

On non-monotonic strategic reasoning.*

Emiliano Catonini[†]

May 2018

Strong- Δ -Rationalizability (Battigalli 2003, Battigalli and Siniscalchi 2003) is a prominent and widely applied solution concept that introduces first-order belief restrictions in forward induction reasoning. In absence of restrictions, it coincides with Strong Rationalizability (Battigalli and Siniscalchi 2002). These solution concepts are based on the notion of Strong Belief (Battigalli and Siniscalchi, 2002). The non-monotonicity of Strong Belief implies that the predictions of Strong- Δ -Rationalizability under given restrictions can have empty intersection with the predictions of Strong Rationalizability. Here we show that the set of outcomes predicted by Strong- Δ -Rationalizability actually shrinks as long as (stricter and stricter) restrictions have no bite off-path. So, Strong- Δ -Rationalizability yields a subset of strongly rationalizable outcomes when the restrictions correspond to the belief in a particular path of play. Moreover, under such restrictions, the epistemic priority between belief in rationality and beliefs in the restrictions (Catonini, 2017) is irrelevant for the predicted outcomes: the predictions of Strong- Δ -Rationalizability and Selective Rationalizability (Catonini 2017) coincide. The workhorse lemma behind these results allows to show also the order independence of the "iterated elimination of never sequential best replies" (of which Strong Rationalizability is the maximal elimination order), and that Strong Rationalizability refines Backward Induction. The outcome equivalence of Strong Rationalizability and Backward Induction in perfect information games with no relevant ties (Battigalli 1997) follows.

Keywords: Strong- Δ -Rationalizability, Strong Rationalizability, First-order Belief Restrictions, Epistemic Priority, Order Independence, Backward Induction.

*The main results of this paper appeared already in "Non-binding Agreements and Forward Induction Reasoning" (Catonini 2015, mimeo).

[†]Higher School of Economics, ICEF, emiliano.catonini@gmail.com

1 Introduction

Strong Rationalizability (Battigalli and Siniscalchi [5]) is a form of extensive-form rationalizability (Pearce [15]) based on the notion of *Strong Belief*.¹ Concretely, it is the maximal iterated deletion of "never sequential best replies" under belief systems that assign probability 1, as long as possible, to opponents' strategies that survive the previous step of the procedure.² Strong- Δ -Rationalizability (Battigalli [3], Battigalli and Siniscalchi [6]) introduces first-order belief restrictions in the same reasoning scheme: only belief systems in an exogenously given set are allowed at all steps.

It is well-known that the introduction of first-order belief restrictions can let the elimination procedure depart completely from Strong Rationalizability. This is due to the non-monotonicity of strong belief: strong belief in a smaller event does not imply strong belief in a larger event. So, for instance, even in a perfect information game without relevant ties, the introduction of first-order belief restrictions can induce completely different outcomes with respect to the only strongly rationalizable one (see, e.g., the introductory example of Catonini [7]). Are there interesting conditions under which the introduction of first-order belief restrictions refines the set of strongly rationalizable outcomes? When such conditions are satisfied, the predictions are reassuringly robust to "restricted" and "unrestricted" forward induction reasoning, as captured, respectively, by Strong- Δ -Rationalizability and Strong Rationalizability.

It turns out that in all games with observable actions (i.e. games where, allowing for simultaneous moves, every player knows the current history of the game) the set of outcomes predicted by Strong- Δ -Rationalizability shrinks as more and more restrictions that "never bite off-path" are imposed. With this, I refer to restrictions that exclude belief systems only based on what they predict about opponents' behavior along the paths that survive all steps of Strong- Δ -Rationalizability. So, off-the-path restrictions are responsible for the general non-monotonicity of Strong- Δ -Rationalizability. The reason is the following. Suppose that at some step of reasoning, the behavior of an opponent that is object of a player's restricted beliefs ends up off-path. Then, the player will not verify whether such restricted beliefs are compatible with higher orders of rationality, because the off-path behavior of the opponent is no more refined by the elimination procedure.

Beside the theoretical insight, though, this broad condition for outcome monotonicity is of little practical use: one cannot verify it without actually performing Strong- Δ -Rationalizability. Yet, a very important class of restrictions always satisfies this condition:

¹i.e. belief as long as compatible with the observed behavior.

²The epistemic justification of Strong Rationalizability requires, at each step n , strong belief in all the previous steps of the procedure. For the iterated elimination of strategies, strong belief in step $n - 1$ suffices.

strong belief in a path of play. This class of restrictions is very important both for theory and practise. Pratically, agreements among real players often specify only the outcome to achieve (i.e., a path) and fall through if a player deviates (i.e., they not specify off-the-path behavior). Theoretically, path restrictions can be used to shed light on some forward induction refinements in the equilibrium literature, stemming from the seminal contribution of Kohlberg and Mertens [12], where deviations from an equilibrium path are interpreted as an attempt of the deviator to improve upon her equilibrium payoff. Examples of this are the epistemic justifications of the Iterated Intuitive Criterion (Cho and Kreps [10]) by Battigalli and Siniscalchi [6] and of "equilibrium paths that cannot be upset by a convincing deviation" (Osborne [13]) by Catonini [8]. This approach is generalized in Catonini [8] to capture in a transparent way forward induction reasoning under this interpretation of deviations from a path prescribed by a non-binding agreement among players.

With the same methodology, I prove that under path restrictions the *epistemic priority* between beliefs in the path and beliefs in rationality does not matter for the predicted outcomes. When a player displays behavior that cannot be rational under her belief restrictions, Strong- Δ -Rationalizability assumes that the opponents drop the belief that this player is rational.³ Selective Rationalizability (Catonini [7]), a refinement of Strong Rationalizability with first-order belief restrictions, assumes instead that opponents keep believing that the player is rational (if per se compatible with the observed behavior) and drop the belief that the player has beliefs in the restricted set. The predictions under path restrictions have the further advantage of being robust to this epistemic priority choice.

The same workhorse lemma that yields the main results also yields and provides new insight on the following result, already proven by Perea [17]: the iterated deletion of never sequential best replies (of which Strong Rationalizability is the maximal elimination order) is order independent in terms of predicted outcomes. Chen and Micali [9] characterize Strong Rationalizability with the iterated elimination of *distinguishably dominated* strategies,⁴ and show its order independence in terms of outcomes. Here, like in the recents work of Perea ([18], [17]), I work directly with the iterated deletion of never sequential best replies, thus with strong belief and without dominance characterizations. This is not the only similarity between Perea's methodology and mine. Perea bases his results on the "monotonicity on reachable histories" of the strong belief operator, which shares with my methodology

³Whether the belief that the player has beliefs in the restricted set is kept or not is immaterial for the procedure, thus Strong- Δ -Rationalizability can be characterized epistemically with or without *transparency* of the restrictions: see Battigalli and Prestipino [4] for details.

⁴By showing the equivalence of the iterated elimination of *distinguishable* and *conditionally* dominated strategies, where the latter was already proved by Shimoji and Watson [20] to be equivalent to Extensive Form Rationalizability (Pearce [15]), which is in turn equivalent to Strong Rationalizability.

the following intuition: strong belief in a smaller set of strategies justifies fewer possible behaviors along the paths induced by this set with respect to strong belief in a larger set. Monotonicity on reachable histories further claims that this intuition goes through also when the second set is richer than the first only in terms of behavior along the paths induced by the first if the sets have been generated through an elimination order of the strong belief operator (which coincides with the "iterated elimination of never sequential best replies" defined here). The workhorse lemma of this paper allows instead to compare two elimination procedures with (nested) belief restrictions by showing directly the following: although the more permissive procedure can actually become more restrictive after deviations from the paths induced by the other procedure, this will not induce players to abandon these paths in absence of belief restrictions after each potential deviation.

As in Chen and Micali [9], Strong Rationalizability is also shown to refine backward induction, here captured by Backwards Extensive Form Rationalizability (Penta [16]), which can be seen as a particular, unfinished elimination order of never sequential best replies. The outcome equivalence of backward induction and forward induction in perfect information games without relevant ties (originally proved by Battigalli [1] and then also by Perea [18] and Heifetz and Perea [11] in a more transparent way) follows.

Section 2 introduces the formal framework for the analysis. Section 3 defines elimination procedures and introduces the workhorse lemma. Section 4 presents the results on outcome monotonicity with respect to first-order belief restrictions and outcome equivalence with respect to the epistemic priority order. Section 5 presents the results on order independence and backward induction. Section 6 provides a sketch of the proof of the workhorse lemma, along with an example. The Appendix provides the formal proof.

2 Preliminaries

Primitives of the game.⁵ Let I be the finite set of *players*. For any profile of sets $(X_i)_{i \in I}$ and any subset of players $\emptyset \neq J \subseteq I$, I write $X_J := \times_{j \in J} X_j$, $X := X_I$, $X_{-i} := X_{I \setminus \{i\}}$. Let $(\bar{A}_i)_{i \in I}$ be the finite sets of *actions* potentially available to each player. Let $\bar{H} \subseteq \cup_{t=1, \dots, T} \bar{A}^t \cup \{h^0\}$ be the set of histories, where $h^0 \in \bar{H}$ is the empty, initial history and T is the finite horizon. The set \bar{H} must have the following properties. First property: For any $h = (a^1, \dots, a^t) \in \bar{H}$ and $l < t$, it holds $h' = (a^1, \dots, a^l) \in \bar{H}$, and I write $h' \prec h$.⁶ Let $Z := \{z \in \bar{H} : \nexists h \in \bar{H}, z \prec h\}$ be the set of terminal histories (henceforth, *outcomes*

⁵The main notation is almost entirely taken from Osborne and Rubinstein [14].

⁶Then, \bar{H} endowed with the precedence relation \prec is a tree with root h^0 .

or *paths*)⁷, and $H := \overline{H} \setminus Z$ the set of non-terminal histories (henceforth, just *histories*). Second property: For every $h \in H$, there exists a non-empty set $A_i(h) \subseteq \overline{A}_i$ for each $i \in I$,⁸ such that $(h, a) \in \overline{H}$ if and only if $a \in A_i(h)$. For each $i \in I$, let $u_i : Z \rightarrow \mathbb{R}$ be the *payoff function*. The list $\Gamma = \langle I, \overline{H}, (u_i)_{i \in I} \rangle$ is a *finite game with complete information and observable actions*.

Derived objects. A *strategy* of player i is an element of $\times_{h \in H} A_i(h)$. Let S_i denote the set of all strategies of i . A *strategy profile* $s \in S$ naturally induces a unique outcome $z \in Z$. Let $\zeta : S \rightarrow Z$ be the function that associates each strategy profile with the induced outcome. For any $h \in \overline{H}$, the set of strategies of i compatible with h is:

$$S_i(h) := \{s_i \in S_i : \exists z \succeq h, \exists s_{-i} \in S_{-i}, \zeta(s_i, s_{-i}) = z\}.$$

For any subset of player $J \subseteq I$ and any $\overline{S}_J \subset S_J$, let $\overline{S}_J(h) := S_J(h) \cap \overline{S}_J$. Let $H(\overline{S}_J) := \{h \in H : \overline{S}_J(h) \neq \emptyset\}$ denote the set of histories compatible with \overline{S}_J . For any $h = (h', a) \in \overline{H}$, let $p(h)$ denote the immediate predecessor h' of h .

Since the game has observable actions, each history $h \in H$ is the root of a subgame $\Gamma(h)$. If $h \neq h^0$, all the objects in $\Gamma(h)$ will be denoted with h as superscript, except for each history $h' \succeq h$ and outcome $z \succ h$, which will be identified with the corresponding history or outcome of the whole game, and not redefined as shorter lists of action profiles. For any $h \in H$, $s_i^h \in S_i^h = \times_{h' \succeq h} A_i(h')$, and $\widehat{h} \in H^h = \{h' \in H : h \preceq h'\}$, $s_i^h | \widehat{h}$ will denote the strategy $s_i^{\widehat{h}} \in S_i^{\widehat{h}}$ such that $s_i^{\widehat{h}}(\widetilde{h}) = s_i^h(\widetilde{h})$ for all $\widetilde{h} \succeq \widehat{h}$. For any $\overline{S}_i^h \subseteq S_i^h$, $\overline{S}_i^h | \widehat{h}$ will denote the set of all strategies $s_i^{\widehat{h}} \in S_i^{\widehat{h}}$ such that $s_i^{\widehat{h}} = s_i^h | \widehat{h}$ for some $s_i^h \in \overline{S}_i^h$.

Beliefs. In this dynamic framework, beliefs are modeled as Conditional Probability Systems (Renyi, [19]; henceforth, CPS).

Definition 1 Fix $i \in I$. An array of probability measures $(\mu_i(\cdot|h))_{h \in H}$ over co-players' strategies S_{-i} is a Conditional Probability System if for all $h \in H$, $\mu_i(S_{-i}(h)|h) = 1$, and for all $h' \succ h$ and $\overline{S}_{-i} \subseteq S_{-i}(h')$,

$$\mu_i(\overline{S}_{-i}|h) = \mu_i(S_{-i}(h')|h) \cdot \mu_i(\overline{S}_{-i}|h').$$

The set of all CPS's on S_{-i} 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. For each subset of

⁷"Path" will be used with emphasis on the moves, and "outcome" with emphasis on the end-point of the game.

⁸When player i is not truly active at history h , $A_i(h)$ consists of just one "wait" action.

opponents' strategies $\bar{S}_{-i} \subseteq S_{-i}$, I say that a CPS $\mu_i \in \Delta^H(S_{-i})$ *strongly believes* \bar{S}_{-i} if, for all $h \in H(\bar{S}_{-i})$, $\mu_i(\bar{S}_{-i}|h) = 1$. I fix the following convention: $H(\emptyset) = \emptyset$. With this, the empty set is always strongly believed, because the condition is vacuously satisfied.

Rationality. 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 payoff given the belief about how opponents will play and the expectation to reply rationally again in the continuation of the game. This is equivalent (see Battigalli [2]) to playing a *sequential best reply* to the CPS.

Definition 2 Fix $\mu_i \in \Delta^H(S_{-i})$. A strategy $s_i \in S_i$ is a *sequential best reply* to μ_i if for every $h \in H(s_i)$,⁹ s_i is a *continuation best reply* to $\mu_i(\cdot|h)$, i.e. for every $\tilde{s}_i \in S_i(h)$,

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

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 μ_i is denoted by $\rho(\mu_i)$. For each $h \in H$, the set of continuation best replies to $\mu_i(\cdot|h)$ is denoted by $\hat{r}(\mu_i, h)$. The set of best replies to a conjecture $\nu_i \in \Delta(S_{-i})$ in the normal form of the game is denoted by $r(\nu_i)$.

3 Elimination procedures and the main lemma

I provide a very general notion of elimination procedure for a subgame $\Gamma(h)$, which encompasses all the procedure I am ultimately interested in, or that will be needed for the proofs.

Definition 3 Fix $h \in H$. An *elimination procedure* in $\Gamma(h)$ is a sequence $((S_{i,q}^h)_{i \in I})_{q=0}^\infty$ where, for every $i \in I$,

EP1 $S_{i,0}^h = S_i^h$;

EP2 $S_{i,n-1}^h \supseteq S_{i,n}^h$ for all $n \in \mathbb{N}$;

EP3 for every $s_i^h \in S_{i,\infty}^h = \bigcap_{n \in \mathbb{N}} S_{i,n}^h$, there exists μ_i^h that strongly believes $(S_{-i,q}^h)_{q=0}^\infty$ such that $s_i^h \in \rho(\mu_i^h) \subseteq S_{i,\infty}^h$.

⁹It would be totally immaterial to require s_i to be optimal also at the histories precluded by itself.

Definition 3 allows $S_{i,n}^h = \emptyset$ for some $n \in \mathbb{N}$, which implies $S_{i,m}^h = \emptyset$ for all $m > n$, but does not imply $S_{j,n+1}^h = \emptyset$ for any $j \neq i$: as already established, the empty set is always strongly believed, hence EP3 can be satisfied by a non-empty $S_{j,n+1}^h$. Moreover, EP2 allows an equality for all players also at steps before "convergence". All this allows Definition 3 to encompass the implications in a subgame of an elimination procedure for a larger subgame (which I will call "truncation" of the elimination procedure in the subgame).

Lemma 1 *For every elimination procedure $((S_{i,q}^h)_{i \in I})_{q=0}^\infty$ and every $\hat{h} \succ h$, $((S_{i,q}^h(\hat{h})|\hat{h})_{i \in I})_{q=0}^\infty$ is an elimination procedure.*

Proof. EP1 and EP2 are obvious. To prove EP3, note the following. For every $i \in I$ and $s_i^{\hat{h}} \in S_{i,\infty}^h(\hat{h})|\hat{h}$, there exists $s_i^h \in S_{i,\infty}^h$ such that $s_i^h|\hat{h} = s_i^{\hat{h}}$. By EP3 for $((S_{i,q}^h)_{i \in I})_{q=0}^\infty$, there exists μ_i^h that strongly believes $(S_{-i,q}^h)_{q=0}^\infty$ such that $s_i^h \in \rho(\mu_i^h) \subseteq S_{i,\infty}^h$. Thus, the pushforward $\mu_i^{\hat{h}}$ of $(\mu_i^h(\cdot|\hat{h}))_{\tilde{h} \in H\hat{h}}$ through the map $s_{-i}^h \mapsto s_{-i}^{\hat{h}}|\hat{h}$ strongly believes $(S_{-i,q}^h(\hat{h})|\hat{h})_{q=0}^\infty$. Clearly $s_i^{\hat{h}} \in \rho(\mu_i^{\hat{h}})$. Finally, fix $\bar{s}_i^{\hat{h}} \in \rho(\mu_i^{\hat{h}})$. Define \bar{s}_i^h as $\bar{s}_i^h(\hat{h}) = s_i^{\hat{h}}(\hat{h})$ for all $\tilde{h} \not\preceq \hat{h}$ and $\bar{s}_i^h|\hat{h} = \bar{s}_i^{\hat{h}}$ for all $\tilde{h} \succeq \hat{h}$. Clearly $\bar{s}_i^h \in \rho(\mu_i^h)$. Thus, $\bar{s}_i^h \in S_{i,\infty}^h(\hat{h})|\hat{h}$. ■

For some $j \neq i$, we can have $S_{j,n+1}^h(\hat{h})|\hat{h} \neq \emptyset$ although $S_{i,n}^h(\hat{h})|\hat{h} = \emptyset$. Moreover, we can have $S_{j,n}^h(\hat{h})|\hat{h} = S_{j,n+1}^h(\hat{h})|\hat{h}$ for all $j \in I$, but $S_{j,n+1}^h(\hat{h})|\hat{h} \supset S_{j,n+2}^h(\hat{h})|\hat{h}$. This is because strategies in $((S_{i,q}^h(\hat{h})|\hat{h})_{i \in I})_{q=0}^\infty$ can be eliminated "exogeneously", due to eliminations from $((S_{i,q}^h)_{i \in I})_{q=0}^\infty$ that affect $((S_{i,q}^h(\hat{h})|\hat{h})_{i \in I})_{q=0}^\infty$ at step $n+2$ and not at step $n+1$, and not because they are not sequential best replies to any valid conjecture in the subgame. For this reason, and to encompass elimination procedures with first-order-belief restrictions, an elimination procedure does not impose to save all the strategies that are sequential best replies to some μ_i that strongly believes $(S_{-i,q}^h)_{q=0}^{n-1}$. This makes an elimination procedure more general than an order of elimination of the "strong belief operator", defined in Perea [?]. Like for an order of elimination of the strong belief operator, instead, an elimination procedure allows to "forget" to eliminate strategies that are not sequential best replies to any μ_i that strongly believes $(S_{-i,q}^h)_{q=0}^{n-1}$ at all steps n before convergence: for this reason, EP3 only refers to the final output of the procedure.

The workhorse lemma of the paper claims the outcome inclusion between two elimination procedures, $((\bar{S}_{i,q}^h)_{i \in I})_{q=0}^\infty$ and $((S_{i,q}^h)_{i \in I})_{q=0}^\infty$, with the following feature. Take the final output \bar{S}_∞^h of the first procedure and, for each player i and each strategy $s_i^h \in \bar{S}_{i,\infty}^h$, fix a CPS $\bar{\mu}_i^h(s_i^h)$ that satisfies EP3, i.e., it strongly believes $(\bar{S}_{-i,q}^h)_{q=0}^\infty$ and justifies s_i^h : $s_i^h \in \rho(\bar{\mu}_i^h(s_i^h)) \subseteq \bar{S}_{i,\infty}^h$. Consider now a CPS μ_i^h that, along the paths predicted by the first procedure, that is at every history $\tilde{h} \in H(\bar{S}_\infty^h)$, assigns the same probability as $\bar{\mu}_i^h(s_i^h)$ to the fact that the opponents will play compatibly with each of these paths $z \in \zeta(\bar{S}_\infty^h)$: $\mu_i^h(S_{-i}(z)|\tilde{h}) = \bar{\mu}_i^h(s_i^h)(S_{-i}(z)|\tilde{h})$. Suppose now that, for each $m \in \mathbb{N}$, if μ_i^h strongly believes

$(\bar{S}_{-i,q}^h)_{q=0}^{m-1}$, then $\rho(\mu_i^h) \subseteq \bar{S}_{i,m}^h$, and if μ_i^h strongly believes $(S_{-i,q}^h)_{q=0}^{m-1}$, then $\rho(\mu_i^h) \subseteq S_{i,m}^h$. The lemma claims that then the second procedure predicts a superset $\zeta(S_\infty^h) \supseteq \zeta(\bar{S}_\infty^h)$ of the paths predicted by the first.

Lemma 2 Fix $h \in H$, two elimination procedures $((\bar{S}_{i,q}^h)_{i \in I})_{q=0}^\infty$, $((S_{i,q}^h)_{i \in I})_{q=0}^\infty$, and, for every $i \in I$, a map $\bar{\mu}_i^h : \bar{S}_{i,\infty}^h \rightarrow \Delta_i^{H^h}(S_{-i}^h)$ such that, for each $s_i^h \in \bar{S}_{i,\infty}^h$, $\bar{\mu}_i^h(s_i^h)$ strongly believes $(\bar{S}_{-i,q}^h)_{q=0}^\infty$ and $s_i^h \in \rho(\bar{\mu}_i^h(s_i^h)) \subseteq \bar{S}_{i,\infty}^h$. Suppose that for every $i \in I$, $s_i^h \in \bar{S}_{i,\infty}^h$, $m \in \mathbb{N}$, and μ_i^h that strongly believes $(S_{-i,q}^h)_{q=0}^{m-1}$ (resp., $(\bar{S}_{-i,q}^h)_{q=0}^{m-1}$),

$$\left(\forall \tilde{h} \in H(\bar{S}_\infty^h), \forall Z^{\tilde{h}} \cap \zeta(\bar{S}_\infty^h), \mu_i^h(S_{-i}(z)|\tilde{h}) = \bar{\mu}_i^h(s_i^h)(S_{-i}(z)|\tilde{h}) \right) \Rightarrow \left(\rho(\mu_i^h) \subseteq S_{i,m}^h \right) \\ \text{(resp., } \rho(\mu_i^h) \subseteq \bar{S}_{i,m}^h \text{)}.$$

Then, $\zeta(\bar{S}_\infty^h) \subseteq \zeta(S_\infty^h)$.

Section 6 contains a sketch of the proof of the lemma, while the Appendix contains the formal proof. Now I focus on the implications of the lemma for the elimination procedures of interest.

4 Belief-restrictions and monotonicity

In this section, I am going to consider the following elimination procedures (for the whole game).

Definition 4 An elimination procedure $((S_{i,q})_{i \in I})_{q=0}^\infty$ is "unconstrained" when for every $n \in \mathbb{N}$, $i \in I$, and μ_i that strongly believes $(S_{-i,q})_{q=0}^{n-1}$, $\rho(\mu_i) \subseteq S_{i,n}$.

Definition 5 An elimination procedure $((S_{i,q})_{i \in I})_{q=0}^\infty$ is "maximal" when for every $n \in \mathbb{N}$, $i \in I$, and $s_i \in S_{i,n}$, $s_i \in \rho(\mu_i)$ for some μ_i that strongly believes $(S_{-i,q})_{q=0}^{n-1}$.

Definition 6 Strong Rationalizability (Battigalli and Siniscalchi, [5]) is the unique unconstrained and maximal elimination procedure. Let $((S_i^q)_{i \in I})_{q=0}^\infty$ denote it, and let M be the $n \in \mathbb{N}$ such that $S^{n-1} \neq S^n = S^{n+1}$.

Definition 7 For each $i \in I$, fix $\Delta_i \subseteq \Delta^H(S_{-i}^h)$. Strong- Δ -Rationalizability (Battigalli [3], Battigalli and Siniscalchi [6]) is the elimination procedure $((S_{i,\Delta}^q)_{i \in I})_{q=0}^\infty$ such that, for every $n \in \mathbb{N}$, $i \in I$, and $s_i \in S_i$, $s_i \in S_{i,n}$ if and only if $s_i \in \rho(\mu_i)$ for some $\mu_i \in \Delta_i$ that strongly believes $(S_{-i,\Delta}^q)_{q=0}^{n-1}$.

Definition 8 For each $i \in I$, fix $\Delta_i \subseteq \Delta^H(S_{-i}^h)$. *Selective Rationalizability* (Catonini [7]) is the elimination procedure $((S_{i,R\Delta}^q)_{i \in I})_{q=0}^\infty$ such that $(S_{R\Delta}^q)_{q=0}^M = (S^q)_{q=0}^M$ and for every $n > M$, $i \in I$, and $s_i \in S_i$, $s_i \in S_{i,R\Delta}^n$ if and only if $s_i \in \rho(\mu_i)$ for some $\mu_i \in \Delta_i$ that strongly believes $(S_{-i,R\Delta}^q)_{q=0}^{n-1}$.¹⁰

Consider first-order belief restrictions $(\Delta_i)_{i \in I}$ with the following characteristic: for each player i and CPS μ_i , all that matters to determine whether μ_i belongs to Δ_i are the probabilities assigned by μ_i at each strongly- Δ -rationalizable history $h \in H(S_\Delta^\infty)$ to the fact that opponents will play compatibly with each strongly- Δ -rationalizable path $z \in \zeta(S_\Delta^\infty)$: $\mu_i(S_{-i}(z)|h)$. Then, Strong- Δ -Rationalizability satisfies the hypotheses of Lemma 2 as first elimination procedure, whereas Strong Rationalizability, being an unconstrained procedure, obviously satisfies the hypotheses of Lemma 2 as second elimination procedure. The desired outcome inclusion with respect to belief restrictions that "do not end up off-path" obtains.

Theorem 1 For each $i \in I$, fix a set of CPS's $\Delta_i \subseteq \Delta^H(S_{-i})$. Suppose that for each $i \in I$, $\mu_i \in \Delta_i$, and $\mu'_i \in \Delta^H(S_{-i})$,

$$\left(\forall \tilde{h} \in H(S_\Delta^\infty), \forall z \in \zeta(S_\Delta^\infty), \mu'_i(S_{-i}(z)|\tilde{h}) = \mu_i(S_{-i}(z)|\tilde{h}) \right) \Rightarrow (\mu'_i \in \Delta_i).$$

Then, $\zeta(S_\Delta^\infty) \subseteq \zeta(S^\infty)$.

Proof. For each $i \in I$ and $s_i \in S_{i,\Delta}^\infty$, fix any $\bar{\mu}_i^h(s_i^h) \in \Delta_i$ that strongly believes $(S_{-i,\Delta}^q)_{q=0}^\infty$ such that $s_i \in \rho(\mu_i)$. By hypothesis of this theorem, the hypothesis of Lemma 2 obtains. For every $m \in \mathbb{N}$ and μ_i that strongly believes $(S_{-i}^q)_{q=0}^{m-1}$, $\rho(\mu_i) \in S_i^m$. Thus, by Lemma 2, $\zeta(S_\Delta^\infty) \subseteq \zeta(S^\infty)$. ■

As discussed in the Introduction, Theorem 1 provides insight on what can determine the non monotonicity of predicted outcomes with respect to belief restrictions: the presence of off-the-path restrictions. Yet, it is of little help in determining ex-ante which belief restrictions preserve the behavioral consequences of common strong belief in rationality and which do not. This is because whether restrictions are off-path or not has to be assessed with respect of the final output of Strong- Δ -Rationalizability itself.

Consider now first-order belief restrictions that correspond to the belief in a specific path $z \in Z$ along the path itself. This is what I call "path restrictions". I note preliminarily that this is equivalent to strong belief in $S_j(z)$ for all $j \neq i$ (the proof is in the Appendix).¹¹

¹⁰Selective Rationalizability is defined in [7] under the more restrictive assumption of *independent rationalization*. That is, a valid μ_i is required to strongly believe $(S_{j,R\Delta}^q)_{q=0}^{n-1}$ for all $j \neq i$, in place of just $(S_{-i,R\Delta}^q)_{q=0}^{n-1}$. However, this assumption is immaterial for the result on Selective Rationalizability of this paper (Theorem 3).

¹¹This corresponds to *belief in the (path) agreement* in [7].

The reason is that after a deviation from the path by a player different than j , believing that j would have kept complying with the path is not restrictive for the expected behavior of j after the deviation.

Lemma 3 *Fix $z \in Z$. For each $i \in I$, let Δ_i be the set of all $\mu_i \in \Delta^H(S_{-i})$ such that $\mu_i(S_{-i}(z)|h) = 1$ for all $h \prec z$, and let Δ_i^* be the set of all $\mu_i \in \Delta^H(S_{-i})$ that strongly believe $S_j(z)$ for all $j \neq i$. Then, $S_\Delta^\infty = S_{\Delta^*}^\infty$ and $S_{R\Delta}^\infty = S_{R\Delta^*}^\infty$.*

Path restrictions hold at histories that precede z . Therefore, if such restrictions ends up off-path, it means that some player has abandoned the path, so the opponents cannot believe in $S_{-i}(z)$ from the start anymore, and Strong- Δ -Rationalizability yields the empty set. Otherwise, Theorem 1 can be applied and, via Lemma 3, monotonicity of strategic reasoning under strong belief in a path obtains.

Theorem 2 *Fix $z \in Z$. Let Δ_i^* be the set of all $\mu_i \in \Delta^H(S_{-i})$ that strongly believe $S_j(z)$ for all $j \neq i$. Then $\zeta(S_{\Delta^*}^\infty) \subseteq \zeta(S^\infty)$.*

Proof. For each $i \in I$, let Δ_i be the set of all μ_i 's such that $\mu_i(S_{-i}(z)|h) = 1$ for all $h \prec z$. If $S_\Delta^\infty = \emptyset$, $\zeta(S_\Delta^\infty) \subseteq \zeta(S^\infty)$ is trivially true, so suppose $S_\Delta^\infty \neq \emptyset$. For each $i \in I$, and $s_i \in S_{i,\Delta}^\infty$, $s_i \in \rho(\bar{\mu}_i)$ for some $\bar{\mu}_i \in \Delta_i$. For each $\bar{\mu}_i \in \Delta_i$ and μ_i with $\mu_i(S_{-i}(z)|h) = \bar{\mu}_i(S_{-i}(z)|h)$ for all $h \prec z$, $\mu_i \in \Delta_i$. Thus, the hypotheses of Theorem 1 hold, and $\zeta(S_\Delta^\infty) \subseteq \zeta(S^\infty)$. Then, by Lemma 3, $\zeta(S_{\Delta^*}^\infty) \subseteq \zeta(S^\infty)$. ■

Also Selective Rationalizability eventually saves only strategies that are sequential best replies under strong belief in the path. Therefore, for path restrictions, Lemma 2 holds with Selective Rationalizability and Strong- Δ -Rationalizability regardless of the roles assigned to the two procedures. Then, via Lemma 3, the outcome equivalence of the two procedures under strong belief in a path obtains.

Theorem 3 *Fix $z \in Z$. Let Δ_i^* be the set of all $\mu_i \in \Delta^H(S_{-i})$ that strongly believe $S_j(z)$ for all $j \neq i$. Then $\zeta(S_{\Delta^*}^\infty) = \zeta(S_{R\Delta^*}^\infty)$.*

Proof. For each $i \in I$, let Δ_i be the set of all μ_i 's such that $\mu_i(S_{-i}(z)|h) = 1$ for all $h \prec z$. First I show that $\zeta(S_\Delta^\infty) \subseteq \zeta(S_{R\Delta}^\infty)$. If $S_\Delta^\infty = \emptyset$ it is trivially true, so suppose $S_\Delta^\infty \neq \emptyset$. For each $i \in I$, and $s_i \in S_{i,\Delta}^\infty$, $s_i \in \rho(\bar{\mu}_i)$ for some $\bar{\mu}_i \in \Delta_i$. For each $\bar{\mu}_i \in \Delta_i$ and μ_i with $\mu_i(S_{-i}(z)|h) = \bar{\mu}_i(S_{-i}(z)|h)$ for all $h \prec z$, $\mu_i \in \Delta_i$. Thus, the hypotheses of Lemma 2 hold. So, $\zeta(S_\Delta^\infty) \subseteq \zeta(S_{R\Delta}^\infty)$. The same proof can be repeated for $\zeta(S_\Delta^\infty) \supseteq \zeta(S_{R\Delta}^\infty)$. Hence $\zeta(S_\Delta^\infty) = \zeta(S_{R\Delta}^\infty)$. Then, by Lemma 3, $\zeta(S_{\Delta^*}^\infty) = \zeta(S_{R\Delta^*}^\infty)$. ■

The last two theorems clearly hold with strong belief in $S_{-i}(z)$ instead of $(S_j(z))_{j \neq i}$.

5 Order independence and backward induction

In absence of belief restrictions, that is, for unconstrained elimination procedures, the hypotheses of Theorem 2 clearly hold. An unconstrained elimination procedure is what we referred to in the Introduction as an order of iterated elimination of never sequential best replies, and in Perea [18] it is called an elimination order of the strong belief operator. Thus, using Theorem 2 in both directions with the maximal unconstrained elimination procedure and any non maximal one, the order independence of iterated elimination of never sequential best replies in terms of predicted outcomes obtains.

Theorem 4 *For any unconstrained elimination procedure $((S_{i,q})_{i \in I})_{q=0}^{\infty}$, $\zeta(S_{\infty}) = \zeta(S^{\infty})$.*

Proof. Any two unconstrained elimination procedures, taken in both orders, obviously satisfy the hypotheses of Lemma 2. ■

In games with observable actions, the well-known backward induction procedure for games with perfect information has been generalized as follows (see, for instance, Chen and Micali [9]). Starting from the bottom of game, an action of a player at a history is eliminated when it is not “folding-back optimal” against any conjecture over the surviving actions of the opponents at the same history and at the future histories, that is, it is not optimal given the already computed optimal actions at the future histories. Penta [16] has translated backward induction for games with observable actions in the language of extensive-form rationalizability, i.e., as a procedure of elimination of strategies that are not sequentially optimal for any appropriate conditional probability system. Penta’s Backwards Extensive-Form Rationalizability is simplified here for games with complete information.

Definition 9 *Backwards Extensive-Form Rationalizability is a sequence $((S_{i,BR}^q)_{i \in I})_{q=0}^{\infty}$ where, for every $i \in I$,*

BR1 $S_{i,BR}^0 = S_i$;

BR2 for each $n \in \mathbb{N}$ and $s_i \in S_i$, $s_i \in S_{i,BR}^n$ if and only if $s_i \in S_{i,B}^{n-1}$ and there exists $\mu_i \in \Delta^H(S_{-i})$ such that, for each $h \in H$,

- (i) there is $\tilde{s}_i \in S_i(h)$ such that $\tilde{s}_i(\tilde{h}) = s_i(\tilde{h})$ for each $\tilde{h} \succeq h$ and $\tilde{s}_i \in \hat{r}(\mu_i, \tilde{h})$;
- (ii) for each $\tilde{s}_{-i} \in S_{-i}(h)$ with $\mu_i(\tilde{s}_{-i}|h) > 0$, there is $s_{-i} \in S_{-i,BR}^{n-1}$ such that $\tilde{s}_{-i}(\tilde{h}) = s_{-i}(\tilde{h})$ for each $\tilde{h} \succeq h$.

Condition BR2.(i) requires s_i to be a continuation best reply not only at each $h \in H(s_i)$, as for sequential best replies of Definition 2, but also at each $h \notin H(s_i)$. Then, condition

BR2.(ii) requires to keep refining beliefs also at histories that cannot be reached anymore. So, Backwards Extensive-Form Rationalizability can be stricter, in terms of strategies, than an unconstrained elimination procedure. Yet, given that realization-equivalent classes are all that matters for elimination procedures and that such refinement of beliefs is off-path, it turns out that Backwards Extensive-Form Rationalizability is outcome-equivalent to an *unfinished*, unconstrained elimination procedure.

Lemma 4 *Let N be the smallest n such that $S_{BR}^n = S_{BR}^{n+1}$. There exists an unconstrained elimination procedure $((S_{i,q})_{i \in I})_{q=0}^\infty$ such that for each $n \leq N$,*

$$S_n = \{s \in S : \exists s' \in S_{BR}^n, \forall h \in H(S_{BR}^n), s(h) = s'(h)\}.$$

Proof. Define $((S_{i,n})_{i \in I})_{n=0}^N$ as above, and for each $n > N$ and $i \in I$, let $s_i \in S_{i,n}$ if and only if there exists μ_i that strongly believes $(S_{-i,q})_{q=0}^{n-1}$ such that $s_i \in \rho(\mu_i)$. It is immediate to check that $((S_{i,q})_{i \in I})_{q=0}^\infty$ is an elimination procedure. To show that it is unconstrained, fix $n \leq N$ and suppose by way of induction that for each $m < n$, $i \in I$, and μ_i that strongly believes $(S_{-i,q})_{q=0}^{m-1}$, we have $\rho(\mu_i) \subseteq S_{i,m}$ (it is vacuously true for $m = 0$). Fix μ_i that strongly believes $(S_{-i,q})_{q=0}^{n-1}$. I show that $\rho(\mu_i) \subseteq S_{i,n}$. By definition of $S_{-i,n-1}$, I can construct μ'_i that satisfies BR2.(ii) such that for all $h \in H(S_{n-1})$ and $z \in \zeta(S_{n-1})$, $\mu_i(S_{-i}(z)|h) = \mu'_i(S_{-i}(z)|h)$. For each $s'_i \in \rho(\mu'_i)$, there is a realization equivalent s''_i that satisfies BR2.(i), so that $s''_i \in S_{i,BR}^n \subseteq S_{i,BR}^{n-1}$. For each $s_i \in \rho(\mu_i)$, by the induction hypothesis we have $s_i \in S_{i,n-1}$. Then, we have $\zeta(\rho(\mu'_i) \times S_{-i,n-1}), \zeta(\rho(\mu_i) \times S_{-i,n-1}) \subseteq \zeta(S_{n-1}) = \zeta(S_{BR}^{n-1})$. Thus, for each $s_i \in \rho(\mu_i)$, there is $s'_i \in \rho(\mu'_i)$ such that $s_i(h) = s'_i(h)$ for all $h \in H(S_{n-1})$. Since there is $s''_i \in S_{i,BR}^n$ realization equivalent to s'_i , so that $s''_i(h) = s'_i(h) = s_i(h)$ for all $h \in H(S_{n-1}) \cap H(s_i)$, by definition of $S_{i,n}$ we have $s_i \in S_{i,n}$. Thus, $\rho(\mu_i) \subseteq S_{i,n}$. ■

Being outcome-equivalent to an unfinished, unconstrained elimination procedure, Backwards Extensive Rationalizability predicts a superset of the outcomes predicted by Strong Rationalizability.

Theorem 5 *Every strongly rationalizable outcome is also a backwards extensive-form rationalizable outcome.*

Proof. Immediate from Lemma 4 and Theorem 4. ■

Since in perfect information games without relevant ties the backward induction outcome is unique, the following obtains.

Corollary 6 (Battigalli, [1]) *In every perfect information game without relevant ties, Strong Rationalizability and backward induction yield the same unique outcome.*

6 Proof of the main lemma.

The rough intuition for the proof of the main lemma is the following. Take the paths induced by the first procedures. If the game had only these paths, they would survive also the second procedure, for the following two reasons. First, the fact that they survive the first procedure indicates that for every player there are beliefs over these paths that justify allowing each of them. Second, all these beliefs are allowed also under the second procedure by assumption. Then, the only way one of these paths can be eliminated along the second procedure is that at some step n , some player finds a deviation *outside* of these paths more profitable, however she believes the opponents will react to the deviation. Since the opponents may be surprised by the deviation (all the paths survived until step $n - 1$), they can react with any continuation plan that survives until step n . So, we have that both the deviator and the opponents can play any sequential best reply to any belief in the subgame that follows the deviation. This allows to generate an auxiliary elimination procedure for the subgame that refines the continuation plans that survive the second procedure for the whole game until step n , and terminates with a non-empty set. Take the subpaths induced by this auxiliary procedure. We want to show that they would have survived also the truncation of the first procedure for the whole game, which implies that the subgame is reached at the end of the procedure and contradicts that the subgame follows a deviation from the paths predicted by the procedure. If the subgame is a static game, i.e., it has depth 1, this is easy to see: the auxiliary procedure generates a best response set where all the beliefs the deviator can have induce a deviation from the original path, which is actually available until the end of the first procedure (it can be sustained by other surviving actions of the opponents in the subgame). If the subgame has depth higher than 1, then suppose by induction that the lemma is true in games of that depth, so that the truncation of the first procedure in the subgame induces a superset of the paths induced by the auxiliary procedure, which leads to the same contradiction.

I refine now this intuition briefly illustrating with mathematical notation the precise arguments of the formal proof. For simplicity, assume that there are two players, i and j ; the argument extends immediately to games with more than 2 players. I argue by induction that for every $\bar{s}_i^h \in \bar{S}_{i,\infty}^h$ and $n \in \mathbb{N}$, there are: (1) $\bar{\mu}_i^h(\bar{s}_i^h)$ that strongly believes $(S_{j,q}^h)_{q=0}^{n-1}$ and assigns the same probabilities to the (opponents playing compatibly with) the paths induced by \bar{S}_∞^h (henceforth, just "paths") as some $\bar{\mu}_i^h(\bar{s}_i^h)$ that justifies $\bar{s}_i^h \in \bar{S}_{i,\infty}^h$; (2) $s_i^h \in \rho(\bar{\mu}_i^h(\bar{s}_i^h))$ that mimicks \bar{s}_i^h along the paths. Then, by the assumption on $(S_{i,q}^h, S_{j,q}^h)_{q=0}^\infty$, we have $s_i^h \in S_{i,n}^h$. All such s_i^h 's allow to construct at step $n + 1$ a CPS $\bar{\mu}_j^h(\bar{s}_j^h)$ as in (1) for each $\bar{s}_j^h \in \bar{S}_{j,\infty}^h$.

Now, suppose by contradiction that for some $\bar{s}_j^h \in \bar{S}_{j,\infty}^h$, every such $\bar{\mu}_j^h(\bar{s}_j^h)$ does not

justify any strategy s_j^h that mimicks \bar{s}_j^h along the paths. For each history \hat{h} that immediately follows a unilateral deviation of player j from the paths, that is, that follows a history along the paths where i takes an action compatible with some of the paths and j does not, consider the most pessimistic belief of j over $S_{i,n}^h(\hat{h})|\hat{h}$. For each $\bar{s}_i^h \in \bar{S}_{i,\infty}^h$, by induction hypothesis there is $s_i^h \in S_{i,n}^h$ that mimicks \bar{s}_i^h along the paths and is a sequential best reply to a belief $\bar{\mu}_i^h(\bar{s}_i^h)$ as in (1), thus which assigns probability zero to each deviation of j until it occurs. Then, the beliefs specified by $\bar{\mu}_i^h(\bar{s}_i^h)$ along the paths can be combined with any beliefs after j 's deviations in a new CPS μ_i^h that satisfies (1). Clearly, there is a sequential best reply to μ_i^h which mimicks s_i^h along the paths and prescribes any sequentially rational reaction to the chosen beliefs after j 's deviations. This is proved by Lemma 6. So, at step $n+1$, player j can have a belief μ_j^h that mimicks $\bar{\mu}_j^h(\bar{s}_j^h)$ along the paths (in the sense of (1)) and, at the same time, features the most pessimistic belief after each deviation. By the initial assumption of this paragraph, player j will still deviate under μ_j^h and, calling \hat{h} a history that immediately follows this deviation, also under $\tilde{\mu}_j^h$ constructed like μ_j^h except for a less pessimistic belief over $S_{i,n}^h(\hat{h})|\hat{h}$. This is proved by Lemma 7. So, $S_{j,n+1}^h(\hat{h})|\hat{h}$ features all the sequential best replies to CPS's μ_j^h that strongly believe $(S_{i,q}^h)_{q=0}^n$, thus $S_{j,n}^h(\hat{h})|\hat{h} \supseteq S_{j,n+1}^h(\hat{h})|\hat{h}$ as well. The same holds with i in place of j , for the argument exposed above.

Refine $S_n^h(\hat{h})|\hat{h}$ by iteratively eliminating strategies that are not sequential best replies to any μ_k^h , $k = i, j$, that strongly believes in the previous steps. Then, we obtain an elimination procedure $((\bar{S}_{k,q}^h)_{k \in I})_{q=0}^\infty$ with non-empty \bar{S}_∞^h that satisfies the assumption of the main lemma. That the deviation is profitable against all beliefs over $\bar{S}_{i,\infty}^h$ with respect to remaining on the paths under $\bar{\mu}_j^h(s_j^h)$ implies that also the elimination procedure $((\hat{S}_{k,q}^h)_{k \in I})_{q=0}^\infty := ((\bar{S}_{k,q}^h(\hat{h})|\hat{h})_{k \in I})_{q=0}^\infty$ satisfies the assumption of the main lemma. Note the inversion of the roles of the two procedures with respect to the original procedures from which they have been derived. If the lemma holds in the subgame $\Gamma(\hat{h})$, we have the desired contradiction: $S_\infty^h = \bar{S}_{k,\infty}^h(\hat{h})|\hat{h}$ is non-empty too, hence $\hat{h} \in H(\bar{S}_\infty^h)$, but \hat{h} follows a deviation from the paths induced by \bar{S}_∞^h . Proceeding by induction on the depth of subgames and observing that the lemma clearly holds for subgames of depth 1, the proof is complete.

Finally, I am going to follow the sketch above on an example. Consider the following game.

$A \backslash B$	W	E	\rightarrow	$A \backslash B$	L	C	R
N	2, 2	·-		U	1, 1	1, 0	0, 0
S	0, 0	2, 2		M	0, 0	0, 1	1, 0
				D	0, 0	0, 0	0, 3

Take Strong Rationalizability, $((S_i^q)_{i \in I})_{q=0}^\infty$, as second procedure in the statement of the main lemma. At the first step, Ann eliminates $N.D$. At the second step, Bob eliminates

E.R. At the third step, Ann eliminates *N.M.* At the fourth step, Bob eliminates *E.C.* The final output is $S_{i,\infty}^{h^0} = (S, N.U) \times (W, E.L)$. Strong Rationalizability trivially satisfies the assumption of the lemma.

For each player $i = A, B$, let Δ_i be the set of CPS's that strongly believe in opponents' strategies that comply with the path $z := (N, W)$:

$$\Delta_i := \{ \mu_i \in \Delta_i^H(S_{-i}) : \mu_i(S_{-i}(z)|h^0) = 1 \}, \quad i = A, B.$$

Take Strong- Δ -Rationalizability, $((S_{i,\Delta}^q)_{i \in I})_{q=0}^\infty$, as first procedure in the statement of the main lemma. At the first step, Ann eliminates *S* and *N.D.*, and Bob eliminates *E.L* and *E.C.* At the second step, Ann eliminates *N.U* and Bob eliminates *E.R.* The final output is: $S_\Delta^\infty = \{(N.M, W)\}$.

Let $\bar{s}_A = N.M$, $\bar{s}_B = W$, $\bar{\mu}_A^h(\bar{s}_A) = (\delta_W, \delta_{E.R})$, and $\bar{\mu}_B^h(\bar{s}_B) = (\delta_{N.M}, \delta_{N.M})$, where δ_s indicates a Dirac measure on s . For every $n \in \mathbb{N}$, $i = A, B$, and μ_i that strongly believes $S_{-i,\Delta}^{n-1}, \dots, S_{-i,\Delta}^0$ with $\mu_i(S_{-i}(z)|h^0) = \bar{\mu}_i^h(S_{-i}(z)|h^0) = 1$, we have $\rho(\mu_i) \subseteq S_{i,\Delta}^n$. So, $((S_{i,\Delta}^q)_{i \in I})_{q=0}^\infty$ satisfies the assumption of the lemma. Indeed, $\zeta(S_\Delta^\infty) = \{z\} \subseteq \zeta(S^\infty)$ (although $S_{A,\Delta}^\infty \cap S_A^\infty = \emptyset$).

Now we follow the sketch above. Fix $n \in \mathbb{N}$ and suppose to have shown that for each $i \in A, B$, there exist:

1. $\bar{\mu}_i(\bar{s}_i)$ that strongly believes $S_{-i}^{n-1}, \dots, S_{-i}^0$ with $\bar{\mu}_i(\bar{s}_i)(S_{-i}(z)|h^0) = \bar{\mu}_i^h(S_{-i}(z)|h^0) = 1$;
2. $s_i \in \rho(\bar{\mu}_i(\bar{s}_i)) \subseteq S_i^n$ with $s_i(h^0) = N$ for $i = A$, $s_i(h^0) = W$ for $i = B$.

Suppose by contradiction the following:

(♠) For every μ_B that strongly believes S_A^n, \dots, S_A^0 with $\mu_B(S_A(z)|h^0) = \bar{\mu}_B^h(S_A(z)|h^0) = 1$, $\rho(\mu_B) \cap S_B(z) = \emptyset$.¹²

Let $\hat{h} := (N, E)$. For each $a \in S_A^n(\hat{h})|\hat{h}$ (non-empty by the induction hypothesis), fix $s_A \in S_A^n(z)$ with $s_A|\hat{h} = a$;¹³ there exists μ_B that strongly believes S_A^n, \dots, S_A^0 with $\mu_B(s_A|h^0) = 1$, so by (♠) $\rho(\mu_B) \subseteq S_B^n(\hat{h})$. For each $b \in S_B^n(\hat{h})|\hat{h}$, fix $s_B \in S_B^n(\hat{h})$ with $s_B|\hat{h} = b$; there exists μ_A that strongly believes S_B^n, \dots, S_B^0 with $\mu_A(W|h^0) = \bar{\mu}_A(\bar{s}_A)(W|h^0) = 1$ ($W \in S_B^n$ by the induction hypothesis) and $\mu_A(s_B|\hat{h}) = 1$, so $\rho(\mu_A) \subseteq S_A^n(\hat{h})$. Hence, $S^n(\hat{h})|\hat{h}$ features all best replies to beliefs in the set.

Let $(\bar{S}_q^{\hat{h}})_{q=0}^\infty = ((S^q(\hat{h})|\hat{h})_{q=0}^n, (\bar{S}_q^{\hat{h}})_{q=n+1}^\infty)$, where for each $m \geq n+1$, $i = A, B$, and μ_i that strongly believes $(\bar{S}_{-i,q}^{\hat{h}})_{q=0}^{m-1}$, $\rho(\mu_i) \subseteq \bar{S}_{i,m}^{\hat{h}}$. Since $S^n(\hat{h})|\hat{h}$ features all best replies to beliefs in the set, $\bar{S}_n^{\hat{h}} \supseteq \bar{S}_{n+1}^{\hat{h}}$ and thus $(\bar{S}_q^{\hat{h}})_{q=0}^\infty$ is an elimination procedure with $\bar{S}_\infty^{\hat{h}} \neq \emptyset$. Let

¹²For Ann, it is obvious that this cannot hold, as *S* is not optimal against *W*.

¹³It obviously exists because any strategy of Ann that allows (N, E) must prescribe *N*.

$(S_{\Delta}^{\widehat{h}})_{q=0}^{\infty} = ((S_{\Delta}^q(\widehat{h})|\widehat{h})_{q=0}^{\infty})$. For each $a \in \overline{S}_{A,\infty}^{\widehat{h}}$ and $q \in \mathbb{N}$, if $a \in S_{A,q}^{\widehat{h}}$, there is $s_A \in S_{A,\Delta}^q(z)$ with $s_A|\widehat{h} = a$; thus, there exists μ_B that strongly believes $S_{B,\Delta}^q, \dots, S_{B,\Delta}^0$ with $\mu_B(s_A|h^0) = 1$, and by the incentives given by (\spadesuit) , $\rho(\mu_B) \subseteq S_{B,\Delta}^n(\widehat{h})$. So, the best replies to a are in $S_{B,q+1}^{\widehat{h}}$. For each $b \in \overline{S}_{B,\infty}^{\widehat{h}}$ and $q \in \mathbb{N}$, if $b \in S_{B,q}^{\widehat{h}}$, there is $s_B \in S_{B,\Delta}^q(\widehat{h})$ with $s_B|\widehat{h} = b$; thus, there exists μ_A that strongly believes $S_{A,\Delta}^q, \dots, S_{A,\Delta}^0$ with $\mu_A(\cdot|h^0) = \bar{\mu}_i(\bar{s}_A)(\cdot|h^0)$ and $\mu_A(s_B|\widehat{h}) = 1$, and $\rho(\mu_A) \subseteq S_{A,\Delta}^n(\widehat{h})$. So, the best replies to b are in $S_{A,q+1}^{\widehat{h}}$. Then, since $\overline{S}_{\infty}^{\widehat{h}}$ is a set with the best reply property, $\emptyset \neq \overline{S}_{\infty}^{\widehat{h}} \subseteq S_{\infty}^{\widehat{h}}$, which contradicts $S_{\Delta}^{\infty}(\widehat{h}) = \emptyset$.

7 Appendix

Proof of Lemma 3. Fix $n \geq 0$ and suppose to have shown that for each $m \leq n$, $S_{\Delta}^m = S_{\Delta^*}^m$ ($S_{\Delta}^0 = S_{\Delta^*}^0$.trivially holds). If $S_{\Delta}^n = \emptyset$, $S_{\Delta}^{n+1} = S_{\Delta^*}^{n+1} = \emptyset$. Else, for each $i \in I$, there exists $\bar{\mu}_i \in \Delta_i$ that strongly believes $(S_{-i,\Delta}^q)_{q=0}^{n-1}$ such that $\rho(\bar{\mu}_i) \cap S_i(z) \neq \emptyset$. Fix $i \in I$ and $s_i \in S_i \setminus S_i(z)$. Let $m := \max \left\{ q \leq n : s_i \in S_{i,\Delta}^q \right\}$. If $m > 0$, there exists $\mu_i \in \Delta_i$ that strongly believes $(S_{-i,\Delta}^q)_{q=0}^{m-1}$ such that $s_i \in \rho(\mu_i)$. Fix $\mu_i^* \in \Delta_i$ that strongly believes $(S_{-i,\Delta}^q)_{q=0}^{m-1}$ such that $\mu_i^*(\cdot|h) = \bar{\mu}_i(\cdot|h)$ for all $h \prec z$ and $\mu_i^*(\cdot|\tilde{h}) = \mu_i(\cdot|\tilde{h})$ for all $\tilde{h} \in H(S_i(z)) \setminus H(S_{-i}(z))$ (it is compatible with CPS-3 because $\bar{\mu}_i(S_{-i}(\tilde{h})|h) = 0$ for all $h \prec z$ and $\tilde{h} \in H(S_i(z)) \setminus H(S_{-i}(z))$). Then, there exists $s_i^* \in \rho(\mu_i^*)(z) \subseteq S_{i,\Delta}^m$ such that for all $\tilde{h} \in H(s_i) \cap H(S_i(z)) \setminus H(S_{-i}(z))$, $s_i^*(\tilde{h}) = s_i(\tilde{h})$. If $m = 0$, fix the unique $s_i^* \in S_i(z)$ such that for all $\tilde{h} \not\prec z$, $s_i^*(\tilde{h}) = s_i(\tilde{h})$. For each $h \in H(S_i(z))$, let $\eta^h(s_i) := s_i^*$. For each $h \notin H(S_i(z))$, let $\eta^h(s_i) := s_i$. For all $s_i \in S_i(z)$ and $h \in H$, let $\eta^h(s_i) := s_i$.

Fix now $i \in I$ and $\mu_i \in \Delta_i$ that strongly believes $(S_{-i,\Delta}^q)_{q=0}^n$. Note that for each $s_i \in S_i$ and $h \in H$, if $s_i \in S_i(h)$, $\eta^h(s_i) \in S_i(h)$, and if $h \in H(S_i(z))$, $\eta^h(s_i) \in S_i(z)$. Thus, I can construct $\mu_i^* \in \Delta_i^*$ that strongly believes $(S_{-i,\Delta}^q)_{q=0}^n = (S_{-i,\Delta^*}^q)_{q=0}^n$ as, for all $h \in H$, $\mu_i^*((s_j)_{j \neq i}|h) = \mu_i(((\eta^h)^{-1}(s_j))_{j \neq i}|h)$. For each $h \prec z$, since $\mu_i(S_{-i}(z)|h) = 1$, $\mu_i^*(\cdot|h) = \mu_i(\cdot|h)$, and for each $h \not\prec z$ and $\tilde{z} \succ h$, by construction, $\mu_i^*(S_{-i}(\tilde{z})|h) = \mu_i(S_{-i}(\tilde{z})|h)$. Hence, $\rho(\mu_i) = \rho(\mu_i^*)$. So, $S_{\Delta}^{n+1} \subseteq S_{\Delta^*}^{n+1}$. By $\Delta_i^* \subseteq \Delta_i$ and $(S_{-i,\Delta}^q)_{q=0}^n = (S_{-i,\Delta^*}^q)_{q=0}^n$, $S_{\Delta^*}^{n+1} \subseteq S_{\Delta}^{n+1}$.

The proof can be repeated for Selective Rationalizability with $n \geq M$ in place of $n \geq 0$, where $(S_{R\Delta}^q)_{q=0}^M = (S_{R\Delta^*}^q)_{q=0}^M$ holds by definition. \blacksquare

Formal proof of Lemma 2.

We need additional notation. For any $h \in H$, $\widehat{h} \succeq h$, $(s_j^h)_{j \in I} \in S^h$, $(\widehat{s}_j^{\widehat{h}})_{j \in I} \in S^{\widehat{h}}$, $\mu_i^h \in \Delta^{H^h}(S_{-i}^h)$, $\mu_i^{\widehat{h}} \in \Delta^{H^{\widehat{h}}}(S_{-i}^{\widehat{h}})$, $\widehat{Z} \subseteq Z^{\widehat{h}}$, and $J \subseteq I$, let:

- $s_j^h = \widehat{Z} s_j^{\widehat{h}}$ if for each $z \in \widehat{Z}$ and $\widehat{h} \preceq \tilde{h} \prec z$, $s_j^h(\tilde{h}) = s_j^{\widehat{h}}(\tilde{h})$;

- $\mu_i^h =_{\widehat{Z}} \mu_i^{\widehat{h}}$ if for each $z \in \widehat{Z}$ and $\widehat{h} \preceq \widetilde{h} \prec z$, $\mu_i^h(S_{-i}^h(z)|\widetilde{h}) = \mu_i^{\widehat{h}}(S_{-i}^{\widehat{h}}(z)|\widetilde{h})$;
- $s_J^h =_{\widehat{h}} s_J^{\widehat{h}}$ and $\mu_i^h =_{\widehat{h}} \mu_i^{\widehat{h}}$ if, respectively, $s_J^h =_{Z^{\widehat{h}}} s_J^{\widehat{h}}$ and $\mu_i^h =_{Z^{\widehat{h}}} \mu_i^{\widehat{h}}$;
- $\widehat{r}(\mu_i^h, \widehat{h})$ is the set of continuation best replies to $\mu_i^h(\cdot|\widehat{h})$.

Moreover, for any $\overline{S}^h = \times_{i \in I} \overline{S}_i^h \subseteq S^h$, define the set of histories that follow a unilateral deviation by player i from the paths induced by \overline{S}^h as:

- $D_i(\overline{S}^h) := \{\widetilde{h} \in H \setminus H(\overline{S}^h) : p(\widetilde{h}) \in H(\overline{S}^h) \wedge \widetilde{h} \in H(\overline{S}_{-i}^h)\}$.

The first two 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 $\widehat{s}_i^h, \overline{s}_i^h \in S_{i,n}^h$ and $\widehat{h}, \overline{h} \in H(\widehat{s}_i^h) \cap H(\overline{s}_i^h)$ such that $\overline{h} \not\prec \widehat{h} \not\prec \overline{h}$: there needs not exist $s_i^h \in S_{i,n}^h(\widehat{h}) \cap S_{i,n}^h(\overline{h})$ such that $s_i^h|\widehat{h} = \overline{s}_i^h|\widehat{h}$ and $s_i^h|\overline{h} = \widehat{s}_i^h|\overline{h}$. The intuitive reason is the following: player i may allow \widehat{h} and \overline{h} either because she is confident that \widehat{h} will be reached and she has appropriate expectations after \widehat{h} , or because she is confident that \overline{h} will be reached and she has appropriate expectations after \overline{h} . If \widehat{s}_i^h is best reply to the first conjecture and \overline{s}_i^h is best reply to the second conjecture, $\widehat{s}_i^h|\overline{h}$ and $\overline{s}_i^h|\widehat{h}$ may be "emergency plans" for an unpredicted contingency, after which the expectations would not have justified the choice to allow \overline{h} and \widehat{h} in the first place. Here is an example. The following is a simplified version of the game in Figure 4 in Battigalli [1], provided by Gul and Reny. The payoffs are of player 1.

		2	← o →	1		
				↓ i		
1	← a →	1	← l →	2	- r →	1 - a' → 1
		↓ b				b' ↓
0	← c →	3				3 - c' → 0
		↓ d				d' ↓
		3				3

Player 1 can rationally play $i.a.b'$ (if she expects r and d' but not d), $i.b.a'$ (if she expects l and d but not d'), but not $i.a.a'$. If one starts from $i.a.b'$, you cannot modify b' into a' because $i.a.b'$ is a sequential best reply only to CPS's that assign initial positive probability to r , therefore the belief at (i,r) cannot be modified without modifying the initial belief, hence the previous choices. Instead, $i.a.b'$ can be modified into $i.b.b'$ because $i.a.b'$ is rational under zero probability to l .

Lemma 5 Fix an elimination procedure $((S_{i,q}^h)_{i \in I})_{q \geq 0}$, $n \in \mathbb{N}$, $i \in I$, $\hat{h} \in H^h$, and μ_i^h that strongly believes $(S_{-i,q}^h)_{q=0}^{n-1}$ such that $\mu_i^h(S_{-i}^h(\hat{h})|p(\hat{h})) = 0$. Fix $s_i^h \in \rho(\mu_i^h) \cap S_i^h(\hat{h})$, $\mu_i^{\hat{h}}$ that strongly believes $(S_{-i,q}^{\hat{h}}(\hat{h})|_{\hat{h}})_{q=0}^{n-1}$, and $\hat{s}_i^h \in \rho(\mu_i^{\hat{h}})$.

Consider the unique $\tilde{s}_i^h =^{\hat{h}} s_i^{\hat{h}}$ such that for every $\tilde{h} \notin H^{\hat{h}}$, $\tilde{s}_i^h(\tilde{h}) = s_i^{\hat{h}}(\tilde{h})$.

There exists $\tilde{\mu}_i^h =^{\hat{h}} \mu_i^{\hat{h}}$ that strongly believes $(S_{-i,q}^{\hat{h}})_{q=0}^{n-1}$ such that $\tilde{\mu}_i^h(\cdot|\tilde{h}) = \mu_i^{\hat{h}}(\cdot|\tilde{h})$ for all $\tilde{h} \notin H^{\hat{h}}$, and $\tilde{s}_i^h \in \rho(\tilde{\mu}_i^h)$ (so, $\rho(\mu_i^h)(\hat{h}) \neq \emptyset$ implies $\rho(\tilde{\mu}_i^h)(\hat{h}) \neq \emptyset$).

Proof.

Fix a map $\varsigma : S_{-i}^{\hat{h}} \rightarrow S_{-i}^h$ such that for each $\hat{s}_{-i}^h \in S_{-i}^{\hat{h}}$, $\varsigma(\hat{s}_{-i}^h) =^{\hat{h}} \hat{s}_{-i}^h$ and $\varsigma(\hat{s}_{-i}^h) \in S_{-i,m}^h(\hat{h})$ for all $m \geq 0$ with $\hat{s}_{-i}^h \in S_{-i,m}^{\hat{h}}(\hat{h})|\hat{h}$. Since ς is injective, we can construct an array of probability measures $\tilde{\mu}_i^h = (\tilde{\mu}_i^h(\cdot|\tilde{h}))_{\tilde{h} \in H^h}$ on S_{-i}^h as $\tilde{\mu}_i^h(\cdot|\tilde{h}) = \mu_i^h(\cdot|\tilde{h})$ for all $\tilde{h} \notin H^{\hat{h}}$ and $\tilde{\mu}_i^h(\varsigma(\hat{s}_{-i}^h)|\tilde{h}) = \mu_i^{\hat{h}}(\hat{s}_{-i}^h|\tilde{h})$ for all $\tilde{h} \in H^{\hat{h}}$ and $\hat{s}_{-i}^h \in S_{-i}^{\hat{h}}$. Thus, $\tilde{\mu}_i^h$ satisfies CPS-1. It is immediate to verify that $\tilde{\mu}_i^h$ satisfies CPS-2, strongly believes $(S_{-i,q}^h)_{q=0}^{n-1}$, $\tilde{\mu}_i^h =^{\hat{h}} \mu_i^{\hat{h}}$. Finally, since $\tilde{\mu}_i^h(S_{-i}^h(\hat{h})|p(\hat{h})) = 0$, $\tilde{\mu}_i^h$ satisfies CPS-3.

Fix $\tilde{h} \in H(\tilde{s}_i^h) \setminus H^{\hat{h}} = H(s_i^h) \setminus H^{\hat{h}}$. If $\tilde{h} \prec \hat{h}$, by $\mu_i^h(S_{-i}^h(\hat{h})|p(\hat{h})) = 0$ and CPS-3, $\mu_i^h(S_{-i}^h(\hat{h})|\tilde{h}) = 0$, and for every $s_{-i}^h \notin S_{-i}^h(\hat{h})$, $\zeta(s_i^h, s_{-i}^h) = \zeta(\tilde{s}_i^h, s_{-i}^h)$. If $\tilde{h} \not\prec \hat{h}$, for every $s_{-i}^h \in S_{-i}^h(\tilde{h})$, $\hat{h} \notin H(s_i^h, s_{-i}^h)$, so $\zeta(s_i^h, s_{-i}^h) = \zeta(\tilde{s}_i^h, s_{-i}^h)$. Hence $s_i^h \in \hat{r}(\mu_i^h, \tilde{h})$ implies $\tilde{s}_i^h \in \hat{r}(\mu_i^h, \tilde{h}) = \hat{r}(\tilde{\mu}_i^h, \tilde{h})$. Fix $\tilde{h} \in H(\tilde{s}_i^h) \cap H^{\hat{h}} = H(\hat{s}_i^h)$. For every $\hat{s}_{-i}^h \in S_{-i}^{\hat{h}}$, $\mu_i^{\hat{h}}(\varsigma(\hat{s}_{-i}^h)|\tilde{h}) = \mu_i^{\hat{h}}(\hat{s}_{-i}^h|\tilde{h})$. For every $\hat{s}_i^h \in S_i^h(\hat{h})$, $\zeta(\hat{s}_i^h|\hat{h}, \hat{s}_{-i}^h) = \zeta(\hat{s}_i^h, \varsigma(\hat{s}_{-i}^h))$. So, $\hat{s}_i^h|\hat{h} = \hat{s}_i^h \in \hat{r}(\mu_i^{\hat{h}}, \tilde{h})$ implies $\tilde{s}_i^h \in \hat{r}(\tilde{\mu}_i^h, \tilde{h})$. ■

Lemma 6 Fix an elimination 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 $Z^S := \zeta(\bar{S}^h)$. For every $i \in I$, suppose that there exists a map $\bar{\mu}_i^h : \bar{S}_i^h \rightarrow \Delta^{H^h}(S_{-i}^h)$ such that for all $s_i^h \in \bar{S}_i^h$, $\bar{\mu}_i^h(s_i^h)$ strongly believes \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 all $s_i^h \in \bar{S}_i^h$, $\bar{\mu}_i^h(s_i^h) =^{Z^S} \bar{\mu}_i^h(s_i^h)$ strongly believes $(\tilde{S}_{-i,q}^h)_{q=0}^{m-1}$ and $\rho(\bar{\mu}_i^h(s_i^h)) \ni \bar{s}_i^h(s_i^h) =^{Z^S} s_i^h$;

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

Fix $l \in I$ and $s_l^h \in \bar{S}_l^h$. Let $D^S := D_l(\bar{S}^h)$. For every $\hat{h} \in D^S$, fix $\tilde{\mu}_l^{\hat{h}}$ that strongly believes $(\tilde{S}_{-l,q}^h(\hat{h})|\hat{h})_{q=0}^m$.

There exists $\tilde{\mu}_l^h =^{Z^S} \bar{\mu}_l^h(s_l^h)$ that strongly believes $(\tilde{S}_{-l,q}^h)_{q=0}^m$ such that $\tilde{\mu}_l^h =^{\hat{h}} \tilde{\mu}_l^{\hat{h}}$ for all $\hat{h} \in D^S$.

Proof.

We show that for every $i \neq l$ and $s_i^h \in \bar{S}_i^h$, and for every map $\varsigma : \hat{h} \in D^S \mapsto s_i^{\hat{h}} \in \tilde{S}_{i,m}^h(\hat{h})|\hat{h}$, there exists $\tilde{s}_i^h \in \tilde{S}_{i,m}^h$ such that $\tilde{s}_i^h =^{Z^S} \bar{s}_i^h(s_i^h)$ and $\tilde{s}_i^h =^{\hat{h}} \varsigma(\hat{h})$ for all $\hat{h} \in D^S$.

The map ς is well defined because for each $\widehat{h} \in D^S$, by A1 $\widehat{h} \in H(\overline{s}_i^h(\widehat{s}_i^h))$ for some $\widehat{s}_i^h \in \overline{S}_i^h$, and by A2, $\overline{s}_i^h(\widehat{s}_i^h) \in \widetilde{S}_{i,m}^h$. Using all such \widehat{s}_i^h 's, it is easy to construct the desired $\widetilde{\mu}_i^h$.

By A1, there exists $\overline{\mu}_i^h(s_i^h) =^{Z^S} \overline{\mu}_i^h(s_i^h)$ that strongly believes $(\widetilde{S}_{-i,q}^h)_{q=0}^{m-1}$ such that $\overline{s}_i^h(s_i^h) \in \rho(\overline{\mu}_i^h(s_i^h))$. Fix $\widehat{h} \in D^S \cap H(s_i^h)$. Since $\overline{\mu}_i^h(s_i^h) =^{Z^S} \overline{\mu}_i^h(s_i^h)$ and $\overline{\mu}_i^h(s_i^h)$ strongly believes \overline{S}_{-i}^h , $\overline{\mu}_i^h(s_i^h)(S_{-i}^h(\widehat{h})|p(\widehat{h})) = 0$. Since $\varsigma(\widehat{h}) \in \widetilde{S}_{i,m}^h(\widehat{h})|\widehat{h}$, there exists μ_i^h that strongly believes $(\widetilde{S}_{-i,q}^h(\widehat{h})|\widehat{h})_{q=0}^{m-1}$ such that $\varsigma(\widehat{h}) \in \rho(\mu_i^h)$. Thus, by Lemma 5, there exist $\widetilde{\mu}_i^h =^{\widehat{h}} \mu_i^h$ that strongly believes $(\widetilde{S}_{-i,q}^h)_{q=0}^{m-1}$ such that $\widetilde{\mu}_i^h(\cdot|\widehat{h}) = \overline{\mu}_i^h(s_i^h)(\cdot|\widehat{h})$ for all $\widehat{h} \notin H^{\widehat{h}}$, and $\widetilde{s}_i^h \in \rho(\widetilde{\mu}_i^h)$ such that $\widetilde{s}_i^h =^{\widehat{h}} \varsigma(\widehat{h})$ and $\widetilde{s}_i^h(\widehat{h}) = \overline{s}_i^h(s_i^h)(\widehat{h})$ for all $\widehat{h} \notin H^{\widehat{h}}$. Iterating for each $\widehat{h} \in D^S$, we obtain $\widetilde{\mu}_i^h =^{Z^S} \overline{\mu}_i^h(s_i^h)$ that strongly believes $(\widetilde{S}_{-i,q}^h)_{q=0}^{m-1}$ such that $\widetilde{\mu}_i^h =^{\widehat{h}} \mu_i^h$ for all $\widehat{h} \in D^S$, and $\widetilde{s}_i^h \in \rho(\widetilde{\mu}_i^h)$ such that $\widetilde{s}_i^h =^{Z^S} s_i^h$ and $\widetilde{s}_i^h =^{\widehat{h}} \varsigma(\widehat{h})$ for all $\widehat{h} \in D^S$. By A2, $\widetilde{s}_i^h \in \widetilde{S}_{i,m}^h$. ■

Lemma 7 Fix two elimination procedures $((\overline{S}_{i,q}^h)_{i \in I, q \geq 0})$ and $((S_{i,q}^h)_{i \in I, q \geq 0})$. For every $i \in I$, let $\overline{S}_i^h := \overline{S}_{i,\infty}^h$ and let $\overline{\mu}_i^h : \overline{S}_i^h \rightarrow \Delta^{H^h}(S_{-i}^h)$ be a map such that for every $s_i^h \in \overline{S}_i^h$, $\overline{\mu}_i^h(s_i^h)$ strongly believes $(\overline{S}_{-i,q}^h)_{q=0}^\infty$ and $s_i^h \in \rho(\overline{\mu}_i^h(s_i^h))$. Let $Z^S := \zeta(\overline{S}^h)$. Fix $n \in \mathbb{N}$, $l \in I$, and $\widehat{s}_l^h \in \overline{S}_l^h$ such that.¹⁴

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 $m \in \mathbb{N}$, $(\overline{S}_q^h)_{q \geq 0}$ satisfies A2;

A6 for every $s_l^h =^{Z^S} \widehat{s}_l^h$ and $\mu_l^h =^{Z^S} \overline{\mu}_l^h(\widehat{s}_l^h)$ that strongly believes $(S_{-l,q}^h)_{q=0}^n$, $s_l^h \notin \rho(\mu_l^h)$.

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 all $\widehat{\mu}_l^h$ that strongly believe $(S_{-l,q}^h(\widehat{h})|\widehat{h})_{q=0}^m$ (resp., $(\overline{S}_{-l,q}^h(\widehat{h})|\widehat{h})_{q=0}^m$) for which there exists $\widehat{\mu}_l^h$ that strongly believes $(S_{-l,q}^h(\widehat{h})|\widehat{h})_{q=0}^n$ such that $\mu_l^h(S_{-i}(z)|\widehat{h}) = \widehat{\mu}_l^h(S_{-i}(z)|\widehat{h})$ for all $z \in \zeta(\widehat{r}(\widehat{\mu}_l^h, \widehat{h}) \times \text{Supp} \widehat{\mu}_l^h(\cdot|\widehat{h}))$.¹⁵

Thus, there exists $\widehat{h} \in D^S$ such that:

1. for every $m \leq n$ and $\widehat{\mu}_l^h \in M_m^{\widehat{h}}$, there exists $\mu_l^h =^{Z^S} \overline{\mu}_l^h(\widehat{s}_l^h)$ that strongly believes $(S_{-l,q}^h)_{q=0}^m$ such that $\mu_l^h =^{\widehat{h}} \widehat{\mu}_l^h$ and $\rho(\mu_l^h)(\widehat{h}) \neq \emptyset$;
2. for every $p \in \mathbb{N}$ and $\widehat{\mu}_l^h \in \overline{M}_p^{\widehat{h}}$, there exists $\widetilde{\mu}_l^h =^{Z^S} \overline{\mu}_l^h(\widehat{s}_l^h)$ that strongly believes $(\overline{S}_{-l,q}^h)_{q=0}^p$ such that $\widetilde{\mu}_l^h =^{\widehat{h}} \widehat{\mu}_l^h$ and $\rho(\widetilde{\mu}_l^h)(\widehat{h}) \neq \emptyset$.¹⁶

¹⁴ A3, A4 and A5 need not hold for $i = l$ to recall Lemma 6 and prove this lemma. However, l has been included to reuse A3, A4 and A5 in the final proof of Lemma 2.

¹⁵ Note: $\widehat{\mu}_l^h$ refers to the second procedure even when μ_l^h refers to the first.

¹⁶ Since $\widehat{h} \notin H^S$, the statement must hold vacuously for some $p \in \mathbb{N}$ (i.e. $\overline{M}_p^{\widehat{h}} = \emptyset$).

Proof.

Suppose by contraposition that there is a partition (D, \bar{D}) of D^S such that for every $\hat{h} \in D$, there exist $m(\hat{h}) \leq n$ and $\mu_l^{\hat{h}} \in M_{m(\hat{h})}^{\hat{h}}$ that violate 1, and for every $\hat{h} \in \bar{D}$ there exist $m(\hat{h}) \in \mathbb{N}$ and $\mu_l^{\hat{h}} \in \bar{M}_{m(\hat{h})}^{\hat{h}}$ that violate 2. For each $\hat{h} \in D^S$, fix corresponding $\hat{\mu}_l^{\hat{h}}$. Let $\bar{\mu}_l^{\hat{h}} := \bar{\mu}_l^{\hat{h}}(\hat{s}_l^{\hat{h}})$. By Lemma 6, there exists $\tilde{\mu}_l^{\hat{h}} =^{Z^S} \bar{\mu}_l^{\hat{h}}$ that strongly believes $(S_{-l,q}^{\hat{h}})_{q=0}^n$ such that for every $\hat{h} \in D^S$, $\tilde{\mu}_l^{\hat{h}} =^{\hat{h}} \hat{\mu}_l^{\hat{h}}$. We want to show that there exists $s_l^{\hat{h}} \in \rho(\tilde{\mu}_l^{\hat{h}})$ such that $s_l^{\hat{h}} =^{Z^S} \hat{s}_l^{\hat{h}}$, violating A6.

Fix $\hat{h} \in D$. Substitute $\hat{\mu}_l^{\hat{h}}$ with $\mu_l^{\hat{h}}$ in the construction of $\tilde{\mu}_l^{\hat{h}}$ and obtain a new $\mu_l^{\hat{h}} =^{\hat{h}} \mu_l^{\hat{h}}$ that strongly believes $(S_{-l,q}^{\hat{h}})_{q=0}^{m(\hat{h})}$ with $\mu_l^{\hat{h}}(S_{-l}(z)|\tilde{h}) = \tilde{\mu}_l^{\hat{h}}(S_{-l}(z)|\tilde{h})$ for all $\tilde{h} \notin H^{\hat{h}}$ and $z \notin Z^{\hat{h}}$. By definition of $M_m^{\hat{h}}$, player l expects a non higher payoff against $\hat{\mu}_l^{\hat{h}}$ than against $\mu_l^{\hat{h}}$. Thus, $\rho(\mu_l^{\hat{h}})(\hat{h}) \neq \emptyset$ (by the contrapositive hypothesis) implies $\rho(\tilde{\mu}_l^{\hat{h}})(\hat{h}) \neq \emptyset$. So, $H(\rho(\tilde{\mu}_l^{\hat{h}})) \cap D = \emptyset$.

Write $\bar{D} = \{h^1, \dots, h^k\}$ where $m(h^1) \geq \dots \geq m(h^k)$. Note that $(\bar{S}_q^{\hat{h}})_{q \geq 0}$ satisfies A1 with $\bar{\mu}_i^{\hat{h}}(\cdot) = \bar{\mu}_i^{\hat{h}}(\cdot)$ and the identity function for $\bar{s}_i^{\hat{h}}(\cdot)$. Then, by Lemma 6,¹⁷ for each $j = 1, \dots, k$, there exists $\mu_{l,j}^{\hat{h}} =^{Z^h \setminus \cup_{t=1}^j Z^{h^t}} \bar{\mu}_l^{\hat{h}}$ that strongly believes $(\bar{S}_{-l,q}^{\hat{h}})_{q=0}^{m(h^j)}$ such that $\mu_{l,j}^{\hat{h}} =^{h^t} \mu_l^{\hat{h}}$ for all $1 \leq t \leq j$. Let $\mu_{l,0}^{\hat{h}} := \bar{\mu}_l^{\hat{h}}$. Fix $j = 1, \dots, k$ and suppose to have shown that $\rho(\mu_{l,j-1}^{\hat{h}}) = \rho(\bar{\mu}_l^{\hat{h}})$. Then $\rho(\mu_{l,j-1}^{\hat{h}}) \cap S_l^{\hat{h}}(h^j) = \emptyset$. By the contrapositive hypothesis, $\rho(\mu_{l,j}^{\hat{h}}) \cap S_l^{\hat{h}}(h^j) = \emptyset$. For all $\tilde{h} \notin H^{h^j}$ and $z \notin Z^{h^j}$, $\mu_{l,j}^{\hat{h}}(S_{-l}(z)|\tilde{h}) = \mu_{l,j-1}^{\hat{h}}(S_{-l}(z)|\tilde{h})$. Then, $\rho(\mu_{l,j}^{\hat{h}}) = \rho(\mu_{l,j-1}^{\hat{h}})$. Inductively, $\rho(\mu_{l,k}^{\hat{h}}) = \rho(\bar{\mu}_l^{\hat{h}}) \ni \hat{s}_l^{\hat{h}}$.

Fix $\tilde{h} \in H(\hat{s}_l^{\hat{h}}) \cap H^S \cap H(\rho(\tilde{\mu}_l^{\hat{h}}))$. By $\tilde{\mu}_l^{\hat{h}} =^{Z^S} \bar{\mu}_l^{\hat{h}} =^{Z^S} \mu_{l,k}^{\hat{h}}$, $\tilde{\mu}_l^{\hat{h}}(S_{-l}(z)|\tilde{h}) = \mu_{l,k}^{\hat{h}}(S_{-l}(z)|\tilde{h})$ for all $z \in Z^{\tilde{h}} \cap Z^S$. Then, since $\bar{\mu}_l^{\hat{h}}$ strongly believes $\bar{S}_{-l}^{\hat{h}}$, $\hat{s}_l^{\hat{h}}$, as well as any other $\hat{s}_l^{\hat{h}} \in S_l^{\hat{h}}$ with $H(\hat{s}_l^{\hat{h}}) \cap D^S = \emptyset$, induces the same outcome distribution against $\tilde{\mu}_l^{\hat{h}}(\cdot|\tilde{h})$ and $\mu_{l,k}^{\hat{h}}(\cdot|\tilde{h})$. Moreover, $H(\rho(\tilde{\mu}_l^{\hat{h}})) \cap D = \emptyset$. Finally, for all $\hat{h} \in \bar{D}$, by definition of $\bar{M}_m^{\hat{h}}$, player l expects a non higher payoff against $\hat{\mu}_l^{\hat{h}}$ than against $\mu_l^{\hat{h}}$, and recall that $\tilde{\mu}_l^{\hat{h}} =^{\hat{h}} \hat{\mu}_l^{\hat{h}}$ and $\mu_{l,k}^{\hat{h}} =^{\hat{h}} \mu_l^{\hat{h}}$. So, $\hat{s}_l^{\hat{h}} \in \hat{r}(\mu_{l,k}^{\hat{h}}, \tilde{h})$ implies $\hat{s}_l^{\hat{h}} \in \hat{r}(\tilde{\mu}_l^{\hat{h}}, \tilde{h})$. Proceeding from the root of the game, this implies $H(\hat{s}_l^{\hat{h}}) \cap H^S \subseteq H(\rho(\tilde{\mu}_l^{\hat{h}})) \cap H^S$. Thus, there exists $s_l^{\hat{h}} \in \rho(\tilde{\mu}_l^{\hat{h}})$ such that $s_l^{\hat{h}}(\tilde{h}) = \hat{s}_l^{\hat{h}}(\tilde{h})$ for all $\tilde{h} \in H^S$. ■

Proof of Lemma 2.

Recall that the depth of a game is the length of the longest terminal history of the game. Suppose that $\Gamma(h)$ has depth $k \in \mathbb{N}$ and, if $k > 1$, that the lemma holds for games of depth $1, \dots, k-1$. Let $\bar{S}_{\infty}^{\hat{h}} \neq \emptyset$, otherwise the lemma trivially holds.

We prove by induction that $\zeta(\bar{S}_{\infty}^{\hat{h}}) \subseteq \zeta(S_{\infty}^{\hat{h}})$. Note first that A4 and A5 hold by hypothesis of the lemma.

¹⁷Using the identity function for $\bar{s}_i^{\hat{h}}(\cdot)$ in the proof of the lemma and without iterating at histories $\hat{h} \in D^S \setminus \{h^1, \dots, h^j\}$, the constructed $\mu_{l,j}^{\hat{h}}$ clearly has the desired features.

Induction Hypothesis (n): $(S_q^h)_{q=0}^\infty$ satisfies A3 at n (so by A4 $\zeta(S_n^h) \supseteq \zeta(\overline{S}_\infty^h)$).

Basis step (1): for all $i \in I$, the Inductive Hypothesis holds with $\overline{\mu}_i^h(\cdot) = \overline{\mu}_i^h(\cdot)$.

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 $\widehat{s}_l^h \in \overline{S}_{l,\infty}^h$. Lemma 7 yields $\widehat{h} \in D_l(\overline{S}_\infty^h)$. If $\Gamma(\widehat{h})$ has depth 1, $D_l(\overline{S}_\infty^h) = \emptyset$, so we have the desired contradiction. Else, define $((\overline{S}_{i,q}^{\widehat{h}})_{i \in I})_{q \geq 0}$ as follows: for every $i \in I$ and $m \leq n$, $\overline{S}_{i,m}^{\widehat{h}} = S_{i,m}^h(\widehat{h})|\widehat{h}$; for every $m > n$, $\widehat{s}_i^h \in \overline{S}_{i,m}^{\widehat{h}}$ if and only if there exists $\mu_i^{\widehat{h}}$ that strongly believes $(\overline{S}_{-i,q}^{\widehat{h}})_{q=0}^{m-1}$ such that $\widehat{s}_i^h \in \rho(\mu_i^{\widehat{h}})$.

For every $i \neq l$, since $\widehat{h} \in D_l(\overline{S}_\infty^h)$, $\overline{S}_{i,\infty}^{\widehat{h}}(\widehat{h}) \neq \emptyset$. So, fix $\widehat{s}_i^h \in \overline{S}_{i,\infty}^{\widehat{h}}(\widehat{h})$. For every $m \leq n$, the Induction Hypothesis provides $\overline{s}_i^h(\widehat{s}_i^h) \in S_{i,m}^h(\widehat{h}) \neq \emptyset$ and $\overline{\mu}_i^h(\widehat{s}_i^h) = {}^{\zeta(\overline{S}_\infty^h)} \overline{\mu}_i^h(\widehat{s}_i^h)$ that strongly believes $(S_{-i,q}^h)_{q=0}^{m-1}$ such that $\overline{\mu}_i^h(\widehat{s}_i^h)(S_{-i}^h(\widehat{h})|p(\widehat{h})) = 0$. Hence, by Lemma 5, for every $\mu_i^{\widehat{h}}$ that strongly believes $(\overline{S}_{-i,q}^{\widehat{h}})_{q=0}^{m-1}$, there exists $\mu_i^h = \widehat{h} \mu_i^{\widehat{h}}$ that strongly believes $(S_{-i,q}^h)_{q=0}^{m-1}$ such that $\mu_i^h = {}^{\zeta(\overline{S}_\infty^h)} \overline{\mu}_i^h(\widehat{s}_i^h)$ and $\rho(\mu_i^h)(\widehat{h}) \neq \emptyset$. By A4, $\rho(\mu_i^h) \subseteq S_{i,m}^h$. So, $\rho(\mu_i^{\widehat{h}}) \subseteq \overline{S}_{i,m}^{\widehat{h}}$.

Fix $\mu_l^{\widehat{h}}$ that strongly believes $(\overline{S}_{-l,q}^{\widehat{h}})_{q=0}^n$: trivially $\mu_l^{\widehat{h}} \in M_n^{\widehat{h}}$. Hence, by Lemma 7.(1), there exists $\widetilde{\mu}_l^h = {}^{\zeta(\overline{S}_\infty^h)} \overline{\mu}_l^h(\widehat{s}_l^h)$ that strongly believes $(S_{-l,n}^h)_{q=0}^n$ such that $\widetilde{\mu}_l^h = \widehat{h} \mu_l^{\widehat{h}}$ and $\rho(\widetilde{\mu}_l^h)(\widehat{h}) \neq \emptyset$. By A4, $\rho(\widetilde{\mu}_l^h) \subseteq S_{l,n}^h$. So $\rho(\mu_l^{\widehat{h}}) \subseteq \overline{S}_{l,n}^{\widehat{h}} \neq \emptyset$.

Hence, for every $i \in I$ and $\mu_i^{\widehat{h}}$ that strongly believes $(\overline{S}_{-i,q}^{\widehat{h}})_{q=0}^n$, $\rho(\mu_i^{\widehat{h}}) \subseteq \overline{S}_{i,n}^{\widehat{h}} \neq \emptyset$. So, $\overline{S}_{i,n}^{\widehat{h}} \supseteq \overline{S}_{i,n+1}^{\widehat{h}}$ and $((\overline{S}_{i,q}^{\widehat{h}})_{i \in I})_{q \geq 0}$ is an elimination procedure with $\overline{S}_\infty^{\widehat{h}} \neq \emptyset$.

For every $m \leq n$, $\mu_l^{\widehat{h}}$ that strongly believes $(\overline{S}_{-l,q}^{\widehat{h}})_{q=0}^\infty$, and $\mu_l^{\widehat{h}} = {}^{\zeta(\overline{S}_\infty^h)} \mu_l^{\widehat{h}}$ that strongly believes $(\overline{S}_{-l,q}^{\widehat{h}})_{q=0}^{m-1}$, $\mu_l^{\widehat{h}} \in M_m^{\widehat{h}}$.¹⁸ Thus, by Lemma 7.(1) there exists $\widetilde{\mu}_l^h = {}^{\zeta(\overline{S}_\infty^h)} \overline{\mu}_l^h(\widehat{s}_l^h)$ that strongly believes $(S_{-l,q}^h)_{q=0}^{m-1}$ such that $\widetilde{\mu}_l^h = \widehat{h} \mu_l^{\widehat{h}}$ and $\rho(\widetilde{\mu}_l^h)(\widehat{h}) \neq \emptyset$. By A4, $\rho(\widetilde{\mu}_l^h) \subseteq S_{l,m}^h$. So $\rho(\mu_l^{\widehat{h}}) \subseteq \overline{S}_{l,m}^{\widehat{h}}$.

Then, for every $m \in \mathbb{N}$, $i \in I$, $\mu_i^{\widehat{h}}$ that strongly believes $(\overline{S}_{-i,q}^{\widehat{h}})_{q=0}^\infty$ and $\mu_i^{\widehat{h}} = {}^{\zeta(\overline{S}_\infty^h)} \mu_i^{\widehat{h}}$ that strongly believes $(\overline{S}_{-i,q}^{\widehat{h}})_{q=0}^{m-1}$, $\rho(\mu_i^{\widehat{h}}) \subseteq \overline{S}_{i,m}^{\widehat{h}}$. Thus, $((\overline{S}_{i,q}^{\widehat{h}})_{i \in I})_{q \geq 0}$ satisfies the hypothesis of Lemma 2.

Define now $((S_{i,q}^{\widehat{h}})_{i \in I})_{q \geq 0}$ as $((\overline{S}_{i,q}^{\widehat{h}}(\widehat{h})|\widehat{h})_{i \in I})_{q \geq 0}$. By Remark 1 it is an elimination procedure.

For every $i \neq l$, $m \in \mathbb{N}$, and $\mu_i^{\widehat{h}}$ that strongly believes $(S_{-i,q}^{\widehat{h}})_{q=0}^{m-1}$, by Lemma 5 there exists $\widetilde{\mu}_i^h = \widehat{h} \mu_i^{\widehat{h}}$ that strongly believes $(\overline{S}_{-i,q}^{\widehat{h}})_{q=0}^{m-1}$ such that for every $\widetilde{h} \notin H^{\widehat{h}}$, $\widetilde{\mu}_i^h(\cdot|\widetilde{h}) = \overline{\mu}_i^h(\widehat{s}_i^h)(\cdot|\widetilde{h})$ and $\rho(\widetilde{\mu}_i^h)(\widehat{h}) \neq \emptyset$. By A5, $\rho(\widetilde{\mu}_i^h) \subseteq \overline{S}_{i,m}^h$.

¹⁸Note that $\mu_l^{\widehat{h}}$ strongly believes $(\overline{S}_{-l,q}^{\widehat{h}})_{q=0}^n = (S_{-l,q}^h(\widehat{h})|\widehat{h})_{q=0}^n$, and that $\rho(\mu_l^{\widehat{h}}) \times \overline{S}_{-l,\infty}^{\widehat{h}} \subseteq \overline{S}_\infty^{\widehat{h}}$, so $\mu_l^{\widehat{h}} = {}^{\zeta(\overline{S}_\infty^h)} \mu_l^{\widehat{h}}$ verifies the definition of $M_m^{\widehat{h}}$ in the statement of Lemma 7.

For every $m \in \mathbb{N}$, $\hat{\mu}_l^{\hat{h}}$ that strongly believes $(\bar{S}_{-l,q}^{\hat{h}})_{q=0}^{\infty}$, and $\mu_l^{\hat{h}} = \zeta(\bar{S}_{\infty}^{\hat{h}}) \hat{\mu}_l^{\hat{h}}$ that strongly believes $(S_{-l,q}^{\hat{h}})_{q=0}^{m-1}$, $\mu_l^{\hat{h}} \in \bar{M}_m^{\hat{h}}$.¹⁹ Thus, by Lemma 7.(2) there exists $\tilde{\mu}_l^{\hat{h}} = \zeta(\bar{S}_{\infty}^{\hat{h}}) \bar{\mu}_l^{\hat{h}}(\hat{s}_l^{\hat{h}})$ that strongly believes $(\bar{S}_{-l,q}^{\hat{h}})_{q=0}^{m-1}$ such that $\tilde{\mu}_l^{\hat{h}} = \hat{\mu}_l^{\hat{h}}$ and $\rho(\tilde{\mu}_l^{\hat{h}})(\hat{h}) \neq \emptyset$. By A5 $\rho(\tilde{\mu}_l^{\hat{h}}) \subseteq \bar{S}_{l,m}^{\hat{h}}$.

Then, for every $m \in \mathbb{N}$, $i \in I$, $\hat{\mu}_i^{\hat{h}}$ that strongly believes $(\bar{S}_{-i,q}^{\hat{h}})_{q=0}^{\infty}$ and $\mu_i^{\hat{h}} = \zeta(\bar{S}_{\infty}^{\hat{h}}) \hat{\mu}_i^{\hat{h}}$ that strongly believes $(S_{-i,q}^{\hat{h}})_{q=0}^{m-1}$, $\rho(\mu_i^{\hat{h}}) \subseteq S_{i,m}^{\hat{h}}$. Thus, $((S_{i,q}^{\hat{h}})_{i \in I})_{q \geq 0}$ satisfies the hypothesis of Lemma 2

Since $\Gamma(\hat{h})$ has strictly lower depth than $\Gamma(h)$, Lemma 2 holds. Hence, $\zeta(\bar{S}_{\infty}^{\hat{h}}) \supseteq \zeta(\bar{S}_{\infty}^{\hat{h}}) \neq \emptyset$. But this contradicts $\hat{h} \in D_l(\bar{S}_{\infty}^{\hat{h}})$. ■

References

- [1] Battigalli, P., “On rationalizability in extensive games”, *Journal of Economic Theory*, **74**, 1997, 40-61.
- [2] Battigalli, P., “Dynamic Consistency and Imperfect Recall”, *Games and Economic Behavior*, **20(1)**, 1997, 31-50.
- [3] Battigalli, P., “Rationalizability in Infinite, Dynamic Games of Incomplete Information,” *Research in Economics*, **57**, 2003, 1-38.
- [4] 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.
- [5] Battigalli, P. and M. Siniscalchi, “Strong Belief and Forward Induction Reasoning,” *Journal of Economic Theory*, **106**, 2002, 356-391.
- [6] Battigalli P. and M. Siniscalchi, “Rationalization and Incomplete Information,” *The B.E. Journal of Theoretical Economics*, **3(1)**, 2003, 1-46.
- [7] Catonini, E., “Rationalizability and Epistemic Priority Orderings”, working paper, 2017.
- [8] Catonini, E., “Self-Enforcing Agreements and Forward Induction Reasoning”, working paper, 2017.
- [9] Chen, J., and S. Micali, “The order independence of iterated dominance in extensive games”, *Theoretical Economics*, **8**, 2013, 125-163.

¹⁹See the previous footnote with $\bar{M}_m^{\hat{h}}$ in place of $M_m^{\hat{h}}$.

- [10] Cho I.K. and D. Kreps, “Signaling Games and Stable Equilibria”, *Quarterly Journal of Economics*, **102**, 1987, 179-222.
- [11] 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.
- [12] Kohlberg, E. and J.F. Mertens, “On the Strategic Stability of Equilibria”, *Econometrica*, **54**, 1986, 1003-1038.
- [13] Osborne, M., “Signaling, Forward Induction, and Stability in Finitely Repeated Games”, *Journal of Economic Theory*, **50**, 1990, 22-36.
- [14] Osborne, M. J. and A. Rubinstein, “A Course in Game Theory”, 1994, Cambridge, Mass.: MIT Press.
- [15] Pearce, D., “Rational Strategic Behavior and the Problem of Perfection”, *Econometrica*, **52**, 1984, 1029-1050.
- [16] Penta, A., “Backward Induction Reasoning in Games with Incomplete Information”, 2011, working paper.
- [17] Perea, A., “Order Independence in Dynamic Games”, working paper, 2017.
- [18] Perea, A., “Why Forward Induction leads to the Backward Induction outcome: a new proof for Battigalli’s theorem”, *Games and Economic Behavior*, **110**, 2018, 120–138.
- [19] Renyi, A., “On a New Axiomatic Theory of Probability”, *Acta Mathematica Academiae Scientiarum Hungaricae*, **6**, 1955, 285-335.
- [20] Shimoji, M. and J. Watson, “Conditional dominance, rationalizability, and game forms”, *Journal of Economic Theory*, **83(2)**, 1998, 161-195.