold true, it is crucial that Cleo does not reveal what she is going to play: for each matrix, Bob has a safe option that dominates *R*. If Ann plays *I*, either she is irrational or she has not believed in the threat. Under both interpretations, she may play *U* or *M* and Cleo may react with *M1* or *M2*. Thus, it is credible that Bob will react with *R*.

The definition of implementability can be strengthened with the path agreement credibility test. First, the outcome should pass the test. Second, the agreement that implements the outcome should be compatible with forward induction reasoning based on the beliefs in the path. A way to obtain this is to derive the agreement from the conjectures under which players remain on path at the end of the test. Yet, such conjectures may not be convertible into an agreement, because of: (correlated) randomizations; strategic incompatibility between the threats of different players; without observable actions, the disagreement of two deviators about the reaction to the deviation of a third player who cannot recognize who the deviator is. Strategic incompatibilities may be solved with an incomplete agreement. All the outcomes implemented in the examples are credible, and the threats are compatible with forward induction reasoning about the path.

If one insists on SPE outcomes, under priority to rationality an important preliminary question arises: is there always the support of a SPE outcome distribution among strongly rationalizable outcomes? The following theorem states that it is ac-



tually so<sup>41</sup> for the wide class of games with observable actions:<sup>42</sup> subgame perfection and strong rationalizability never give completely disjoint predictions.

**Theorem 6** *Consider the set of strongly rationalizable strategy profiles  $S^\infty$ . There exists a SPE  $(\sigma_i)_{i \in I} \in \times_{i \in I} \Delta(S_i)$  and an equilibrium  $(\tilde{\sigma}_i)_{i \in I} \in \times_{i \in I} \Delta(S_i^\infty)$  such that for every  $z \in Z$ ,  $\prod_{i \in I} \sigma_i(S_i(z)) = \prod_{i \in I} \tilde{\sigma}_i(S_i(z))$ .*<sup>43</sup>

## 5 Directions for future research

The issue of compliance with agreements could be interestingly analyzed in psychological games (see for instance Battigalli and Dufwenberg [6]) or games with ambiguity averse players. Both ambiguity aversion and belief-dependent payoffs, like in the case of guilt-averse players, could sustain the self-enforceability a wider range of agreements, including path ones. Other psychological considerations could motivate formally the preference for path agreements, which do not involve the discussion of what to do in case someone defects.

The methodology developed here can be applied to a wide range of economic problems, from non-binding commitment of governments (see for instance Bassetto [1] on capital taxation) to dynamic collusion in oligopolies. On this topic, Harrington [22] documents instances of mutual partial understanding among firms whose implications can be understood with the methodology of this paper. Another application which is close in spirit to this work can be found in Gossner [19]. In his study of incomplete codes, he calls self-enforcing the equilibrium paths that can be sustained by any choice of equilibria in the off-the-path subgames. As already argued, such paths are at least credible in the sense of this paper (but, so far, using the finite horizon hypothesis).

Although the focus of the paper has been kept on complete information for interpretative and notational easiness, the methodology can be applied also to dynamic games with incomplete information. In an incomplete information environment, players reaching an agreement at the interim stage may reveal their types in the bargaining process or not. In the first case, first-order-beliefs restrictions could be extended to

---

<sup>41</sup>And these outcomes are all induced by a mixed rationalizable Nash, an important feature if one wants to go on proving implementability with probabilistic agreements.

<sup>42</sup>I conjecture that the result holds for all extensive form games with perfect recall.

<sup>43</sup>While  $\times$  denotes a cartesian product,  $\prod$  will always denote an algebraic product.

payoff-relevant types. In the second case, what the agreement suggests about co-players types would be embodied in the rationalizability procedure.

Finally, there are two complementary issues, already investigated in the literature, to which this analysis could be profitably connected. The first is the pre-play bargaining issue. This paper sheds light on which agreements are self-enforcing and which outcomes are implementable through agreements. How will players ultimately choose among them? Welfare considerations and a theory of bargaining could refine the answer. Dufwenberg, Servátka and Vadovic [16] propose an interesting approach to the pre-play bargaining issue. Miller and Watson [25], instead, analyze a problem of bargaining between players who can reach an agreement at every stage of a repeated game. When players can communicate during the game, the second issue is renegotiation proofness (see, for instance, [17]). Past moves and consequent forward induction considerations could influence the bargaining power of players in renegotiation.

## 6 Appendix

**Additional notation (for a subgame  $\Gamma(h)$  of  $\Gamma = \langle I, \bar{H}, (u_i)_{i \in I} \rangle$ ).**

Histories of  $\Gamma(h)$  are identified for convenience by the histories of  $\Gamma$  following  $h$ , and not redefined as shorter lists of action profiles occurring *after*  $h$ .

- For any  $h \in H$  and  $i \in I$ :

- $H^h := \{\tilde{h} \in H : \tilde{h} \succeq h\}$ ,  $Z^h := \{z \in Z : z \succ h\}$ ;
- $S_i^h$  is the set of all strategies  $s_i^h : \tilde{h} \in H^h \mapsto a_i \in A_i(\tilde{h})$  of  $\Gamma(h)$ .

- for any  $J \subseteq I$ ,  $z \in H^h \cup Z^h$ ,  $\hat{h} \in H^h$  and subsets of strategies  $(\bar{S}_j^h)_{j \in I}$ :

- $\bar{S}_J^h(z)$ ,  $\bar{S}_J^h[\hat{h}]$ ,  $H(\bar{S}_J^h)$ ,  $\Delta^{H^h}(S_{-i}^h)$ ,  $\zeta(\cdot)$ ,  $r(\cdot)$  and  $\rho(\cdot)$  are defined like in  $\Gamma$ ;
- $D_i(\bar{S}^h) := \{\tilde{h} \in H : \exists \bar{h} \in H(\bar{S}^h), \exists a_i \in A_i(\bar{h}) \setminus \bar{S}_i^h[\bar{h}], \exists a_{-i} \in \bar{S}_{-i}^h[\bar{h}], \tilde{h} = (\bar{h}, (a_i, a_{-i}))\}$ ; <sup>44</sup>

- for any  $(s_j^h)_{j \in I} \in S^h$ ,  $(\hat{s}_j^h)_{j \in I} \in \hat{S}^h$ ,  $\mu_i^h \in \Delta^{H^h}(S_{-i}^h)$ ,  $\tilde{\mu}_i^{\hat{h}} \in \Delta^{H^{\hat{h}}}(S_{-i}^{\hat{h}})$  and  $\hat{H} \subseteq H^{\hat{h}}$ :

- $s_i^h|_{\hat{h}}$  is the restriction of  $s_i^h$  to  $H^{\hat{h}}$ ,  $s_J^h|_{\hat{h}} := (s_j^h|_{\hat{h}})_{j \in J}$ ,  $\bar{S}_J^h|_{\hat{h}} := (\tilde{s}_j^h|_{\hat{h}})_{\tilde{s}_j^h \in \bar{S}_J^h}$ ;

---

<sup>44</sup>Set of histories that follow a unilateral deviation by player  $i$  from the histories induced by  $\bar{S}^h$ .

- $\widehat{r}(\mu_i^h, \widehat{h})$  is the set of continuation best replies to  $\mu_i^h(\cdot|\widehat{h})$ ;
- $s_J^h =_{\widehat{H}} \widehat{s}_J^h$  if for every  $\widetilde{h} \in \widehat{H}$ ,  $s_J^h(\widetilde{h}) = \widehat{s}_J^h(\widetilde{h})$ ;  $S_{J, \widehat{H}}^{h, \widehat{s}_J^h} := \{\widetilde{s}_J^h \in S_J^h : \widetilde{s}_J^h =_{\widehat{H}} \widehat{s}_J^h\}$ ;
- $\mu_i^h =_{\widehat{H}} \widehat{\mu}_i^h$  if for every  $\widetilde{h} \in \widehat{H}$  and  $\widetilde{z} \succ p(\widetilde{z}) \in \widehat{H} \cap H^{\widetilde{h}}$ ,  $\mu_i^h(S_{-i}^h(\widetilde{z})|\widetilde{h}) = \widehat{\mu}_i^h(S_{-i}^h(\widetilde{z})|\widetilde{h})$ ;
- $s_J^h =_{\widehat{h}} \widehat{s}_J^h$  and  $\mu_i^h =_{\widehat{h}} \widehat{\mu}_i^h$  if, respectively,  $s_J^h =_{H^{\widehat{h}}} \widehat{s}_J^h$  and  $\mu_i^h =_{H^{\widehat{h}}} \widehat{\mu}_i^h$ .

I will often use the fact that  $=_{\widehat{H}}$  is transitive and that if  $\mu_i^h =_{\widehat{H}} \widehat{\mu}_i^h$  and  $\widetilde{H} \subseteq \widehat{H}$ , then  $\mu_i^h =_{\widetilde{H}} \widehat{\mu}_i^h$ .<sup>45</sup> Moreover, note that by CPS-3 (with  $\widehat{h}$  as "intermediate" history):<sup>46</sup>

- ♠ if  $\widetilde{H} \cap H^{\widehat{h}} = \emptyset$ ,  $p(\widehat{h}) \in \widetilde{H}$ ,  $\mu_i^h =_{\widetilde{H}} \widetilde{\mu}_i^h$ ,  $\mu_i^h =_{\widehat{H}} \widehat{\mu}_i^h$  and  $\widetilde{\mu}_i^h =_{\widehat{H}} \widehat{\mu}_i^h$ , then  $\mu_i^h =_{\widetilde{H} \cup \widehat{H}} \widehat{\mu}_i^h$ .

Rewrite selective rationalizability as  $(S_i^0, \dots, S_i^M, S_{i, R\Delta^e}^0, S_{i, R\Delta^e}^1, \dots)_{i \in I}$ , where  $((S_i^q)_{i \in I})_{q \geq 0}$  is strong rationalizability and  $M$  is the smallest  $q \in \mathbb{N}$  such that  $S^q = S^{q+1} = S^\infty$  (it exists by finiteness of the game).

In a subgame, substrategies can be eliminated "exogenously" and not because they are not sequential best replies to any valid conjecture in the subgame. On the other hand, substrategies can survive even if the opponents do not reach the subgame anymore. Thus I define the following.

**Definition 11** Fix  $h \in H$ . A reduction procedure is a sequence  $((S_{i,q}^h)_{i \in I})_{q \geq 0}$  where:

- for every  $i \in I$ ,  $S_{i,0}^h = S_i^h$ ;
- for every  $n > 0$  and  $s_i^h \in S_{i,n}^h$ ,  $s_i^h \in S_{i,n-1}^h$  and there exists  $\mu_i^h \in \Delta^{H^h}(S_{-i}^h)$  s. t.:
  - $s_i^h \in \rho(\mu_i^h)$ ;
  - for every  $q < n$  and  $\widetilde{h} \in H^h$ , if  $S_{-i,q}^h(\widetilde{h}) \neq \emptyset$ , then  $\mu_i^h(S_{-i,q}^h|\widetilde{h}) = 1$ .

Let  $M$  be the smallest  $q \in \mathbb{N}$  such that  $S_{i,q}^h = S_{i,q+1}^h$  and let  $S_{i,\infty}^h := S_{i,M}^h$  (it exists by finiteness). Note that even if  $S_{-i,n-1}^h$  is empty,  $S_{i,n}^h$  needs not be empty. Still, for simplicity I will say that  $s_i^h \in S_{i,n}^h$  is justified by a  $\mu_i^h$  that strongly believes (henceforth, t.s.b.)  $(S_{-i,q}^h)_{q=0}^{n-1}$ . The definition of reduction procedure encompasses

<sup>45</sup>These properties clearly hold also when  $=_{\widehat{H}}$  applies between strategies.

<sup>46</sup>By CPS-3, for every  $\widetilde{h} \in \widetilde{H}$  and  $\widetilde{z} \succ p(\widetilde{z}) \in \widehat{H} \cap H^{\widetilde{h}}$ ,  $\mu_i^h(S_{-i}^h(\widetilde{z})|\widetilde{h}) = \mu_i^h(S_{-i}^h(\widehat{h})|\widetilde{h})\mu_i^h(S_{-i}^h(\widetilde{z})|\widehat{h}) = \widetilde{\mu}_i^h(S_{-i}^h(\widehat{h})|\widetilde{h})\widehat{\mu}_i^h(S_{-i}^h(\widetilde{z})|\widehat{h}) = \widetilde{\mu}_i^h(S_{-i}^h(\widehat{h})|\widetilde{h})\widetilde{\mu}_i^h(S_{-i}^h(\widetilde{z})|\widehat{h}) = \widetilde{\mu}_i^h(S_{-i}^h(\widetilde{z})|\widetilde{h})$ .

both strong- $\Delta$ -rationalizability and selective rationalizability (rewritten as before) and their implications in the subgames. Indeed, if  $((S_{i,q}^h)_{i \in I})_{q \geq 0}$  is a reduction procedure,  $((S_{i,q}^h(\widehat{h})|\widehat{h})_{i \in I})_{q \geq 0}$  is a reduction procedure for any  $\widehat{h} \in H^h$ .

Some lemmata claim the survival of strategies, or conjectures over such strategies, which combine substrategies that have survived by assumption. The reason why such lemmata are needed is merely the following. Fix  $\widehat{s}_i^h, \bar{s}_i^h \in S_{i,n}^h$  and  $\widehat{h}, \bar{h} \in H(\widehat{s}_i^h) \cap H(\bar{s}_i^h)$  such that  $\bar{h} \not\preceq \widehat{h} \not\preceq \bar{h}$ : there needs not exist  $s_i^h \in S_{i,n}^h(\widehat{h}) \cap S_{i,n}^h(\bar{h})$  such that  $s_i^h|\widehat{h} = \bar{s}_i^h|\widehat{h}$  and  $s_i^h|\bar{h} = \widehat{s}_i^h|\bar{h}$ . The intuitive reason is the following: player  $i$  may allow  $\widehat{h}$  and  $\bar{h}$  either because she is confident that  $\widehat{h}$  will be reached and she has certain expectations after  $\widehat{h}$ , or because she is confident that  $\bar{h}$  will be reached and she has certain expectations after  $\bar{h}$ . If  $\widehat{s}_i^h$  is best reply to the first conjecture and  $\bar{s}_i^h$  is best reply to the second conjecture,  $\widehat{s}_i^h|\bar{h}$  and  $\bar{s}_i^h|\widehat{h}$  may be "emergency plans" for an unpredicted contingency, after which the expectations would not have justified the choice to allow  $\bar{h}$  and  $\widehat{h}$  in the first place.

I start with a lemma about the combination of substrategies. Consider a player who may be surprised by history  $\widehat{h}$ , in the sense that she may play a strategy that allows  $\widehat{h}$  while believing that the co-players do not until  $\widehat{h}$  is actually reached. This player can keep the same beliefs and the same strategy out of  $\Gamma(\widehat{h})$ , whatever she believes the co-players would do and hence however she may play after  $\widehat{h}$ . Being a rather intuitive result, the proof is relegated to the Online Appendix.

**Lemma 1** *Fix a reduction procedure  $((S_{i,q}^h)_{i \in I})_{q \geq 0}$ ,  $i \in I$ ,  $n \in \mathbb{N}$ ,  $\widehat{h} \in H^h$  and  $\mu_i^h \in \Delta^{H^h}(S_{-i}^h)$  t.s.b.  $(S_{-i,q}^h)_{q=0}^{n-1}$  such that  $\mu_i^h(S_{-i}^h(\widehat{h})|p(\widehat{h})) = 0$ . Fix  $s_i^h \in \rho(\mu_i^h)$ ,  $\mu_i^{\widehat{h}}$  t.s.b.  $(S_{-i,q}^h(\widehat{h})|\widehat{h})_{q=0}^{n-1}$  and  $\widehat{s}_i^h \in \rho(\mu_i^{\widehat{h}})$ .*

*Consider the unique  $\widetilde{s}_i^h =_{\widehat{h}} \widehat{s}_i^h$  such that for every  $\widetilde{h} \notin H^h$ ,  $\widetilde{s}_i^h(\widetilde{h}) = s_i^h(\widetilde{h})$ .*

*There exists  $\widetilde{\mu}_i^h =_{\widehat{h}} \mu_i^{\widehat{h}}$  t.s.b.  $(S_{-i,q}^h)_{q=0}^{n-1}$  such that for every  $\widetilde{h} \notin H^h$ ,  $\widetilde{\mu}_i^h(\cdot|\widetilde{h}) = \mu_i^h(\cdot|\widetilde{h})$ , and  $\widetilde{s}_i^h \in \rho(\widetilde{\mu}_i^h)$  (so that  $\rho(\mu_i^h)(\widehat{h}) \neq \emptyset$  implies  $\rho(\widetilde{\mu}_i^h)(\widehat{h}) \neq \emptyset$ ).*

### PROOF OF THEOREMS 1 AND 3

Fix a set of outcomes and call  $\widehat{H}$  the set of histories of  $\Gamma(h)$  that precede them. Suppose that  $\mu_i^h =_{\widehat{H}} \bar{\mu}_i^h$  and that the co-players are not expected to deviate outside of  $\widehat{H}$ , i.e. they are expected to play compatibly with some of those outcomes at every history in  $\widehat{H}$ . If also the sequential best replies to  $\mu_i^h$  and the sequential best reply to  $\bar{\mu}_i^h$  do not deviate outside of  $\widehat{H}$ , then they mimic each other inside of  $\widehat{H}$ .

**Lemma 2** Fix  $l \in I$ ,  $h \in H$  and  $\widehat{Z} \subseteq Z^h$ . Let  $\widehat{H} := \{\tilde{h} \in H^h : \exists z \in \widehat{Z}, \tilde{h} \prec z\}$ . Fix  $\mu_l^h =^{\widehat{H}} \bar{\mu}_l^h$  such that for every  $\tilde{\mu}_l^h \in \{\mu_l^h, \bar{\mu}_l^h\}$ :

$$\rho(\tilde{\mu}_l^h) \subseteq \{\tilde{s}_l^h \in S_l^h : \forall \tilde{h} \in \widehat{H} \cap H(\tilde{s}_l^h), \forall \tilde{s}_{-l}^h \in \text{Supp}\tilde{\mu}_l^h(\cdot|\tilde{h}), H(\tilde{s}_l^h, \tilde{s}_{-l}^h) \subseteq \widehat{H}\}. \quad (\text{A0})$$

Then for every  $\tilde{s}_l^h \in \rho(\bar{\mu}_l^h)$  there exists  $s_l^h \in \rho(\mu_l^h)$  such that  $s_l^h =^{\widehat{H}} \tilde{s}_l^h$ .

**Proof.** Fix  $\tilde{s}_l^h \in \rho(\bar{\mu}_l^h)$ . For every  $\tilde{h} \notin \widehat{H}$  and  $\bar{h} \succ \tilde{h}$ ,  $\bar{h} \notin \widehat{H}$ . Hence there exists  $s_l^h =^{\widehat{H}} \tilde{s}_l^h$  such that for every  $\tilde{h} \notin \widehat{H}$ ,  $s_l^h \in \widehat{r}(\mu_l^h, \tilde{h})$ .

Fix  $\tilde{h} \in H(s_l^h) \cap H(\rho(\mu_l^h)) \cap \widehat{H}$  and  $\tilde{s}_l^h \in \rho(\mu_l^h)(\tilde{h})$ . For every  $h \preceq \bar{h} \prec \tilde{h}$ ,  $\bar{h} \in \widehat{H}$ . Then, since  $s_l^h =^{\widehat{H}} \tilde{s}_l^h$ ,  $\tilde{h} \in H(\tilde{s}_l^h)$ . Hence  $\tilde{s}_l^h \in \widehat{r}(\bar{\mu}_l^h, \tilde{h})$  and for every  $\tilde{s}_{-l}^h \in \text{Supp}\bar{\mu}_l^h(\cdot|\tilde{h})$ , by A0  $H(\tilde{s}_l^h, \tilde{s}_{-l}^h) \subseteq \widehat{H}$ ; so, since  $s_l^h =^{\widehat{H}} \tilde{s}_l^h$ ,  $\zeta(s_l^h, \tilde{s}_{-l}^h) = \zeta(\tilde{s}_l^h, \tilde{s}_{-l}^h)$ . Thus  $s_l^h \in \widehat{r}(\bar{\mu}_l^h, \tilde{h})$  too and for every  $z \in \zeta(\{s_l^h\} \times \text{Supp}\bar{\mu}_l^h(\cdot|\tilde{h})) =: \bar{Z}$ ,  $p(z) \in \widehat{H}$ ; so by  $\mu_l^h =^{\widehat{H}} \bar{\mu}_l^h$ ,  $\mu_l^h(S_{-l}^h(z)|\tilde{h}) = \bar{\mu}_l^h(S_{-l}^h(z)|\tilde{h}) =: \nu(z)$ . Then  $s_l^h$  induces the same probability distribution  $(\nu(z))_{z \in \bar{Z}}$ <sup>47</sup> with  $\bar{\mu}_l^h(\cdot|\tilde{h})$  and  $\mu_l^h(\cdot|\tilde{h})$  over  $Z^h$ . For every  $\tilde{s}_{-l}^h \in \text{Supp}\mu_l^h(\cdot|\tilde{h})$ , by A0  $H(\tilde{s}_l^h, \tilde{s}_{-l}^h) \subseteq \widehat{H}$ . Thus for every  $z \in \zeta(\{\tilde{s}_l^h\} \times \text{Supp}\mu_l^h(\cdot|\tilde{h})) =: \tilde{Z}$ ,  $p(z) \in \widehat{H}$ ; so by  $\mu_l^h =^{\widehat{H}} \bar{\mu}_l^h$ ,  $\bar{\mu}_l^h(S_{-l}^h(z)|\tilde{h}) = \mu_l^h(S_{-l}^h(z)|\tilde{h}) =: \eta(z)$ . Then  $\tilde{s}_l^h$  induces the same probability distribution  $(\eta(z))_{z \in \tilde{Z}}$  with  $\mu_l^h(\cdot|\tilde{h})$  and  $\bar{\mu}_l^h(\cdot|\tilde{h})$  over  $Z^h$ . Hence, since  $s_l^h \in \widehat{r}(\bar{\mu}_l^h, \tilde{h})$  and  $\tilde{s}_l^h \in \widehat{r}(\mu_l^h, \tilde{h})$ ,  $\tilde{s}_l^h \in \widehat{r}(\bar{\mu}_l^h, \tilde{h})$  and  $s_l^h \in \widehat{r}(\mu_l^h, \tilde{h})$  too.

So for every  $\tilde{h} \in H(s_l^h) \cap H(\rho(\mu_l^h))$ ,  $s_l^h \in \widehat{r}(\mu_l^h, \tilde{h})$ . I show that  $H(s_l^h) \setminus H(\rho(\mu_l^h)) = \emptyset$ , so that  $s_l^h \in \rho(\mu_l^h)$ . Suppose not. Then there exists  $\hat{h} \in H(s_l^h) \setminus H(\rho(\mu_l^h))$  such that  $p(\hat{h}) \in H(s_l^h) \cap H(\rho(\mu_l^h))$ . But since  $s_l^h \in \widehat{r}(\mu_l^h, p(\hat{h}))$ , from the observation about the relationship between continuation and sequential best replies in Section 3, I can deduce that  $\hat{h} \in H(\rho(\mu_l^h))$  too. ■

From now on, whenever  $\mu_l^h =^{\widehat{H}} \bar{\mu}_l^h$ , there will exist a set of outcomes such that  $\widehat{H}$  is the set of histories of  $\Gamma(h)$  that precede them.

Lemma 3 deals with a similar situation as in Lemma 1, but from the perspective of the deviator. If the co-players may be surprised by different deviations from the same predicted behavior, the deviator can expect any combination of reactions. The Lemma is less general to target the particular setting in which it will be used. The Lemma is rather intuitive and the proof is mostly a tedious book-keeping exercise. Therefore, it is relegated to the Online Appendix.

<sup>47</sup>The probability distribution over  $Z$  induced by a strategy  $s_i^h \in S_i^h(\tilde{h})$  with  $\mu_i^h(\cdot|\tilde{h})$  is  $(\mu_i^h(S_{-i}^h(z)|\tilde{h}))_{z \in \zeta(\{s_i^h\} \times S_{-i}^h(\tilde{h}))}$  (and probability 0 to every  $z \notin \zeta(\{s_i^h\} \times S_{-i}^h(\tilde{h}))$ ).

**Lemma 3** Fix a red. procedure  $((\tilde{S}_{i,q}^h)_{i \in I})_{q \geq 0}$ , subsets of strategies  $(\bar{S}_i^h)_{i \in I}$ ,  $m \in \mathbb{N}$  and  $l \in I$ . Let  $H^S := H(\bar{S}^h)$  and  $D^S := D_l(\bar{S}^h)$ . For every  $i \neq l$ , suppose that there exists a map  $\bar{\mu}_i^h : \bar{S}_i^h \rightarrow \Delta^{H^h}(S_{-i}^h)$  such that for every  $s_i^h \in \bar{S}_i^h$ ,  $\bar{\mu}_i^h(s_i^h)$  s.bel.  $\bar{S}_{-i}^h$ , and:

A1 there exist maps  $\bar{\mu}_i^h : \bar{S}_i^h \rightarrow \Delta^{H^h}(S_{-i}^h)$  and  $\bar{s}_i^h : \bar{S}_i^h \rightarrow S_i^h$  such that for every  $s_i^h \in \bar{S}_i^h$ ,  $\bar{\mu}_i^h(s_i^h) =^{H^S} \bar{\mu}_i^h(s_i^h)$  strongly bel.  $(\tilde{S}_{-i,q}^h)_{q=0}^{m-1}$  and  $\rho(\bar{\mu}_i^h(s_i^h)) \ni \bar{s}_i^h(s_i^h) =^{H^S} s_i^h$ ;

A2 for every  $s_i^h \in \bar{S}_i^h$  and  $\mu_i^h =^{H^S} \bar{\mu}_i^h(s_i^h)$  t.s.b.  $(\tilde{S}_{-i,q}^h)_{q=0}^{m-1}$ ,  $\rho(\mu_i^h) \subseteq \tilde{S}_{i,m}^h$ .

Fix  $\mu_l^h$  t.s.b.  $(\tilde{S}_{-l,q}^h)_{q=0}^m$  and  $\times_{i \neq l} (\bar{s}_i^h(s_i^h))_{s_i^h \in \bar{S}_i^h}$ . Fix  $\tilde{D} \subseteq D^S$  and for every  $\hat{h} \in \tilde{D}$ , fix  $\tilde{\mu}_l^{\hat{h}}$  t.s.b.  $(\tilde{S}_{-l,q}^h(\hat{h})|\hat{h})_{q=0}^m$ . Let  $H^* := H^h \setminus \cup_{\hat{h} \in \tilde{D}} H^{\hat{h}}$ .

There exists  $\tilde{\mu}_l^h =^{H^*} \mu_l^h$  t.s.b.  $(\tilde{S}_{-l,q}^h)_{q=0}^m$  such that for every  $\hat{h} \in \tilde{D}$ ,  $\tilde{\mu}_l^h =^{\hat{h}} \tilde{\mu}_l^{\hat{h}}$ .

For future reference, note that every  $\bar{\mu}_i^h$  t.s.b.  $(\tilde{S}_{-i,q}^h)_{q=0}^\infty$  and  $H(\tilde{S}_\infty^h)$  satisfy A0.

Fix a set of strategy profiles  $\bar{S}^h$  delivered by a reduction procedure. Suppose that until step  $n$ , each player  $i$  is willing to play strategies that mimic those in  $\bar{S}_i^h$  along the paths induced by  $\bar{S}^h$  while expecting the co-players to do the same. At step  $n+1$ , instead, some player  $l$  stops playing any strategy of hers that mimics a strategy  $\hat{s}_l^h$  in  $\bar{S}_l^h$ . Since at  $n$  the co-players may be surprised by any deviation, player  $l$  might expect them to play any combination of substrategies that survive  $n$  steps after the potential deviations. Hence, there must exist one particular deviation that player  $l$  prefers to mimicking  $\hat{s}_l^h$  whatever she may conjecture thereafter.

**Lemma 4** Fix reduction procedures  $((\bar{S}_{i,q}^h)_{i \in I})_{q \geq 0}$  and  $((S_{i,q}^h)_{i \in I})_{q \geq 0}$ . For every  $i \in I$  call  $\bar{S}_i^h := \bar{S}_{i,\infty}^h$  and let  $\bar{\mu}_i^h : \bar{S}_i^h \rightarrow \Delta^{H^h}(S_{-i}^h)$  be a map such that for every  $s_i^h \in \bar{S}_i^h$ ,  $\bar{\mu}_i^h(s_i^h)$  strongly believes  $(\bar{S}_{-i,q}^h)_{q=0}^\infty$  and  $s_i^h \in \rho(\bar{\mu}_i^h(s_i^h))$  (so that  $(\bar{S}_n^h)_{n \geq 0}$  satisfies A1 with  $\bar{\mu}_i^h(\cdot) = \bar{\mu}_i^h(\cdot)$ ). Let  $H^S := H(\bar{S}^h)$ . Fix  $n \in \mathbb{N}$ ,  $l \in I$  and  $\hat{s}_l^h \in \bar{S}_l^h$  such that:<sup>48</sup>

A3 for every  $i \in I$  and  $m \leq n$ ,  $(S_q^h)_{q \geq 0}$  satisfies A1;

A4 for every  $i \in I$  and  $m \in \mathbb{N}$ ,  $(S_q^h)_{q \geq 0}$  satisfies A2;

A5 for every  $i \in I$  and  $m \in \mathbb{N}$ ,  $(\bar{S}_q^h)_{q \geq 0}$  satisfies A2;

A6 for every  $s_l^h =^{H^S} \hat{s}_l^h$  and  $\mu_l^h =^{H^S} \bar{\mu}_l^h(\hat{s}_l^h)$  t.s.b.  $(S_{-l,q}^h)_{q=0}^n$ ,  $s_l^h \notin \rho(\mu_l^h)$ .

<sup>48</sup>A3, A4 and A5 need not hold for  $i = l$  to claim Lemma 3 and prove this Lemma. However,  $l$  has been included to reuse A3, A4 and A5 in the final proof of the theorems.

Let  $D^S := D_l(\overline{S}^h)$ . For every  $\widehat{h} \in D^S$  and  $m \in \mathbb{N}$  call  $M_m^{\widehat{h}}$  (resp.  $\overline{M}_m^{\widehat{h}}$ ) the set of  $\mu_l^{\widehat{h}}$  t.s.b.  $(S_{-l,q}^h(\widehat{h})|\widehat{h})_{q=0}^m$  (resp.  $(\overline{S}_{-l,q}^h(\widehat{h})|\widehat{h})_{q=0}^m$ ) such that there exist  $H_{\widehat{h}} \subseteq H^{\widehat{h}}$  and  $\widehat{\mu}_l^{\widehat{h}} =^{H_{\widehat{h}}} \mu_l^{\widehat{h}}$  t.s.b.  $(S_{-l,q}^h(\widehat{h})|\widehat{h})_{q=0}^n$  which satisfy A0.<sup>49</sup>

Then there exists  $\widehat{h} \in D^S$  such that for every  $m \leq n$ ,  $p \in \mathbb{N}$  and  $(\mu_l^{\widehat{h}}, \widetilde{\mu}_l^{\widehat{h}}) \in M_m^{\widehat{h}} \times \overline{M}_p^{\widehat{h}}$ , there exist (1)  $\mu_l^{\widehat{h}} =^{H^S} \overline{\mu}_l^{\widehat{h}}(\widehat{s}^h)$  t.s.b.  $(S_{-l,q}^h)_{q=0}^m$  and (2)  $\widetilde{\mu}_l^{\widehat{h}} =^{H^S} \overline{\mu}_l^{\widehat{h}}(\widehat{s}^h)$  t.s.b.  $(\overline{S}_{-l,q}^h)_{q=0}^p$  such that  $\mu_l^{\widehat{h}} =^{\widehat{h}} \mu_l^{\widehat{h}}$ ,  $\widetilde{\mu}_l^{\widehat{h}} =^{\widehat{h}} \widetilde{\mu}_l^{\widehat{h}}$  and  $\rho(\mu_l^{\widehat{h}})(\widehat{h}) \neq \emptyset \neq \rho(\widetilde{\mu}_l^{\widehat{h}})(\widehat{h})$ .

**Proof.** Suppose by contraposition that there is a partition  $(D, \overline{D})$  of  $D^S$  such that for every  $\widehat{h} \in D$  there exist  $m(\widehat{h}) \leq n$  and  $\mu_l^{\widehat{h}} \in M_{m(\widehat{h})}^{\widehat{h}}$  that violate 1, and for every  $\widehat{h} \in \overline{D}$  there exist  $m(\widehat{h}) \in \mathbb{N}$  and  $\mu_l^{\widehat{h}} \in \overline{M}_{m(\widehat{h})}^{\widehat{h}}$  that violate 2. For every  $\widehat{h} \in D^S$  fix corresponding  $\widehat{\mu}_l^{\widehat{h}}$  and  $H_{\widehat{h}}$ . Write  $\overline{D} = \{h^1, \dots, h^k\}$  where  $m(h^1) \geq \dots \geq m(h^k)$ . Let  $\overline{\mu}_l^{\widehat{h}} := \overline{\mu}_l^{\widehat{h}}(\widehat{s}^h)$ .

Fix any  $\mu_l^{\widehat{h}} =^{H^S} \overline{\mu}_l^{\widehat{h}}$  t.s.b.  $(S_{-l,q}^h)_{q=0}^n$  and  $\times_{i \neq l} (\overline{s}_i^h(s_i^h))_{s_i^h \in \overline{S}_i^h}$  (one clearly exists); let  $H^* := H^h \setminus \cup_{\widehat{h} \in D^S} H_{\widehat{h}}$ ; by Lemma 3 there exist:

- $\widetilde{\mu}_l^{\widehat{h}} =^{H^*} \mu_l^{\widehat{h}} =^{H^S} \overline{\mu}_l^{\widehat{h}}$  t.s.b.  $(S_{-l,q}^h)_{q=0}^n$  such that for every  $\widehat{h} \in D^S$ ,  $\widetilde{\mu}_l^{\widehat{h}} =^{\widehat{h}} \widehat{\mu}_l^{\widehat{h}}$ ;
- for every  $\widehat{h} \in D$ ,  $\widetilde{\mu}_{l,\widehat{h}}^{\widehat{h}} =^{H^*} \mu_l^{\widehat{h}} =^{H^S} \overline{\mu}_l^{\widehat{h}}$  t.s.b.  $(S_{-l,q}^h)_{q=0}^{m(\widehat{h})}$  such that  $\widetilde{\mu}_{l,\widehat{h}}^{\widehat{h}} =^{\widehat{h}} \mu_l^{\widehat{h}} =^{H_{\widehat{h}}} \widehat{\mu}_l^{\widehat{h}}$  and for every  $\widehat{h} \in D^S \setminus \{\widehat{h}\}$ ,  $\widetilde{\mu}_{l,\widehat{h}}^{\widehat{h}} =^{\widehat{h}} \widetilde{\mu}_l^{\widehat{h}}$ , so that by  $\spadesuit$ ,  $\widetilde{\mu}_{l,\widehat{h}}^{\widehat{h}} =^{H^h \setminus (H_{\widehat{h}} \setminus H_{\widehat{h}})} \widetilde{\mu}_l^{\widehat{h}}$ ;

Let  $\widetilde{\mu}_{l,0}^{\widehat{h}} := \overline{\mu}_l^{\widehat{h}} =: \mu_l^{\widehat{h}}$  and  $\widetilde{H} := \cup_{\widehat{h} \in \overline{D}} H_{\widehat{h}} \cup H^S$ ; by Lemma 3 there exist:

- for every  $1 \leq j \leq k$ ,  $\widetilde{\mu}_{l,j}^{\widehat{h}} =^{H^h \setminus \cup_{i=1}^j H_{h^i}} \overline{\mu}_l^{\widehat{h}}$  t.s.b.  $(\overline{S}_{-l,q}^h)_{q=0}^{m(h^j)}$  such that for every  $1 \leq t \leq j$ ,  $\widetilde{\mu}_{l,j}^{\widehat{h}} =^{h^t} \mu_l^{h^t} =^{H_{h^t}} \widetilde{\mu}_l^{h^t}$ ; so by  $\spadesuit$ ,  $\widetilde{\mu}_{l,j}^{\widehat{h}} =^{H^h \setminus H_{h^j}} \widetilde{\mu}_{l,j-1}^{\widehat{h}}$  and  $\widetilde{\mu}_{l,k}^{\widehat{h}} =^{\widetilde{H}} \widetilde{\mu}_l^{\widehat{h}}$ .

Fix  $\widehat{h} \in D$  and let  $\widehat{H} := H^h \setminus (H_{\widehat{h}} \setminus H_{\widehat{h}})$ :  $\widetilde{\mu}_l^{\widehat{h}}$  and  $\widehat{H}$  satisfy A0 because  $\widetilde{\mu}_l^{\widehat{h}} =^{\widehat{h}} \widehat{\mu}_l^{\widehat{h}}$ ,<sup>50</sup>  $\widetilde{\mu}_{l,\widehat{h}}^{\widehat{h}}$  and  $\widehat{H}$  satisfy A0 because by the contrapositive hypothesis  $\rho(\widetilde{\mu}_{l,\widehat{h}}^{\widehat{h}})(\widehat{h}) = \emptyset$ . Then by Lemma 2  $\rho(\widetilde{\mu}_l^{\widehat{h}})(\widehat{h}) = \emptyset$  too.<sup>51</sup> So  $H(\rho(\widetilde{\mu}_l^{\widehat{h}})) \cap D = \emptyset$ .

Now I show inductively that  $H(\rho(\widetilde{\mu}_{l,k}^{\widehat{h}})) \cap D^S = \emptyset$ . Fix  $1 \leq j \leq k$  and suppose that  $H(\rho(\widetilde{\mu}_{l,j-1}^{\widehat{h}})) \cap D^S = \emptyset$ , which is true for  $j = 1$  because  $\rho(\overline{\mu}_l^{\widehat{h}}) \subseteq \overline{S}_l^h$ . Hence  $\widetilde{\mu}_{l,j-1}^{\widehat{h}}$

<sup>49</sup>Note that:  $\mu_l^{\widehat{h}}$  refers to the second procedure even when  $\mu_l^{\widehat{h}}$  refers to the first;  $\mu_l^{\widehat{h}}$  and  $H_{\widehat{h}}$  need not satisfy A0.

<sup>50</sup>The sequential best replies to  $\widetilde{\mu}_l^{\widehat{h}}$  that reach  $\widehat{h}$ , obviously mimic those to  $\widehat{\mu}_l^{\widehat{h}}$  after  $\widehat{h}$ .

<sup>51</sup>Player  $l$  expects the same payoff against  $\widetilde{\mu}_{l,\widehat{h}}^{\widehat{h}}$  and  $\widetilde{\mu}_l^{\widehat{h}}$  without reaching  $\widehat{h}$  and a non-higher payoff against  $\widetilde{\mu}_l^{\widehat{h}}$  after reaching  $\widehat{h}$ . Hence, since she does not want to deviate under  $\widetilde{\mu}_{l,\widehat{h}}^{\widehat{h}}$ , she does not want to deviate under  $\widetilde{\mu}_l^{\widehat{h}}$ .

and  $H^h \setminus H^{hj}$  satisfy A0;  $\tilde{\mu}_{l,j}^h$  and  $H^h \setminus H^{hj}$  satisfy A0 because by the contrapositive hypothesis  $\rho(\tilde{\mu}_{l,j}^h)(h^j) = \emptyset$ ; then by Lemma 2  $H(\rho(\tilde{\mu}_{l,j}^h)) \cap D^S = \emptyset$  too. Inductively,  $H(\rho(\tilde{\mu}_{l,k}^h)) \cap D^S = \emptyset$ .

Hence,  $\tilde{\mu}_{l,k}^h$  and  $\tilde{H}$  satisfy A0;  $\tilde{\mu}_l^h$  and  $\tilde{H}$  satisfy A0 because  $H(\rho(\tilde{\mu}_l^h)) \cap D = \emptyset$  and for every  $\hat{h} \in \overline{D}$ ,  $\tilde{\mu}_l^h = \hat{h} \hat{\mu}_l^h$ ; then by Lemma 2  $H(\rho(\tilde{\mu}_l^h)) \cap D^S = \emptyset$  too. Hence  $\tilde{\mu}_l^h$  and  $H^S$  satisfy A0;  $\overline{\mu}_l^h$  and  $H^S$  satisfy A0. So by Lemma 2 there exists  $s_l^h \in \rho(\tilde{\mu}_l^h)$  such that  $s_l^h = {}^{H^S} \hat{s}_l^h$ , violating A6. ■

### Proof of Theorem 1 and 3.

The idea is the following. Take the set delivered by strong- $\Delta$ -rationalizability or selective rationalizability for a given path agreement. Suppose that at some step  $n$  of the same procedure for a looser agreement or of the alternative procedure for the same agreement, a player  $l$  excludes to mimic a strategy of the set for some profitable deviation outside the histories induced by the set. More precisely, at step  $n$  of the second procedure A3 and A6 are satisfied. The absence of off-the-path restrictions allows to claim A4 and A5 and apply Lemma 4 (first statement), which together with Lemma 1 implies that both deviator and co-players will play at step  $n$  of the second procedure any sequential best reply to any CPS in the subgame. But then, with a reduction procedure, I can find an set of strategy profiles of the subgame which should have survived also in the first procedure, a contradiction. To prove this last statement, I can apply the same logics, using as second procedure the implications in the subgame of the first procedure and as first procedure the new reduction procedure. For the latter, A5 holds by construction, for the former A4 holds by the previous application of Lemma 4 (second statement). Thus, the problem is recursive and ends when the post-deviation subgame, by finiteness, is just a simultaneous moves game.

The same reasoning applies when the first procedure refers to a looser agreement and the second procedure to the path agreement in case the first procedure implements the path, so that the set it yields is compatible with the path restrictions.

Technically, the proof shows that A3, A4, A5 and A6 cannot hold together also for player  $l$ , in this particular structure, because otherwise the second statement of Lemma 4 would hold non vacuously for every  $p \in \mathbb{N}$ , which contradicts  $\hat{h} \in D^S$ .

So, let  $(\overline{S}_q)_{q=0}^\infty$  be strong- $\Delta$ -rationalizability or selective rationalizability for a path agreement  $e$  on  $z \in Z$ . Let  $(S_q)_{q=0}^\infty$  be any of the two for  $e$  or a looser agreement  $e'$ . I show that  $\zeta(S_\infty) \supseteq \zeta(\overline{S}_\infty)$ . Thus,  $\zeta(S_{R\Delta}^\infty) \subseteq \zeta(S_{\Delta}^\infty)$  and  $\zeta(S_{R\Delta}^\infty) \supseteq \zeta(S_{\Delta}^\infty)$  (Theorem 1) and  $\zeta(S_e^\infty) \supseteq \zeta(S_e^\infty)$ . Vice versa, let  $(\overline{S}_q)_{q=0}^\infty$  apply to  $e'$  and  $(S_q)_{q=0}^\infty$  to  $e$ . I show



that if  $\zeta(\bar{S}_\infty) = z$ , still  $\zeta(S_\infty) \supseteq \zeta(\bar{S}_\infty)$ , so  $\zeta(S_{e'}^\infty) \subseteq \zeta(S_e^\infty)$  and Theorem 3 holds.

The proof of  $\zeta(S_\infty) \supseteq \zeta(\bar{S}_\infty)$  is common to all cases (see the next two footnotes). To better understand the proof, I suggest to read it first for  $k = 0$ , second for  $k = 1$ , and third for  $k = 2$ : for  $k > 2$  the logics are the same as for  $k = 2$ .

Let  $\bar{S}_\infty \neq \emptyset$ , otherwise the inclusion is trivially verified. Set  $k = 0$ . For notational convenience,  $h^k$  will be substituted by just  $k$  in subscripts and superscripts.

### RECURSIVE STEP $k$

If  $k = 0$ ,  $h^0$  is the root of the game. If  $k > 0$ ,  $h^k$  is defined in Recursive Step  $k - 1$ .

If  $k = 0$ , let  $(\bar{S}_q^0)_{q=0}^\infty := (\bar{S}_q)_{q=0}^\infty$ ; else,  $(\bar{S}_q^k)_{q=0}^\infty$  is defined in Recursive Step  $k - 1$ .

If  $k = 0$ , let  $(S_q^0)_{q=0}^\infty := (S_q)_{q=0}^\infty$ ; else, let  $(S_q^k)_{q=0}^\infty := (\bar{S}_q^{k-1}(h^k)|h^k)_{q=0}^\infty$ .

For every  $i \in I$ , let  $\bar{\mu}_i^k : \bar{S}_{i,\infty}^k \rightarrow \Delta^{H^k}(S_{-i}^k)$  be a map such that for every  $s_i^k \in \bar{S}_{i,\infty}^k$ ,  $\bar{\mu}_i^k(s_i^k)$  strongly believes  $(\bar{S}_{-i,q}^k)_{q=0}^\infty$ ,  $s_i^k \in \rho(\bar{\mu}_i^k(s_i^k))$ , and, if  $k = 0$ ,  $\bar{\mu}_i^0(s_i^0) \in \Delta_i^e$ .<sup>52</sup>

I prove by induction that  $\zeta(\bar{S}_\infty^k) \subseteq \zeta(S_\infty^k)$ .

**Premise (k=0):** for every  $n \in \mathbb{N}$ ,  $i \in I$ ,  $s_i \in \bar{S}_{i,\infty}$  and  $\mu_i = {}^{H(\bar{S}_\infty)}\bar{\mu}_i^0(s_i)$  t.s.b.  $(S_{-i,q})_{q=0}^n$  or  $(\bar{S}_{-i,q})_{q=0}^n$ , by  $\bar{\mu}_i^0(s_i) \in \Delta_i^e$ ,  $\bar{\mu}_i^0(s_i)(S_{-i}(z)|h^0) = 1$ , so  $\mu_i(S_{-i}(z)|h^0) = 1$ . Thus  $\mu_i \in \Delta_i^e \subseteq \Delta_i^{e'}$ . So A4 and A5 hold.

**Premise (k>0):** A4 holds by Fact  $k - 1$ . A5 holds by Claim  $k - 1$ . (To be read in Recursive step  $k - 1$  with  $k$  and  $k - 1$  in place of  $k + 1$  and  $k$ )

**Inductive Hypothesis (n):**  $(S_q^k)_{q=0}^\infty$  satisfies A3 at  $n$  (so by A4  $\zeta(S_n^k) \supseteq \zeta(\bar{S}_\infty^k)$ ).

**Basis step (1):** for every  $i \in I$ , the I.H. holds with  $\bar{\mu}_i^k(\cdot) = \bar{\mu}_i^k(\cdot)$ .

### Inductive step (n+1).

Suppose by contradiction that the Inductive Hypothesis does not hold at  $n + 1$ . Then A6 holds for some  $l \in I$  and  $\hat{s}_l^k \in \bar{S}_{l,\infty}^k$ . Lemma 4 yields  $h^{k+1} \in D_l(\bar{S}_\infty^k)$ . If  $\Gamma(h^k)$  has depth<sup>53</sup> 1,  $h^{k+1}$  cannot exist and we have the desired contradiction (**Exit Rule**). Else, define  $((\bar{S}_{i,q}^{k+1})_{i \in I})_{q \geq 0}$  as follows: for every  $i \in I$  and  $m \leq n$ ,  $\bar{S}_{i,m}^{k+1} = S_{i,m}^k(h^{k+1})|h^{k+1}$ ; for every  $m > n$ ,  $s_i^{k+1} \in \bar{S}_{i,m}^{k+1}$  if and only if there exists  $\mu_i^{k+1}$  t.s.b.  $(\bar{S}_{-i,q}^{k+1})_{q=0}^{m-1}$  such that  $s_i^{k+1} \in \rho(\mu_i^{k+1})$ .

For every  $i \neq l$ , since  $h^{k+1} \in D_l(\bar{S}_\infty^k)$ ,  $\emptyset \neq \bar{S}_{i,\infty}^k(h^{k+1}) \ni \hat{s}_i^k$ . For every  $m \leq n$ , the I. H. provides  $\bar{s}_i^k(\hat{s}_i^k) \in S_{i,m}^k(h^{k+1}) \neq \emptyset$  and  $\bar{\mu}_i^k(\hat{s}_i^k) = {}^{H(\bar{S}_\infty^k)}\bar{\mu}_i^k(\hat{s}_i^k)$  t.s.b.  $(S_{-i,q}^k)_{q=0}^{m-1}$  such that  $\bar{\mu}_i^k(\hat{s}_i^k)(S_{-i}^k(h^{k+1})|p(h^{k+1})) = 0$ . Hence, by Lemma 1, for every  $\mu_i^{k+1}$  t.s.b.

<sup>52</sup>In case  $(\bar{S}_q)_{q=0}^\infty = (S_{e'}^q)_{q=0}^\infty$ , I assumed  $\bar{S}_{-i,\infty} \subseteq S_{-i}(z)$ ; hence, for every  $\mu_i$  t.s.b.  $\bar{S}_{-i,\infty}$ ,  $\mu_i \in \Delta_i^e$ .

<sup>53</sup>The difference between the length  $T$  of the longest path  $(a^1, \dots, a^T) \succ h$  and the length of  $h$ .

$(\overline{S}_{-i,q}^{k+1})_{q=0}^{m-1}$ , there exists  $\mu_i^k =^{k+1} \mu_i^{k+1}$  t.s.b.  $(S_{-i,q}^k)_{q=0}^{m-1}$  such that  $\mu_i^k =_{H(\overline{S}_\infty^k)} \overline{\mu}_i^k(\widehat{s}_i^k)$  and  $\rho(\mu_i^k)(h^{k+1}) \neq \emptyset$ . By A4  $\rho(\mu_i^k) \subseteq S_{i,m}^k$ . So  $\rho(\mu_i^{k+1}) \subseteq \overline{S}_{i,m}^{k+1}$ .

Fix  $\mu_l^{k+1}$  t.s.b.  $(\overline{S}_{-l,q}^{k+1})_{q=0}^n$ : trivially  $\mu_l^{k+1} \in M_n^{k+1}$ . Hence by Lemma 4 (1) there exists  $\widetilde{\mu}_l^k =_{H(\overline{S}_\infty^k)} \overline{\mu}_l^k(\widehat{s}_l^k)$  t.s.b.  $(S_{-l,n}^k)_{q=0}^n$  such that  $\widetilde{\mu}_l^k =^{k+1} \mu_l^{k+1}$  and  $\rho(\widetilde{\mu}_l^k)(h^{k+1}) \neq \emptyset$ . By A4,  $\rho(\mu_l^{k+1}) \subseteq S_{l,n}^k$ . So  $\rho(\mu_l^{k+1}) \subseteq \overline{S}_{l,n}^{k+1} \neq \emptyset$ .

Hence for every  $i \in I$  and  $\mu_i^{k+1}$  t.s.b.  $(\overline{S}_{-i,q}^{k+1})_{q=0}^n$ ,  $\rho(\mu_i^{k+1}) \subseteq \overline{S}_{i,n}^{k+1} \neq \emptyset$ , so that  $\overline{S}_{i,n}^{k+1} \supseteq \overline{S}_{i,n+1}^{k+1}$  and  $((\overline{S}_{i,q}^{k+1})_{i \in I})_{q \geq 0}$  is a reduction procedure with  $\overline{S}_\infty^{k+1} \neq \emptyset$ .

For every  $m \leq n$  and  $\widehat{\mu}_l^{k+1}$  t.s.b.  $(\overline{S}_{-l,q}^{k+1})_{q=0}^\infty$ , fix  $\mu_l^{k+1} =_{H(\overline{S}_\infty^{k+1})} \widehat{\mu}_l^{k+1}$  t.s.b.  $(\overline{S}_{-l,q}^{k+1})_{q=0}^{m-1}$ . Since  $\widehat{\mu}_l^{k+1}$  and  $H(\overline{S}_\infty^{k+1})$  satisfy A0, by Lemma 4 (1) there exists  $\widetilde{\mu}_l^k =_{H(\overline{S}_\infty^k)} \overline{\mu}_l^k(\widehat{s}_l^k)$  t.s.b.  $(S_{-l,q}^k)_{q=0}^{m-1}$  such that  $\widetilde{\mu}_l^k =^{k+1} \mu_l^{k+1}$  and  $\rho(\widetilde{\mu}_l^k)(h^{k+1}) \neq \emptyset$ . By A4,  $\rho(\mu_l^{k+1}) \subseteq S_{l,m}^k$ . So  $\rho(\mu_l^{k+1}) \subseteq \overline{S}_{l,m}^{k+1}$ .

Then, for every  $m \in \mathbb{N}$ ,  $i \in I$ ,  $\widehat{\mu}_i^{k+1}$  t.s.b.  $(\overline{S}_{-i,q}^{k+1})_{q=0}^\infty$  and  $\mu_i^{k+1} =_{H(\overline{S}_\infty^{k+1})} \widehat{\mu}_i^{k+1}$  t.s.b.  $(\overline{S}_{-i,q}^{k+1})_{q=0}^{m-1}$ ,  $\rho(\mu_i^{k+1}) \subseteq \overline{S}_{i,m}^{k+1}$ . (**Claim k**)

If  $\zeta(\overline{S}_\infty^{k+1}) \subseteq \zeta(\overline{S}_\infty^k)$ , then  $h^{k+1} \in H(\overline{S}_\infty^k)$ , the desired contradiction to claim that the I. H. holds for  $n + 1$ . The goal of Recursive Step  $k + 1$  is precisely to prove  $\zeta(\overline{S}_\infty^{k+1}) \subseteq \zeta(\overline{S}_\infty^k)$ . Before, observe what follows.

For every  $i \neq l$ ,  $m \in \mathbb{N}$  and  $\mu_i^{k+1}$  t.s.b.  $(\overline{S}_{-i,q}^k(h^{k+1})|h^{k+1})_{q=0}^{m-1}$ , by Lemma 1 there exists  $\widetilde{\mu}_i^k =^{k+1} \mu_i^{k+1}$  t.s.b.  $(\overline{S}_{-i,q}^k)_{q=0}^{m-1}$  such that for every  $\widetilde{h} \notin H^{k+1}$ ,  $\widetilde{\mu}_i^k(\cdot|\widetilde{h}) = \overline{\mu}_i^k(\widehat{s}_i^k)(\cdot|\widetilde{h})$  and  $\rho(\widetilde{\mu}_i^k)(h^{k+1}) \neq \emptyset$ .<sup>54</sup> By A5,  $\rho(\mu_i^{k+1}) \subseteq \overline{S}_{i,m}^k$ .

For every  $m \in \mathbb{N}$  and  $\widehat{\mu}_l^{k+1}$  t.s.b.  $(\overline{S}_{-l,q}^{k+1})_{q=0}^\infty$ , fix  $\mu_l^{k+1} =_{H(\overline{S}_\infty^{k+1})} \widehat{\mu}_l^{k+1}$  t.s.b.  $(\overline{S}_{-l,q}^k(h^{k+1})|h^{k+1})_{q=0}^{m-1}$ . Since  $\widehat{\mu}_l^{k+1}$  and  $H(\overline{S}_\infty^{k+1})$  satisfy A0, by Lemma 4 (2) there exists  $\widetilde{\mu}_l^k =_{H(\overline{S}_\infty^k)} \overline{\mu}_l^k(\widehat{s}_l^k)$  t.s.b.  $(\overline{S}_{-l,q}^k)_{q=0}^{m-1}$  such that  $\widetilde{\mu}_l^k =^{k+1} \mu_l^{k+1}$  and  $\rho(\widetilde{\mu}_l^k)(h^{k+1}) \neq \emptyset$ . By A5  $\rho(\mu_l^{k+1}) \subseteq \overline{S}_{l,m}^k$ .

Then, for every  $m \in \mathbb{N}$ ,  $i \in I$ ,  $\widehat{\mu}_i^{k+1}$  t.s.b.  $(\overline{S}_{-i,q}^{k+1})_{q=0}^\infty$  and  $\mu_i^{k+1} =_{H(\overline{S}_\infty^{k+1})} \widehat{\mu}_i^{k+1}$  t.s.b.  $(\overline{S}_{-l,q}^k(h^{k+1})|h^{k+1})_{q=0}^{m-1}$ ,  $\rho(\mu_i^{k+1}) \subseteq \overline{S}_{i,m}^k(h^{k+1})|h^{k+1}$ . (**Fact k**)

Note that the Exit Rule must be verified at some  $k$ , because the game has finite depth. Then  $\zeta(\overline{S}_\infty^k) \subseteq \zeta(S_\infty^k)$ , the desired contradiction to claim that  $\zeta(\overline{S}_\infty^{k-1}) \subseteq \zeta(S_\infty^{k-1})$ . Proceeding backwards, we obtain  $\zeta(\overline{S}_\infty) \subseteq \zeta(S_\infty)$ . ■

The proof can be employed also to show an interesting relationship between strong rationalizability, which can be seen as strong- $\Delta$ -rationalizability under the loosest possible agreement, and strong- $\Delta$ -rationalizability. The latter does not in general deliver a subset of the former in terms of strategy profiles, nor in terms of outcomes. Yet, it does deliver a subset of outcomes when at the end of the procedure, there is no

<sup>54</sup>Recall that:  $\overline{\mu}_i^k(s_i^k)$  s.b.  $(\overline{S}_{-i,q}^k)_{q=0}^\infty$ ;  $\overline{\mu}_i^k(s_i^k)(S_{-i}^k(h^{k+1})|p(h^{k+1})) = 0$ ;  $s_i^k \in \rho(\overline{\mu}_i^k(s_i^k))(h^{k+1}) \neq \emptyset$ .

restriction left off the induced paths. For coherency with the framework, restrictions are expressed here in terms of an agreement, but the result holds through for more general restrictions (such as those deriving from an outcome distribution).

**Remark 7** *Fix an agreement  $e = (e_i)_{i \in I}$ . If for every  $h \notin H(S_{\Delta^e}^\infty)$  and  $i \in I$ ,  $e_i(h) = \emptyset$ , then  $\zeta(S_{\Delta^e}^\infty) \subseteq \zeta(S^\infty)$ .*

### PROOF OF THEOREMS 2 AND 6.

The proofs of theorems 2 and 6 are applications of Lemma 10, which is based on the following idea. For simplicity, think of SPE and equilibria in pure strategies. When I say that a SPE survives/is eliminated I mean that an equilibrium which induces the SPE path survives/is eliminated. Fix a reduction procedure and a collection of unordered subgames, each with an associated SPE, such that until step  $m$ : (A1) the associated SPE survive; (A2) the subgames follow unilateral deviations from equilibria which survive  $m - 1$  steps (used in next paragraph); (A3) players save the sequential best replies to CPS's against which they yield in expectation at least the payoff of a SPE in the class of SPE that induce in the subgames the associated SPE. I show that a SPE in the restricted class survives until step  $m$ . Suppose not. Take one of the SPE in the class that survive more steps, say  $n - 1$ . There is a unilateral deviation from its path that the deviator takes at  $n$  whatever she conjectures thereafter. This is proved in Lemma 8, which resembles Lemma 4. The same holds for the co-players since they can be surprised by the deviation. But then (see next paragraph), in the new post-deviation subgame some SPE in the class of SPE that induce in the smaller old subgames the associated SPE survives until step  $n$ . Since the new subgame cannot coincide with or be contained in an old subgame (also shown in Lemma 8), substituting the smaller old subgames with the new subgame, the new class of SPE is smaller than the first. Then, iterating,  $m$  weakly decreases. So A1,2,3 still hold. Since the subgames keep enlarging, I obtain a contradiction.

I still need to prove that a SPE of the new subgame in the restricted class survives until step  $n$ . Suppose not. However the surviving substrategies feature an equilibrium of the subgame which induces the equilibria of the smaller old subgames (I prove this inside Lemma 10 using A1, A2 and the fact that players save all the sequential best replies to CPS's at step  $n$ ). There is a unilateral deviation from the equilibrium path which the deviator takes if in the post-deviation subgame she expects at least the payoff of a SPE in the class of SPE that induce in the smaller old subgames the

associated SPE. This is proved in Lemma 9 which is the dual of Lemma 8. The same holds for the co-players since they can be surprised by the deviation. This yields A3 for the new post-deviation subgame: suppose that the survival of a SPE of the new subgame in the restricted class has been proved inductively on the depth of subgames. Since the new subgame cannot coincide with or be contained in an old subgame (also shown in Lemma 9), substituting the smaller old subgames with the new subgame, the new class of SPE is smaller than the first. Then, iterating, the contrapositive hypothesis still holds. Since the subgames keep enlarging, I obtain a contradiction.

In the proof of Lemma 10, the two iterative procedures described above constitute the outer and the inner recursive proofs for the inductive step. I will always refer to uncorrelated CPS's<sup>55</sup> and distributions.

**Additional notation:**

Fix  $h \preceq \hat{h}$ ,  $i \in J \subseteq I$ ,  $\mu_i^h \in \Delta^{H^h}(S_{-i}^h)$ ,  $(\tilde{\sigma}_j^h)_{j \in I} \in \times_{j \in I} \Delta(S_j^h)$ ,  $(\sigma_j^{\hat{h}})_{j \in I} \in \times_{j \in I} \Delta(S_j^{\hat{h}})$ :

- $\tilde{\sigma}_j^h(\cdot) \in \Delta(S_j^h)$  is the product of the marginal distributions  $(\tilde{\sigma}_j^h)_{j \in J}$ ;<sup>56</sup>
- $H(\tilde{\sigma}_j^h) := H(\text{Supp} \tilde{\sigma}_j^h)$ ,  $\tilde{\sigma}_j^h[\hat{h}] := (\text{Supp} \tilde{\sigma}_j^h)[\hat{h}]$ ,  $D_i(\tilde{\sigma}^h) := D_i(\text{Supp} \tilde{\sigma}^h)$ ;
- $\tilde{\sigma}_j^h[\hat{h}]$  is the product of  $(\tilde{\sigma}_j^h[\hat{h}])_{j \in J}$  and  $\tilde{\sigma}_i^h[\hat{h}] \in \Delta(S_i^{\hat{h}})$  is def. for every  $s_i^{\hat{h}} \in S_i^{\hat{h}}$  as
  - $(\tilde{\sigma}_i^h[\hat{h}])(s_i^{\hat{h}}) = \tilde{\sigma}_i^h(\{s_i^h \in S_i^h(\hat{h}) : s_i^h[\hat{h}] = s_i^{\hat{h}}\}) / \tilde{\sigma}_i^h(S_i^h(\hat{h}))$  if  $\hat{h} \in H(\tilde{\sigma}_i^h)$ ,
  - $(\tilde{\sigma}_i^h[\hat{h}])(s_i^{\hat{h}}) = \tilde{\sigma}_i^h(\{s_i^h \in S_i^h : s_i^h[\hat{h}] = s_i^{\hat{h}}\})$  else;
- $\tilde{\sigma}_j^h =^* \sigma_j^{\hat{h}}$  if for every  $z \succ \hat{h}$  with  $p(z) \in H(\sigma_j^{\hat{h}})$  and  $j \in J$ ,  $(\tilde{\sigma}_j^h[\hat{h}])(S_j^{\hat{h}}(z)) = \sigma_j^{\hat{h}}(S_j^{\hat{h}}(z))$ ;<sup>57</sup>
- $\tilde{\sigma}_j^h =^{\hat{h}} \sigma_j^{\hat{h}}$  if for every  $z \succ \hat{h}$  and  $j \in J$ ,  $(\tilde{\sigma}_j^h[\hat{h}])(S_j^{\hat{h}}(z)) = \sigma_j^{\hat{h}}(S_j^{\hat{h}}(z))$ ;
- $\mu_i^h =^* \sigma_{-i}^{\hat{h}}$  if  $\mu_i^h(\cdot|\hat{h}) =^* \sigma_{-i}^{\hat{h}}$ ;  $\mu_i^h =^{\hat{h}} \sigma_{-i}^{\hat{h}}$  if  $\mu_i^h(\cdot|\hat{h}) =^{\hat{h}} \sigma_{-i}^{\hat{h}}$ ;
- $\pi_i(\tilde{\sigma}^h)$  is  $i$ 's exp. payoff under  $\tilde{\sigma}^h$ ;  $\pi(\tilde{\sigma}_{-i}^h) := \max_{s_i^h \in S_i^h} \sum_{s_{-i}^h \in \text{Supp} \tilde{\sigma}_{-i}^h} u_i(\zeta(\tilde{s}_i^h, s_{-i}^h)) \tilde{\sigma}_{-i}^h(s_{-i}^h)$ ;
- $\tilde{\sigma}_{-i}^h[\hat{h}] \geq \sigma_{-i}^{\hat{h}}$  if  $\pi(\tilde{\sigma}_{-i}^h[\hat{h}]) \geq \pi(\sigma_{-i}^{\hat{h}})$ ;  $\mu_i^h \geq \tilde{\sigma}_{-i}^h$  if  $\mu_i^h(\cdot|h) \geq \tilde{\sigma}_{-i}^h$ ;

<sup>55</sup>A CPS  $\mu_i^h$  is uncorrelated if for every  $\tilde{h} \in H^h$ ,  $\mu_i^h(\cdot|\tilde{h}) = \times_{j \neq i} \text{Marg}_{S_j^h} \mu_i^h(\cdot|\tilde{h})$

<sup>56</sup>This is an exception to the rule of subscripts:  $\tilde{\sigma}_j^h$  is not a (sub-)profile of distributions but an uncorrelated joint distribution. Equilibria  $(\tilde{\sigma}_j^h)_{j \in I}$  will be represented as the joint uncorrelated distribution  $\tilde{\sigma}^h$  they induce, and then  $\tilde{\sigma}_j^h := \text{Marg}_{S_j^h} \tilde{\sigma}^h$ .

<sup>57</sup>Notice that  $z$  is not necessarily a terminal history. Notice also that only the histories whose predecessor is induced with positive probability by  $\sigma^{\hat{h}}$  and not all those compatible with  $\sigma_j^{\hat{h}}$  matter.

- $\tilde{\sigma}^h$  is a SPE of  $\Gamma(h)$  if for every  $\tilde{h} \in H^h$ ,  $\tilde{\sigma}^h|\tilde{h}$  is an equilibrium of  $\Gamma(\tilde{h})$ ;
- for any set of unordered<sup>58</sup> non-terminal histories  $\tilde{H} \subseteq H$  and any set of SPE  $\Sigma^{\tilde{H}} = (\tilde{\sigma}^h)_{\tilde{h} \in \tilde{H}}$  of the corresponding subgames,  $E^h(\Sigma^{\tilde{H}})$  is the set of SPE of  $\Gamma(h)$   $\sigma^h$  such that for every  $\tilde{h} \in \tilde{H} \cap H^h$ ,  $\sigma^h|\tilde{h} = \tilde{\sigma}^h$ .

I will often use the fact that  $=^*$  and  $=^{\hat{h}}$  are transitive and that  $=^{\hat{h}}$  implies  $=^*$ .<sup>59</sup> Moreover note that when  $\tilde{\sigma}_i^h =^* \sigma_i^{\hat{h}}$ :

- ♡ for every  $\tilde{h} \succ \hat{h}$  with  $p(\tilde{h}) \in H(\sigma^{\hat{h}})$  and  $\hat{h} \preceq \bar{h} \prec \tilde{h}$ ,  $(\tilde{\sigma}_i^h|\tilde{h})(S_i^{\bar{h}}(\tilde{h})) = (\sigma_i^{\hat{h}}|\bar{h})(S_i^{\bar{h}}(\tilde{h}))$ ;<sup>60</sup>
- ♠ if for every  $j \neq i$  and  $\tilde{h} \in D_j(\sigma^{\hat{h}})$ ,  $\tilde{\sigma}_i^h|\tilde{h} = \sigma_i^{\hat{h}}|\tilde{h}$ , then  $\tilde{\sigma}_i^h =^{\hat{h}} \sigma_i^{\hat{h}}$ .<sup>61</sup>

When  $\sigma^{\hat{h}}$  is an equilibrium and  $\tilde{\sigma}^h =^{\hat{h}} \sigma^{\hat{h}}$ , I will often use the fact that for every  $\tilde{h} \in H(\sigma^{\hat{h}})$ ,  $\tilde{\sigma}^h|\tilde{h}$  is an equilibrium. Moreover:

- ♦ if  $\tilde{\sigma}_{-i}^{\hat{h}} =^* \sigma_{-i}^{\hat{h}}$ ,  $\pi(\tilde{\sigma}_{-i}^{\hat{h}}) \geq \pi(\sigma_{-i}^{\hat{h}}) = \pi_i(\sigma^{\hat{h}})$  and if  $\tilde{\sigma}^{\hat{h}}$  is an equil.,  $\pi(\tilde{\sigma}_{-i}^{\hat{h}}) = \pi(\sigma_{-i}^{\hat{h}})$ ;<sup>62</sup>
- ♣ if  $\mu_i^{\hat{h}} =^{\hat{h}} \sigma_{-i}^{\hat{h}}$ , for every  $s_i^{\hat{h}} \in \text{Supp}\sigma_i^{\hat{h}}$ , there is  $\tilde{s}_i^{\hat{h}} \in \rho(\mu_i^{\hat{h}})$  such that  $\tilde{s}_i^{\hat{h}} =^{H(\sigma^{\hat{h}})} s_i^{\hat{h}}$ .<sup>63</sup>

For the next five Lemmata, fix  $n \in \mathbb{N}$ ,  $h \in H$ , a red. procedure  $(S_q^h)_{q \geq 0}$  and  $\tilde{\sigma}^h \in \Delta(S_{n-1}^h)$ . Let  $H^\sigma := H(\tilde{\sigma}^h)$ . For every  $i \in I$ , let  $D_i := D_i(\tilde{\sigma}^h)$  and  $D_{-i} := \cup_{j \neq i} D_j$ .

The first Lemma, like in the proof of Lemma 3, combines different reactions of player  $i$  to unexpected deviations from  $H^\sigma$ . Its proof is in the Online Appendix too.

**Lemma 5** Fix  $v \leq n$  and  $i \in I$  such that  $\tilde{\sigma}_i^h \in \Delta(r(\tilde{\sigma}_{-i}^h))$  and for every  $\mu_i^h =^* \tilde{\sigma}_{-i}^h$  t.s.b.  $(S_{-i,q}^h)_{q=0}^{v-1}$ ,  $\rho(\mu_i^h) \subseteq S_{i,v}^h$ . For every  $\tilde{h} \in D_{-i}$ , fix  $\tilde{\sigma}_i^{\tilde{h}} \in \Delta(S_{i,v}^h(\tilde{h})|\tilde{h})$ . There exists  $\hat{\sigma}_i^h \in \Delta(S_{i,v}^h)$  such that  $\hat{\sigma}_i^h =^* \tilde{\sigma}_i^h$  and for every  $\tilde{h} \in D_{-i}$ ,  $\hat{\sigma}_i^h|\tilde{h} = \tilde{\sigma}_i^{\tilde{h}}$ .

<sup>58</sup>For every two histories  $h, h'$  in the set,  $h \not\preceq h'$  and  $h' \not\preceq h$ .

<sup>59</sup>This fact resembles the set monotonicity of  $=^{\tilde{H}}$  for CPS's.

<sup>60</sup>The equivalent condition for CPS's is incorporated in the definition of  $=^{\hat{h}}$ . Here it can be derived from the conjectures at  $\hat{h}$  because  $\bar{h}$  is reached with positive probability.

<sup>61</sup>This is the analogous of ♠ for CPS's. Since  $\tilde{\sigma}_i^h =^* \sigma_i^{\hat{h}}$ ,  $(\tilde{\sigma}_i^h|\tilde{h})(S_i^{\bar{h}}(\tilde{h})) = \sigma_i^{\hat{h}}|\tilde{h})(S_i^{\bar{h}}(\tilde{h}))$ . For every  $z \succ \tilde{h}$ , since  $\tilde{\sigma}_i^h|\tilde{h} = \sigma_i^{\hat{h}}|\tilde{h}$ ,  $(\tilde{\sigma}_i^h|\tilde{h})(S_i^{\bar{h}}(z)) = (\sigma_i^{\hat{h}}|\tilde{h})(S_i^{\bar{h}}(z))$ . Together,  $(\tilde{\sigma}_i^h|\tilde{h})(S_i^{\bar{h}}(z)) = \sigma_i^{\hat{h}}|\tilde{h})(S_i^{\bar{h}}(z))$ .

<sup>62</sup>The equivalent payoff relation for CPS's was implied by  $\mu_i^{\hat{h}} =^{H_{\hat{h}}} \hat{\mu}_i^{\hat{h}}$  where  $\hat{\mu}_i^{\hat{h}}$  and  $H_{\hat{h}}$  satisfied A0, but the latter condition instead of its implication for payoffs was employed via Lemma 2.

<sup>63</sup>From the observation in Section 3 about the relationship between continuation and sequential best replies. It is the analogous of Lemma 2 for equilibria, net of indifferences: if the CPS has an initial equilibrium conjecture, there are sequential best replies that mimic the equilibrium strategies within the histories that precede the equilibrium outcomes.

For the next four lemmata suppose that  $\tilde{\sigma}^h$  is an equilibrium and:

A0 for every  $v \leq n$ ,  $i \in I$  and  $\mu_i^h =^* \tilde{\sigma}_{-i}^h$  t.s.b.  $(S_{-i,q}^h)_{q=0}^{v-1}$ ,  $\rho(\mu_i^h) \subseteq S_{i,v}^h$ ,

so that every  $i \in I$  satisfies the hypotheses of Lemma 5.

Lemma 6 is a characterization of equilibrium which will turn out to be useful. Since the arguments for it are standard, the proof is omitted.

**Lemma 6** Fix  $(\hat{\sigma}_i^h)_{i \in I} \in \times_{i \in I} \Delta(S_i^h)$ :  $\hat{\sigma}^h$  is an equilibrium if and only if for every  $\bar{h} \in H(\hat{\sigma}^h)$ ,  $i \in I$  and  $a_i \in A_i(\bar{h}) \setminus \hat{\sigma}_i^h[\bar{h}]$ , calling  $H_{a_i}^{\bar{h}} := (\bar{h}, (a_i, a_{-i}))_{a_{-i} \in \hat{\sigma}_{-i}^h[\bar{h}]}$ ,

$$\sum_{\tilde{h} \in H_{a_i}^{\bar{h}} \setminus Z} \pi(\hat{\sigma}_{-i}^h|\tilde{h}) \cdot (\hat{\sigma}_{-i}^h|\bar{h})(S_{-i}^{\bar{h}}(\tilde{h})) + \sum_{z \in H_{a_i}^{\bar{h}} \cap Z} u_i(z) \cdot (\hat{\sigma}_{-i}^h|\bar{h})(S_{-i}^{\bar{h}}(z)) \leq \pi_i(\hat{\sigma}^h|\bar{h}). \quad (\star)$$

Lemma 7 converts a condition on CPS's into  $\star$  for some related conjectures.

**Lemma 7** Fix  $\hat{\sigma}^h =^* \tilde{\sigma}^h$ ,  $\bar{h} \in H^\sigma$ ,  $i \in I$ ,  $a_i \in A_i(\bar{h}) \setminus \hat{\sigma}_i^h[\bar{h}]$ ,  $\hat{h} \in H_{a_i}^{\bar{h}} \setminus Z$ ,  $v \leq n$  and  $\tilde{\mu}_i^{\hat{h}}$  t.s.b.  $(S_{-i,q}^h(\hat{h})|\hat{h})_{q=0}^v$  such that (i)  $\hat{\sigma}_{-i}^h|\hat{h} \leq \tilde{\mu}_i^{\hat{h}}$ , (ii) for every  $\tilde{h} \in H_{a_i}^{\bar{h}} \setminus (Z \cup \{\hat{h}\})$ ,  $\hat{\sigma}_{-i}^h|\tilde{h} \leq \bar{\sigma}_{-i}^{\tilde{h}}$  for some  $\bar{\sigma}_{-i}^{\tilde{h}} \in \Delta(S_{-i,v}^h(\tilde{h})|\tilde{h})$ , and (iii) for every  $\mu_i^h =^* \tilde{\sigma}_{-i}^h$  t.s.b.  $(S_{-i,q}^h)_{q=0}^{v-1}$ , if  $\mu_i^h =^{\hat{h}} \tilde{\mu}_i^{\hat{h}}$ , then  $\rho(\mu_i^h)(\hat{h}) = \emptyset$ . Then  $\star$  holds.

**Proof.** Let  $\bar{\sigma}_{-i}^{\hat{h}} := \tilde{\mu}_i^{\hat{h}}(\cdot|\hat{h})$ . By Lemma 5 there exists  $\bar{\sigma}_{-i}^h \in \Delta(S_v^h)$  such that  $\bar{\sigma}_{-i}^h =^* \tilde{\sigma}_{-i}^h$ , for every  $\tilde{h} \in H_{a_i}^{\bar{h}} \setminus Z$ ,  $\bar{\sigma}_{-i}^h|\tilde{h} = \bar{\sigma}_{-i}^{\tilde{h}}$  and for every  $\tilde{h} \in D_i \setminus H_{a_i}^{\bar{h}}$ ,  $\bar{\sigma}_{-i}^h|\tilde{h} = \tilde{\sigma}_{-i}^h|\tilde{h}$ . Fix  $\mu_i^h =^h \bar{\sigma}_{-i}^h$  t.s.b.  $(S_{-i,q}^h)_{q=0}^v$  such that  $\mu_i^h =^{\hat{h}} \tilde{\mu}_i^{\hat{h}}$  (one exists because  $\mu_i^h(\cdot|\hat{h}) = \bar{\sigma}_{-i}^h$  implies  $\mu_i^h(\cdot|\hat{h})|\hat{h} = \bar{\sigma}_{-i}^h|\hat{h} = \tilde{\mu}_i^{\hat{h}}(\cdot|\hat{h})$ ).

For every  $\tilde{h} \in D_i \setminus H_{a_i}^{\bar{h}}$  and  $z \in Z$  such that  $z \succ \tilde{h}$ , by the argument used for  $\spadesuit$ ,  $\bar{\sigma}_{-i}^h(S_{-i}^h(z)) = \tilde{\sigma}_{-i}^h(S_{-i}^h(z))$ . For every  $z \in \varsigma(\text{Supp}\tilde{\sigma}^h)$ , by  $\bar{\sigma}^h =^* \tilde{\sigma}^h$ ,  $\bar{\sigma}_{-i}^h(S_{-i}^h(z)) = \tilde{\sigma}_{-i}^h(S_{-i}^h(z))$ . Hence every  $s_i^h \notin S_i^h(\hat{h}) = \cup_{\tilde{h} \in H_{a_i}^{\bar{h}}} S_i^h(\tilde{h})$  induces with  $\bar{\sigma}_{-i}^h$  and  $\tilde{\sigma}_{-i}^h$  (1) and with  $\bar{\sigma}_{-i}^h|\bar{h}$  and  $\tilde{\sigma}_{-i}^h|\bar{h}$  (2) the same distribution over outcomes. By 1 and  $r(\tilde{\sigma}_{-i}^h)(\bar{h}) \neq \emptyset$ ,  $r(\bar{\sigma}_{-i}^h)(\bar{h}) \cup r(\bar{\sigma}_{-i}^h)(\hat{h}) \neq \emptyset$ ; by  $\bar{h} \prec \hat{h}$ ,  $r(\bar{\sigma}_{-i}^h)(\bar{h}) \neq \emptyset$ ; by  $\bar{h} \in H(\bar{\sigma}_{-i}^h)$ ,  $\rho(\mu_i^h)(\bar{h}) \neq \emptyset$ ;<sup>64</sup> by  $\rho(\mu_i^h)(\hat{h}) = \emptyset$ ,  $\hat{r}(\mu_i^h, \bar{h})(\hat{h}) = \emptyset$ ; by  $\mu_i^h(\cdot|\bar{h})|\bar{h} = \bar{\sigma}_{-i}^h|\bar{h}$ ,  $r(\bar{\sigma}_{-i}^h|\bar{h})(\hat{h}) = \emptyset$ ; by 2,  $\pi(\bar{\sigma}_{-i}^h|\bar{h}) = \pi(\tilde{\sigma}_{-i}^h|\bar{h}) = \pi_i(\tilde{\sigma}^h|\bar{h})$ ; thus,

$$\sum_{\tilde{h} \in H_{a_i}^{\bar{h}} \setminus Z} \pi(\bar{\sigma}_{-i}^h|\tilde{h}) \cdot (\bar{\sigma}_{-i}^h|\bar{h})(S_{-i}^{\bar{h}}(\tilde{h})) + \sum_{z \in H_{a_i}^{\bar{h}} \cap Z} u_i(z) \cdot (\bar{\sigma}_{-i}^h|\bar{h})(S_{-i}^{\bar{h}}(z)) \leq \pi_i(\tilde{\sigma}^h|\bar{h}).$$

<sup>64</sup>See the relationship between continuation and sequential best replies in Section 3.

By  $\tilde{\sigma}^h =^* \hat{\sigma}^h$ ,  $\pi_i(\tilde{\sigma}^h|\bar{h}) = \pi_i(\hat{\sigma}^h|\bar{h})$ . By  $\bar{\sigma}_{-i}^h =^* \tilde{\sigma}_{-i}^h =^* \hat{\sigma}_{-i}^h$  and  $\heartsuit$ , for every  $\tilde{h} \in H_{a_i}^{\bar{h}}$ ,  $(\bar{\sigma}_{-i}^h|\bar{h})(S_{-i}^{\bar{h}}(\tilde{h})) = (\hat{\sigma}_{-i}^h|\bar{h})(S_{-i}^{\bar{h}}(\tilde{h}))$ . For each  $\tilde{h} \in H_{a_i}^{\bar{h}} \setminus Z$ ,  $\bar{\sigma}_{-i}^h \geq \hat{\sigma}_{-i}^h|\bar{h}$ . So  $\star$  holds.  $\blacksquare$

For the next two Lemmata, fix a set of unordered histories  $\hat{H} \subseteq H^h$  and a set of SPE  $\Sigma^{\hat{H}} = (\sigma^{\tilde{h}})_{\tilde{h} \in \hat{H}}$  such that:

A1 for every  $\tilde{h} \in \hat{H}$ , there exists an equilibrium  $\tilde{\sigma}^{\tilde{h}} =^* \sigma^{\tilde{h}}$  such that  $\tilde{\sigma}^{\tilde{h}} \in \Delta(S_n^h(\tilde{h})|\tilde{h})$ .

Lemma 8 augments Lemma 4 for the simpler case in which  $\bar{S}^h$  is the support of a SPE. If only until step  $n - 1$  an equilibrium that mimics the SPE survives, at step  $n$  one deviation is always *strictly* preferred to continuing as the equilibrium prescribes.

**Lemma 8** *Suppose that there exists  $\sigma^h \in E^h(\Sigma^{\hat{H}})$  such that  $\tilde{\sigma}^h =^* \sigma^h$  but there is no equilibrium  $\hat{\sigma}^h \in \Delta(S_n^h)$  such that  $\hat{\sigma}^h =^* \sigma^h$ . Then there exist  $l \in I$  and  $\hat{h} \in D_l \setminus (\cup_{\tilde{h} \in \hat{H}} H^{\tilde{h}})$  such that for every  $\hat{\sigma}_{-l}^{\hat{h}} \in \Delta(S_{-l,n}^h(\hat{h})|\hat{h})$ ,  $v \leq n$  and  $\tilde{\mu}_l^{\hat{h}} \geq \hat{\sigma}_{-l}^{\hat{h}}$  t.s.b.  $(S_{-l,q}^h(\hat{h})|\hat{h})_{q=0}^v$ , there exists  $\tilde{\mu}_l^h =^* \hat{\sigma}_{-l}^h$  t.s.b.  $(S_{-l,q}^h)_{q=0}^{v-1}$  such that  $\tilde{\mu}_l^h =^{\hat{h}} \tilde{\mu}_l^{\hat{h}}$  and  $\rho(\tilde{\mu}_l^h)(\hat{h}) \neq \emptyset$  (so by A0  $\rho(\tilde{\mu}_l^h) \subseteq S_{l,v}^h$ ).*

**Proof.** Suppose not. For every  $i \in I$  and  $\hat{h} \in D_i \setminus (\cup_{\tilde{h} \in \hat{H}} H^{\tilde{h}}) =: \bar{D}_i$ , fix  $\hat{\sigma}_{-i}^{\hat{h}}$ ,  $v(\hat{h})$  and  $\tilde{\mu}_i^{\hat{h}}$  that violate the statement and let  $\tilde{\sigma}_{-i}^{\hat{h}} := \hat{\sigma}_{-i}^{\hat{h}}$ .

By Lemma 5 there exists  $\hat{\sigma}^h \in \Delta(S_n^h)$  such that  $\hat{\sigma}^h =^* \tilde{\sigma}^h =^* \sigma^h$  and for every  $i \in I$ ,  $\tilde{h} \in \hat{H} \cup \bar{D}_i$  and  $\hat{h} \in D_i \cap H^{\tilde{h}}$ ,  $\hat{\sigma}_{-i}^h|\hat{h} = \tilde{\sigma}_{-i}^{\tilde{h}}|\hat{h}$ . I show that  $\hat{\sigma}^h$  is an equilibrium, a contradiction.

Fix  $\bar{h} \in H^\sigma$ ,  $i \in I$  and  $a_i \in A_i(\bar{h}) \setminus \tilde{\sigma}_i^h|\bar{h}$ . If there exists  $\tilde{h} \preceq \bar{h}$  such that  $\tilde{h} \in \hat{H}$ ,  $\hat{\sigma}^h =^* \sigma^h =^{\tilde{h}} \tilde{\sigma}^h =^* \tilde{\sigma}^h$ , so by  $\spadesuit$   $\hat{\sigma}^h =^{\tilde{h}} \tilde{\sigma}^h$ ; then  $\hat{\sigma}^h|\bar{h}$  is an equilibrium, so by Lemma 6 ("only if")  $\star$  holds. If  $H_{a_i}^{\bar{h}} \setminus Z \subseteq \hat{H}$ , for every  $\hat{h} \in H_{a_i}^{\bar{h}} \setminus Z$ ,  $\hat{\sigma}_{-i}^h|\hat{h} = \tilde{\sigma}_{-i}^{\hat{h}}$  and by  $\diamond$ ,  $\pi(\tilde{\sigma}_{-i}^{\hat{h}}) = \pi(\sigma_{-i}^h|\hat{h})$ ; by  $\hat{\sigma}^h =^* \sigma^h$ ,  $\pi_i(\hat{\sigma}^h|\bar{h}) = \pi_i(\sigma^h|\bar{h})$  and by  $\heartsuit$ ,  $(\hat{\sigma}_{-i}^h|\bar{h})(S_{-i}^{\bar{h}}(\hat{h})) = (\sigma_{-i}^h|\bar{h})(S_{-i}^{\bar{h}}(\hat{h}))$ ; so since  $\sigma^h|\bar{h}$  is an equilibrium by Lemma 6 ("only if")  $\star$  holds. If  $H_{a_i}^{\bar{h}} \cap \bar{D}_i \neq \emptyset$ , fix  $v := \min_{\tilde{h} \in H_{a_i}^{\bar{h}} \cap \bar{D}_i} v(\tilde{h})$  and  $\hat{h} := \arg \min_{\tilde{h} \in H_{a_i}^{\bar{h}} \cap \bar{D}_i} v(\tilde{h})$ : for every  $\tilde{h} \in H_{a_i}^{\bar{h}} \cap \bar{D}_i$ ,  $\hat{\sigma}_{-i}^h|\tilde{h} \leq \tilde{\mu}_i^{\tilde{h}}(\cdot|\tilde{h}) \in \Delta(S_{-i,v}^h(\tilde{h})|\tilde{h}) \neq \emptyset^{65}$  and for every  $\tilde{h} \in H_{a_i}^{\bar{h}} \setminus (Z \cup \bar{D}_i)$ ,  $\hat{\sigma}_{-i}^h|\tilde{h} \in \Delta(S_{-i,v}^h(\tilde{h})|\tilde{h})$ ; therefore by Lemma 7  $\star$  holds. Thus by Lemma 6 ("if"),  $\hat{\sigma}^h$  is an equilibrium.  $\blacksquare$

Lemma 9 is the "dual" of Lemma 8: if an equilibrium has survived  $n$  steps but it does not mimic a SPE (within a subset), then there is a deviation from one of the

<sup>65</sup>For every  $j \neq i$ , there is  $\mu_j^h =^h \tilde{\sigma}_{-j}^h$  t.s.b.  $(S_{-l,q}^h)_{q=0}^{n-1}$ : by  $\clubsuit$ ,  $\rho(\mu_j^h)(\tilde{h}) \neq \emptyset$ ; by A0  $\rho(\mu_j^h) \subseteq S_{j,n}^h$ .

equilibrium paths that the deviator could take whenever thereafter she expects at least the payoff of a SPE of the subgame (within a subset).

**Lemma 9** *Suppose that  $\tilde{\sigma}^h \in \Delta(S_n^h)$  and for every  $\tilde{h} \in \widehat{H}$ ,  $\tilde{\sigma}^h =_{\tilde{h}} \tilde{\sigma}^{\tilde{h}}$ , but there is no  $\sigma^h \in E^h(\Sigma^{\widehat{H}})$  with  $\tilde{\sigma}^h =^* \sigma^h$ . Then there exist  $p \in I$  and  $\bar{h} \in D_p(\tilde{\sigma}^h) \setminus (\cup_{\tilde{h} \in \widehat{H}} H^{\tilde{h}})$  such that for every  $\sigma^{\bar{h}} \in E^{\bar{h}}(\Sigma^{\widehat{H}})$ ,  $v \leq n$  and  $\tilde{\mu}_p^{\bar{h}} \geq \sigma_{-p}^{\bar{h}}$  t.s.b.  $(S_{-p,q}^h(\bar{h})|_{\bar{h}})_{q=0}^v$ , there exists  $\tilde{\mu}_p^h =^* \tilde{\sigma}_{-p}^h$  t.s.b.  $(S_{-p,q}^h)_{q=0}^{v-1}$  such that  $\tilde{\mu}_p^h =_{\bar{h}} \tilde{\mu}_p^{\bar{h}}$  and  $\rho(\tilde{\mu}_p^h)(\bar{h}) \neq \emptyset$  (so by A0  $\rho(\tilde{\mu}_p^h) \subseteq S_{p,v}^h$ ).*

**Proof.** Suppose not. For every  $i \in I$  and  $\hat{h} \in D_i \setminus (\cup_{\tilde{h} \in \widehat{H}} H^{\tilde{h}}) =: \bar{D}_i$ , fix  $\sigma^{\hat{h}}$ ,  $v(\hat{h})$  and  $\tilde{\mu}_i^{\hat{h}}$  that violate the statement.

Construct  $\hat{\sigma}^h =^* \tilde{\sigma}^h$  such that for every  $\tilde{h} \in \cup_{i \in I} \bar{D}_i \cup \widehat{H}$  and  $\hat{h} \in \cup_{i \in I} D_i \cap H^{\tilde{h}}$ ,  $\hat{\sigma}^h|_{\hat{h}} = \sigma^{\tilde{h}}|_{\hat{h}} \in E^{\tilde{h}}(\Sigma^{\widehat{H}})$ , and for every  $\tilde{h} \notin \cup_{i \in I} \bar{D}_i$  such that  $p(\tilde{h}) \in H^\sigma$ ,  $\hat{\sigma}^h|_{\tilde{h}} \in E^{\tilde{h}}(\Sigma^{\widehat{H}})$ . I show that  $\hat{\sigma}^h$  is an equilibrium such that for every  $\tilde{h} \in \widehat{H} \cap H^\sigma$ ,  $\hat{\sigma} =_{\tilde{h}} \sigma^{\tilde{h}} \in E^{\tilde{h}}(\Sigma^{\widehat{H}})$ . Then, by Lemma 6 ("only if") for every  $\tilde{h} \in H^\sigma$ ,  $\hat{\sigma}^h|_{\tilde{h}}$  is an equilibrium, and so  $\hat{\sigma}^h \in E^h(\Sigma^{\widehat{H}})$ , a contradiction.

Fix  $\bar{h} \in H^\sigma$ ,  $i \in I$  and  $a_i \in A_i(\bar{h}) \setminus \tilde{\sigma}_i^h[\bar{h}]$ . If there exists  $\tilde{h} \preceq \bar{h}$  such that  $\tilde{h} \in \widehat{H}$ ,  $\tilde{\sigma}^h =^* \tilde{\sigma}^{\tilde{h}} =_{\tilde{h}} \tilde{\sigma}^{\tilde{h}} =^* \sigma^{\tilde{h}}$ , so by  $\spadesuit$   $\hat{\sigma} =_{\tilde{h}} \sigma^{\tilde{h}}$ ; then  $\hat{\sigma}^h|_{\bar{h}}$  is an equilibrium, so by Lemma 6 ("only if")  $\star$  holds. If  $H_{a_i}^{\bar{h}} \setminus Z \subseteq \widehat{H}$ , for every  $\hat{h} \in H_{a_i}^{\bar{h}} \setminus Z$ ,  $\hat{\sigma}_{-i}^h|_{\hat{h}} = \sigma_{-i}^{\hat{h}}$  and by  $\blacklozenge$ ,  $\pi(\sigma_{-i}^{\hat{h}}) = \pi(\hat{\sigma}_{-i}^h|_{\hat{h}})$ ; by  $\hat{\sigma}^h =^* \tilde{\sigma}^h$ ,  $\pi_i(\hat{\sigma}^h|_{\bar{h}}) = \pi_i(\tilde{\sigma}^h|_{\bar{h}})$  and by  $\heartsuit$ ,  $(\hat{\sigma}_{-i}^h|_{\bar{h}})(S_{-i}^{\bar{h}}(\hat{h})) = (\tilde{\sigma}_{-i}^h|_{\bar{h}})(S_{-i}^{\bar{h}}(\hat{h}))$ ; so since  $\tilde{\sigma}^h|_{\bar{h}}$  is an equilibrium by Lemma 6 ("only if")  $\star$  holds. If  $H_{a_i}^{\bar{h}} \cap \bar{D}_i \neq \emptyset$ , fix  $v := \min_{\tilde{h} \in H_{a_i}^{\bar{h}} \cap \bar{D}_i} v(\tilde{h})$  and  $\hat{h} := \arg \min_{\tilde{h} \in H_{a_i}^{\bar{h}} \cap \bar{D}_i} v(\tilde{h})$ : for every  $\tilde{h} \in H_{a_i}^{\bar{h}} \cap \bar{D}_i$ ,  $\hat{\sigma}_{-i}^h|_{\tilde{h}} \leq \tilde{\mu}_i^{\tilde{h}}(\cdot|_{\tilde{h}}) \in \Delta(S_{-i,v}^h(\tilde{h})|_{\tilde{h}}) \neq \emptyset^{66}$  and for every  $\tilde{h} \in H_{a_i}^{\bar{h}} \setminus (Z \cup \bar{D}_i)$ ,  $\hat{\sigma}_{-i}^h|_{\tilde{h}} \leq \tilde{\sigma}^{\tilde{h}} \in \Delta(S_{-i,v}^h(\tilde{h})|_{\tilde{h}})$ ; hence by Lemma 7  $\star$  holds. Thus by Lemma 6 ("if")  $\hat{\sigma}^h$  is an equilibrium.  $\blacksquare$

Now I can prove the main Lemma.

**Lemma 10** *Fix  $h \in H$ ,  $m \in \mathbb{N}$ , a red. proced.  $(S_q^h)_{q \geq 0}$ , a set of unordered histories  $\widehat{H} = \{h^1, \dots, h^w\} \subseteq H^h$  and a set of SPE  $\Sigma^{\widehat{H}} = (\sigma^{\tilde{h}})_{\tilde{h} \in \widehat{H}}$  s.t. A1 holds for  $n = m$  and:*

A2 *for every  $v \leq w$ , there exists an equilibrium  $\hat{\sigma}^{h,v} \in \Delta(S_{m-1}^h)$  such that  $h^v \in \cup_{i \in I} D_i(\hat{\sigma}^{h,v})$  and for every  $q < v$ , if  $h^q \in H(\hat{\sigma}^{h,v})$ ,  $\hat{\sigma}^{h,v} =^* \tilde{\sigma}^{h^q}$ ;*

A3 *for each  $i \in I$ ,  $n \leq m$ ,  $\sigma^h \in E^h(\Sigma^{\widehat{H}})$  and  $\mu_i^h \geq \sigma_{-i}^h$  t.s.b.  $(S_{-i}^h)_{q=0}^{n-1}$ ,  $\rho(\mu_i^h) \subseteq S_{i,n}^h$ .*

<sup>66</sup>See the previous footnote.



Then there exist  $\sigma^h \in E^h(\Sigma^{\widehat{H}})$  and an equilibrium  $\tilde{\sigma}^h \in \Delta(S_m^h)$  such that  $\tilde{\sigma}^h =^* \sigma^h$ .

**Proof.** The proof is by induction on the depth of  $\Gamma(h)$ .

**Inductive hypothesis (d)**

The Lemma holds for every  $h \in H$  such that  $\Gamma(h)$  has depth not bigger than  $d$ .

**Basis step (1)** For every  $i \in I$ ,  $n \leq m$  and equilibrium of  $\Gamma(h)$   $\sigma^h$  such that  $\text{Supp}\sigma^h \subseteq S_{n-1}^h$ , by A3  $r(\sigma_{-i}^h) \subseteq S_{i,n}^h$ . Inductively,  $\text{Supp}\sigma^h \subseteq S_m^h$ .

**INDUCTIVE STEP (d+1)** Suppose not. I will find a contradiction through a recursive procedure. Set  $k = 0$  and  $\overline{H}^0 := \widehat{H}$ .

**Recursive step (k)**

If  $k > 0$ ,  $\overline{H}^k$  and  $\Sigma^{\overline{H}^k}$  are defined in step  $k - 1$ . Let  $n \leq m$  be the greatest  $q \in \mathbb{N}$  such that there exist  $\sigma^{h,k} \in E^h(\Sigma^{\overline{H}^k})$  and an equilibrium  $\tilde{\sigma}^{h,k} \in \Delta(S_{q-1}^h)$  with  $\tilde{\sigma}^{h,k} =^* \sigma^{h,k}$ . If  $k > 0$ , by the last remark of the previous steps,  $E^h(\Sigma^{\overline{H}^k}) \subseteq \dots \subseteq E^h(\Sigma^{\overline{H}^0})$ . Then by  $\blacklozenge$  A3 implies A0. Moreover,  $n$  weakly decreases with  $k$ . Then  $\Sigma^{\overline{H}^k}$  satisfies A1 and A2 with  $n$  in place of  $m$ .<sup>67</sup> Lemma 8 yields  $l \in I$  and  $\widehat{h} \in D_l(\tilde{\sigma}^{h,k}) \setminus \cup_{\bar{h} \in \overline{H}^k} H^{\bar{h}}$ .

Define the reduction procedure  $((S_{i,q}^{\widehat{h}})_{i \in I})_{q=0}^\infty := ((S_{i,q}^h(\widehat{h})|_{\widehat{h}})_{i \in I})_{q=0}^\infty$ . Fix  $i \neq l$  and  $v \leq n$ . For every  $\tilde{\mu}_i^{\widehat{h}}$  t.s.b.  $(S_{-i,q}^{\widehat{h}})_{q=0}^{v-1}$ , since  $\widehat{h} \in D_l(\tilde{\sigma}^{h,k})$ , by Lemma 1 there exists  $\tilde{\mu}_i^h =^h \tilde{\sigma}_{-i}^{h,k}$  t.s.b.  $(S_{-i,q}^h)_{q=0}^{v-1}$  such that  $\tilde{\mu}_i^h =^{\widehat{h}} \tilde{\mu}_i^{\widehat{h}}$ . By A0,  $\rho(\tilde{\mu}_i^h) \subseteq S_{i,v}^h$ ; by  $\clubsuit$ ,  $\rho(\tilde{\mu}_i^h)(\widehat{h}) \neq \emptyset$ . Hence, together with Lemma 8, for every  $i \in I$ ,  $v \leq n$ ,  $\widehat{\sigma}_{-i}^{\widehat{h}} \in \Delta(S_{-i,n}^{\widehat{h}})$  and  $\tilde{\mu}_i^{\widehat{h}} \geq \widehat{\sigma}_{-i}^{\widehat{h}}$  t.s.b.  $(S_{-i,q}^{\widehat{h}})_{q=0}^v$ ,  $\rho(\tilde{\mu}_i^{\widehat{h}}) \subseteq S_{i,v}^{\widehat{h}} \neq \emptyset$  (**F**).

Let  $\widetilde{H}^0 := \overline{H}^k \cap H^{\widehat{h}}$ : I show that there exist  $\sigma^{\widehat{h}} \in E^{\widehat{h}}(\Sigma^{\widetilde{H}^0})$  and an equilibrium  $\tilde{\sigma}^{\widehat{h}} =^* \sigma^{\widehat{h}}$  such that  $\text{Supp}\tilde{\sigma}^{\widehat{h}} \subseteq S_n^{\widehat{h}}$ . Suppose not (**G**). I will find a contradiction through a recursive procedure. For every  $q \leq w + k$  such that  $h^q \in \widetilde{H}^0$ , let  $\overline{h}^q := h^q$  and  $\widehat{\sigma}^{\widehat{h},q} := \widehat{\sigma}^{h^q}|_{\widehat{h}}$ , which is an equilibrium because  $\cup_{\bar{h} \in \overline{H}^k} H^{\bar{h}} \not\ni \widehat{h} \prec \overline{h}^q \in \cup_{i \in I} D_i(\widehat{\sigma}^{h^q})$ , so  $\widehat{h} \in H(\widehat{\sigma}^{h^q})$ . Set  $t = 0$ .

**Recursive step (t)** If  $t > 0$ ,  $\widetilde{H}^t$  and  $\Sigma^{\widetilde{H}^t}$  are defined in step  $t - 1$ , and satisfy A1 and A2 with  $n$  in place of  $m$  and  $\widehat{h}$  in place of  $h$ .<sup>68</sup> For every  $i \in I$ , let  $\Sigma_i^{\widehat{h},t}$  be the set of  $\widehat{\sigma}_i^{\widehat{h}} \in \Delta(S_{i,n}^{\widehat{h}})$  such that for every  $\tilde{h} \in \widetilde{H}^t \cap H(\widehat{\sigma}_i^{\widehat{h}})$ ,  $\widehat{\sigma}_i^{\widehat{h}} =^{\tilde{h}} \tilde{\sigma}_i^{\tilde{h}}$ .

First I show that  $\Sigma^{\widehat{h},t}$  is non-empty and features an equilibrium of  $\Gamma(\widehat{h})$ . Let

<sup>67</sup>For every  $h^q \in \overline{H}^k$ ,  $\widehat{\sigma}^{h^q} \in \Delta(S_n^h(h^q)|_{h^q})$  and  $\widehat{\sigma}^{h^q} \in \Delta(S_{n-1}^h)$  come from A1 and A2 if  $h^q \in \widehat{H}$ , from some previous step if  $h^q \notin \widehat{H}$ .

<sup>68</sup>For every  $\overline{h}^q \in \widetilde{H}^t$ ,  $\widehat{\sigma}^{\overline{h}^q} \in \Delta(S_n^h(\overline{h}^q)|_{\overline{h}^q})$  and  $\widehat{\sigma}^{\overline{h}^q} \in \Delta(S_{n-1}^h)$  come from the outer recursive step if  $\overline{h}^q \in \overline{H}^k$ , from some previous step if  $\overline{h}^q \notin \overline{H}^k$ .

$\tau := w + k + t$ . Note that for every  $i \in I$ ,

$$\Sigma_i^{\hat{h},t} = \cap_{\tilde{h} \in \tilde{H}^t} \cap_{z \in Z^{\tilde{h}} \cup H^{\tilde{h}}} \{ \hat{\sigma}_i^{\tilde{h}} \in \Delta(S_i^{\tilde{h}}) : \hat{\sigma}_i^{\tilde{h}}(S_i^{\tilde{h}}(z)) = \hat{\sigma}_i^{\tilde{h}}(S_i^{\tilde{h}}(\tilde{h})) \cdot \tilde{\sigma}_i^{\tilde{h}}(S_i^{\tilde{h}}(z)) \} \cap \Delta(S_{i,n}^{\hat{h}}),$$

an intersection of convex and compact sets.<sup>69</sup> Hence  $\Sigma_i^{\hat{h},t}$  is convex and compact. Then, since expected utility is linear, the reduced game with strategy sets  $(\Sigma_i^{\hat{h},t})_{i \in I}$ , if non-empty, features an equilibrium  $\tilde{\sigma}^{\hat{h},t}$ . Fix  $i \in I$  and  $\mu_i^{\hat{h}} = \tilde{h} \tilde{\sigma}_{-i}^{\hat{h},t}$  t.s.b.  $(S_{-i,q}^{\hat{h}})_{q=0}^n$ . There exists  $\hat{\sigma}_{-i}^{\hat{h},\tau+1} \in \Delta(\rho(\mu_i^{\hat{h}})) \subseteq \Delta(r(\tilde{\sigma}_{-i}^{\hat{h},t}))$  such that for every  $\tilde{h} \in \tilde{H}^t \cap H(\tilde{\sigma}_{-i}^{\hat{h},t})$ , since  $\tilde{\sigma}_{-i}^{\hat{h},t} = \tilde{h} \tilde{\sigma}_{-i}^{\hat{h}}$ , by  $\clubsuit$   $\hat{\sigma}_{-i}^{\hat{h},\tau+1} = \tilde{\sigma}_{-i}^{\hat{h}}$ . By F,  $\hat{\sigma}_{-i}^{\hat{h},\tau+1} \in \Delta(S_{i,n}^{\hat{h}})$ . If  $t = 0$  and  $\tilde{H}^0 = \emptyset$ ,  $\Sigma_i^{\hat{h},0} = \Delta(S_i^{\hat{h}}) \neq \emptyset$  (by F) and  $\hat{\sigma}_{-i}^{\hat{h},\tau+1} \in \Sigma_i^{\hat{h},0}$ , so that  $\tilde{\sigma}_{-i}^{\hat{h},t} \in \Delta(r(\tilde{\sigma}_{-i}^{\hat{h},t}))$  too. Else, for notational convenience let  $\hat{\sigma}_{-i}^{\hat{h},\tau+1} := \hat{\sigma}_{-i}^{\hat{h},t}$  and proceed as follows.

For every  $\tilde{h} \in H^{\hat{h}}$  and  $q \leq \tau$ , let  $Q_q^{\tilde{h}} := \{g \leq q : \bar{h}^g \in \tilde{H}^t \cap H^{\tilde{h}}\}$  and  $Q_{\tau+1}^{\tilde{h}} := Q_{\tau}^{\tilde{h}}$ . I show that for every  $q \in Q_{\tau}^{\tilde{h}} \cup \{\tau+1\}$ , there exists  $\bar{\sigma}_i^{\tilde{h},q} \in \Delta(S_{i,n}^{\hat{h}})$  such that  $\bar{\sigma}_i^{\tilde{h},q} = \tilde{\sigma}_i^{\tilde{h},q}$  and for every  $g \in Q_q^{\tilde{h}}$ ,  $\bar{\sigma}_i^{\tilde{h},q} = \bar{h}^g \tilde{\sigma}_i^{\tilde{h},g}$ . Fix  $q \in Q_{\tau}^{\tilde{h}} \cup \{\tau+1\}$  and suppose to have shown it already for every  $g \in Q_q^{\tilde{h}} \setminus \{q\}$ . Since  $\tilde{\sigma}_i^{\tilde{h},q} \in \Delta(r(\tilde{\sigma}_{-i}^{\tilde{h},q}))$ , by  $\hat{\sigma}_i^{\tilde{h},q} \in \Delta(S_n^{\hat{h}})$  and F, Lemma 5 yields  $\bar{\sigma}_i^{\tilde{h},q} \in \Delta(S_{i,n}^{\hat{h}})$  such that  $\bar{\sigma}_i^{\tilde{h},q} = \tilde{\sigma}_i^{\tilde{h},q}$  and for every  $\tilde{h} \in D_{-i}(\tilde{\sigma}_i^{\tilde{h},q})$ : if  $\tilde{h} \in \tilde{H}^t$ ,  $\bar{\sigma}_i^{\tilde{h},q} | \tilde{h} = \tilde{\sigma}_i^{\tilde{h}}$ ; if  $\tilde{h} \notin \tilde{H}^t$  but  $Q_q^{\tilde{h}} \neq \emptyset$ ,  $\bar{\sigma}_i^{\tilde{h},q} | \tilde{h} = \bar{\sigma}_i^{\tilde{h},\max Q_q^{\tilde{h}}} | \tilde{h}$  (where  $\max Q_q^{\tilde{h}} < q$  because  $\bar{h}^q \in \tilde{H}^t \cap (\cup_{j \in I} D_j(\tilde{\sigma}_j^{\tilde{h},q}))$  and  $\tilde{H}^t \cap (\cup_{j \in I} D_j(\tilde{\sigma}_j^{\tilde{h},q})) \cap H^{\tilde{h}} = \emptyset$ ); if  $\tilde{h} \succ \bar{h}$  for some  $\bar{h} \in \tilde{H}^t$ ,  $\bar{\sigma}_i^{\tilde{h},q} | \tilde{h} = \tilde{\sigma}_i^{\tilde{h}} | \tilde{h}$ , so that by  $\bar{\sigma}_i^{\tilde{h},q} = \tilde{\sigma}_i^{\tilde{h},q} = \tilde{\sigma}_i^{\tilde{h}}$  and  $\spadesuit$ ,  $\bar{\sigma}_i^{\tilde{h},q} = \bar{h} \tilde{\sigma}_i^{\tilde{h}}$ . Then  $\bar{\sigma}_i^{\tilde{h},\tau} \in \Sigma_i^{\hat{h},t} \neq \emptyset$ .<sup>70</sup> So  $\tilde{\sigma}^{\hat{h},t}$  and  $\hat{\sigma}^{\hat{h},\tau+1}$  exist and  $\bar{\sigma}_i^{\hat{h},\tau+1} \in \Sigma_i^{\hat{h},t}$ . Since  $\bar{\sigma}_i^{\hat{h},\tau+1} = \tilde{\sigma}_i^{\hat{h},\tau+1}$ ,  $\bar{\sigma}_i^{\hat{h},\tau+1} \in \Delta(r(\tilde{\sigma}_{-i}^{\hat{h},t}))$ . Then  $\tilde{\sigma}_i^{\hat{h},t} \in \Delta(r(\tilde{\sigma}_{-i}^{\hat{h},t}))$  too.

By  $\blacklozenge$ , F implies A0 with  $\hat{h}$  in place of  $h$ ;  $\tilde{H}^t$  satisfies A1 with  $n$  in place of  $m$ . By the last remark of the previous steps  $E^h(\Sigma^{\tilde{H}^t}) \subseteq E^h(\Sigma^{\tilde{H}^0})$ , so, by G,  $\tilde{\sigma}^{\hat{h},t}$  satisfies the hypotheses of Lemma 9.<sup>71</sup> Lemma 9 yields  $p \in I$  and  $\bar{h} \in D_p(\tilde{\sigma}^{\hat{h},t}) \setminus \cup_{\tilde{h} \in \tilde{H}^t} H^{\tilde{h}}$ .

Define the reduction procedure  $((S_{i,q}^{\bar{h}})_{i \in I})_{q=0}^{\infty} := ((S_{i,q}^{\hat{h}}(\bar{h}) | \bar{h})_{i \in I})_{q=0}^{\infty}$ . Fix  $i \neq p$  and  $v \leq n$ . For every  $\tilde{\mu}_i^{\bar{h}}$  t.s.b.  $(S_{-i,q}^{\bar{h}})_{q=0}^v$ , since  $\bar{h} \in D_p(\tilde{\sigma}^{\hat{h},t})$ , by Lemma 1 there exists  $\tilde{\mu}_i^{\bar{h}} = \bar{h} \tilde{\sigma}_{-i}^{\hat{h},t}$  t.s.b.  $(S_{-i,q}^{\bar{h}})_{q=0}^{v-1}$  such that  $\tilde{\mu}_i^{\bar{h}} = \bar{h} \tilde{\mu}_i^{\bar{h}}$ . By A0,  $\rho(\tilde{\mu}_i^{\bar{h}}) \subseteq S_{i,v}^{\bar{h}}$ ; by  $\clubsuit$ ,  $\rho(\tilde{\mu}_i^{\bar{h}})(\bar{h}) \neq \emptyset$ . Hence, together with Lemma 9, for every  $i \in I$ ,  $v \leq n$ ,  $\sigma^{\bar{h}} \in E^{\bar{h}}(\Sigma^{\tilde{H}^t})$  and  $\tilde{\mu}_i^{\bar{h}} \geq \sigma_{-i}^{\bar{h}}$  t.s.b.  $(S_{-i,q}^{\bar{h}})_{q=0}^v$ ,  $\rho(\mu_i^{\bar{h}}) \subseteq S_{i,v}^{\bar{h}}$ . Moreover, for every  $q \in Q_{\tau}^{\bar{h}}$ , since

<sup>69</sup>Clearly,  $\Delta(S_{i,n}^{\hat{h}})$  is convex and compact. Each set of the kind  $\{ \hat{\sigma}_i^{\tilde{h}} \in \Delta(S_i^{\tilde{h}}) : (\hat{\sigma}_i^{\tilde{h}})(S_i^{\tilde{h}}(z)) = (\hat{\sigma}_i^{\tilde{h}})(S_i^{\tilde{h}}(\tilde{h})) \cdot c \}$ , where  $c$  is a constant, is clearly convex and compact too. Notice that if  $\tilde{h} \notin H(\tilde{\sigma}_i^{\tilde{h}})$ ,  $\hat{\sigma}_i^{\tilde{h}}$  satisfies  $(\hat{\sigma}_i^{\tilde{h}})(S_i^{\tilde{h}}(z)) = (\hat{\sigma}_i^{\tilde{h}})(S_i^{\tilde{h}}(\tilde{h})) \cdot c$  as  $0 = 0$ .

<sup>70</sup>By recursive step  $t-1$ ,  $\max Q_{\tau}^{\tilde{h}} = \tau$ .

<sup>71</sup>Without loss of generality assume that for every  $i \in I$  and  $\tilde{h} \in \tilde{H}^t \setminus H(\tilde{\sigma}_i^{\hat{h},t})$ ,  $\bar{\sigma}_i^{\hat{h},t} = \tilde{h} \tilde{\sigma}_i^{\hat{h}}$ .

$\bar{h} \notin \tilde{H}^t$ ,  $\bar{h}^q \succ \bar{h} \in H(\hat{\sigma}^{\bar{h},q})$ , so set  $\hat{\sigma}^{\bar{h},q} := \hat{\sigma}^{\bar{h},q}|_{\bar{h}}$ . Then A3, A2 and A1 are satisfied with  $\tilde{H}^t \cap H^{\bar{h}}$  in place of  $\Sigma^{\hat{H}}$ ,  $n$  in place of  $m$  and  $\bar{h}$  in place of  $h$ . So by the I.H. there exist  $\sigma^{\bar{h}} \in E^{\bar{h}}(\Sigma^{\tilde{H}^t})$  and an equilibrium  $\tilde{\sigma}^{\bar{h}} = * \sigma^{\bar{h}}$  such that  $\text{Supp} \tilde{\sigma}^{\bar{h}} \subseteq S_n^{\bar{h}} = S_n^{\hat{h}}(\bar{h})|_{\bar{h}}$ .

Since  $\bar{h} \notin \cup_{\tilde{h} \in \tilde{H}^t} H^{\tilde{h}}$  and  $H^{\hat{h}}$  is finite,  $\tilde{H}^{t+1} := \{\tilde{h} \in \tilde{H}^t : \tilde{h} \not\succeq \bar{h}\} \cup \{\bar{h}\}$  is a set of unordered histories that keep shortening with  $t$ , until a contradiction is obtained. Before, let  $\Sigma^{\tilde{H}^{t+1}} := \Sigma^{\tilde{H}^t} \cup \{\sigma^{\bar{h}}\} \setminus (\sigma^{\bar{h}})_{\tilde{h} \in \tilde{H}^t \setminus \tilde{H}^{t+1}}$ ,  $\bar{h}^{\tau+1} := \bar{h}$ ,  $\hat{\sigma}^{\bar{h},\tau+1} := \hat{\sigma}^{\bar{h},t}$ .<sup>72</sup> Then, increase  $t$  by 1 and run again noting what follows: for every  $\tilde{h} \in \tilde{H}^t$  such that  $\tilde{h} \succ \bar{h}$ ,  $\sigma^{\bar{h}}|_{\tilde{h}} = \sigma^{\tilde{h}}$ , so that  $E^h(\Sigma^{\tilde{H}^{t+1}}) \subseteq E^h(\Sigma^{\tilde{H}^t})$ .

Since  $\hat{h} \notin \cup_{\bar{h} \in \bar{H}^k} H^{\bar{h}}$  and  $H^h$  is finite,  $\bar{H}^{k+1} := \{\bar{h} \in \bar{H}^k : \bar{h} \not\succeq \hat{h}\} \cup \{\hat{h}\}$  is a set of unordered histories that keep shortening with  $k$ , until a contradiction is obtained. Before, let  $\Sigma^{\bar{H}^{k+1}} := \Sigma^{\bar{H}^k} \cup \{\sigma^{\hat{h}}\} \setminus (\sigma^{\hat{h}})_{\bar{h} \in \bar{H}^k \setminus \bar{H}^{k+1}}$ ,  $h^{w+k+1} := \hat{h}$ ,  $\hat{\sigma}^{h,w+k+1} := \hat{\sigma}^{h,k}$ .<sup>73</sup> Increase  $k$  by 1 and run again noting what follows: for every  $\bar{h} \in \bar{H}^k$  such that  $\bar{h} \succ \hat{h}$ ,  $\sigma^{\hat{h}}|_{\bar{h}} = \sigma^{\bar{h}}$ , so that  $E^h(\Sigma^{\bar{H}^{k+1}}) \subseteq E^h(\Sigma^{\bar{H}^k})$ . ■

**Proof of theorem 2.** Fix  $j \in I$ ,  $h \in D_j(S_{\Delta^e}^\infty)$  and  $(S_n^h)_{n=0}^\infty = (S_{\Delta^e}^n(h)|_h)_{n=0}^\infty$ . Fix  $m \in \mathbb{N}$ ,  $i \neq j$  and  $\mu_i^h$  t.s.b.  $(S_{-i,n}^h)_{n=0}^{m-1}$ . By self-enforceability ( $\zeta(S_{\Delta^e}^\infty) = \{z\}$ ) there exists  $\mu_i$  t.s.b.  $(S_{-i,\Delta^e}^n)_{n=0}^{m-1}$  such that  $\mu_i(S_{-i}(h)|_p(h)) = 0$  and  $\rho(\mu_i)(h) \neq \emptyset$ . Hence, by Lemma 1 there exists  $\tilde{\mu}_i =^h \mu_i^h$  t.s.b.  $(S_{-i,\Delta^e}^n)_{n=0}^{m-1}$  such that  $\rho(\tilde{\mu}_i)(h) \neq \emptyset$ . So A3 holds for every  $i \neq j$ ; A1 and A2 hold with  $\hat{H} = \emptyset$ . By  $\zeta(S_{\Delta^e}^\infty) = \{z\}$  there is  $m \in \mathbb{N}$  such that  $S_m^h = \emptyset$ , so Lemma 10 cannot hold. Thus A3 must be violated for  $j$  and some  $v \leq m$ , SPE  $\sigma^h$  of  $\Gamma(h)$  and  $\mu_j^h \geq \sigma_{-j}^h$  t.s.b.  $(S_{-j,v}^h)_{q=0}^v$ . Fix  $s_{-j} \in S_{-j,\Delta^e}^\infty \subseteq S_{-j}(z)$ : for every  $s_{-j}^h \in \text{Supp} \mu_j^h(\cdot|h)$ , by Lemma 1 I can construct  $\tilde{s}_{-j} \in S_{-j,\Delta^e}^v(z)$  such that  $\tilde{s}_{-j} =^h s_{-j}^h$  and  $\tilde{s}_{-j} =^{H \setminus \{h\}} s_{-j}$ . Using all such  $\tilde{s}_{-j}$ 's I can construct  $\mu_j =^h \mu_j^h$  t.s.b.  $(S_{-j,v}^h)_{q=0}^v$  and  $S_{-j}(z)$  so that  $\rho(\mu_j)(h) = \emptyset$  and  $\rho(\mu_j)(z) \neq \emptyset$ . Then  $j$  prefers  $z$  to the distribution over outcomes induced by  $\sigma^h$ . Hence,  $z$  is a SPE path. ■

**Proof of theorem 6.** Lemma 10 can be applied with strong rationalizability as reduction procedure,  $h := h^0$ , empty  $\hat{H}$ , and  $m$  after convergence. ■

## References

- [1] Bassetto, M., ‘‘Equilibrium and government commitment’’, *Journal of Economic Theory*, **124**(1), 2005, 79-105.

<sup>72</sup>Overwrite the previous, temporary choice.

<sup>73</sup>Note that for each  $\bar{h} \in \bar{H}^k \cap H(\hat{\sigma}^{h,k})$ ,  $\tilde{\sigma}^{h,k} = * \sigma^{h,k} = * \sigma^{\bar{h}} = * \tilde{\sigma}^{\bar{h}}$ , so  $\tilde{\sigma}^{h,k}$  satisfies A2 with  $\bar{H}^{k+1}$ .

- [2] Battigalli, P., “Comportamento razionale ed equilibrio nei giochi e nelle situazioni sociali”, 1987, undergraduate dissertation, Universita’ Bocconi, Milano.
- [3] Battigalli, P., “Dynamic Consistency and Imperfect Recall”, *Games and Economic Behavior*, **20(1)**, 1997, 31-50.
- [4] Battigalli, P., “On Rationalizability in Extensive Games”, *Journal of Economic Theory*, **74**, 1997, 40-61.
- [5] Battigalli, P., “Rationalizability in Infinite, Dynamic Games of Incomplete Information”, *Research in Economics*, **57**, 2003, 1-38.
- [6] Battigalli, P. and M. Dufwenberg, “Dynamic psychological games”, *Journal of Economic Theory*, **144(1)**, 2009, 1-35.
- [7] Battigalli, P. and A. Friedenberg, “Forward induction reasoning revisited”, *Theoretical Economics*, **7**, 2012, 57-98.
- [8] Battigalli, P. and A. Prestipino, “Transparent Restrictions on Beliefs and Forward Induction Reasoning in Games with Asymmetric Information”, *The B.E. Journal of Theoretical Economics (Contributions)*, **13**, 2013, Issue 1.
- [9] Battigalli, P. and M. Siniscalchi, “Interactive beliefs, epistemic independence and strong rationalizability”, *Research in economics*, **53**, 1999, 247-273.
- [10] Battigalli, P. and M. Siniscalchi, “Strong Belief and Forward Induction Reasoning”, *Journal of Economic Theory*, **106**, 2002, 356-391.
- [11] Battigalli P. and M. Siniscalchi, “Rationalization and Incomplete Information,” *The B.E. Journal of Theoretical Economics*, **3(1)**, 2003, 1-46.
- [12] Battigalli, P. and M. Siniscalchi, “Interactive Epistemology in Games with Payoff Uncertainty”, *Research in Economics*, **61**, 2007, 165-184.
- [13] Catonini, E., “Selecting strongly rationalizable strategies”, working paper, 2015.
- [14] Chen, J., and S. Micali, “The order independence of iterated dominance in extensive games”, *Theoretical Economics*, **8**, 2013, 125-163.
- [15] Cho I.K. and D. Kreps, “Signaling Games and Stable Equilibria”, *Quarterly Journal of Economics*, **102**, 1987, 179-222.

- [16] Dufwenberg, M., Servátka, M., Vadovic, R., “ABC on Deals”, Working Papers in Economics 12/14, University of Canterbury, Dpt. of Economics and Finance.
- [17] Farrell, J. P., and Maskin, E., “Renegotiation in repeated games”, *Games and Economic Behavior*, **1(4)**, 1989, 327-360.
- [18] Fudenberg, D., and D. Levine, “Self-confirming equilibrium”, *Econometrica*, **61**, 1993, 523–546.
- [19] Gossner, O., “The robustness of incomplete penal codes in repeated interactions”, working paper, 2014.
- [20] Govindan, S., Wilson, R., “On forward induction,” *Econometrica*, **77**, 2009, 1-28.
- [21] Greenberg, J., Gupta, S., Luo, X., “Mutually acceptable courses of action”, *Economic Theory*, **40**, 2009, 91-112.
- [22] Harrington, J. “A Theory of Collusion with Partial Mutual Understanding”, working paper, 2015.
- [23] Heifetz, A., and A. Perea, “On the Outcome Equivalence of Backward Induction and Extensive Form Rationalizability”, *International Journal of Game Theory*, **44**, 2015, 37–59.
- [24] Kohlberg, E. and J.F. Mertens, “On the Strategic Stability of Equilibria”, *Econometrica*, **54**, 1986, 1003-1038.
- [25] Miller, D., and J. Watson, “A Theory of Disagreement in Repeated Games with Bargaining”, *Econometrica*, **81(6)**, 2013, 2303–2350.
- [26] Osborne, M., “Signaling, Forward Induction, and Stability in Finitely Repeated Games”, *Journal of Economic Theory*, **50**, 1990, 22-36.
- [27] Pearce, D., “Rational Strategic Behavior and the Problem of Perfection”, *Econometrica*, **52**, 1984, 1029-1050.
- [28] Reny, P., “Backward Induction, Normal Form Perfection and Explicable Equilibria”, *Econometrica*, **60(3)**, 1992, 627-49.
- [29] Renyi, A., “On a New Axiomatic Theory of Probability”, *Acta Mathematica Academiae Scientiarum Hungaricae*, **6**, 1955, 285-335.