

# Online Appendices to Renegotiation and Coordination with Private Values (Yuval Heller and Christoph Kuzmics, September 2019)

## B More on Properties of Strategies

In this appendix we demonstrate that no single one of the three properties (mutual-preference consistency, coordination, and binary communication) is implied by the other two. Clearly a strategy that has binary communication and is coordinated must be an equilibrium. No other combination of two of the three properties implies that a strategy is an equilibrium. Finally, we also define what it means for a strategy to be ordinal preference-revealing and show that this is implied by it being mutual-preference consistent.

Consider the following strategy  $\sigma = (\mu, \xi)$  in the game with communication with a message set  $M$  that contains at least three elements. Let  $m_L^1, m_L^2, m_R \in M$ , let

$$\mu(u) = \begin{cases} m_L^1 & \text{if } u \leq \frac{1}{4} \\ m_L^2 & \text{if } \frac{1}{4} < u \leq \frac{1}{2} \\ m_R & \text{if } u > \frac{1}{2} \end{cases},$$

and let  $\xi$  be such that  $\xi(m_L^i, m_L^j) = L$  for all  $i, j \in \{1, 2\}$ ,  $\xi(m_R, m_R) = R$ ,  $\xi(m_L^1, m_R) = \xi(m_R, m_L^1) = R$ , and  $\xi(m_L^2, m_R) = \xi(m_R, m_L^2) = L$ . This strategy is mutual-preference consistent and coordinated but does not have binary communication. It is not an equilibrium as types  $u \leq 1/4$  would strictly prefer to send message  $m_L^2$ .

Consider the following strategy  $\sigma = (\mu, \xi)$  in the game with communication with a message set  $M$  that contains at least two elements. Let  $m_L, m_R \in M$ , let

$$\mu(u) = \begin{cases} m_L & \text{if } u \leq \frac{1}{2} \\ m_R & \text{if } u > \frac{1}{2} \end{cases},$$

and let  $\xi$  be such that  $\xi(m_L, m_L) = L$ ,  $\xi(m_R, m_R) = R$ ,  $\xi(m_L, m_R) = 1/4$ , and  $\xi(m_R, m_L) = 3/4$ . This strategy is mutual-preference consistent, has binary communication, but is not coordinated. For almost all type distributions  $F$  this is not an equilibrium: it is only an equilibrium if  $F$  satisfies  $(F(3/4) - F(1/2)) / (1 - F(1/2)) = 1/4$  and  $F(1/4) / F(1/2) = 3/4$ .

Finally, for a strategy that has binary communication and is coordinated but not mutual-preference consistent, consider the equilibrium strategy that always leads to coordination on action  $L$  for any pair of messages.

Note also that an equilibrium does not necessarily satisfy any of the three properties. The interior

cutoff babbling equilibria mentioned in Section 3 are not coordinated and not mutual-preference consistent. The equilibrium of Example 1 does not have binary communication.

Call a strategy  $\sigma = (\mu, \xi) \in \Sigma$  *ordinal preference-revealing* if there exist two nonempty, disjoint, and exhaustive subsets of  $\text{supp}(\bar{\mu})$  denoted by  $M_L$  and  $M_R$  (i.e.,  $\text{supp}(\bar{\mu}) = M_L \dot{\cup} M_R$ ) such that if  $u < 1/2$ , then  $\mu_u(m) = 0$  for each  $m \in M_R$ , and if  $u > 1/2$ , then  $\mu_u(m) = 0$  for each  $m \in M_L$ . With an ordinal preference-revealing strategy a player indicates her ordinal preferences.

A strategy  $\sigma$  that is mutual-preference consistent is also ordinal preference-revealing (but not vice versa). Suppose not. Then there is a message  $m$  and two types  $u < 1/2$  and  $v > 1/2$  such that  $\mu_u(m), \mu_v(m) > 0$ . But then no matter how we specify  $\xi(m, m)$  we get either that if two types  $u$  meet they do not coordinate on  $L$  with probability one or if two types  $v$  meet they do not coordinate on  $R$  with probability one.

## C Non-binary Communication Equilibrium

We here formally present the example, which is discussed informally at the end of Section 3, of an equilibrium in which agents reveal some information about the cardinality of their preferences.

Suppose that  $|M| \geq 4$  and consider the game with two rounds of communication. Let  $e, m_L, m_R \in M$  and let  $\sigma = (\mu_1, \mu_2, \xi)$ , with  $\mu_1 : U \rightarrow \Delta(M)$ ,  $\mu_2 : M \times M \times U \rightarrow \Delta(M)$ , and  $\xi : (M \times M)^2 \rightarrow U$  be as follows. For the first round of messages there is an  $x \in [0, 1]$  such that

$$\mu_1(u) = \begin{cases} e & \text{if } u \leq x \text{ or } u > 1 - x \\ m_L & \text{if } x < u \leq \frac{1}{2} \\ m_R & \text{if } \frac{1}{2} < u \leq 1 - x. \end{cases}$$

The second round of messages depends on the outcome of the first round and is best described in the following table.

	$e$	$m_L$	$m_R$
$e$	$e$	$\mu^*$	$\mu^*$
$m_L$	$m_L$	$m_L$	$\mu_C$
$m_R$	$m_R$	$\mu_C$	$m_R$

Each entry in this table describes the message function that a player follows if her first-stage message is the one indicated on the left and her opponent's first-stage message is the one indicated at the top. The message function  $\mu^*$  (after for instance a message pair of  $(e, m_L)$ ) is just as in the definition of  $\sigma_L$  and  $\sigma_R$  (in Section 3). The message function  $\mu_C$  is as in the definition of  $\sigma_C$  with an appropriate relabeling of four messages in  $M$ .

The action function is also best given in table form as a function of the result of the first round of

communication (or the second round when so indicated).

	$e$	$m_L$	$m_R$
$e$	$\begin{cases} L & \text{if } u \leq \frac{1}{2} \\ R & \text{if } u > \frac{1}{2} \end{cases}$	$\begin{cases} L & \text{if } u \leq \frac{1}{2} \\ R & \text{if } u > \frac{1}{2} \end{cases}$	$\begin{cases} L & \text{if } u \leq \frac{1}{2} \\ R & \text{if } u > \frac{1}{2} \end{cases}$
$m_L$	$\begin{cases} L & \text{if } \mu_2 = (m_L, m_L) \\ R & \text{otherwise} \end{cases}$	$L$	$\xi_C$
$m_R$	$\begin{cases} R & \text{if } \mu_2 = (m_R, m_R) \\ L & \text{otherwise} \end{cases}$	$\xi_C$	$R$

Action function  $\xi_C$  is as defined for  $\sigma_C$  applied to the second round of communication only.

We can complete the description of this strategy by requiring that all other messages in  $M$  be treated exactly the same as one of the messages  $e, m_L, m_R$ .

**Proposition 3.** *Let  $F$  be a nondegenerate symmetric distribution around  $1/2$ , i.e.,  $F(x) = 1 - F(1 - x)$  for all  $x \in [0, 1]$ . Then there is an  $x \in (0, 1/2)$  such that the above-defined strategy of the coordination game, with two rounds of communication and with  $|M| \geq 4$ , is a Nash equilibrium.*

*Proof.* Consider the given strategy for some arbitrary  $x \in (0, 1/2)$ . First note that whenever messages lead the two players to coordinate their action then clearly both players are best responding to each other with their action choice. This is so in all cases except when both players send message  $e$  in the first round. In this case players choose action  $L$  if their type  $u \leq 1/2$  and  $R$  otherwise. Each player, in this case, faces an opponent that either has  $u \leq x$  or  $u > 1 - x$ . In the first case the opponent plays  $L$ ; in the second case,  $R$ . Given that  $F(x) = 1 - F(1 - x)$  the two cases are equally likely. Given this, players' action choices are indeed best responses.

Thus, all action choices are best responses to the given strategy. We now turn to message choices. Consider the second round. After moderate messages in the first round messages in the second round either do not affect play at all (after  $(m_L, m_L)$  and  $(m_R, m_R)$ ) or do so as in strategy  $\sigma_C$ . In either case players are indifferent between all messages. After message pairs  $(m_L, e)$  and  $(m_R, e)$  the sender of the moderate message has a strict incentive to send the same message again, while the sender of the extreme message has a strict incentive to send  $m_L$  if her type  $u < 1/2$  or to send  $m_R$  if her type  $u > 1/2$  as this induces coordination on her preferred outcome. After both players send message  $e$ , play will not depend on messages in the second round either and so both players will be indifferent between all messages. Thus, the behavior in the second round of communication is a best response to the given strategy.

Finally, we need to consider the incentives to send messages in the first round. It is obvious that any type  $u < 1/2$  prefers sending message  $m_L$  to sending message  $m_R$  and vice versa for types  $u > 1/2$ . The only remaining thing to show is that types  $u \leq x$  and  $u > 1 - x$  and only these weakly prefer

to send message  $e$  in the first round. Given the symmetry it is without loss of generality to consider a type  $u \leq 1/2$ . Given the strategy, sending message  $e$  yields to this type a payoff of

$$F(x)(1-u) + \left(F\left(\frac{1}{2}\right) - F(x)\right)(1-u) + \left(F(1-x) - F\left(\frac{1}{2}\right)\right)(1-u) + (1 - F(1-x))0,$$

where  $F(x)$  is the probability that her opponent is an extreme left type,  $(F(1/2) - F(x))$  is the probability that the opponent is a moderate left type,  $(F(1-x) - F(1/2))$  is the probability that the opponent is a moderate right type, in all of which cases both players eventually play  $L$ , and where  $(1 - F(1-x))$  is the probability that her opponent is an extreme right type, in which case the two players miscoordinate. Sending message  $m_L$  yields a payoff of

$$F(x)(1-u) + \left(F\left(\frac{1}{2}\right) - F(x)\right)(1-u) + \left(F(1-x) - F\left(\frac{1}{2}\right)\right)\frac{1}{2} + (1 - F(1-x))u.$$

A type  $u \leq 1/2$  therefore weakly prefers sending message  $e$  to sending message  $m_L$  if and only if

$$D_x(u) \equiv \left(F(1-x) - F\left(\frac{1}{2}\right)\right)(1-u) - \left(F(1-x) - F\left(\frac{1}{2}\right)\right)\frac{1}{2} - (1 - F(1-x))u \geq 0.$$

Using the symmetry of  $F$  we can rewrite  $D(x)$  as

$$D_x(u) = \left(\frac{1}{2} - F(x)\right)\left(\frac{1}{2} - u\right) - F(x)u.$$

Note that  $D_x(u)$  is linear and downward sloping in  $u$  if  $x \in (0, 1/2)$ . In an equilibrium we then must have that  $D_x(x) = 0$ . This implies

$$D_x(x) = \left(\frac{1}{2} - F(x)\right)\left(\frac{1}{2} - x\right) - F(x)x = 0,$$

or, equivalently,

$$D_x(x) = \frac{1}{4} - \frac{1}{2}F(x) - \frac{1}{2}x = 0.$$

As  $D_0(0) = 1/4 > 0$ ,  $D_{1/2}(1/2) = -1/4 < 0$ , and  $D_x(x)$  is a continuous function in  $x$ , there is an  $x \in (0, 1/2)$  such that  $D_x(x) = 0$ . For this  $x$  the given strategy is thus an equilibrium.  $\square$

## D Evolutionary Stability Analysis (Section 7)

In this appendix we analyze the stability properties of strategies  $\sigma_L$  and  $\sigma_R$  (the results can be extended to other renegotiation-proof equilibria, but we omit the details here for brevity). Specifically, we show that  $\sigma_L$  and  $\sigma_R$  satisfy three properties that imply robustness to various perturbations: neutral stability (à la [Maynard Smith and Price, 1973](#)); weakly dominant first-stage behavior (given the second-stage behavior); and neighborhood-invader second-stage behavior ([Apaloo, 1997](#); [Cressman, 2010](#)).

## D.1 Evolutionary/Neutral Stability

We say that two strategies are almost surely realization equivalent (abbr., equivalent) if they induce the same behavior in almost all types (regardless of the opponent's behavior).

**Definition 4.** A condition holds for *almost all types* if the set of types that satisfy the condition  $\tilde{U} \subseteq U$  has mass one (as measured by the distribution  $f$ ), i.e.,

$$\int_{u \in \tilde{U}} f(u) du = 1.$$

**Definition 5.** Strategies  $\sigma = (\mu, \xi)$  and  $\tilde{\sigma} = (\tilde{\mu}, \tilde{\xi})$  are *almost surely realization equivalent* (abbr., *equivalent*) if for almost all types  $u \in [0, 1]$ :  $\mu_u(m) = \tilde{\mu}_u(m)$  for every message  $m \in M$ , and  $F_m(\xi(m, m')) = F_m(\tilde{\xi}(m, m'))$  for all messages  $m, m' \in \text{supp}(\bar{\mu})$ .

If  $\sigma$  and  $\tilde{\sigma}$  are equivalent strategies we denote this by  $\sigma \approx \tilde{\sigma}$ . It is immediate that equivalent strategies always obtain the same ex-ante expected payoff.

An equilibrium strategy  $\sigma$  is neutrally (evolutionarily) stable if it achieves a weakly (strictly) higher ex-ante expected payoff against any (non-equivalent) best-reply strategy, relative to the payoff that the best-reply strategy achieves against itself.

**Definition 6** (adaptation of [Maynard Smith and Price, 1973](#)). Equilibrium strategy  $\sigma \in \mathcal{E}$  is neutrally stable if for any nonequivalent strategy  $\tilde{\sigma} \not\approx \sigma$ ,

$$\pi(\tilde{\sigma}, \sigma) = \pi(\sigma, \sigma) \Rightarrow \pi(\sigma, \tilde{\sigma}) \geq \pi(\tilde{\sigma}, \tilde{\sigma}).$$

It is evolutionarily stable if this last weak inequality is replaced by a strict one.

The refinement of neutral stability is arguably a necessary requirement for an equilibrium to be a stable convention in a population (see, e.g., [Banerjee and Weibull, 2000](#)). If  $\sigma$  is an equilibrium strategy that is not neutrally stable, then a few experimenting agents who play a best-reply strategy  $\sigma'$  can invade the population. These experimenting agents would fare the same against the incumbents, whereas they would outperform the incumbents when being matched with other experimenting agents. This implies that, on average, these experimenting agents would be more successful than the incumbents, and their frequency in the population would increase in any payoff-monotone learning dynamics. This, in turn, implies that the population will move away from  $\sigma$ .

Our first result shows that both  $\sigma_L$  and  $\sigma_R$  are neutrally stable, and, moreover, they are evolutionarily stable if there are two feasible messages.

**Proposition 4.** *Strategies  $\sigma_L$  and  $\sigma_R$  are neutrally stable strategies of the coordination game with communication  $\langle \Gamma, M \rangle$ . Moreover, if  $|M| = 2$ , then  $\sigma_L$  and  $\sigma_R$  are evolutionarily stable strategies.*

*Proof.* We here prove this result for  $\sigma_L$ . The proof for  $\sigma_R$  proceeds analogously and is omitted. In order to prove this result we first characterize all strategies  $\sigma$  that are best responses to  $\sigma_L$  and thus satisfy  $\pi(\sigma, \sigma_L) = \pi(\sigma_L, \sigma_L)$ .

Consider the message and action choice of a type  $u$  when her opponent uses strategy  $\sigma_L$ . If our type  $u$  chooses any message other than  $m_R$ , her opponent, sending message  $m_L$  or  $m_R$ , plays action  $L$  in either case. Our type  $u$  could then choose action  $L$  (as prescribed by  $\sigma_L$ ), which provides a payoff of  $1 - u$ , or action  $R$ , which provides a payoff of zero. Thus all types  $u < 1$  are strictly better off choosing action  $L$  in this case. Also note that sending any message other than  $m_L$  leads to a best possible payoff of  $1 - u$ .

If our type  $u$  chooses to send message  $m_R$  then there are two cases. First, suppose that her opponent sends message  $m_L$ , in which case her opponent chooses action  $L$ . Our type  $u$  could then choose action  $L$  (as prescribed by  $\sigma_L$ ), which provides a payoff of  $1 - u$ , or action  $R$  which provides a payoff of zero. Thus all types  $u < 1$  are strictly better off choosing action  $L$  in this case. Second, suppose that her opponent sends message  $m_R$ , in which case her opponent chooses action  $R$ . Our type  $u$  could then choose action  $R$  (as prescribed by  $\sigma_L$ ), which provides a payoff of  $u$ , or action  $L$  which provides a payoff of zero. Thus all types  $u > 0$  are strictly better off choosing action  $R$  in this case. Note that sending message  $m_R$  thus provides a best possible payoff of  $F(1/2)(1 - u) + (1 - F(1/2))u$ .

For our type  $u$  it is therefore a strict best response to send message  $m_R$  if  $F(1/2)(1 - u) + (1 - F(1/2))u > 1 - u$ , which is the case if and only if  $u > 1/2$  (as  $F(1/2) \in (0, 1)$  by assumption). For the case of  $|M| = 2$  we then have that any best response to  $\sigma_L$  is equivalent to  $\sigma_L$  as only three possible types have an alternative best reply: types  $u = 0$ ,  $u = 1/2$ , and  $u = 1$  (all zero measures under the assumption of an atomless distribution  $F$ ). Any strategy that differs from  $\sigma_L$  for a positive measure of types yields a strictly inferior payoff against  $\sigma_L$  than  $\sigma_L$  does. This proves that  $\sigma_L$  is evolutionarily stable in the case of  $|M| = 2$  simply by virtue of the fact that there are no nonequivalent strategies  $\sigma$  that satisfy  $\pi(\sigma, \sigma_L) = \pi(\sigma_L, \sigma_L)$ .

Suppose from now on that  $|M| > 2$ . Our type  $u$  then has a choice of messages  $m \neq m_R$  when  $u < 1/2$ . All of these messages can at best lead to a payoff of  $u$  (from playing  $L, L$ ) and, therefore, all of them are equally good when playing against  $\sigma_L$ . As her opponent never chooses any message other than  $m_L$  or  $m_R$  (each has probability zero under  $\sigma_L$ ) our type  $u < 1/2$  when best responding can play anything after any message pair  $(m, m')$  when both  $m, m' \notin \{m_L, m_R\}$ . Let  $\sigma$  be a strategy that satisfies all the previous restrictions, where all types  $u$  play a (in most cases unique and strict) best response against  $\sigma_L$ . Then we have that  $\pi_u(\sigma, \sigma) = \pi_u(\sigma_L, \sigma)$  for all  $u \geq 1/2$  (as the behavior under  $\sigma$  for types  $u \geq 1/2$  (except for possibly types  $u = 1/2$  and  $u = 1$ ) is identical to that under  $\sigma_L$ ),  $\pi_u(\sigma, \sigma) \leq 1 - u$  for all  $u < 1/2$  (since this type can achieve at best  $1 - u$ ), and  $\pi_u(\sigma_L, \sigma) = 1 - u$  (since  $\sigma$  similarly to  $\sigma_L$  prescribes playing  $L$  in this case). We thus have for any such  $\sigma$  by construction  $\pi(\sigma, \sigma_L) = \pi(\sigma_L, \sigma_L)$ . We also have that  $\pi(\sigma_L, \sigma) \geq \pi(\sigma, \sigma_L)$  for any such

$\sigma$ . Finally any best-reply strategy to strategy  $\sigma_L$  must be equivalent to some such strategy  $\sigma$  and thus  $\sigma_L$  is neutrally stable.  $\square$

## D.2 Message Function is Dominant

In this subsection we show that the behavior in the first stage induced by strategy  $\sigma_L$  (resp.,  $\sigma_R$ ), namely, the message function  $\mu^*$ , is a weakly dominant message function (and strictly dominant when  $|M| = 2$ ), when taking as given that the behavior in the second stage is according to the action function  $\xi_L$  (resp.,  $\xi_R$ ). This suggests that the behavior in the first stage that is induced by  $\sigma_L$  (resp., by  $\sigma_R$ ) is robust to any perturbation that keeps the behavior in the second stage unchanged. Specifically, it implies that even if the message function used by the population is perturbed in an arbitrary (and possibly significant) way, then the original function  $\mu^*$  yields a weakly higher payoff than any other message function, which suggests that the behavior in the first stage would converge back to play  $\mu^*$  under any payoff-monotone learning dynamics.

Proposition 5 shows that message function  $\mu^*$  yields a weakly higher payoff relative to any other message function when the action function is given by  $\xi_L$  or  $\xi_R$ . Moreover, the inequality is strict whenever the alternative message function is essentially different from  $\mu^*$  in the sense of inducing low types to play  $m_R$  or inducing high types to play  $m \neq m_R$ .

**Proposition 5.** *Let  $\mu', \mu''$  be arbitrary message functions. Then for,  $\xi \in \{\xi_L, \xi_R\}$  and for any type  $u \neq 1/2$ ,*

$$\pi_u((\mu^*, \xi), (\mu', \xi)) \geq \pi_u((\mu'', \xi), (\mu', \xi)).$$

*This inequality is strict for  $\xi = \xi_L$  if  $\mu'_u(m_R) > 0$  for a positive measure of types  $u$  and, either  $\mu''_u(m_R) > 0$  for a positive measure of types  $u < 1/2$ , or  $\mu''_u(m_R) < 1$  for a positive measure of types  $u > 1/2$ . This inequality is strict for  $\xi = \xi_R$  if  $\mu'_u(m_L) > 0$  for a positive measure of types  $u$  and, either  $\mu''_u(m_L) > 0$  for a positive measure of types  $u > 1/2$ , or  $\mu''_u(m_L) < 1$  for a positive measure of types  $u < 1/2$ .*

*Proof.* Consider the case of  $\xi = \xi_L$  (the other case is proven analogously). Let  $\gamma$  denote the probability that a player following strategy  $(\mu', \xi_L)$  sends message  $m_R$ . Then sending any message other than  $m_R$  when the partner sends  $(\mu', \xi_L)$  yields a payoff of  $1 - u$ , and sending message  $m_R$  yields a payoff of  $\gamma u + (1 - \gamma)(1 - u)$ . Thus any type  $u > 1/2$  weakly prefers sending message  $m_R$  to sending any other message (and strictly prefers this if  $\gamma > 0$ ), while any type  $u < 1/2$  weakly prefers sending any message other than  $m_R$  to sending message  $m_R$  (and strictly prefers this if  $\gamma < 1$ ). Thus, for every message function  $\mu'$  of the opponent,  $\mu^*$  optimizes the message choice for every type  $u$  universally.  $\square$

### D.3 Action Function is a Neighborhood Invader Strategy

In the induced second-stage game  $\Gamma(F_m, F_{m'})$  (the game played after players observe a pair of messages  $(m, m')$ ), players choose a cutoff to determine whether to play action  $L$  (if their type is below or equal to that cutoff) or action  $R$  (otherwise). Thus, players essentially choose a number (their cutoff) from the unit interval. Note also that this induced game is asymmetric whenever the message profile is asymmetric, i.e., when  $m \neq m'$ . As argued by [Eshel and Motro \(1981\)](#) and [Eshel \(1983\)](#), when the set of strategies is a continuum, a stable convention should be robust to perturbations that slightly change the strategy played by all agents in the population. [Cressman \(2010\)](#) formalizes this requirement using the notion of neighborhood invader strategy (adapting the related notion of [Apaloo, 1997](#)). In what follows we show that the action function induced by  $\sigma_L$  and  $\sigma_R$  is a neighborhood invader strategy in any induced game  $\Gamma(F_m, F_{m'})$  on the path of play.

Fix a message function  $\mu$  and a pair of messages  $m_1, m_2 \in \text{supp}(\bar{\mu})$ . We identify a strategy in the induced game  $\Gamma(F_{m_1}, F_{m_2})$  with thresholds  $x_i$ , which is interpreted as the maximal type for which player  $i \in \{1, 2\}$  plays  $L$ . We say that strategy  $x_i$  of player  $i$  is equivalent to  $x'_i$  (denoted by  $x_i \approx x'_i$ ) in the induced game  $\Gamma(F_{m_1}, F_{m_2})$ , if  $F_{m_i}(x_i) = F_{m_i}(x'_i)$ , which implies that both thresholds induce the same behavior with probability one. Let  $\pi^{m_1, m_2}(x_1, x_2)$  denote the expected payoff of an agent with a random type sampled from  $f_{m_1}$  who uses threshold  $x_1$  when facing a partner with a random unknown type sampled from  $f_{m_2}$  who uses threshold  $x_2$ .

A strategy profile  $(x_1, x_2)$  is a strict equilibrium in the induced game  $\Gamma(F_{m_1}, F_{m_2})$ , if any best reply to  $x_j$  is equivalent to  $x_i$ , i.e.,  $\pi^{m_1, m_2}(x'_1, x_2) \geq \pi^{m_1, m_2}(x_1, x_2) \Rightarrow x'_1 \approx x_1$ , and  $\pi^{m_2, m_1}(x_2, x'_1) \geq \pi^{m_2, m_1}(x_2, x_1) \Rightarrow x'_1 \approx x_1$ .

We say that the strict equilibrium  $(x_1, x_2)$  is a neighborhood invader strategy in the induced game  $\Gamma(F_{m_1}, F_{m_2})$  if the population converges to  $(x_1, x_2)$  from any nonequivalent nearby strategy profile  $(x'_1, x'_2)$  in two steps: (1) strategy  $x_i$  yields a strictly higher payoff against  $x_j$  relative to the payoff of  $x'_i$  against  $x_j$  (which implies convergence from  $(x'_i, x'_j)$  to  $(x_i, x'_j)$ ), and (2) due to  $(x_1, x_2)$  being a strict equilibrium, strategy  $x_j$  yields a strictly higher payoff against  $x_i$  relative to the payoff of  $x'_j$  against  $x_i$  (which implies the convergence from  $(x_i, x'_j)$  to  $(x_i, x_j)$ ).

**Definition 7** (Adaptation of [Cressman, 2010](#), Def. 5). Fix a message function  $\mu$  and a pair of messages  $m_1, m_2 \in \text{supp}(\bar{\mu})$ . A strict Nash equilibrium  $(x_1, x_2)$  is a *neighborhood invader strategy profile* in the induced game  $\Gamma(F_{m_1}, F_{m_2})$  if there exists  $\epsilon > 0$ , such that for each  $(x'_1, x'_2)$  satisfying  $x'_1 \not\approx x_1$ ,  $x'_2 \not\approx x_2$ ,  $|x'_1 - x_1| < \epsilon$  and  $|x'_2 - x_2| < \epsilon$ , then either  $\pi^{m_1, m_2}(x_1, x'_2) > \pi^{m_1, m_2}(x'_1, x'_2)$  or  $\pi^{m_2, m_1}(x_2, x'_1) > \pi^{m_2, m_1}(x'_2, x'_1)$ .

[Proposition 6](#) shows that the profile of action functions induced by  $\sigma_L$  (or, similarly, by  $\sigma_R$ ) is a neighborhood invader strategy in any induced game.

**Proposition 6.** *Let  $m_1, m_2 \in \text{supp}(\bar{\mu}^*)$ . Then strategy profiles  $(\xi_L(m_1, m_2), \xi_L(m_2, m_1))$  and*



$(\xi_R(m_1, m_2), \xi_R(m_2, m_1))$  are strict equilibria and neighborhood invader strategy profiles in the induced game  $\Gamma_{F_{m_1}, F_{m_2}}$ .

*Proof.* We present the proof for  $(\xi_L(m_1, m_2), \xi_L(m_2, m_1))$  (the proof for  $(\xi_R(m_1, m_2), \xi_R(m_2, m_1))$  is analogous). Observe that  $m_1, m_2 \in \text{supp}(\bar{\mu}^*)$  implies one of three cases:  $m_1 = m_2 = m_L$ ,  $m_1 = m_2 = m_R$ , or  $m_1 = m_R, m_2 = m_L$ . We analyze each case as follows.

Suppose first that  $m_1 = m_2 = m_L$ . This implies that  $\xi_L(m_1, m_2) = \xi_L(m_2, m_1) = 1$  and  $F_{m_1}(1/2) = F_{m_2}(1/2) = 1$ . Let  $\bar{x} < 1/2$  be sufficiently close to  $1/2$  such that  $F_{m_1}(\bar{x}), F_{m_2}(\bar{x}) > 1/2$ . Observe that  $\pi^{m_1, m_2}(1, x) > \pi^{m_1, m_2}(y, x)$  for any  $x > \bar{x}$  and any  $y \neq 1$ . This proves that  $(\xi_L(m_1, m_2), \xi_L(m_2, m_1))$  is a strict equilibrium and a neighborhood invader strategy profile.

Now suppose that  $m_1 = m_2 = m_R$ . This implies that  $\xi_L(m_1, m_2) = \xi_L(m_2, m_1) = 0$  and  $F_{m_1}(1/2) = F_{m_2}(1/2) = 0$ . Let  $\bar{x} > 1/2$  be sufficiently close to  $1/2$  such that  $F_{m_1}(\bar{x}), F_{m_2}(\bar{x}) < 1/2$ . Observe that  $\pi^{m_1, m_2}(0, x) > \pi^{m_1, m_2}(y, x)$  for any  $x < \bar{x}$  and any  $y \neq 0$ . This proves that  $(\xi_L(m_1, m_2), \xi_L(m_2, m_1))$  is a strict equilibrium and a neighborhood invader strategy profile.

Suppose finally that  $m_1 = m_R, m_2 = m_L$ . This implies that  $\xi_L(m_1, m_2) = \xi_L(m_2, m_1) = 1$ ,  $F_{m_1}(1/2) = 0$ , and  $F_{m_2}(1/2) = 1$ . Observe that  $\pi^{m_1, m_2}(1, 1) > \pi^{m_1, m_2}(x, 1)$  for any  $x \neq 1$  and  $\pi^{m_2, m_1}(1, 1) > \pi^{m_2, m_1}(x, 1)$  for any  $x \neq 1$ , which implies that  $(\xi_L(m_1, m_2), \xi_L(m_2, m_1))$  is a strict equilibrium. Let  $\bar{x} > 1/2$  be sufficiently close to  $1/2$  such that  $F_{m_1}(\bar{x}) < 1/2$ . Observe that  $\pi^{m_2, m_1}(1, x) > \pi^{m_1, m_2}(y, x)$  for any  $x < \bar{x}$  and any  $y \neq 1$ . This proves that strategy profile  $(\xi_L(m_1, m_2), \xi_L(m_2, m_1))$  is a neighborhood invader strategy profile.  $\square$

#### D.4 Remark on Evolutionary Robustness

[Oechssler and Riedel \(2002\)](#) present a strong notion of stability, called evolutionary robustness, that refines both evolutionary stability and the neighborhood invader strategy. An evolutionary robust strategy  $\sigma^*$  is required to be robust against small perturbation in the strategy played by the population, which may comprise both (1) a few experimenting agents who follow arbitrary strategies, and (2) many agents who follow strategies that are only slightly different than  $\sigma^*$ . Specifically, if  $\sigma$  is a distribution of strategies that is sufficiently close to  $\sigma^*$  (in the  $L_1$  norm induced by the weak topology), evolutionary robustness à la [Oechssler and Riedel](#) requires that  $\pi(\sigma^*, \sigma) > (\sigma, \sigma)$ .

One can show that  $\sigma_L$  and  $\sigma_R$  do not satisfy this condition (and, we conjecture, that no strategy can satisfy such a strong condition in our setup). However, we conjecture that one can show that  $\sigma_L$  and  $\sigma_R$  satisfy a somewhat weaker notion of evolutionary robustness: for each strategy distribution  $\sigma$  sufficiently close to  $\sigma_L$  ( $\sigma_R$ ), there exists a finite sequence of strategy distributions  $\sigma_1, \sigma_2, \dots, \sigma_k$ , such that  $\pi(\sigma_1, \sigma) \geq (\sigma, \sigma)$ ,  $\pi(\sigma_2, \sigma_1) \geq (\sigma_1, \sigma_1)$ ,  $\dots$ ,  $\pi(\sigma_k, \sigma_{k-1}) \geq (\sigma_{k-1}, \sigma_{k-1})$ ,

and  $\pi(\sigma_L, \sigma_1) \geq (\sigma_1, \sigma_1)$  (resp.,  $\pi(\sigma_R, \sigma_1) \geq (\sigma_1, \sigma_1)$ ), with strict inequalities if  $|M| = 2$  and  $\sigma$  is not realization equivalent to  $\sigma_L$  ( $\sigma_R$ ).

## E Analysis of Extensions (Section 8)

In this appendix we formally analyze the six extensions presented informally in Section 8. Formal proofs are postponed to Section E.7.

### E.1 Multiple Rounds of Communication

Consider a variant of the coordination game with communication in which players have a fixed and finite number  $T \geq 1$  of rounds of communication. In each such round of this communication phase players simultaneously send messages from the set of messages  $M$ . Players observe messages after each round and can, thus, condition their message choice and then their final action choice on the history of observed message pairs up to the point in time where they take their message or action decision. Renegotiation then possibly takes place once at the end of this communication phase but before the final action choices are made. Let  $\mathcal{M} = \bigcup_{t=0}^{T-1} (M \times M)^t$ , where  $(M \times M)^0 = \emptyset$ .

A (pure) *message protocol* is a function  $\mathbf{m} : \mathcal{M} \rightarrow M$  that describes the message sent by an agent as a deterministic function of the message profiles observed in the previous rounds of communication. Let  $\mathfrak{M}$  be the set of all message protocols. A strategy  $\sigma = (\mu, \xi)$  is a pair where  $\mu : U \rightarrow \Delta(\mathfrak{M})$  denotes the *message function*, prescribing a (possibly random) message protocol for each type, and  $\xi : (M \times M)^T \rightarrow U$  denotes the *action function* by means of describing the cutoff (the highest possible value of  $u$ ) for the two players to choose action  $L$  after observing the final message history. Renegotiation is modeled, as in the main text, as a possibility for the two players to play an equilibrium of a new game with another round of communication after all messages are sent, possibly using a different message set.

Next, we adapt the notion of binary communication to fit multiple rounds of communication. For any message protocol  $\mathbf{m} \in \mathfrak{M}$ , let  $\beta^\sigma(\mathbf{m})$  denote the expected probability of a player's opponent playing  $L$  conditional on the player following message protocol  $\mathbf{m} \in \mathfrak{M}$  and the opponent following strategy  $\sigma = (\mu, \xi) \in \Sigma$ . We say that strategy  $\sigma$  has *binary communication* if there are two numbers  $0 \leq \underline{\beta}^\sigma \leq \overline{\beta}^\sigma \leq 1$  such that for all message protocols  $\mathbf{m} \in \mathfrak{M}$  we have  $\beta^\sigma(\mathbf{m}) \in [\underline{\beta}^\sigma, \overline{\beta}^\sigma]$ , for all message protocols  $\mathbf{m} \in \mathfrak{M}$  such that there is a type  $u < 1/2$  with  $\mu_u(\mathbf{m}) > 0$  we have  $\beta^\sigma(\mathbf{m}) = \overline{\beta}^\sigma$ , and for all message protocols  $\mathbf{m} \in \mathfrak{M}$  such that there is a type  $u > 1/2$  with  $\mu_u(\mathbf{m}) > 0$  we have  $\beta^\sigma(\mathbf{m}) = \underline{\beta}^\sigma$ . That is, binary communication implies that players use just two kinds of message protocols: any message protocol used by types  $u < 1/2$  induces the consequence of maximizing the probability of the opponent to play  $L$ , and any message protocol used by types  $u > 1/2$  induces the

opposite consequence of maximizing the probability of the opponent to play  $R$ .

Theorem 1, together with Propositions 1 and 2, holds in this setting with minor adaptations to the proof (omitted for brevity). Thus, regardless of the length of the pre-play communication, agents can reveal only their preferred outcome (but not the strength of their preference), and, regardless of having access to additional rounds of communication, they cannot improve the ex-ante expected payoff relative to the payoff induced by a single round of communication with a binary message.

## E.2 MultiDimensional Sets of Types

In our model we made the simplifying assumption that miscoordination provides the same payoff (normalized to zero) to both players. This is not completely innocuous. In this section we explore which results are still true in this more general setting.

Consider the following multidimensional set of types. Let  $\hat{U}$ , a subset of  $\mathbb{R}^4$ , be the set of payoff matrices of binary coordination games, with  $u_{ab}$  being the payoff if a player chooses action  $a \in \{L, R\}$  while her opponent chooses action  $b \in \{L, R\}$ :

$$\hat{U} = \{(u_{LL}, u_{LR}, u_{RL}, u_{RR}) \mid u_{LL} > u_{RL} \text{ and } u_{RR} > u_{LR}\}.$$

Thus, all types strictly prefer coordination on the same action as the partner to miscoordination. Note that any affine transformation of all payoffs neither changes the player's incentives nor changes how she compares any two outcome distributions  $\in \Delta(\{L, R\})$ . We can thus subtract  $\min\{u_{RL}, u_{LR}\}$  from all payoffs and then divide all payoffs by some number such that the sum of the diagonal entries is equal to one. This leaves two parameters to describe a payoff vector in  $\hat{U}$ . This means that for our purposes the set  $\hat{U}$  is two-dimensional. Let  $\hat{\Gamma} = \hat{\Gamma}(G)$  denote the coordination game with the two-dimensional type space  $\hat{U}$ , endowed with an atomless CDF  $G$  over  $\hat{U}$  with a density  $g$ . Similarly, let  $(\hat{\Gamma}, M)$  be the corresponding game with communication.

Given a type  $u = (u_{LL}, u_{LR}, u_{RL}, u_{RR})$ , let  $\varphi_u \in [0, 1]$  denote type  $u$ 's *indifference threshold*, which is the probability of the opponent playing  $L$  that induces an agent of type  $u$  to be indifferent between the two actions:

$$\varphi_u = \frac{u_{RR} - u_{LR}}{u_{LL} - u_{RL} + u_{RR} - u_{LR}}.$$

Observe that an agent with indifference threshold  $\varphi_u$ , where  $\varphi_u$  is a number always between 0 and 1, prefers to play  $L$  ( $R$ ) if her partner plays  $L$  with probability larger (smaller) than  $\varphi_u$ . In other words, for a given probability of her partner playing  $L$ , a type  $u$  prefers to play  $L$  if and only if  $\varphi_u$  is less than that probability. Thus, the indifference threshold  $\varphi_u$  replaces what we denoted by  $u$  in the main model. In particular, in this setting we can also restrict attention to cutoff action functions. These are now applied to  $\varphi_u$  instead of to  $u$ . Thus, under a strategy  $\sigma = (\mu, \xi)$  a player

plays action  $L$  after observing a message pair  $(m, m')$  if and only if  $\varphi_u \leq \xi(m, m')$ .

Recall, that action  $L$  is *risk-dominant* (Harsanyi and Selten, 1988) if it is a best reply against the opponent randomizing equally over the two actions, i.e., if

$$\varphi_u \leq \frac{1}{2} \Leftrightarrow u_{LL} - u_{LR} \geq u_{RR} - u_{RL}.$$

The crucial assumption that we implicitly made in our (one-dimensional) main model is that for any type of player the action that she prefers to coordinate on is also risk-dominant for her.

**Definition 8.** An atomless probability distribution over the payoff space  $U$  with density function  $g : U \rightarrow \mathbb{R}$  satisfies *unambiguous coordination preferences* if for any  $u \in U$  with  $g(u) > 0$  we have

$$u_{LL} \geq u_{RR} \Leftrightarrow \varphi_u \leq \frac{1}{2}.$$

Under a probability distribution over types with unambiguous coordination preferences, every type in its support prefers coordinating on action  $L$  if and only if that type finds action  $L$  risk dominant. Under the assumption that the probability distribution satisfies unambiguous coordination preferences, Theorem 1 goes through unchanged if we set

$$F(\varphi) = \int_{\{u \in U : \varphi_u \leq \varphi\}} g(u) du$$

to be the implied distribution over the players' indifference threshold induced by density  $g$ . As in the baseline model, we assume that  $F(\varphi)$  has full support on the interval  $[0, 1]$ .

**Theorem 2** (Theorem 1 adapted to a multidimensional set of types). *A strategy  $\sigma$  of a game  $(\hat{\Gamma}, M)$  that satisfies unambiguous coordination preferences is a renegotiation-proof equilibrium strategy if and only if it is mutual-preference consistent, coordinated, and has binary communication.*

The proof is presented in Appendix E.7.1. The intuition is the same as in Theorem 1. The adaptation of Lemma 2, to the current setup relies on having unambiguous coordination preferences.

While we cannot say that the restriction of unambiguous coordination preferences is necessary for the result, we present an example that suggests that if this restriction is not satisfied, then equilibria with miscoordination may be renegotiation-proof.

**Example 2.** There are four possible preference types as follows:

$u_{L1}$	L	R	$u_{L2}$	L	R	$u_{R1}$	L	R	$u_{R2}$	L	R
L	2	$\frac{1}{2}$	L	2	-8	L	1	0	L	1	0
R	0	1	R	0	1	R	$\frac{1}{2}$	2	R	-8	2

Define the distribution of types  $F$  such that<sup>24</sup>  $P(u_{L1}) = P(u_{R1}) = 1/8$  and  $P(u_{L2}) = P(u_{R2}) = 3/8$ . Let  $M = \{m_L, m_R\}$  and let  $\sigma = (\mu, \xi)$  be such that  $\mu(u_{L1}) = \mu(u_{L2}) = m_L$  and  $\mu(u_{R1}) = \mu(u_{R2}) = m_R$  (making  $\sigma$  mutual-preference consistent), and  $\xi(m_L, m_L) = L$ ,  $\xi(m_R, m_R) = R$ ,  $\xi(u_{L1}, m_L, m_R) = \xi(u_{R2}, m_R, m_L) = L$ , and  $\xi(u_{L2}, m_L, m_R) = \xi(u_{R1}, m_R, m_L) = R$ .

It is straightforward to verify that  $\sigma$  is an equilibrium strategy.<sup>25</sup> Note that it is mutual-preference consistent and has binary communication, but it is not coordinated. We now show that  $\sigma$  is not Pareto-dominated by any coordinated equilibrium strategy in any induced post-communication game. To see this note that the non-coordinated equilibrium following  $(m_L, m_R)$  is not Pareto-dominated by any coordinated equilibrium  $\sigma_\alpha$  with left tendency of  $\alpha$  (with additional communication): under  $\sigma_\alpha$  the expected payoff of an agent, conditional on observing message pair  $(m_L, m_R)$ , is given by  $(3/4) \cdot 2 + (1/4) \cdot (1/2) = 1 + 5/8$  for types  $u_{L1}$  and  $u_{R1}$ , and equal to  $1/4$  for types  $u_{L2}$  and  $u_{R2}$ . In order to induce a payoff of at least  $1 + 5/8$  to type  $u_{L1}$  with any coordinated equilibrium strategy  $\sigma_\alpha$ , it must be that  $\alpha \geq 5/8$ , while in order to induce a payoff of at least  $1 + 5/8$  to type  $u_{L2}$ , it must be that  $\alpha \leq 3/8$ . Thus, there is no  $\alpha$  that satisfies both requirements.

Note, however, that any strategy that is coordinated and mutual-preference consistent and has binary communication is renegotiation-proof also in the general setting, and that only the other direction may fail without the assumption of unambiguous coordination preferences. There may be additional renegotiation-proof equilibria in the general setting. One can show that any such renegotiation-proof equilibrium strategy must satisfy mutual-preference consistency, but need not satisfy the other two properties (namely, coordination and binary communication).

### E.3 More Than Two Players

Consider a variant of the coordination game in which there are  $n \geq 2$  players who play a symmetric coordination game (with private values) with pre-play communication. The action set is  $\{L, R\}$  for every player and the payoff to player  $i$  is equal to  $u_i$  if every player chooses action  $R$ , equal to  $1 - u_i$  if every player chooses  $L$ , and equal to zero otherwise. The payoff to type  $u_i$  is independent and identically drawn from some given distribution  $F$  with support in the unit interval  $[0, 1]$ . Before players choose actions, they simultaneously send messages from a finite set of messages  $M$  and observe all these messages. Let  $\langle \Gamma_n, M \rangle$  denote this  $n$ -player coordination game with pre-play communication.

In this setting the appropriate version of Theorem 1 still holds.

<sup>24</sup>Note that this distribution is discrete, but could easily be modified to a nearby atomless distribution without changing the result.

<sup>25</sup>The expected payoffs are:  $\frac{1}{2} \cdot 2 + \frac{1}{2} \cdot (\frac{3}{4} \cdot 2 + \frac{1}{4} \cdot \frac{1}{2}) = 1 + \frac{13}{16}$  for types  $u_{L1}$  and  $u_{R1}$ , and  $\frac{1}{2} \cdot 2 + \frac{1}{2} \cdot \frac{1}{4} \cdot 1 = 1 + \frac{1}{8}$  for types  $u_{L2}$  and  $u_{R2}$ . This implies that no type wants to misreport her preferred outcome in round one. In particular, a misreporting type  $u_{L2}$  will get a payoff of  $\frac{1}{2} \cdot 1 + \frac{1}{2} \cdot \frac{3}{4} \cdot 1 = \frac{7}{8} < 1 + \frac{1}{8}$ .

**Theorem 3** (Theorem 1 adapted to more than two players). *A strategy  $\sigma$  of the  $n$ -player coordination game  $\langle \Gamma_n, M \rangle$  is a renegotiation-proof equilibrium strategy if and only if it is mutual-preference consistent, coordinated, and has binary communication.*

*Sketch of proof; for the formal proof see Appendix E.7.2.* The proof of the “only if” direction has to be adapted (the proof of the “if” direction remains, essentially, the same). In this setting it is not generally true that any play that involves miscoordination is post-communication Pareto-dominated by  $\sigma_L$ ,  $\sigma_R$ , or  $\sigma_C$ . The proof instead first establishes that miscoordination after all players send the *same* message must be Pareto-dominated by either  $\sigma_L$  or  $\sigma_R$  (Lemma 7). This is then used to show that a renegotiation-proof equilibrium strategy must be mutual-preference consistent (Lemma 8). Then one can show that a renegotiation-proof equilibrium strategy must be coordinated and must have binary communication (Lemma 9).  $\square$

Proposition 1 and Corollary 1 also hold in the multi-player setting: renegotiation-proof equilibrium strategies are interim (pre-communication) Pareto-undominated and are Pareto-improving relative to all symmetric equilibria of the game without communication. By contrast, Proposition 2 does not extend to this setting: with three players, for instance, for some distributions of values  $F$ , the strategy that determines the fallback option by majority vote (in the case of messages that indicate different preferred actions) is an ex-ante payoff improvement over a simple fallback norm of choosing, say, action  $L$  in every case of disagreement.

## E.4 Asymmetric Coordination Games

**Adapted Model** Consider a setup similar to our baseline model except that the distributions of the types of the two players’ positions differ: the type of player 1 is distributed according to  $F_1$  and the type of player 2 is distributed according to  $F_2$ . As in the baseline model, both distributions are assumed to be atomless with full support in  $[0, 1]$ . Let  $\langle \Gamma(F_1, F_2), M \rangle$  denote the asymmetric coordination game with communication (to ease notation, we assume that both players have the same set of messages at their disposal). Let  $\Sigma^i$  denote the set of all strategies of player  $i \in \{1, 2\}$ . We let  $i$  denote the index of one player and  $j$  denote the index of the opponent.

*Remark 6.* The game  $\langle \Gamma(F, F), M \rangle$  in which both players have the same distribution of types corresponds to a setup, in which the payoff-irrelevant position of player 1 or player 2 is identifiable, and the players can condition their play on their positions.

Given a strategy profile  $(\sigma_1, \sigma_2)$ , let  $\pi_u^i(\sigma_1, \sigma_2)$  denote the (interim) payoff of type  $u$  of player  $i \in \{1, 2\}$ , and let  $\pi^i(\sigma_1, \sigma_2) = \mathbb{E}_{u \sim F_i} [\pi_u^i(\sigma_1, \sigma_2)]$  denote the ex-ante payoff of player  $i \in \{1, 2\}$ . A strategy profile  $(\sigma_1, \sigma_2)$  is a (Bayes Nash) *equilibrium* if  $\pi_u^1(\sigma_1, \sigma_2) \geq \pi_u^1(\sigma'_1, \sigma_2)$  for each strategy

$\sigma'_1 \in \Sigma^1$  and for each type  $u$  of player 1, and  $\pi_u^2(\sigma_1, \sigma_2) \geq \pi_u^2(\sigma_1, \sigma'_2)$  for each strategy  $\sigma'_2 \in \Sigma^2$  and for each type  $u$  of player 2.

**Adapted Properties** We adapt the three key properties of Section 3 as follows. Let  $\mu_u^i(m_i)$  denote the probability, given message function  $\mu^i$ , that player  $i$  sends message  $m_i$  if she is of type  $u_i$ . Let  $\mu^i(m_i) = \mathbb{E}_{u \sim F_i}[\mu_u^i(m_i)]$  be the average (ex-ante) probability of player  $i$  sending message  $m_i$ . A strategy profile  $(\sigma_1, \sigma_2)$  is *mutual-preference consistent* if whenever  $u_1, u_2 < 1/2$  then  $\xi_1(m_1, m_2) = \xi_2(m_1, m_2) = L$  for all  $m_1 \in \text{supp}(\mu_u^1)$  and  $m_2 \in \text{supp}(\mu_u^2)$ , and whenever  $u_1, u_2 > 1/2$  then  $\xi_1(m_1, m_2) = \xi_2(m_1, m_2) = R$  for all  $m_1 \in \text{supp}(\mu_u^1)$  and  $m_2 \in \text{supp}(\mu_u^2)$ .

A strategy profile  $(\sigma_1, \sigma_2)$  is *coordinated* if  $\xi_1(m_1, m_2) = \xi_2(m_1, m_2) \in \{L, R\}$  for each pair of messages  $m_1 \in \text{supp}(\mu^1)$  and  $m_2 \in \text{supp}(\mu^2)$ .

For any strategy profile  $\sigma = ((\mu^1, \xi_1), (\mu^2, \xi_2)) \in \Sigma^1 \times \Sigma^2$  and any message  $m_j \in M$ , define

$$\beta_i^\sigma(m_j) = E_{u \sim F_i} \left[ \sum_{m_i \in M} \mu_u^i(m_i) \mathbf{1}_{\{u \leq \xi_i(m_i, m_j)\}} \right]$$

as the expected probability of player  $i$  playing  $L$  conditional on player  $j$  sending message  $m_j \in M$ . We say that strategy profile  $\sigma = (\sigma_1, \sigma_2)$  has (*essentially*) *binary communication* if there are two pairs of numbers  $0 \leq \underline{\beta}_1^\sigma \leq \overline{\beta}_1^\sigma \leq 1$  and  $0 \leq \underline{\beta}_2^\sigma \leq \overline{\beta}_2^\sigma \leq 1$  such that for all messages  $m \in M$  and each player  $i \in \{1, 2\}$  we have  $\beta_i^\sigma(m) \in [\underline{\beta}_i^\sigma, \overline{\beta}_i^\sigma]$ ; for all messages  $m \in M$  such that there is a type  $u < 1/2$  with  $\mu_u^j(m) > 0$  we have  $\beta_i^\sigma(m) = \overline{\beta}_i^\sigma$ ; and for all messages  $m \in M$  such that there is a type  $u > 1/2$  with  $\mu_u(m) > 0$  we have  $\beta_i^\sigma(m) = \underline{\beta}_i^\sigma$ .

Consider a strategy profile  $\sigma = (\sigma_1, \sigma_2)$  that is coordinated and mutual-preference consistent and has binary communication. Then there are  $\alpha_1^\sigma, \alpha_2^\sigma \in [0, 1]$  such that, for each  $i \in \{1, 2\}$ ,

$$\underline{\beta}_i^\sigma = \left(1 - F_j\left(\frac{1}{2}\right)\right) \alpha_i^\sigma \text{ and } \overline{\beta}_i^\sigma = F_j\left(\frac{1}{2}\right) + \left(1 - F_j\left(\frac{1}{2}\right)\right) \alpha_i^\sigma,$$

where  $\alpha_i^\sigma$  is the probability of coordination on  $L$  conditional on player  $i$  having type  $u_i < 1/2$  and player  $j$  having type  $u_i > 1/2$ . We refer to  $\alpha^\sigma = (\alpha_1^\sigma, \alpha_2^\sigma)$  as the *left-tendency profile* of a strategy profile  $\sigma$  that is coordinated and mutual-preference consistent and has binary communication. It is simple to see that the set of strategies satisfying the above three properties (coordination, mutual-preference consistency, and binary communication) is essentially two-dimensional because the left-tendency profile  $\alpha^\sigma = (\alpha_1^\sigma, \alpha_2^\sigma)$  of such a strategy profile  $\sigma$  describes all payoff-relevant aspects. Two such strategy profiles  $\sigma$  and  $\sigma'$  with the same left-tendency profile (i.e., with  $\alpha^\sigma = \alpha^{\sigma'}$ ) can only differ in the way in which the players implement the joint lottery when they have different preferred outcomes, but these implementation differences are not payoff-relevant, as the probability of the joint lottery inducing the players to play  $L$  remains the same.

**Adaptation of Renegotiation-proofness** Given a strategy profile of the game  $\langle \Gamma, M \rangle$  we denote



the induced “renegotiation” game after a positive probability message pair  $m_1, m_2 \in M$  is sent by  $\langle \Gamma(F_{m_1}, F_{m_2}), \tilde{M} \rangle$ . For a strategy profile  $\sigma'$  of such a renegotiation game  $\langle \Gamma(G_1, G_2), \tilde{M} \rangle$ , define the *post-communication* expected payoffs for a player  $i$  of type  $u$  by

$$\pi_u^{i, G_2}(\sigma') = \mathbb{E}_{v \sim G_2} [\pi_{u,v}^i(\sigma')] \equiv \int_{v=0}^1 \pi_{u,v}^i(\sigma') g_2(v) dv.$$

Define  $\mathcal{E}(G_1, G_2)$  as the set of all (possibly asymmetric) equilibrium profiles of the coordination game with communication  $\langle \Gamma(G_1, G_2), \tilde{M} \rangle$  for some finite message set  $\tilde{M}$ .

We say that a strategy profile  $\sigma$  is *post-communication equilibrium Pareto-dominated* if there is a pair of messages  $m_1 \in \text{supp}(\mu^1)$  and  $m_2 \in \text{supp}(\mu^2)$  and an equilibrium profile  $\tilde{\sigma} \in \mathcal{E}(F_{m_1}, F_{m_2})$  such that  $\pi_u^{i, F_{m_i}}(\sigma) \leq \pi_u^{i, F_{m_i}}(\tilde{\sigma})$  for each player  $i \in \{1, 2\}$  and all  $u \in \text{supp}(F_{m_i})$  with strict inequality for some  $u \in \text{supp}(F_{m_i})$  of some player  $i \in \{1, 2\}$ . A strategy profile  $\sigma$  is *renegotiation-proof* if it is not post-communication equilibrium Pareto-dominated.

**Adapted Results** Our main result remains the same in the setup of asymmetric coordination games. The proof, which is analogous to the proof of Theorem 1, is omitted for brevity.

**Theorem 4** (Theorem 1 adapted to asymmetric coordination games). *A strategy profile  $\sigma$  of  $\langle \Gamma(F_1, F_2), M \rangle$  is a renegotiation-proof equilibrium strategy if and only if it is mutual-preference consistent, coordinated, and has binary communication.*

Propositions 1 and 2 as well as Corollary 1 can be adapted to the current setup analogously and the proofs are omitted for brevity.

## E.5 Coordination Games with More Than 2 Actions

In this subsection we extend our main model to coordination games with more than two actions. We now consider a coordination game with two players in which the two players first send one message from a finite message set  $M$  and then, after observing the message pair, choose one action from the ordered set  $A = (a_1, \dots, a_k)$  with  $2 < k < \infty$ .

A player’s type is now a vector  $u = (u_1, \dots, u_k) \in [0, 1]^k$ , where we interpret the  $i$ -th component  $u_i$  as the payoff of the agent if both players choose action  $a_i$ . If the players choose different actions (miscoordinate), then they both get a payoff of zero. We assume that the distribution of types  $F$  is a continuous (atomless) distribution with full support in  $[0, 1]^k$ . For each action  $a_i$ , let  $p_i$  be the probability that the preferred action of a random type is  $a_i$  (i.e., the probability that  $u_i = \max(\{u_1, \dots, u_k\})$ ). Let  $\langle \Gamma_A, M \rangle = \langle \Gamma_A(F), M \rangle$  be the coordination game with set of actions  $A$  and pre-play communication.

A player’s (ex-ante) strategy is a pair  $\sigma = (\mu, \xi)$ , where  $\mu : U \rightarrow \Delta(M)$  is a *message function* that



describes which (possibly random) message is sent for each possible realization of the player's type, and  $\xi : M \times M \times U \rightarrow \Delta(A)$  is an *action function* that describes the distribution of actions chosen as a measurable function of the player's type and the observed message profile. That is, when a player of type  $u$  who follows strategy  $(\mu, \xi)$  observes a message profile  $(m, m')$ , then this player plays action  $a_i$  with probability  $\xi_u(m, m')(a_i)$ .

This game, just as like the main model, has many equilibria. For every action there is a babbling equilibrium in which players of all types after observing any message pair play this action. For every pair of actions  $a_i, a_j$  there are also equilibria in which players send only one of two messages, one message indicating a preference for  $a_i$  and another for not  $a_i$  with play coordinated on  $a_i$  if both players send the appropriate message and play coordinated on  $a_j$  otherwise. None of these equilibria are renegotiation-proof as they are not mutual-preference consistent and mutual-preference consistency is a necessary condition for a strategy to be a renegotiation-proof equilibrium strategy also in the present context, as we shall see below.

It is more difficult to find equilibria that are mutual-preference consistent, that is equilibria in which each player indicates her most preferred action out of all  $k$  actions and play is coordinated on that action if both players indicate a preference for it. Simple adaptations of  $\sigma_L$  and  $\sigma_R$  are not equilibria in the present context. To see this, consider a strategy in which there is a "fallback" action, say action  $a_1$ , in which players indicate their most preferred action (the action with the highest  $u_i$ ), and in which play is coordinated on either action  $a_i$  if both players indicate a preference for it, or coordinated on action  $a_1$  otherwise. Suppose that the distribution of types is such that there are two actions  $a_i$  and  $a_j$  (unequal to each other and unequal to  $a_1$ ) with  $p_j > p_i$ . But then there is a player type  $u = (u_1, u_2, \dots, u_k)$  with  $u_i = \max\{u_1, \dots, u_k\}$ ,  $u_j$  very close to  $u_i$ , and  $u_1 < u_j$ , who would prefer to indicate a preference for action  $a_j$ . Indicating a preference for action  $a_i$ , under the given strategy, provides her with a payoff of  $p_i u_i + (1 - p_i) u_1$ . Indicating a preference for  $a_j$  yields a payoff of  $p_j u_j + (1 - p_j) u_1$ . But then for a suitably chosen vector  $u = (u_1, u_2, \dots, u_k)$  the latter expression is greater than the former, which contradicts the supposition that the given strategy is an equilibrium.

Next we show that a simple adaptation of  $\sigma_C$  remains a renegotiation-proof equilibrium strategy also when there are more than two actions. Fix  $2k$  distinct messages  $m_1^0, m_2^0, \dots, m_k^0, m_1^1, m_2^1, \dots, m_k^1 \in M$ , where the index  $i$  of message  $m_i^b$  is interpreted as denoting that the agent's preferred outcome is the  $i$ -th outcome, and the index  $b \in \{0, 1\}$  is interpreted as a random binary number. Let  $\sigma_C = (\mu_C, \xi_C)$  be extended to the current setup as follows. Define

$$\mu_C(u) = \frac{1}{2} m_i^0 \oplus \frac{1}{2} m_i^1,$$

where  $i = \operatorname{argmin}_j \{u_j \mid u_j = \max\{u_1, \dots, u_k\}\}$ . Thus, the message function  $\mu_C$  induces each agent to reveal her preferred outcome, and to uniformly choose a binary number (either, zero or one). In the second stage, if both agents share the same preferred outcome they play it. Otherwise,

they coordinate on the preferred action with the smaller index if both agents have chosen the same random number, and they coordinate on the preferred outcome with the larger index if both agents have chosen different random numbers, i.e.,

$$\xi_C(m_i^b, m_j^c) = \begin{cases} a_i & (i \leq j \text{ and } b = c) \text{ OR } (i \geq j \text{ and } b \neq c) \\ a_j & \text{otherwise.} \end{cases}$$

We then have the following proposition.

**Proposition 7.** *Strategy  $\sigma_C$  is a renegotiation-proof equilibrium strategy in the game  $\langle \Gamma_A, M \rangle$ .*

*Proof.* Observe that an agent who sends message  $m_i^b$  obtains an expected payoff of

$$\frac{1}{2}u_i + \frac{1}{2}\sum_{j=1}^k p_j u_j$$

when facing a partner who follows strategy  $\sigma_C$ . As the second term in this sum is the same for all messages, an agent of type  $u$  sends this message only if  $u_i = \max\{u_1, \dots, u_k\}$ , as required. The remaining arguments as to why the second-stage behavior is a best reply and why  $\sigma_C$  is renegotiation-proof are analogous to the proof of the “if” part of Theorem 1 and are omitted for brevity.  $\square$

A strategy  $\sigma = (\mu, \xi)$  is *same-message coordinated* if for all messages  $m \in \text{supp}(\bar{\mu})$  there is an action  $a_i$  such that for all  $u$  with  $\mu_u(m) > 0$  we have  $\xi(u, m, m) = a_i$ . In what follows we show that a necessary condition for a strategy to be a renegotiation-proof equilibrium strategy is that this strategy is same-message coordinated and mutual-preference consistent.

**Proposition 8.** *If strategy  $\sigma$  of the game  $\langle \Gamma_A, M \rangle$  with action set  $A$  is renegotiation-proof, then it is same-message coordinated and mutual-preference consistent.*

*Sketch of proof; for the formal proof see Appendix E.7.3.* To show that a renegotiation-proof equilibrium strategy is same-message coordinated we cannot, in fact, use the proof of the main theorem because Lemma 2 crucially depends on the game having only two actions. Instead, we suppose to the contrary that there is a renegotiation-proof equilibrium strategy in which play is not coordinated after both players have sent the same message  $m$ . This strategy thus induces some nondegenerate probability distribution over actions after both players send message  $m$ . We then construct a post-communication Pareto-dominating equilibrium of the induced game that is fully coordinated and has a probability of coordination on every action exactly equal to the probability of this action being played under the original strategy conditional on observing  $(m, m)$ , which contradicts the supposition. The construction is achieved by players sending random messages in such a way that

they are indifferent between all messages and this joint lottery is implemented.<sup>26</sup> The proof that a renegotiation-proof equilibrium strategy must be mutual-preference consistent is then achieved straightforward (by a simple adaptation of Lemma 8).  $\square$

We are able neither to show nor to provide a counterexample that a renegotiation-proof strategy must be coordinated after the agents observe a pair of different messages, and that it must have binary communication.

## E.6 Extreme Types with Dominant Actions

In this subsection we show how to extend our analysis to a setup in which some types have an extreme preference for one of the actions such that it becomes a dominant action for them.

Let  $a < 0$  and  $b > 1$ . We extend the set of types to be the interval  $[a, b]$ . Observe that action  $L$  ( $R$ ) is a dominant action for any type  $u < 0$  ( $u > 1$ ) as coordinating on  $R$  ( $L$ ) yields to such a type a negative payoff of  $u < 0$  ( $1 - u < 0$ ). We call types with a dominant action (i.e.,  $u < 0$  or  $u > 1$ ) *extreme*, and types that do not have a strictly dominant action (i.e.,  $u \in [0, 1]$ ) *moderate*. We assume that the cumulative distribution of types  $F$  is continuous (atomless) and has full support in the interval  $[a, b]$ .

We further assume that the extreme types are a minority both among the agents who prefer action  $R$  and among the agents who prefer action  $L$ , i.e.,

$$F(0) < \frac{1}{2}F\left(\frac{1}{2}\right) \text{ and } 1 - F(1) < \frac{1}{2}\left(1 - F\left(\frac{1}{2}\right)\right).$$

Next, we adapt the definitions of coordination and binary communication to the current setup. The original definition of coordination is too strong in the current setup, as, clearly, when extreme types with different preferred outcomes meet they must miscoordinate. Thus, we present a mild notion of weak coordination. A strategy is *weakly coordinated* if whenever two moderate types meet they never miscoordinate. Note that the definition does not impose any restriction on what happens when an extreme type meets a moderate type.

The original definition of binariness is too weak in the current setup. This is because coordinated strategies must allow for some miscoordination between extreme types, which implies that an agent cares not only about the average probability of the opponent playing left (i.e.,  $\beta^\sigma(m)$ ), but also about the total probability of miscoordination. Thus, we strengthen binariness (and combine it with

---

<sup>26</sup>Using simultaneous communication to implement a jointly controlled lottery was introduced in [Aumann and Maschler \(1968\)](#) (see also [Heller, 2010](#) for a recent implementation, which is robust to joint deviations). The original implementation works perfectly if the probabilities of different actions are rational numbers. In [Appendix E.7.3](#) we present a more elaborate implementation that allows to deal also with irrational numbers in the current setup.

ordinal-preference revelation) by requiring that there exist two distributions of messages, which are used by all types below  $1/2$  and all types above  $1/2$ , respectively. Formally, a strategy  $\sigma = (\mu, \xi)$  has *strongly binary communication* if  $\mu(u) = \mu(u')$  if either  $u, u' \leq 1/2$  or  $u, u' > 1/2$ . It is easy to see that the strategies  $\sigma_L, \sigma_R, \sigma_C$  defined in Section 3 all satisfy strongly binary communication. Moreover, one can show, for any  $\alpha \in [0, 1]$ , that if there exists a strategy  $\sigma$  that is coordinated, mutual-preference consistent, and has binary communication with left tendency  $\alpha$ , then there also exists strategy  $\tilde{\sigma}$  with the same properties that is strongly binary communication.

Our next result shows that there exists, essentially, a unique renegotiation-proof equilibrium strategy that is coordinated, mutual-preference consistent, and has strongly binary communication.

**Proposition 9.** *In a coordination game with communication and with dominant action types, a strategy  $\sigma$  that is coordinated, mutual-preference consistent, and has strongly binary communication is a renegotiation-proof equilibrium strategy if and only if it has a left tendency of*

$$\alpha = \frac{F(0)}{F(0) + (1 - F(1))}.$$

The formal proof is presented in Appendix E.7.4. The key intuition behind this proposition is that given the frequency of dominant action types  $F(0) > 0$  (of  $L$ -dominant action types) and  $1 - F(1) > 0$  (of  $R$ -dominant action types) to make the agent of type  $u = 1/2$  indifferent between signaling a lower than half or higher than half type we must have a strategy that counterbalances these frequencies of dominant action types. To see this, consider a straightforward adaptation of  $\sigma_L = (\mu^*, \xi_L)$  to this setting by having extreme types follow their dominant action in the second stage (and moderate types play in the same way as in the baseline model). Note that  $\sigma_L$  is no longer an equilibrium strategy with extreme types. Observe that having a moderate type send message  $m_R$  leads to coordination with probability one (sometimes on  $R$  and sometimes on  $L$  depending on the opponent's message), while having a moderate type send message  $m_L$  leads to coordination (on  $L$  only) with probability  $F(1) < 1$ . This implies that agents of type  $u < 1/2$  sufficiently close to  $1/2$  strictly prefer sending message  $m_R$  to sending message  $m_L$  (as the former induces a higher probability of coordination on the same action as the partner), which contradicts the supposition that  $\sigma_L$  is an equilibrium strategy.

Appendix E.7.4 also shows that a left tendency  $\alpha$  renegotiation-proof strategy that is coordinated and has strongly binary communication can be implemented whenever  $\alpha$  is a rational number and the set of messages  $M$  is sufficiently large (and irrational  $\alpha$ -s can be approximately implemented by  $\epsilon$ -equilibria).

Observe that in the symmetric case ( $F(0) = 1 - F(1)$ ), the essentially unique renegotiation-proof strategy with the above two properties is  $\sigma_C$ . Further observe that in the asymmetric case, the moderate types gain if the extreme types with the same preferred outcome are more frequent

than the extreme types of the opposite preferred outcome. Specifically, if there are more extreme “leftists” than extreme “rightists” (i.e.,  $F(0) > 1 - F(1)$ ), then the essentially unique renegotiation-proof strategy with properties of coordination and strongly binary communication induces higher probability to coordinate on action  $L$  (rather than on action  $R$ ) whenever two moderate agents with different preferred outcomes meet.

## E.7 Formal Proofs of Extensions

### E.7.1 Proof of Theorem 2 (MultiDimensional types, Section E.2)

The proof of Theorem 2 mimics the proof of Theorem 1 except that Lemma 2 has to be adapted somewhat as follows (and this is the only place where one needs to use the assumption of unambiguous coordination preferences).

**Lemma 6.** *Let  $U = \{(u_{LL}, u_{LR}, u_{RL}, u_{RR}) \mid u_{LL} > u_{RL} \text{ and } u_{RR} > u_{LR}\}$  and let the atomless distribution of types have unambiguous coordination preferences. Let  $\sigma = (\mu, \xi)$  be a renegotiation-proof equilibrium strategy. Then it is coordinated.*

*Proof.* We need to show that for any message pair  $m, m' \in \text{supp}(\bar{\mu})$ ,

$$\text{either } \xi(m, m') \geq \sup \{\varphi_u \mid \mu_u(m) > 0\} \text{ or } \xi(m, m') \leq \inf \{\varphi_u \mid \mu_u(m) > 0\}.$$

Let  $m, m' \in \text{supp}(\bar{\mu})$  and assume to the contrary that

$$\inf \{\varphi_u \mid \mu_u(m) > 0\} < \xi(m, m') < \sup \{\varphi_u \mid \mu_u(m) > 0\}.$$

As  $\sigma$  is an equilibrium, we must have

$$\inf \{\varphi_u \mid \mu_u(m') > 0\} < \xi(m', m) < \sup \{\varphi_u \mid \mu_u(m') > 0\}.$$

(Otherwise the  $m'$  message sender would play  $L$  with probability one or  $R$  with probability one, in which case the  $m$  message sender’s best response would be to play  $L$  (or  $R$ ) regardless of her type). Let  $x = \xi(m, m')$  and  $x' = \xi(m', m)$ . In what follows we will show that the equilibrium  $(x, x')$  of the game without communication  $\Gamma(F_m, F_{m'})$  is Pareto-dominated by either  $\sigma_L$ ,  $\sigma_R$ , or  $\sigma_C$  (all based on  $\varphi_u$  instead of  $u$ ).

There are three cases to be considered. Case 1: Suppose that  $x, x' \leq 1/2$ . We now show that in this case the equilibrium  $(x, x')$  is Pareto-dominated by  $\sigma_R$ . Consider the player who sent message  $m$ .

Case 1a: Consider a type  $u$  with  $\varphi_u \leq x$ . Then we have

$$\begin{aligned} u_{LL}F_{m'}(x') + (1 - F_{m'}(x')) u_{LR} &\leq u_{LL}F_{m'}(\tfrac{1}{2}) + u_{LR} \left(1 - F_{m'}(\tfrac{1}{2})\right) \\ &\leq u_{LL}F_{m'}(\tfrac{1}{2}) + u_{RR} \left(1 - F_{m'}(\tfrac{1}{2})\right), \end{aligned}$$

where the first expression is the type  $u$  agent's payoff under strategy profile  $(x, x')$  and the last expression is her payoff under strategy profile  $\sigma_R$ . The first inequality follows from  $u_{LL} \geq u_{LR}$  and  $F_{m'}(1/2) \geq F_{m'}(x')$  by the fact that  $F_{m'}$  is nondecreasing (as it is a cumulative distribution function), and the second inequality follows from  $u_{RR} \geq u_{LR}$ . This inequality is strict when  $u_{LL} > u_{LR}$  and  $F_{m'}(1/2) > F_{m'}(x')$  or when  $u_{RR} > u_{LR}$ .

Case 1b: Now consider a type  $u$  with  $x < \varphi_u \leq 1/2$ . Then we have

$$\begin{aligned} u_{RL}F_{m'}(x') + u_{RR}(1 - F_{m'}(x')) &\leq u_{LL}F_{m'}(x') + u_{RR}(1 - F_{m'}(x')) \\ &\leq u_{LL}F_{m'}(\tfrac{1}{2}) + u_{RR} \left(1 - F_{m'}(\tfrac{1}{2})\right), \end{aligned}$$

where the first expression is the type  $u$  agent's payoff under strategy profile  $(x, x')$  and the last expression is her payoff under strategy profile  $\sigma_R$ . The first inequality follows from  $u_{LL} \geq u_{RL}$  and the second one from  $F_{m'}(1/2) \geq F_{m'}(x')$  and  $u_{LL} \geq u_{RR}$ . Note also that the second inequality follows from the assumption of unambiguous coordination preferences and  $\varphi_u \leq 1/2$ . This inequality is strict when  $u_{LL} > u_{RL}$  or when  $F_{m'}(1/2) > F_{m'}(x')$  and  $u_{LL} > u_{RR}$ .

Case 1c: Finally, consider a type  $u$  with  $\varphi_u > 1/2$ . Then we have

$$u_{RR} > u_{RL}F_{m'}(x') + u_{RR}(1 - F_{m'}(x')),$$

where the right-hand side is the type  $u$  agent's payoff under strategy profile  $(x, x')$  and the left-hand side is her payoff under strategy profile  $\sigma_R$ . The inequality follows from the observation that  $u_{RR} > u_{RL}$  because  $u_{RR} > u_{LL}$  by the assumption of unambiguous coordination preferences, and  $u_{LL} \geq u_{RL}$  by the fact that it is a coordination game.

The analysis for the player who sent message  $m'$  is analogous.

Case 2: Suppose that  $x, x' \geq 1/2$ . The analysis is analogous to Case 1 if we replace  $\sigma_R$  with  $\sigma_L$ .

Case 3: Suppose, without loss of generality for the remaining cases, that  $x \leq 1/2 \leq x'$ . We now show that the equilibrium  $(x, x')$  in this case is Pareto-dominated by  $\sigma_C$ . Consider the player who sent message  $m$ .

Case 3a: Consider a type  $u$  such that  $\varphi_u \leq x$ . Then we have

$$u_{LL} \left[ F_{m'}(\tfrac{1}{2}) + \tfrac{1}{2} \left(1 - F_{m'}(\tfrac{1}{2})\right) \right] + u_{RR} \tfrac{1}{2} \left(1 - F_{m'}(\tfrac{1}{2})\right) > u_{LL}F_{m'}(x') + u_{LR}(1 - F_{m'}(x')),$$

where the right-hand side is the type  $u$  agent's payoff under strategy profile  $(x, x')$  and the left-hand side is her payoff under strategy profile  $\sigma_C$ . The inequality follows from the observation that  $u_{RR} \geq u_{LR}$  and  $F_{m'}(x') \leq 1/2$  by the fact that  $F_{m'}(x') = x$  when  $(x, x')$  is an equilibrium.

Case 3b: Now consider a type  $u$  with  $x < \varphi_u \leq 1/2$ . Then we have

$$\begin{aligned} u_{RL}F_{m'}(x') + u_{RR}(1 - F_{m'}(x')) &\leq u_{LL}F_{m'}(x') + u_{RR}(1 - F_{m'}(x')) \\ &\leq u_{LL}\left[\frac{1}{2} + \frac{1}{2}F_{m'}(x')\right] + u_{RR}\frac{1}{2}\left(1 - F_{m'}\left(\frac{1}{2}\right)\right), \end{aligned}$$

where the first expression is the type  $u$  agent's payoff under strategy profile  $(x, x')$  and the last expression is her payoff under strategy profile  $\sigma_C$ . The first inequality follows from  $u_{LL} \geq u_{RL}$  and the second one from  $u_{LL} \geq u_{RR}$  by the assumption of unambiguous coordination preferences given  $\varphi_u \leq 1/2$  and  $F_{m'}(x') = x$  by  $(x, x')$  being an equilibrium and  $x < 1/2$ . The inequality is strict if  $u_{LL} > u_{RL}$  or  $u_{LL} > u_{RR}$ .

Case 3c: Finally, consider a type  $u$  with  $\varphi_u > 1/2$ . Then we have

$$\begin{aligned} u_{RL}F_{m'}(x') + u_{RR}(1 - F_{m'}(x')) &< u_{RL}\frac{1}{2}F_{m'}\left(\frac{1}{2}\right) + u_{RR}\left[\left(1 - F_{m'}\left(\frac{1}{2}\right)\right) + \frac{1}{2}F_{m'}\left(\frac{1}{2}\right)\right] \\ &\leq u_{LL}\frac{1}{2}F_{m'}\left(\frac{1}{2}\right) + u_{RR}\left[\left(1 - F_{m'}\left(\frac{1}{2}\right)\right) + \frac{1}{2}F_{m'}\left(\frac{1}{2}\right)\right], \end{aligned}$$

where the first expression is the type  $u$  agent's payoff under strategy profile  $(x, x')$  and the last expression is her payoff under strategy profile  $\sigma_C$ . The first inequality follows from  $u_{RR} > u_{LL} \geq u_{RL}$  by the assumption of unambiguous coordination preferences and from  $(1 - F_{m'}(1/2)) \geq (1 - F_{m'}(x'))$  as  $F_{m'}$  is nondecreasing.

The analysis for the player who sent message  $m'$  is analogous. □

## E.7.2 Proof of Theorem 3 (Multiple Players, Section E.3)

The “if” part is analogous to the proof of the “if” part of Theorem 1.

The proof of the “only if” part does not extend directly and has to be adapted as follows. The following lemma states that play is coordinated whenever all players send the same message. Note that it does not (yet) claim that play is coordinated after any message pair is observed.

**Lemma 7.** *Let  $\sigma = (\mu, \xi)$  be a renegotiation-proof equilibrium strategy. Let  $m \in \text{supp}(\bar{\mu})$  and let  $\mathbf{m} = (m, \dots, m)$  be the vector with  $n$  identical entries of  $m$ , which represents the case of all  $n$  players sending message  $m$ . Then either  $\xi(\mathbf{m}) \geq \sup\{u \mid \mu_u(m) > 0\}$  or  $\xi(\mathbf{m}) \leq \inf\{u \mid \mu_u(m) > 0\}$ .*

*Proof.* Let  $m \in \text{supp}(\bar{\mu})$  and assume to the contrary that

$$\inf \{u \mid \mu_u(m) > 0\} < \xi(\mathbf{m}) < \sup \{u \mid \mu_u(m) > 0\}.$$

Let  $x = \xi(\mathbf{m})$ . In what follows we will show that the symmetric equilibrium in which all players use cutoff  $x$  after sending the identical message  $m$ , denoted by  $\mathbf{x} = (x, \dots, x)$ , is equilibrium Pareto-dominated by either  $\sigma_L$  or  $\sigma_R$ .

There are two cases to be considered. Case 1: Suppose that  $x \leq 1/2$ . We now show that in this case the equilibrium  $\mathbf{x}$  is Pareto-dominated by  $\sigma_R$ .

Case 1a: Consider a type  $u \leq x$ . Then we have

$$(1 - u) \left(F_m\left(\frac{1}{2}\right)\right)^{n-1} + u \left(1 - \left(F_m\left(\frac{1}{2}\right)\right)^{n-1}\right) \geq (1 - u) (F_m(x))^{n-1},$$

where the left-hand side is the type  $u$  agent's payoff under strategy profile  $\sigma_R$  and the right-hand side is her payoff under strategy profile  $\mathbf{x}$ . The inequality follows from the fact that  $u \left(1 - \left(F_m(1/2)\right)^{n-1}\right) \geq 0$  and  $F_m(1/2) \geq F_m(x)$  by the fact that  $F_m$  is nondecreasing (as it is a cumulative distribution function). Note also that this inequality is strict for all  $u$  except for  $u = 0$  in the case of  $x = 1/2$ .

Case 1b: Now consider a type  $u$  with  $x < u \leq 1/2$ . Then we have

$$(1 - u) \left(F_m\left(\frac{1}{2}\right)\right)^{n-1} + u \left(1 - \left(F_m\left(\frac{1}{2}\right)\right)^{n-1}\right) \geq u \left(1 - \left(F_m(x)\right)^{n-1}\right),$$

where the left-hand side is the type  $u$  agent's payoff under strategy profile  $\sigma_R$  and the right-hand side is her payoff under strategy profile  $\mathbf{x}$ . The inequality follows from the fact that given  $u \leq 1/2$  we have that  $1 - u \geq u$  and therefore

$$(1 - u) \left(F_m\left(\frac{1}{2}\right)\right)^{n-1} + u \left(1 - \left(F_m\left(\frac{1}{2}\right)\right)^{n-1}\right) \geq u.$$

Note that this inequality actually holds strictly for all  $u$ .

Case 1c: Finally, consider a type  $u > 1/2$ . Then we have

$$u > u (1 - F_m(x))^{n-1},$$

where the left-hand side is the type  $u$  agent's payoff under strategy profile  $\sigma_R$  and the right-hand side is her payoff under strategy profile  $\mathbf{x}$ .

Case 2: Suppose that  $x \geq 1/2$ . The analysis is analogous to Case 1 if we replace  $\sigma_R$  with  $\sigma_L$ .  $\square$

**Lemma 8.** *Every renegotiation-proof equilibrium strategy  $\sigma = (\mu, \xi)$  is mutual-preference consistent.*



*Proof.* The proof of this lemma involves two steps. In the first step we show that a renegotiation-proof equilibrium strategy  $\sigma$  is ordinal preference revealing, i.e., such that for any message  $m \in \text{supp}(\bar{\mu})$ ,  $F_m(1/2) \in \{0, 1\}$ . We then use this to show that  $\sigma$  is mutual-preference consistent.

Assume, first, that  $\sigma$  is a renegotiation-proof equilibrium strategy but not ordinal preference-revealing. That is, suppose to the contrary that  $F_m(1/2) \in (0, 1)$ . Then there are types  $u < 1/2$  as well as types  $u > 1/2$  who both send message  $m$  with positive probability. By Lemma 7 play after message pair  $(m, m)$  must be either  $L$  or  $R$ . If it is  $L$  then the equilibrium strategy  $\sigma_L$  Pareto-dominates playing  $L$ , with a strict payoff improvement for all types  $u > 1/2$  (and unchanged payoffs for all types  $u \leq 1/2$ ). If it is  $R$  then the equilibrium strategy  $\sigma_R$  Pareto-dominates playing  $R$ , with a strict payoff improvement for all types  $u < 1/2$  (and unchanged payoffs for all types  $u \geq 1/2$ ).

Given a renegotiation-proof equilibrium strategy  $\sigma = (\mu, \xi)$ , we can classify messages in the support of  $\mu$  into two distinct sets,  $M_L = M_L(\sigma) = \{m \in \text{supp}(\mu) \mid F_m(1/2) = 1\}$  and  $M_R = M_R(\sigma) = \{m \in \text{supp}(\mu) \mid F_m(1/2) = 0\}$ , where  $M_L \cap M_R = \emptyset$  and  $M_L \cup M_R = \text{supp}(\bar{\mu})$ .

To show that a renegotiation-proof equilibrium strategy  $\sigma = (\mu, \xi)$  is mutual-preference consistent, then consider any profile of types  $(u_1, u_2, \dots, u_n)$  such that  $u_i < 1/2$  for all  $i \in \{1, \dots, n\}$ . They must each send a message in  $M_L$ , which we denote by the profile  $\mathbf{m} = (m_1, \dots, m_n)$ . Any play after message profile  $m$  that is not coordinated on  $L$  is now clearly Pareto-dominated (given that all types are at most  $1/2$ ) by playing the equilibrium strategy  $L$ . The case for a profile of types  $u_i > 1/2$  for all  $i \in \{1, \dots, n\}$  is proven analogously.  $\square$

The following lemma shows that, in a renegotiation-proof equilibrium strategy, agents never miscoordinate after observing any message profile.

**Lemma 9.** *Every renegotiation-proof equilibrium strategy  $\sigma$  is coordinated.*

*Proof.* Suppose that  $\sigma = (\mu, \xi)$  is a renegotiation-proof equilibrium strategy. Given Lemmas 7 and 8 it only remains to prove that play under  $\sigma$  is coordinated even after mixed messages are sent, i.e., when there is at least one player who sends a message in  $M_L$  and another player who sends a message in  $M_R$ , where  $M_L$  and  $M_R$  are as defined in the proof of Lemma 8. Suppose that this is the case. Then let  $I \subset \{1, \dots, n\}$  be the set of all players who send a message  $m_i \in M_L$ . Let  $I^c$  denote its complement. By Lemma 8 all  $i \in I^c$  satisfy  $m_i \in M_R$ . Let  $x_i = \xi(m_i, m_{-i})$  be the cutoff used by player  $i$  after observing message profile  $(m_1, m_2, \dots, m_n)$ . Then by Lemma 8 we have  $x_i \leq 1/2$  for all  $i \in I$  and  $x_i \geq 1/2$  for all  $i \in I^c$ . For this profile  $\mathbf{x} = (x_1, \dots, x_n)$  to be an equilibrium after the players observe message profile  $(m_1, m_2, \dots, m_n)$ , we must have that for each  $i = 1, 2, \dots, n$ , the probability that player  $i$ 's opponents coordinate their action on  $L$  conditional on them coordinating (on either  $L$  or  $R$ ) is

$$x_i = \frac{\prod_{j \neq i} F_{m_j}(x_j)}{\prod_{j \neq i} F_{m_j}(x_j) + \prod_{j \neq i} (1 - F_{m_j}(x_j))}.$$

But then, all types of all players, after observing message profile  $(m_1, m_2, \dots, m_n)$ , weakly (and some strictly) prefer to play  $\sigma_C$ , which is a payoff identical in this case to a public fair coin toss to determine whether coordination should be on  $L$  or  $R$ . To see this, consider a player  $i$  who sent a message in  $M_L$  (i.e.,  $u_i \leq 1/2$ , which implies that  $x_i \leq 1/2$ ) and consider the following two cases.

Case 1: Suppose that  $u_i \leq x_i$ . Then, under the given strategy, this type's payoff is  $(1-u_i) \prod_{j \neq i} F_{m_j}(x_j)$  with  $\prod_{j \neq i} F_{m_j}(x_j) \leq 1/2$ . The equilibrium strategy  $\sigma_C$  yields to this type a payoff of  $1/2(1-u_i) + 1/2u_i$ , which exceeds the former payoff, which contradicts the supposition that  $\sigma$  is renegotiation-proof.

Case 2: Suppose that  $x_i < u_i \leq 1/2$ . Then, under the given strategy, this type's payoff is

$$u_i \prod_{j \neq i} (1 - F_{m_j}(x_j)).$$

The equilibrium strategy  $\sigma_C$  yields to this type a payoff of  $1/2(1-u_i) + 1/2u_i$ , which, by virtue of  $1-u_i \geq u_i$ , again exceeds the former payoff, which contradicts the supposition that  $\sigma$  is renegotiation-proof.

The analysis for a player who sent a message in  $M_R$  is proven analogously. □

To complete the proof of Theorem 1 for the case of many players, we need to prove that any renegotiation-proof equilibrium strategy also has binary communication. The proof of this statement is analogous to the proof of Lemma 3 and therefore omitted.

### E.7.3 Proof of Proposition 8 (Multiple Actions, Section E.2)

For the proof, two lemmas about approximating real numbers by rational numbers will be useful. The first lemma shows that any discrete distribution with (some) probabilities being irrational numbers and at least three elements in its support can be approximated from below by a vector of rational numbers, such that the profile of differences (between the irrational exact probability and its rational approximation from below) is roughly uniform in the sense that no difference is larger than the half the sum of all the differences.

**Lemma 10.** *Let  $p \in \Delta(A)$  be a distribution satisfying  $|\text{supp}(p)| \geq 3$ . Then there exists a function  $q : A \rightarrow \mathbb{R}^+$  such that, for each  $1 \leq i \leq k$ ,  $q(a_i)$  is a rational number,  $q(a_i) \leq p(a_i)$ , and*

$$p(a_i) - q(a_i) \leq \frac{1}{2} \sum_{1 \leq j \leq k} (p(a_j) - q(a_j)).$$

*Proof.* Let  $\delta < \min \{p(a)/2 \mid a \in \text{supp}(p)\}$ . As the rational numbers are dense in the reals, for each real number  $\hat{p} > \delta$ , there exists a rational number  $\hat{q} \in (0, \hat{p})$  such that  $\hat{p} - \hat{q} \in ((9/10)\delta, \delta)$ . Call this  $\hat{q}$  a *rational approximation* of  $\hat{p}$ . For each  $a \in \text{supp} p$ , let  $q(a)$  be a rational approximation of  $p(a)$ .

For each  $a \notin \text{supp}(p)$  let  $q(a) = q(a) = 0$ . Then it follows that, for each  $1 \leq i \leq k$ ,  $q(a_i)$  is a rational number, and  $q(a_i) \leq p(a_i)$ . Finally we get, for each  $1 \leq i \leq k$ , that

$$p(a_i) - q(a_i) \leq \delta \leq \frac{1}{2} |\text{supp}(p)| \frac{9}{10} \delta \leq \frac{1}{2} \sum_{1 \leq j \leq k} (p(a_j) - q(a_j)),$$

where the first inequality follows directly from the definition of a rational approximation, the second one follows from the assumption that  $|\text{supp}(p)| \geq 3$ , and the last one from the assumption that, for each  $a_j \in \text{supp}(p)$ ,  $p(a_j) - q(a_j) > (9/10)\delta$  by the definition of a rational approximation.  $\square$

Note that the closer  $\delta$  is to zero, the better the rational approximation constructed in this proof. Note, however, that this does not matter in the proof of Proposition 8 below, which simply uses any (possibly quite rough) rational approximation.

The second lemma is utilized in the proof of Proposition 8 for the case where the distribution in question has exactly two elements in its support.

**Lemma 11.** *Let  $p, q \in (0, 1)$ . Then there exists a rational number  $\alpha \in (0, 1)$  satisfying*

$$\frac{p - q}{1 - q} < \alpha < \frac{p}{q}.$$

*Proof.* Note that, as  $p < 1$ ,

$$0 \leq (p - q)^2 = p^2 - 2pq + q^2 < p - 2pq + q^2.$$

The inequality  $0 < p - 2pq + q^2$  then implies that

$$qp - q^2 < p - pq \Leftrightarrow q(p - q) < p(1 - q) \Leftrightarrow \frac{p - q}{1 - q} < \frac{p}{q}.$$

The result then follows from the fact that  $p - q < 1 - q$ .  $\square$

We now turn to the proof of Proposition 8. Let  $\sigma = (\mu, \xi)$  be a renegotiation-proof equilibrium strategy of  $\langle \Gamma_A, M \rangle$ . We begin by showing that  $\sigma$  is same-message coordinated. Let  $m \in \text{supp}(\bar{\mu})$  let and  $p \in \Delta(A)$  be the distribution of play under  $\sigma$  conditional on message pair  $(m, m)$  being observed. Assume to the contrary that  $p(a) < 1$  for each  $a \in A$  (i.e., that there is miscoordination).

**Case I:** Assume that  $|\text{supp}(p)| \geq 3$ . Let  $q : A \rightarrow \mathbb{R}^+$  be a rational approximation of  $p$  satisfying the requirements of Lemma 10. The fact that all  $q(a_i)$ -s are rational numbers implies that there are  $l_1, \dots, l_k, n \in \mathbb{N}$ , such that  $q(a_i) = l_i/n$  for each  $i$  and  $l_1 + \dots + l_k \leq n$ . Consider the following equilibrium strategy  $\tilde{\sigma} = (\tilde{\mu}, \tilde{\xi})$  of the game induced after players observe message pair  $(m, m)$  with

an additional communication round with the set of messages

$$\tilde{M} = \{m_{B,i,b} \mid 1 \leq B \leq n, 1 \leq i \leq k, b \in \{0,1\}\}.$$

We let  $1 \leq B \leq n$  denote a random integer used for a joint lottery,  $1 \leq i \leq k$  denote the index of the player's preferred outcome (i.e.,  $u_i = \max \{u_j \mid 1 \leq j \leq k\}$ , and  $b \in \{0,1\}$  denotes a random bit). The message function  $\tilde{\mu}$  induces each agent to choose the indexes  $B$  and  $b$  randomly (uniformly, and independently of each other), and to choose  $i$  such that  $u_i = \max \{u_j \mid 1 \leq j \leq k\}$  is her preferred outcome.

The action function  $\tilde{\xi}(m_{B,i,b}, m_{B',i',b'})$  is defined as follows. Let  $\hat{B} = (B + B') \bmod n$  be the sum of the random  $B$ -s sent by the players. Both players play action  $a_j$  if  $l_1 + \dots + l_{j-1} \leq \hat{B} < l_1 + \dots + l_j$ . If  $\hat{B} \geq l_1 + \dots + l_k$ , then both players play the action in  $\{a_i, a_{i'}\}$  with the smaller index if  $b = b'$  and the action with the larger index if  $b \neq b'$ . Strategy  $\tilde{\sigma}$  induces both players to coordinate on a random action with probability  $\bar{q} \equiv q(a_1) + \dots + q(a_k)$  (and, conditional on that, the random action is chosen to be  $a_j$  with probability  $q(a_j)/\bar{q}$ ), and to coordinate on the preferred action of one of the two players (chosen uniformly at random) with probability  $1 - \bar{q}$ , which can be written as  $1 - \bar{q} = \sum_{1 \leq j \leq k} (p(a_j) - q(a_j))$ , by the fact that  $p$  is a probability distribution and, thus,  $\sum_{1 \leq j \leq k} p(a_j) = 1$ .

The proof that  $\tilde{\sigma}$  is an equilibrium strategy is analogous to the proof of Proposition 7 and therefore omitted. The expected payoff that the original equilibrium  $\sigma$  yields to each type  $u = (u_1, \dots, u_k) \in U$  in the game induced after observing  $(m, m)$  is  $\max_j \{p(a_j) u_j\}$ . The expected payoff that  $\tilde{\sigma}$  induces for each type  $u$  is at least

$$\sum_j q(a_j) u_j + \frac{1}{2} (1 - \bar{q}) \max_j \{u_j\} \geq \max_j \{q(a_j) u_j\} + \frac{1}{2} \max_j \{u_j\} \sum_j (p(a_j) - q(a_j)).$$

Thus, the difference between the payoff of  $\tilde{\sigma}$  and the payoff of  $\sigma$  for a type  $u$  is at least

$$\begin{aligned} \frac{1}{2} \max_j \{u_j\} \sum_j (p(a_j) - q(a_j)) - (\max_j \{p(a_j) u_j\} - \max_j \{q(a_j) u_j\}) &\geq \\ \frac{1}{2} \max_j \{u_j\} \sum_j (p(a_j) - q(a_j)) - \frac{1}{2} \sum_{1 \leq j \leq k} (p(a_j) - q(a_j)) u_j &\geq 0, \end{aligned}$$

where the first inequality is due to

$$\max_j \{p(a_j) u_j\} - \max_j \{q(a_j) u_j\} = p(a_l) u_l - \max_j \{q(a_j) u_j\} \leq p(a_l) u_l - q(a_l) u_l,$$

where  $l = \operatorname{argmax}_j \{p(a_j) u_j\}$ ; and the second inequality is due to  $q$  being a rational approximation of  $p$  as given by Lemma 10.

This then implies that  $\sigma$  is post-communication equilibrium Pareto-dominated by  $\tilde{\sigma}$ , which contradicts the supposition that  $\sigma$  is a renegotiation-proof equilibrium.

**Case II:** We are left with the case of  $|\text{supp}(p)| = 2$ . Let  $\text{supp}(p) = \{a_i, a_j\}$ . Let  $q \in (0, 1)$  be the posterior probability of a player having a type  $u$  with  $u_i \geq u_j$ , conditional on sending message  $m$ . Let  $p \equiv p(a_i)$ . By Lemma 11, there exists a rational number  $\alpha \equiv k/n \in (0, 1)$  satisfying

$$\frac{p - q}{1 - q} < \alpha < \frac{p}{q}.$$

Consider the following symmetric equilibrium  $\tilde{\sigma} = (\tilde{\mu}, \tilde{\xi})$  of the game induced after players observe message pair  $(m, m)$  with an additional communication round with the set of messages

$$\tilde{M} = \{i, j\} \times \{1, \dots, n\}.$$

The first component of the message of each player is interpreted as her preferred coordinated outcome of  $a_i$  and  $a_j$ , and the second component is a random number between 1 and  $n$ . When following strategy  $\tilde{\sigma}$  the players send message component  $i$  if and only if  $u_i \geq u_j$ , send a random number between 1 and  $n$  according to the uniform distribution, play  $a_i$  after observing  $((i, a), (i, b))$  for any numbers  $a$  and  $b$ , play  $a_j$  after observing  $((j, a), (j, b))$  for any numbers  $a$  and  $b$ , play  $a_i$  after observing  $((i, a), (j, b))$  if  $a + b < k \pmod n$ , and play  $a_j$  after observing  $((a_i, a), (a_j, b))$  if  $a + b \geq k \pmod n$ .

Observe that  $\tilde{\sigma}$  is indeed an equilibrium of the induced game, because following any pair of messages the players coordinate for sure, each agent with  $u_i > u_j$  (resp.,  $u_j > u_i$ ) strictly prefers to report that her preferred outcome is  $a_i$  (resp.,  $a_j$ ) as this induces her to coordinate on  $a_i$  (resp.,  $a_j$ ) with a high probability of  $q + \alpha(1 - q)$  (resp.,  $1 - q + (1 - \alpha)q$ ) instead of with a low probability of  $\alpha q$  (resp.,  $(1 - q)(1 - \alpha)$ ), and each agent is indifferent between sending any random number, as this has no effect on the probability of coordinating on  $a_i$  (which is equal to  $\alpha = k/n$ ), given that her opponent chooses his random number uniformly as well.

Recall that the payoff of each type who follows  $\sigma$  in the game induced after observing  $(m, m)$  is equal to  $\max\{u_i p, u_j(1 - p)\}$ . The payoff of each type  $u$  with  $u_i \geq u_j$  in equilibrium  $\tilde{\sigma}$  is given by

$$(q + \alpha(1 - q))u_i + (1 - (q + \alpha(1 - q)))u_j > pu_i + (1 - p)u_j \geq \max\{u_i p, u_j(1 - p)\},$$

where the first inequality is implied by  $u_i \geq u_j$  and

$$\frac{p - q}{1 - q} < \alpha \Leftrightarrow p < q + \alpha(1 - q).$$

The payoff of each type  $u$  with  $u_i < u_j$  in equilibrium  $\tilde{\sigma}$  is given by

$$\begin{aligned} & ((1 - q) + (1 - \alpha)q)u_j + (1 - ((1 - q) + (1 - \alpha)q))u_i > \\ & (1 - p)u_j + pu_i \geq \max\{(1 - p)u_j, pu_i\}, \end{aligned}$$

where the first inequality is implied by

$$\frac{q}{p} > \alpha \Leftrightarrow 1 - \alpha > \frac{q-p}{q} \Leftrightarrow (1-q) + (1-\alpha)q > 1-p$$

and  $u_i < u_j$ . This implies that all types obtain a strictly larger payoff in  $\tilde{\sigma}$  (relative to the expected payoff of  $\sigma$  in the game induced after players observe message pair  $m, m$ ), which implies that  $\tilde{\sigma}$  Pareto-dominates  $\sigma$ , which contradicts the supposition that  $\sigma$  is renegotiation-proof.

Next, we show that  $\sigma = (\mu, \xi)$  is mutual-preference consistent. For each  $i \in \{1, \dots, k\}$ , let  $U_i \subset [0, 1]^k$  be the set of types such that  $u_i \geq \max_{j \neq i} u_j$ . Assume, first, that  $\sigma$  is not ordinal preference-revealing. That is, suppose that there is a message  $m \in \text{supp}(\bar{\mu})$  such that there are action indices  $i, j$  with  $i \neq j$  and  $\mu_u(m) > 0$  for some  $u \in U_i$  and some  $u \in U_j$ . We have shown above that the play after players observe message pair  $(m, m)$  must be coordinated on some action  $a_l \in A$ . Now consider the following strategy with new message space  $\tilde{M} = \{m_i, m_{-i}\}$  in which players of type  $u \in U_i$  send message  $m_i$ , while all others send message  $m_{-i}$  and play is  $a_l$  unless both players send message  $m_i$ , in which case it is  $a_i$ . This is an equilibrium strategy of the induced game after players observe message pair  $(m, m)$  and it Pareto-dominates  $\sigma$ , a contradiction. This proves that a renegotiation-proof equilibrium strategy  $\sigma$  must be ordinal preference-revealing.

Given a renegotiation-proof equilibrium strategy  $\sigma = (\mu, \xi)$ , we can classify messages in the support of  $\mu$  into  $k$  distinct sets  $M_i = M_i(\sigma) = \{m \in \text{supp}(\bar{\mu}) \mid u \in U_i\}$  for each  $i = 1, \dots, k$ , where for each  $i, j$  with  $i \neq j$   $M_i \cap M_j = \emptyset$  and  $\bigcup_{i=1}^k M_i = \text{supp}(\bar{\mu})$ .

To show that a renegotiation-proof equilibrium strategy  $\sigma$  is mutual-preference consistent, consider a message pair  $(m, m')$  with  $m, m' \in M_i$  for some  $i = 1, \dots, k$ . Since  $\sigma$  is ordinal preference-revealing, the updated support of types who observe either message  $m$  or  $m'$  is then in  $U_i$ . But then any joint action distribution that the two players could play after  $(m, m')$  is Pareto-dominated by the equilibrium strategy of playing action  $a_i$ .

#### E.7.4 Proof of Proposition 9 (Extreme Types, Section E.6)

*Proof of Proposition 9.* Any strategy  $\sigma$  that is coordinated, mutual-preference consistent, and has strongly binary communication can be characterized by its left tendency (see Section 3) as follows. Under such a mutual-preference consistent strategy players indicate whether their type is below or above  $1/2$ . This means that there are two disjoint sets of messages,  $M_L$  and  $M_R$ , such that players of type  $u \leq 1/2$  send a message in  $M_L$  and players of type  $u > 1/2$  send a message in  $M_R$ . Also whenever two players both send messages in  $M_L$  they then play  $L$  and if both send messages in  $M_R$  they both play  $R$ . The left tendency  $\alpha = \alpha^\sigma$  then describes how moderate players coordinate if one of them sends a message from  $M_L$  and the other sends a message from  $M_R$ . The left tendency  $\alpha$  is then the probability that the two players coordinate on  $L$  (through random message selection

within the respective sets of messages), while  $1 - \alpha$  is then the remaining probability that they coordinate on  $R$ .

To prove the “only if” part of the proposition, consider an arbitrary left tendency of  $\alpha \in [0, 1]$ . Then consider a player of type  $1/2$  who needs to be indifferent between sending a message in  $M_L$  and sending a message in  $M_R$  for this strategy to be an equilibrium strategy. If she sends a message in  $M_L$  she coordinates on  $L$  whenever either her opponent sends a message in  $M_L$  (which happens with probability  $F(1/2)$ ), or her moderate opponent sends a message in  $M_R$  (which happens with probability  $F(1) - F(1/2)$ ) and the joint lottery yields the outcome  $L$  (which happens with probability  $\alpha$ ). By contrast, she coordinates on  $R$  whenever her opponent sends a message in  $M_R$  and the joint lottery yields the outcome  $R$  (which happens with probability  $1 - \alpha$ ). Therefore, her expected payoff from sending a message in  $M_L$  is given by

$$\frac{1}{2}F\left(\frac{1}{2}\right) + \frac{1}{2}\alpha\left(F(1) - F\left(\frac{1}{2}\right)\right) + \frac{1}{2}(1 - \alpha)\left(1 - F\left(\frac{1}{2}\right)\right).$$

Similarly, her expected payoff from sending a message in  $M_R$  is given by

$$\frac{1}{2}\alpha F\left(\frac{1}{2}\right) + \frac{1}{2}(1 - \alpha)\left(F\left(\frac{1}{2}\right) - F(0)\right) + \frac{1}{2}\left(1 - F\left(\frac{1}{2}\right)\right).$$

It is easily verified that her expected payoff from sending a message in  $M_L$  is equal to her expected payoff from sending a message in  $M_R$  if and only if  $\alpha = F(0)/(F(0) + (1 - F(1)))$ , as required. This proves the “only if” direction.

To prove the “if” direction of this proposition, we need to show that a coordinated and mutual-preference consistent strategy with strongly binary communication and with a left tendency of  $\alpha = F(0)/(F(0) + (1 - F(1)))$  is both an equilibrium and a renegotiation-proof strategy. To prove the latter condition the same arguments as in the relevant parts of the proof of the “if” direction of Theorem 1 apply directly. It remains to show that such a strategy is an equilibrium strategy. We have already shown that the message function is a best reply to itself and the action function. All that remains to prove is that the action function is a best reply to the given strategy. It is easy to see that playing  $L$  is the optimal strategy when both players send a message in  $M_L$  and thus are of type  $u < 1/2$ . In doing so, they coordinate on their most preferred outcome with probability one. Similarly, playing  $R$  after two messages in  $M_R$  is clearly optimal. Now suppose that one player sends a message in  $M_L$  and the other player sends a message in  $M_R$ . There are two possibilities. Either they are now supposed to both play  $L$  (unless they are an extreme  $R$  type) or they are now supposed to both play  $R$  (unless they are an extreme  $L$  type). Consider first the person who sends a message in  $M_L$  and therefore be of type  $u < 1/2$ . Suppose that the two players are expected to coordinate on  $R$ . Since her opponent sent a message in  $M_R$ , our  $M_L$  sender expects  $R$  with a probability of one (as all  $M_R$  senders are of type  $u > 1/2$ , which excludes  $L$ -dominant action types). But then our  $M_L$  sender of any type  $u > 0$  has a strict incentive to play  $R$  as well. Now suppose that the two players are expected to coordinate on  $L$ . Then our  $M_L$  sender expects her opponent

to play  $L$  with a probability of  $(F(1) - F(1/2)) / (1 - F(1/2))$ , which is the conditional probability of an  $R$ -type to be moderate, which by assumption is greater than or equal to  $1/2$ . Playing  $L$  in this case is therefore optimal for all  $M_L$  senders. That the  $M_R$  sender has the correct incentives in her choice of action after any mixed-message pair (one in  $M_R$  and one in  $M_L$ ) is proven analogously and requires the assumption that  $F(0)/F(1/2) \leq 1/2$ .  $\square$

We now show when and how one can implement a coordinated, mutual-preference consistent strategy with binary communication that has the required left tendency of  $\alpha = F(0)/(F(0) + (1 - F(1)))$ . This implementation requires two things. First,  $\alpha$  needs to be a rational number, and second, the finite message space needs to be sufficiently large.<sup>27</sup> Note that in the case of a symmetric distribution  $F$  (i.e.,  $F(0) = 1 - F(1)$ ) the required left tendency is exactly  $\alpha = 1/2$  and the required strategy is  $\sigma_C$ , as described in Section 3. More generally, let  $\alpha = k/n$ , and assume that  $|M| \geq 2n$ . Denote  $2n$  distinct messages as  $\{m_{L,1}, \dots, m_{L,n}, m_{R,1}, \dots, m_{R,n}\} \in M$ , where we interpret sending messages  $m_{L,i}$  as expressing a preference for  $L$  and sending messages  $m_{R,i}$  as expressing a preference for  $R$  and choosing at random the number  $i$  from the set of numbers  $\{1, \dots, n\}$  in the joint lottery described below. We arbitrarily interpret any message  $m \in M \setminus \{m_{L,1}, \dots, m_{L,n}, m_{R,1}, \dots, m_{R,n}\}$  as equivalent to  $m_{L,1}$ . Given message  $m \in M$ , let  $i(m)$  denote its associated random number, e.g.,  $i(m_{L,j}) = j$ . Let  $M_R = \{m_{R,1}, \dots, m_{R,n}\}$  and  $M_L = M \setminus M_R$ .

Then  $\sigma_\alpha = (\mu_\alpha, \xi_\alpha)$  can be defined as follows:

$$\mu_\alpha(u) = \begin{cases} \frac{1}{n}m_{L,1} + \dots + \frac{1}{n}m_{L,n} & u \leq \frac{1}{2} \\ \frac{1}{n}m_{R,1} + \dots + \frac{1}{n}m_{R,n} & u > \frac{1}{2}. \end{cases}$$

Thus, the message function  $\mu_\alpha$  induces each agent to reveal whether her preferred outcome is  $L$  or  $R$ , and to uniformly choose a number between 1 and  $n$ . In the second stage, if both agents share the same preferred outcome they play it. Otherwise, moderate types coordinate on  $L$  if the sum of their random numbers modulo  $n$  is at most  $k$ , and coordinate on  $R$  otherwise. Extreme types play their strictly dominant action. Formally:

$$\xi_\alpha(m, m') = \begin{cases} 0 & (m, m') \in M_R \times M_R \text{ or } [(m, m') \notin M_L \times M_L \text{ and } i(m) + i(m') \pmod n > k] \\ 1 & \text{otherwise.} \end{cases}$$

---

<sup>27</sup>The method for implementing a binary joint lottery of  $\alpha$  and  $1 - \alpha$  is based on [Aumann and Maschler \(1968\)](#). In order to deal with irrational  $\alpha$ -s one needs either to slightly weaken the result to show that there exists a renegotiation-proof  $\epsilon$ -equilibrium strategy (in which each type of each player gains at most  $\epsilon$  from deviating) for any  $\epsilon > 0$ , or to allow an infinite set of messages or a continuous “sunspot.”