

8 Supplemental Appendix

8.1 On SPE and self-enforcing agreements

Consider the following game.

$A \setminus B$	W	E
N	6, 6	·-
S	0, 0	2, 2

→

$A \setminus B$	L	C	R
U	9, 0	0, 5	0, 3
M	0, 5	9, 0	0, 3
D	0, 7	0, 7	1, 8

All plans are justifiable, hence they are all rationalizable. The subgame has one pure equilibrium, (D, R) , and no mixed equilibrium: for Ann to be indifferent between U and M , Bob must randomize over $\{L, C\}$, but when he is indifferent between them, he prefers R ; for Ann to be indifferent between U and D or M and D , Bob must randomize over, respectively, $\{L, R\}$ and $\{C, R\}$, but R dominates L over $\{U, D\}$ and C over $\{M, D\}$. So, the game has only one SPE, inducing outcome (S, E) .

The SPE outcome (S, E) is implementable, but differently from the game in the Introduction, only with an agreement that features also the off-the-path threat R by Bob. For instance, the reduced agreement on the realization-strict Nash $(S, E.R)$ is self-enforcing by Theorem 2. Instead, the path agreement on $z = (S, E)$ is not self-enforcing because Ann may rationally deviate and then play U or M , hence Bob could best reply with any action, and not just with R . Formally, we have $S_{A,z}^\infty = S_{A,z}^1 = \{S, N.U, N.M\}$ and $S_{B,z}^\infty = S_{B,z}^1 = \{E.L, E.C, E.R\}$.

Note moreover that if Ann believes in the SPE path, it is not rational for her to deviate and then play D . Thus, if Bob interprets the deviation of Ann as an attempt to increase her payoff with respect to the equilibrium (as implicitly assumed by strategic stability and related refinements, see Section 6), the fact that R is the best reply to D which is the best reply to R itself is of no value: R is a credible reaction of Bob only by virtue of other beliefs he may have.

Players can also implement the outcome (N, W) , and differently from the game in the Introduction, only with an agreement that *does not* feature a threat played with positive probability in an equilibrium of the subgame (here, just D). For instance, the reduced agreement with $e_A^0 = \{N.U, N.M\}$ and $e_B^0 = \{W\}$ is self-enforcing: we have $S_e^1 = \{N.U, N.M, N.D\} \times \{W\}$, thus $S_e^\infty = S_e^1 = S((N, W))$.

To conclude, note that there is no conceptual difference behind the reasons for self-enforceability of the SPE and of the Pareto-superior Nash outcome.

8.2 Another form of agreement incompleteness

Consider the following game.

4, 9, 5					$A \setminus B$	w	e
$\uparrow o$					n	3, 9, 0	0, 8, 2
Ann	5, 0, 1				s	0, 3, 0	1, 5, 2
$\downarrow i$	$u \uparrow$				\uparrow		
Bob	\longrightarrow	$Cleo$	\longrightarrow	a	\longrightarrow	Bob	
$\downarrow d$					\downarrow		
$C \setminus B$	l	c	r		$A \setminus B$	w	e
t	5, 4, 1	5, 6, 0	5, 0, 0		n	3, 9, 0	0, 8, 2
b	5, 4, 0	5, 0, 1	5, 10, 1		s	0, 3, 0	1, 5, 2

All plans are justifiable, hence they are all rationalizable. Players want to implement outcome (o) . As suggested in Section 4, we first look for the sets $S^* = S_A^* \times S_B^* \times S_C^* \subseteq S^\infty = S$ that induce (o) and satisfy Self-Enforceability and Self-Justifiability. Ann's Self-Enforceability requires Bob not to play d and Cleo not to play u . Then, Bob's Self-Justifiability requires that Cleo may play t , and Cleo's Self-Justifiability requires that Bob may play e in a subgame he allows. Hence, calling S_B^w and S_B^e the binary sets of plans of Bob where the last move is w and e respectively, the required sets S^* coincide with those that

satisfy

$$S_A^* = \{o\}, \quad S_B^* \subseteq S_B^w \cup S_B^e, \quad S_B^e \cap S_B^* \neq \emptyset, \quad S_C^* \subseteq \{t.a, b.a\}, \quad t.a \in S_C^*.$$

Does any of these sets satisfy Forward Induction? No. Under belief in S_C^* , it is irrational for Bob to play $d.l$, because both plans in S_B^e guarantee a higher payoff. Yet, it is rational to play $d.c$, because $t.a \in S_C^*$. Therefore, Forward Induction requires Cleo to play b and not t , a contradiction. Thus, there is no SES that implements (o) .

So, we look for a tight agreement e where e^0 satisfies the conditions above and alternative plans of Ann and Bob, e_A^1 and e_B^1 , are introduced to reach all histories (for T2) and restrict their behavior after deviations to i and d . First, observe that we need $e_C^0 = \{t.a\}$. If $b.a \in e_C^0$, then, regardless of e_A^1 , we have $d.r \in \rho_B(\Delta_B^e) \cap S_B((i, d))$, but $d.l \notin \rho_B(\Delta_B^e)$. So, for Bob, T3 imposes $d.l \notin e_B^1 \cap S_B((i, d)) \neq \emptyset$, but then $t.a \notin \rho_C(\Delta_C^e)$, a violation of T3. Still, without restrictions on e_A^1 , we have $d.c \in \rho_B(\Delta_B^e) \cap S_B((i, d))$, so again $d.l \notin e_B^1 \cap S_B((i, d)) \neq \emptyset$ and $t.a \notin \rho_C(\Delta_C^e)$. Hence, we must obtain $d.c \notin \rho_B(\Delta_B^e)$. So, we must impose $i.s.s \notin e_A^1$. If Ann guarantees to play n in a specific subgame, then we have $\rho_i(\Delta_B^e) \subseteq S_B^w$; hence, T3 imposes $e_B^0 \subseteq S_B^w$, a contradiction of the conditions on e_B^0 . So, the only remaining option is $e_A^1 = \{i.n.n, i.n.s, i.s.n\}$. Then, on the one hand there is $\mu_B \in \Delta_B^e$ with $\mu_B(i.n.s|i) = \mu_B(i.s.n|i) = 1/2$ and $\rho_B(\mu_B) = S_B^e$; on the other hand, for every $\mu_B \in \Delta_B^e$, there is $s_B \in \rho_B(\mu_B) \cap (S_B^w \cup S_B^e)$ that gives to Bob an expected payoff of at least 6.5, so $d.c \notin \rho_B(\Delta_B^e) \cap S_B((i, d)) = \emptyset$.

$$\begin{aligned} e_A^0 &= \{o\}, & e_B^0 &= S_B^w \cup S_B^e, & e_C^0 &= \{t.a\}; \\ e_A^1 &= \{i.n.n, i.n.s, i.s.n\}, & e_B^1 &= \{d.l, d.c, d.r\}. \end{aligned}$$

The vagueness of Ann about in which subgame she is going to play n is a kind of agreement incompleteness that, like here, can be necessary to implement an outcome. It can be interpreted as Ann doing the following speech: “I guarantee that I will be prepared to play n in at least one contingency, but I

cannot guarantee that I will be prepared to play n in both.”

This kind of strategic uncertainty also arises naturally from strategic reasoning. The example on page 50 in Battigalli [6] (provided by Gul and Reny) shows that already the set of justifiable plans of a player is not a Cartesian product of sets of actions at different information sets. This is the reason why (selective) rationalizability is defined as an elimination procedure of plans and not of actions at different information sets, and agreements are defined in terms of plans as well.

8.3 Proofs for Sections 5 and 6

For any $h \in \overline{H} \setminus \{h^0\}$, let $p(h) \in H$ be the immediate predecessor of h .

Proof of Proposition 8

Fix an outcome set $P \subseteq Z$ that is implementable under priority to rationality. Then, by Theorem 3, P is implemented by an agreement $e = (e_i)_{i \in I}$ which is tight under priority to rationality. The proof is constructive. Let M be the smallest m such that $S^m = S^\infty$ (it exists by finiteness of the game). For each $i \in I$ and $n = 0, \dots, k_i$, let

$$\bar{e}_i^n := \{s_i \in S_i^\infty : \exists s'_i \in e_i^n, \forall h \in H(s'_i) \cap H(S^\infty), s'_i(h) = s_i(h)\};$$

for each $n = k_i + 1, \dots, k_i + M + 1$, let $\bar{e}_i^n = S_i^{k_i + M + 1 - n}$. Assume without loss of generality that $\bar{e}_i^n \subsetneq \bar{e}_i^{n+1}$ for each $n = 0, \dots, k_i + M$,⁴⁷ so that $\bar{e} = (\bar{e}_i)_{i \in I}$ is an agreement. I am going to show that \bar{e} is tight under priority to the agreement. Indicate with T1^a, T2^a and T3^a the conditions of tightness under priority to the agreement (i.e., with S_i in place of S_i^∞).

First, I show that \bar{e} satisfies T1^a (which is identical to T1). Fix $i \in I$ and μ_i that strongly believes \bar{e}_{-i}^0 . For each $j \in I$ and $s_j \in e_j^0$, there is $s'_j \in \bar{e}_j^0$ such that $s'_j(h) = s_j(h)$ for all $h \in H(S^\infty) \cap H(s'_j)$, and vice versa. Hence, (i) $\zeta(S_i \times \bar{e}_{-i}^0) \cap \zeta(S^\infty) = \zeta(S_i \times e_{-i}^0) \cap \zeta(S^\infty)$, and there exists μ'_i that strongly

⁴⁷If $\bar{e}_i^n = \bar{e}_i^{n+1}$ for some n , \bar{e}_i^{n+1} can simply be eliminated from the chain.

believes e_{-i}^0 such that (ii) $\mu_i(S_{-i}(z)|h) = \mu'_i(S_{-i}(z)|h)$ for all $h \in H(S^\infty)$ and $z \in \zeta(S^\infty)$. Note that $\zeta(\rho_i(\mu_i) \times \bar{e}_{-i}^0), \zeta(\rho_i(\mu'_i) \times \bar{e}_{-i}^0) \subseteq \zeta(S^\infty)$.⁴⁸ Then: by (ii), $\zeta(\rho_i(\mu_i) \times \bar{e}_{-i}^0) = \zeta(\rho_i(\mu'_i) \times \bar{e}_{-i}^0)$; by (i), $\zeta(\rho_i(\mu'_i) \times e_{-i}^0) = \zeta(\rho_i(\mu'_i) \times \bar{e}_{-i}^0)$. So, we obtain

$$\zeta(\rho_i(\mu_i) \times \bar{e}_{-i}^0) = \zeta(\rho_i(\mu'_i) \times \bar{e}_{-i}^0) = \zeta(\rho_i(\mu'_i) \times e_{-i}^0) \subseteq \zeta(e^0) = \zeta(\bar{e}^0),$$

where the inclusion holds by T1 and the last equality by construction.

Moreover, \bar{e} satisfies T2^a by $\bar{e}_i^{k_i+M+1} = S_i$. It remains to show that \bar{e} satisfies T3^a. I will show later that⁴⁹

$$\rho_i(\Delta_i^{\bar{e}}) \cap S_i = \{s_i \in S_i^\infty : \exists s'_i \in \rho_i(\Delta_i^e) \cap S_i^\infty, \forall h \in H(s'_i) \cap H(S^\infty), s'_i(h) = s_i(h)\}. \quad (2)$$

Now, fix $h \in H(\rho_i(\Delta_i^{\bar{e}}) \cap S_i)$. Suppose first that either $h = h^0$ or $p(h) \in H(S^\infty)$. Then, by (2), $h \in H(\rho_i(\Delta_i^e) \cap S_i^\infty)$. By T3, there is n such that $\emptyset \neq e_i^n \cap S_i(h) \subseteq \rho_i(\Delta_i^e) \cap S_i^\infty$. Then, by definition of \bar{e}_i^n , $\bar{e}_i^n \cap S_i(h) \neq \emptyset$, and for each $s_i \in \bar{e}_i^n \cap S_i(h) \subset S_i^\infty$, there is $s'_i \in e_i^n$ such that $s'_i(h) = s_i(h)$ for all $h \in H(s'_i) \cap H(S^\infty)$, thus $s'_i \in e_i^n \cap S_i(h)$. So, $s'_i \in \rho_i(\Delta_i^e) \cap S_i^\infty$. But then, by (2), $s_i \in \rho_i(\Delta_i^{\bar{e}}) \cap S_i$.

Suppose now that $p(h) \notin H(S^\infty)$. Fix the unique $h' \prec h$ such that $h' \notin H(S^\infty)$ but $p(h') \in H(S^\infty)$. As shown, there is n such that $\emptyset \neq \bar{e}_i^n \cap S_i(h') \subseteq \rho_i(\Delta_i^{\bar{e}}) \cap S_i$. So, it suffices to show that $\bar{e}_i^n \cap S_i(h) \neq \emptyset$. Fix $s_i \in \rho_i(\Delta_i^{\bar{e}}) \cap S_i \cap S_i(h)$ and $s'_i \in \bar{e}_i^n \cap S_i(h') \subseteq \rho_i(\Delta_i^{\bar{e}}) \cap S_i$. By (2), $s_i, s'_i \in S_i^\infty$. Fix μ_i, μ'_i that strongly believe $((S_j^q)_{j \neq i})_{q=0}^\infty$ such that $s_i \in \rho_i(\mu_i)$ and $s_i \in \rho_i(\mu'_i)$. Since $h' \notin H(S^\infty)$, $p(h') \in H(S^\infty)$, and $h' \in H(S_i^\infty)$, we have $h' \notin H(S_{-i}^\infty)$, so $\mu'_i(S_{-i}(h')|p(h)) = 0$. Then, I can construct μ''_i that strongly believes $((S_j^q)_{j \neq i})_{q=0}^\infty$ such that $\mu''_i(\cdot|\tilde{h}) = \mu'_i(\cdot|\tilde{h})$ for each $\tilde{h} \not\geq h'$ and $\mu''_i(\cdot|\tilde{h}) = \mu_i(\cdot|\tilde{h})$ for each $\tilde{h} \succeq h'$. So, there is $s''_i \in \rho_i(\mu''_i) \subseteq S_i^\infty$ such that $s''_i(\tilde{h}) = s'_i(\tilde{h})$ for each $\tilde{h} \not\geq h'$ with $\tilde{h} \in H(s_i)$ and $s''_i(\tilde{h}) = s_i(\tilde{h})$ for each $\tilde{h} \succeq h'$ with $\tilde{h} \in H(s'_i)$. Hence, $s''_i \in \bar{e}_i^n$

⁴⁸To see this, fix $\tilde{\mu}_i$ that strongly believes $((S_j^q)_{j \neq i})_{q=0}^\infty$ with $\tilde{\mu}_i(\cdot|h) = \mu_i(\cdot|h)$ for all $h \in H(\bar{e}_{-i}^0)$. So, $\zeta(\rho_i(\tilde{\mu}_i) \times \bar{e}_{-i}^0) = \zeta(\rho_i(\mu_i) \times \bar{e}_{-i}^0)$. By $\rho_i(\tilde{\mu}_i) \subseteq S_i^\infty$ and $\bar{e}_{-i}^0 \subseteq S_{-i}^\infty$, $\zeta(\rho_i(\tilde{\mu}_i) \times \bar{e}_{-i}^0) \subseteq \zeta(S^\infty)$. Hence, $\zeta(\rho_i(\mu_i) \times \bar{e}_{-i}^0) \subseteq \zeta(S^\infty)$ as well.

⁴⁹Of course, the intersection with S_i is superfluous here. It will be substituted with S_i^∞ in the next proof.

and $s'_i \in S_i(h)$.

Finally I prove (2). First I prove “ \subseteq ”. For each $j \neq i$, $n = 0, \dots, k_j$, and $s_j \in \bar{e}_j^n$, there is $s'_j \in e_j^n$ such that $s'_j(h) = s_j(h)$ for all $h \in H(S^\infty) \cap H(s'_j)$. Moreover, by T2, $H(e_j^{k_j}) \supseteq H(S_j^\infty) \supseteq H(S^\infty)$. Hence, for each $\mu_i \in \Delta_i^{\bar{e}}$ (which strongly believes $((S_j^q)_{j \neq i})_{q=0}^\infty$ by construction of \bar{e}) I can construct $\mu'_i \in \Delta_i^{\bar{e}}$ that strongly believes $((S_j^q)_{j \neq i})_{q=0}^\infty$ (which is possible because $e_j^{k_j} \subseteq S_j^\infty$) such that $\mu'_i(S_{-i}(z)|h) = \mu_i(S_{-i}(z)|h)$ for all $h \in H(S^\infty)$ and $z \in \zeta(S^\infty)$. Then, for every $s_i \in \rho_i(\mu_i) \cap S_i \subset S_i^\infty$, there is $s'_i \in \rho_i(\mu'_i) \subset S_i^\infty$ such that $s'_i(h) = s_i(h)$ for all $h \in H(s'_i) \cap H(S^\infty)$.

Now I prove “ \supseteq ”. Fix $s_i \in S_i^\infty$, $\mu'_i \in \Delta_i^{\bar{e}}$, and $s'_i \in \rho_i(\mu'_i) \cap S_i^\infty$ with $s'_i(h) = s_i(h)$ for all $h \in H(s'_i) \cap H(S^\infty)$. Fix μ''_i that strongly believes $((S_j^q)_{j \neq i})_{q=0}^\infty$ such that $s_i \in \rho_i(\mu''_i)$. For each $j \neq i$, $n = 0, \dots, k_j$, $s_j \in e_j^n \subseteq S_j^\infty$, $h \in H(s_j) \setminus H(S^\infty)$ with $p(h) \in H(S^\infty)$, and $s'_j \in S_j^\infty \cap S_j(h)$, fix μ_j, μ'_j that strongly believe $((S_k^q)_{k \neq j})_{q=0}^\infty$ such that $s_j \in \rho_j(\mu_j)$ and $s'_j \in \rho_j(\mu'_j)$. By $\mu_j(S_{-j}(h)|p(h)) = 0$, I can construct μ''_j that strongly believes $((S_k^q)_{k \neq j})_{q=0}^\infty$ such that $\mu''_j(\cdot|\tilde{h}) = \mu_j(\cdot|\tilde{h})$ for each $\tilde{h} \not\geq h$, and $\mu''_j(\cdot|\tilde{h}) = \mu'_j(\cdot|\tilde{h})$ for each $\tilde{h} \succeq h$. So, there is $s''_j \in \rho_j(\mu''_j) \subseteq S_j^\infty$ such that $s''_j(\tilde{h}) = s_j(\tilde{h})$ for each $\tilde{h} \not\geq h$ with $\tilde{h} \in H(s_j)$ and $s''_j(\tilde{h}) = s'_j(\tilde{h})$ for each $\tilde{h} \succeq h$ with $\tilde{h} \in H(s'_j)$. Hence, $s''_j \in \bar{e}_j^n$. With all such s''_j 's, I can construct $\mu_i \in \Delta_i^{\bar{e}}$ such that $\mu_i(S_{-i}(z)|h) = \mu'_i(S_{-i}(z)|h)$ for all $h \in H(S^\infty)$ and $z \in \zeta(S^\infty)$, and $\mu_i(S_{-i}(z)|h) = \mu''_i(S_{-i}(z)|h)$ for all $h' \in H(S_i^\infty) \setminus H(S^\infty)$ with $p(h') \in H(S^\infty)$, $h \succeq h'$, and $z \succeq h$. Hence, $s_i \in \rho_i(\Delta_i^{\bar{e}})$. ■

Proof of Proposition 9. For each $P \subseteq Z$ which is implementable under priority to the path, a tight agreement \bar{e} that implements P under priority to rationality can be constructed exactly like in the proof of Proposition 8, substituting T1, T2, T3 with T1^p, T2^p, T3^p (the requirements of tightness under priority to the path, that is, with S_z^∞ in place of S^∞), T1^a, T2^a, T3^a with T1, T2, T3, S with S^∞ , S^∞ with S_z^∞ , and $((S_i^q)_{i \in I})_{q=0}^\infty$ with $((S_{i,z}^q)_{i \in I})_{q=0}^\infty$. ■

Proposition 10 *Let $\bar{z} = (\bar{a}^1, \dots, \bar{a}^T)$ be a path that can be upset by a convincing deviation. The path agreement on \bar{z} is not credible.*

Proof. Fix a two-players (i and j) static game G with action sets A_i and A_j and payoff function $v_k : A_i \times A_j \rightarrow \mathbb{R}$, $k = i, j$. Let b^k and c^k be the first- and second-ranked stage-outcomes of G for player $k = i, j$. A path $\bar{z} = (\bar{a}^1, \dots, \bar{a}^T)$ of Nash equilibria of the T -fold repetition of G can be upset by a convincing deviation if there exist $\tau \in \{1, \dots, T-1\}$ and $\hat{a}_i \neq \bar{a}_i^\tau$ such that, letting $\bar{T} := T - \tau$,

$$v_i(\hat{a}_i, \bar{a}_j^\tau) + v_i(c^i) + (\bar{T} - 1)v_i(b^i) < \sum_{t=\tau}^T v_i(\bar{a}^t) < v_i(\hat{a}_i, \bar{a}_j^\tau) + \bar{T}v_i(b^i); \quad (\text{I})$$

$$\bar{T}v_j(b^j) > \max_{a_j \in A_j \setminus \{b_j^i\}} v_j(b_j^i, a_j) + (\bar{T} - 1)v_j(b^j). \quad (\text{J})$$

Condition I says that player i benefits from a unilateral deviation at τ only if followed by her preferred subpath.⁵⁰ Condition J says that player j cannot benefit from a unilateral deviation from that subpath even if followed by her preferred subpath.⁵¹

Now I prove the proposition. Let $e_i = (S_i(\bar{z}))$ and $e_j = (S_j(\bar{z}))$. Let $\hat{h} := (\bar{a}^1, \dots, (\hat{a}_i, \bar{a}_j^\tau))$ and $z := (\bar{a}^1, \dots, (\hat{a}_i, \bar{a}_j^\tau), b^i, \dots, b^i)$. Suppose that $S_e^1(\bar{z}) \neq \emptyset$, otherwise $S_e^2 = \emptyset$. Then, for each $k = i, j$, there exists $\bar{\mu}_k$ that strongly believes S_{-k}^∞ and $S_{-k}(\bar{z})$ such that $\rho_k(\bar{\mu}_k) \cap S_k(\bar{z}) \neq \emptyset$.

Fix $n \in \mathbb{N}$ and suppose that $S_i^{n-1}(z) \neq \emptyset$. Fix $s_j \in S_j$ with $\bar{\mu}_i(s_j|h^0) \neq 0$. Since $\bar{\mu}_i$ strongly believes S_j^∞ and $S_j(\bar{z})$, $s_j \in S_j^\infty(\bar{z})$. Fix μ_j that strongly believes $(S_i^q)_{q=0}^\infty$ with $s_j \in \rho_j(\mu_j)$. Since $\bar{\mu}_j$ strongly believes $S_i(\bar{z})$, for each $h \notin H(S_i(\bar{z}))$ with $p(h) \prec \bar{z}$, $\bar{\mu}_j(S_i(h)|p(h)) = 0$. Thus, there exists μ'_j that strongly believes $(S_i^q)_{q=0}^{n-1}$ such that (i) $\mu'_j(\cdot|h^0) = \bar{\mu}_j(\cdot|h^0)$, (ii) $\mu'_j(S_i(z)|\hat{h}) = 1$, and (iii) $\mu'_j(\cdot|h) = \mu_j(\cdot|h)$ for all $h \in H(S_j(\bar{z}))$ with $h \neq \bar{z}$ and $h \not\prec \hat{h}$. Then, there exists $s'_j \in \rho_j(\mu'_j) \subseteq S_j^n$ such that: by $\rho_j(\bar{\mu}_j) \cap S_j(\bar{z}) \neq \emptyset$, $\bar{\mu}_j(S_i(z)|h^0) = 1$, and (i), $s'_j \in S_j(\bar{z}) \subseteq S_j(\hat{h})$; by (ii) and (J), $s'_j \in S_j(z)$; by (iii) and

⁵⁰In the example of Section 5, $i = \text{Ann}$, $j = \text{Bob}$, $(\bar{a}^1, \bar{a}^2) = ((FR, W), (W, FR))$, $b^i = (FR, W)$, $c^i = (W, W)$, $\tau = 1$, $\hat{a}_i = W$, thus $\bar{T} - 1 = 0$. Formally, the first inequality in (I) is not satisfied (equality holds), but this is immaterial because b^i and c^i entail the same action for Bob, against which the best reply of Ann induces b^i .

⁵¹This implies that i 's preferred stage-outcome is Nash, reason why Osborne (1991) refers to coordination games.

$s_j, s'_j \in S_j(\bar{z})$, $s'_j(h) = s_j(h)$ for all $h \in H(S_j(\bar{z}))$ with $h \not\prec \hat{h}$. With these s'_j 's, I can construct μ_i that strongly believes $(S_j^q)_{q=0}^n$ such that $\mu_i(S_j(z)|h^0) = 1$, and $\mu_i(S_j(\tilde{z})|h^0) = \bar{\mu}_i(S_j(\tilde{z})|h^0)$ for all $\tilde{z} \not\prec \hat{h}$. Thus, by $\rho_i(\bar{\mu}_i) \cap S_i(\bar{z}) \neq \emptyset$, $\bar{\mu}_i(S_j(\bar{z})|h^0) = 1$, and (I), $\emptyset \neq \rho_i(\mu_i) \cap S_i(z) \subseteq S_i^{n+1}(z)$. So, by induction, there exists μ_i that strongly believes $(S_j^q)_{q=0}^\infty$ and $S_j(\bar{z})$ such that $\emptyset \neq \rho_i(\mu_i) \cap S_i(z) \subseteq S_{i,e}^1(z)$. On the other hand, for every μ_i that strongly believes $S_j(\bar{z})$, by (I) $\rho_i(\mu_i) \cap S_i(\hat{h}) \subseteq S_i(z)$, so $S_{i,e}^1(\hat{h}) \subseteq S_i(z)$. The two things combined imply that for every μ_j that strongly believes $S_{i,e}^1$ and $S_i(\bar{z})$, $\mu_j(S_i(z)|\hat{h}) = 1$. So, by (J), $S_{j,e}^2(\hat{h}) \subseteq S_j(z)$. Since $S_j(\bar{z}) \subseteq S_j(\hat{h})$, for every μ_i that strongly believes $S_{j,e}^2$ and $S_j(\bar{z})$, $\mu_i(S_j(z)|h^0) = 1$, so by (I) $\rho_i(\mu_i)(\bar{z}) = \emptyset$. Hence $S_{i,e}^3(\bar{z}) = \emptyset$. So, $S_{j,e}^4 = \emptyset$. ■

8.4 On the definition of Selective Rationalizability.

Consider the following, alternative definition of Selective Rationalizability.

Definition 16 Let $((S_i^m)_{i \in I})_{m=0}^\infty$ denote the Rationalizability procedure. Consider the following procedure.

(Step 0) For each $i \in I$, let $\hat{S}_{i,e}^0 = S_i^\infty$.

(Step $n > 0$) For each $i \in I$ and $s_i \in S_i$, let $s_i \in \hat{S}_{i,e}^n$ if and only if there is $\mu_i \in \Delta_i^e$ such that:

S1 $s_i \in \rho_i(\mu_i)$;

S2 μ_i strongly believes $\hat{S}_{j,e}^q$ for all $j \neq i$ and $q < n$;

S3 μ_i strongly believes \hat{S}_j^q for all $j \neq i$ and $q \in \mathbb{N}$.

Finally, let $\hat{S}_{i,e}^\infty = \bigcap_{n \geq 0} \hat{S}_{i,e}^n$. The profiles in \hat{S}_e^∞ are called selectively-rationalizable.

This is the definition of Selective Rationalizability provided and characterized epistemically in [18]. It differs from the definition used in this paper

because of requirement S3 in place of the requirement that $s_i \in S_i^\infty$. Here I argue that the two definitions are equivalent for the analysis of agreements.

The two definitions are equivalent for the *same* agreement whenever the agreed-upon plans are chosen only according to what they prescribe at the rationalizable histories ($H(S^\infty)$).

Proposition 11 *Fix an agreement $e = (e_i)_{i \in I}$ such that, for each $i \in I$, $n = 0, \dots, k_i$, $s_i \in e_i^n$, and $s'_i \in S_i^\infty$, if $s'_i(h) = s_i(h)$ for all $h \in H(S^\infty) \cap H(s'_i)$, then $s'_i \in e_i^n$. Then, $\widehat{S}_e^\infty = S_e^\infty$.*

Proof. By induction.

Induction hypothesis: for each $m \leq n$, $\widehat{S}_e^m = S_e^m$; moreover, unless $\widehat{S}_e^{n+1} = S_e^{n+1} = \emptyset$, for each $i \in I$ and $\bar{h} \notin H(S^\infty)$ with $p(\bar{h}) \in H(S^\infty)$, there exists a map $\eta_{i,n}^{\bar{h}} : S_i(\bar{h}) \rightarrow S_i(\bar{h})$ such that:

a) for each $\bar{s}_i \in S_i(\bar{h}) \setminus S_i^\infty(\bar{h})$, $\eta_{i,n}^{\bar{h}}(s_i) = \bar{s}_i$;

b) for each $\bar{s}_i \in S_i^\infty(\bar{h})$,

(i) $\eta_{i,n}^{\bar{h}}(\bar{s}_i)(h) = \bar{s}_i(h)$ for all $h \in H(\bar{s}_i)$ with $h \succeq \bar{h}$,

(ii) $\eta_{i,n}^{\bar{h}}(\bar{s}_i) \in S_{i,e}^m$ for all $m \leq n$ with $S_{i,e}^m(\bar{h}) \neq \emptyset$,

(iii) if $e_i^q \cap S_i(\bar{h}) \neq \emptyset$ for some $q = 0, \dots, k_i$, $\eta_{i,n}^{\bar{h}}(\bar{s}_i) \in e_i^q$.

Basis step: $S_e^0 = \widehat{S}_e^0 = S^\infty$, and the required maps exist by property of e (in particular, (iii) can always be satisfied).

Inductive step. For $\widehat{S}_e^{m+1} = S_e^{m+1}$, since by the induction hypothesis $\widehat{S}_e^m = S_e^m$ for each $m \leq n$, it suffices to show that for every $i \in I$ and $s_i \in S_{i,e}^{m+1}$, there is $\widehat{\mu}_i \in \Delta_i^e$ that strongly believes $((S_{j,e}^m)_{j \neq i})_{m=0}^n$ and $((S_j^m)_{j \neq i})_{m=0}^\infty$ such that $s_i \in \rho_i(\widehat{\mu}_i)$. So, fix $\mu_i \in \Delta_i^e$ that strongly believes $((S_{j,e}^m)_{j \neq i})_{m=0}^n$ and μ'_i that strongly believes $((S_j^m)_{j \neq i})_{m=0}^\infty$ such that $s_i \in \rho_i(\mu_i) \cap \rho_i(\mu'_i)$. By the induction hypothesis, I can construct $\widehat{\mu}_i$ such that $\widehat{\mu}_i(\cdot|h) = \mu_i(\cdot|h)$ for all $h \in H(S_{-i}^\infty)$ and $\widehat{\mu}_i(s_{-i}|h) = \mu'_i((\times_{j \neq i} \eta_{j,n}^{\bar{h}})^{-1}(s_{-i})|h)$ for all $\bar{h} \notin H(S_{-i}^\infty)$ with $p(\bar{h}) \in H(S^\infty)$, $h \succeq \bar{h}$, and $s_{-i} \in \times_{j \neq i} \eta_{j,n}^{\bar{h}}(S_{-i}(\bar{h}))$. By (iii), $\widehat{\mu}_i \in \Delta_i^e$. By

(ii), $\widehat{\mu}_i$ strongly believes $((S_{j,e}^m)_{j \neq i})_{m=0}^n$ and, by (a), also $((S_j^m)_{j \neq i})_{m=0}^\infty$. By (i), $s_i \in \rho_i(\widehat{\mu}_i)$.

Now fix $\bar{h} \notin H(S^\infty)$ with $p(\bar{h}) \in H(S^\infty)$. If $S_{i,e}^{n+1}(\bar{h}) = \emptyset$, let $\eta_{i,n+1}^{\bar{h}} = \eta_{i,n}^{\bar{h}}$. Else, we need to update $\eta_{i,n}^{\bar{h}}(\bar{s}_i)$ for each $\bar{s}_i \in S_i^\infty(\bar{h})$. Fix $\bar{\mu}_i$ that strongly believes $((S_j^m)_{j \neq i})_{m=0}^\infty$ such that $\bar{s}_i \in \rho_i(\bar{\mu}_i)$. Unless $\widehat{S}_e^{n+1} = S_e^{n+1} = \emptyset$, there exists $\widehat{s}_i \in S_{i,e}^{n+1}(\bar{h}) = \widehat{S}_{i,e}^{n+1}(\bar{h})$ with $\widehat{s}_i \in e_i^q$ for all $q = 0, \dots, k_i$ such that $e_i^q \cap S_i(\bar{h}) \neq \emptyset$, otherwise, for any $j \neq i$, there would not be any $\widehat{\mu}_j \in \Delta_j^e$ that strongly believes $S_{i,e}^{n+1} = \widehat{S}_{i,e}^{n+1}$. Fix $\widehat{\mu}_i \in \Delta_i^e$ that strongly believes $((S_{j,e}^m)_{j \neq i})_{m=0}^n$ such that $\widehat{s}_i \in \rho_i(\widehat{\mu}_i)$. Since $\widehat{\mu}_i$ strongly believes $S_{-i,e}^0 = S_{-i}^\infty$, by the induction hypothesis I can construct μ_i such that $\mu_i(\cdot|h) = \widehat{\mu}_i(\cdot|h)$ for all $h \not\geq \bar{h}$ and $\mu_i(s_{-i}|h) = \bar{\mu}_i((\times_{j \neq i} \eta_{j,n}^{\bar{h}})^{-1}(s_{-i})|h)$ for all $h \succeq \bar{h}$ and $s_{-i} \in \times_{j \neq i} \eta_{j,n}^{\bar{h}}(S_{-i}(\bar{h}))$. By (iii), $\mu_i \in \Delta_i^e$. By (ii), μ_i strongly believes $((S_{j,e}^m)_{j \neq i})_{m=0}^{n-1}$. By (i), there is $s_i \in \rho_i(\mu_i)$ such that $s_i(h) = \widehat{s}_i(h)$ for all $h \in H(s_i)$ with $h \not\geq \bar{h}$ (thus $s_i \in S_i(\bar{h})$) and $s_i(h) = \bar{s}_i(h)$ for all $h \in H(\bar{s}_i)$ with $h \succeq \bar{h}$. So, $\eta_{i,n+1}^{\bar{h}}(\bar{s}_i) = s_i$ satisfies (i). If $s_i \in S_i^\infty$, then $s_i \in S_{i,e}^{n+1}$, satisfying (ii), and by the property of e , $s_i \in e_i^m$ for every m such that $\widehat{s}_i \in e_i^m$, satisfying (iii). So, it only remains to show that $s_i \in S_i^\infty$. Since $\widehat{s}_i \in S_i^\infty$, there is also $\widehat{\mu}_i$ that strongly believes $((S_j^m)_{j \neq i})_{m=0}^\infty$ such that $\widehat{s}_i \in \rho_i(\widehat{\mu}_i)$. Thus, I can construct also μ_i that strongly believes $((S_j^m)_{j \neq i})_{m=0}^\infty$ such that $\mu_i(\cdot|h) = \widehat{\mu}_i(\cdot|h)$ for all $h \not\geq \bar{h}$ and $\mu_i(\cdot|h) = \bar{\mu}_i(\cdot|h)$ for all $h \succeq \bar{h}$, so clearly $s_i \in \rho_i(\mu_i) \subseteq S_i^\infty$. ■

The intuition is the following: under an agreement in this class, all rationalizable plans can always be justified at the non-rationalizable histories under both definitions, while the two definitions do not differ in terms of beliefs they allow at the rationalizable histories. This class of agreements suffices for the implementation of all implementable outcome sets, for the following reason. Restricting behavior at the non-rationalizable histories cannot have a direct effect on the induced paths, which are always rationalizable. It can only have an indirect effect via a player's beliefs by combining an co-player's agreed behavior at rationalizable and non-rationalizable histories in a particular way. But given that the behavior of the co-player before and after our player leaves the rationalizable histories can always be "disentangled" (because the co-player gets surprised by finding herself at the non-rationalizable histories and has to

come up with new beliefs), this indirect effect can also be obtained directly by only restricting her behavior at the rationalizable histories. This can be proven formally with the same arguments of the proof of Proposition 8.

Proposition 12 *Fix a self-enforcing agreement $e^* = (e_i^*)_{i \in I}$. Then, there exists an agreement $\bar{e} = (\bar{e}_i)_{i \in I}$ that satisfies the condition in Definition 11 such that $\zeta(S_{\bar{e}}^\infty) = \zeta(S_{e^*}^\infty)$.*

Proof. By Theorem 2, there exists a tight agreement $e = (e_i)_{i \in I}$ such that $\zeta(S_{e^*}^\infty) = \zeta(S_e^\infty)$. Define an agreement $\bar{e} = (\bar{e}_i)_{i \in I}$ by letting, for each $i \in I$ and $n = 0, \dots, k_i$,

$$\bar{e}_i^n = \{s_i \in S_i^\infty : \exists s'_i \in e_i^n, \forall h \in H(S^\infty) \cap H(s_i), s_i(h) = s'_i(h)\}.$$

I show that also \bar{e} is tight, so that $\zeta(S_{\bar{e}}^\infty) = \zeta(\bar{e}^0) = \zeta(e^0) = \zeta(S_e^\infty) = \zeta(S_{e^*}^\infty)$. T2 is obvious.

To see T1, follow the proof for Proposition 8 that \bar{e} satisfies T1^a (which coincides with T1).

To see T3, fix $i \in I$, $\mu_i \in \Delta_i^{\bar{e}}$, and $\bar{h} \in H(\rho_i(\mu_i) \cap S_i^\infty)$, and follow the proof for Proposition 8 that \bar{e} satisfies T3^a (which coincides with T3 by $\rho_i(\Delta_i^{\bar{e}}) = \rho_i(\Delta_i^{\bar{e}}) \cap S_i^\infty$ in that proof). ■

The same is true if self-enforceability is defined using Definition 16 (and it can be proven in the same way). Therefore, we have the following.

Corollary 5 *The implementable outcome sets under the two definitions of Selective Rationalizability coincide.*

This would not be true if agreements were allowed to feature non rationalizable plans. In this case, some e_i^n could reach a history $h \notin H(S_i^\infty)$ with some plan $s_i \notin S_i^m$, although $h \in H(S_i^m)$, so that no $s_j \in S_j^\infty \cap S_j(h) \neq \emptyset$ is compatible with the belief in e_i^n . This can imply the elimination of a move by j at a rationalizable history (possibly dominant within the rationalizable paths!) under Definition 7, whereas the agreement would not be credible at all under Definition 16.