

# Language and Coordination Games

Pei-yu Lo  
Brown University

March, 2007

## Abstract

Intuitively, if players can communicate, they should be able to reach coordinated play in a coordination game. However, simply adding a communication stage before the play of the game does not render coordination as a unique prediction. To further refine the set of equilibria, Farrell suggested that a self-committing cheap talk statement about one's planned behavior should be believed and thus would for sure lead to coordinated play. Aumann, however, argued that the statement has to be both self-committing and self-signaling for it to guarantee coordination. In this paper, the concept of common knowledge of language is formally incorporated into the cheap talk extension game. This paper shows that, if the stage game satisfies both the self-committing and the self-signaling condition, then every iteratively admissible outcome in the language game constitutes a coordinated play and gives the Sender her Stackelberg outcome. On the other hand, this paper identifies a class of generic games that violate self-signaling condition where every strategy profile of the stage game is an iteratively admissible outcome of the language game. This result supports Aumann's argument that the self-signaling condition is necessary for coordinated play to be guaranteed by one-sided communication.

# 1 Introduction

This paper applies the idea of common knowledge of language to complete-information games with one-sided communication. There is a debate in the literature over what criterion for a cheap talk statement makes it credible. Farrell (1988) argues that a cheap talk statement about one's planned behavior is credible if it is *self-committing*, that is, if the speaker believes that the statement will be believed, she will have the incentive to carry it out. A self-committing statement should be believed because, if the speaker is sure that it will be believed, the speaker will indeed carry it out. Aumann (1990), on the other hand, argues that self-committing criterion is not enough; a credible cheap talk statement about one's planned behavior has to be *self-signaling* as well, that is, the speaker would want it to be believed only if she indeed plans to carry it out.

The difficulty in formalizing the credibility criterion lies in how to incorporate the strategy of the hypothetical speaker who intends to *not* carry out her statement into the analysis. Baliga and Morris (2002) tackle this problem by expanding the original game into one in which the Sender has private information. In this expanded game, each action that the Sender may take in the original game is the dominant action of one Sender type. Given any claim about planned behavior, every type of Sender whose dominant action is not equal to this claimed action represents a hypothetical speaker who intends *not* to carry out her claim. This transforms the question of when the Sender could credibly transmit information about her intended action into the question of when a fully-separating perfect Bayesian equilibrium exists, i.e. an equilibrium where the informed player fully reveals her type. Since the common prior puts positive weight on every Sender type, the strategy of every Sender type has to be taken into consideration by the Receiver in a perfect Bayesian equilibrium. Baliga and Morris (2002) show that the self-committing condition alone is not sufficient for establishing a credible Sender claim by demonstrating that there is no communication in a class of games which are self-committing but not self-signaling. In this class of games, the Receiver has only two actions. The self-signaling condition is violated in this class of games because the Sender's preference of the Receiver's actions is independent of her own actions.

We notice that every Sender action in the original stage game that is not strictly dominated is associated with a belief about the Receiver's actions, and that every rationalizable Sender action is a best response to a possibly mixed Receiver action. In addition, if the Receiver puts positive weight on every belief that the Sender holds, the Receiver has to take into account the strategy

of every hypothetical Sender with different intentions. Iterative admissibility is a solution concept with this property.

In this expanded game, an instruction is actually a recommendation for the Receiver to take an action in a specified subset. Thus, an opposite instruction is then a recommendation of actions in the complement of that subset. An instruction is more precise if the recommended subset of Receiver actions is smaller. We assume that the language is rich enough to contain every possible sequence of instructions with increasing precision. Two such sequences may share the first several instructions. So, we can think of the common instructions as the common ancestor of the original sequences. Roughly speaking, if the common ancestor of a pair of such messages contains the common ancestor of another pair as a subsequence, we say that the former pair is more similar to each other than the latter pair. With this relationship, we apply two conditions that define the set of strategies consistent with language: (1) literal meaning condition, i.e.: if the Receiver reacts to a message with a specific action, then he reacts with the same action to the related messages that literally recommends that specific action; (2) convexity condition, i.e.: if the Receiver takes the same action after receiving two different messages, then he takes that same action after any message that may have been delivered with some component of the original message. Our language assumption combined with weak dominance enables messages to convey some information about the Sender's preference regarding the actions of the Receiver.

We focus on stage games where the best response correspondences are functions. The stage game is self-signaling when the Sender always prefers the Receiver to take his best response, and the stage game is self-committing if the Receiver should take an action whenever one is recommended, given that the Sender believes that the recommendation will be followed. With these definitions in mind, we find that if the stage game is self-committing and self-signalling, there is a unique iterative admissible outcome of the language game which gives the Sender her Stackelberg payoff. On the other hand, if the stage game is self-committing, but the Sender's preference over the Receiver's actions does not depend her own action, every rationalizable action profile is the outcome of an iterative admissible strategy profile in the language game.

The rest of this paper is structured as follows. Section 2 provides three simple examples to illustrate the role of the self-signalling condition and to motivate the language assumption. Section 3 describes the model and the language assumptions. Section 4 presents our main results described above. Section 4.3 briefly reviews the main results in Baliga and Morris (2002) and

		Receiver's actions	
		Opera	Club
Sender's actions	Opera	2,1	0,0
	Club	0,0	1,2

Table 1: Battle of Sex Game

	“opera”	“club”
<i>Always Opera</i>	Opera	Opera
<i>Always Club</i>	Club	Club
<i>Literal</i>	Opera	Club
<i>Perverse</i>	Club	Opera

Table 2: Receiver's Strategies in Battle-of-the-Sex Game

compares theirs with ours. Section 5 concludes.

## 2 Motivating Examples

The main idea of this paper is best understood through examples. The battle-of-the-sex game example in section 2.1 illustrates that self-signalling is sufficient to guarantee Stackleberg payoff for the speaker. The investment game example in section 2.2 shows that a severe violation of the self-signalling criterion makes communication ineffective. The partial-common-interest game in section 2.3 motivates the hierarchical messages and language assumptions formally described in section 3.2.

### 2.1 Coordination without positive spillovers

In the Battle-of-the-Sexes game in table 2.1, there are two Nash equilibria: both go to the Opera and both go to the Club. The Sender prefers the first equilibrium and the Receiver prefers the second. The promise “I will go to the opera” is self-committing because if the Sender believes that the Receiver will believe this statement and play his best response *Opera*, the Sender would prefer to go to the *Opera* and carry out her promise. The promise is self-signalling as well because had the Sender not intended to go to the *Opera*, i.e. had she intended to go to the *Club*, she would prefer the Receiver to go to the *Club* instead of the *Opera* and hence she would not want the Receiver to believe the promise “I will go to the *Opera*”.

Suppose  $M = \{\text{“opera”}, \text{“club”}\}$ . It can be interpreted as a promise to carry out a certain action, or a recommended action for the Receiver. The Sender

sends a message  $m \in M$ , and then plays an action  $a^S \in A^S$  in the stage game. The set of strategies for the Sender is thus

$$S^S := \left\{ \begin{array}{l} (\textit{“Opera”}, \textit{Opera}), (\textit{“Club”}, \textit{Club}), \\ (\textit{“Opera”}, \textit{Club}), (\textit{“Club”}, \textit{Opera}) \end{array} \right\},$$

while the set of strategies for the Receiver is listed in table 2.1. Both the *Always Opera* and *Always Club* strategies ignore the messages completely. *Literal* strategy and *Opposite* strategy both respond to a message by going to the *Opera* and the other message by going to the *Club*, and hence are essentially the same up to their renaming.

In addition, We can see from table 2.1 that message “*opera*” and message “*club*” are complete symmetric in the sense that if we swap the names of these two messages, we end up with exactly the same strategy set  $S^R$  as in table 2.1. This should not be surprising because in traditional economic models of communication, messages have no inherent meanings — the meaning is determined by the equilibrium.

However, the idea that messages have no inherent meanings is counter-intuitive. If the Receiver does respond differently to the two messages “*opera*” and “*club*,” it’s generally common knowledge how he is going to respond. Suppose, the Sender says the messages in a very sincere and literal way, it is natural that if the Receiver responds differently to different messages, he will use the *Literal* strategy, not the *Opposite* strategy. Suppose, to the contrary, the Sender says “You’d better go to the Opera” in a sarcastic way. If this sarcasm is commonly understood by the Sender and the Receiver, possibly through the tone or the gesture, then it is natural that there is common knowledge that the Receiver would use the *Opposite* strategy if he decides to respond differently to the two different messages.

If we assume that the two players are both native English speakers and come from the same cultural background, and thus they perfectly understand the meaning that the other person tries to convey from the words uttered, the tone, and the body language, then it is without loss of generality to consider only the sincere tone. Suppose it is common knowledge that the Receiver follows the convention of language and never uses *Opposite*. We are thus describing a different game which I call the language game  $G_L$ , where the set of strategies for the Receiver is

$$S_L^R := \{ \textit{Always Opera}, \textit{Always Club}, \textit{Literal} \}.$$

We will show that the unique outcome that survive three rounds of deletion of weakly dominated strategies is for both players to go to the *Opera*.

In the first round of deletion of weakly dominated strategies, sending the message “opera” and going to the *Club* is weakly dominated for the Sender by sending the message “club” and going to the *Club*. This is because if the Sender is going to the *Club*, she prefers the Receiver to go to the *Club*. If what the Sender says affect what the Receiver does, she gets her preferred action only if she says “club.” Likewise, the strategy (“{*Club*”}, *Opera*) is weakly dominated for the Sender by the strategy (“{*Club*”}, *Opera*).

Therefore, in the second round of deletion of weakly dominated strategies, the strategy *Always Opera* is weakly dominated by the strategy *Literal*, and the strategy *Always Club* is weakly dominated by *Literal* strategy. The only Receiver strategy that survives the second round is thus the *Literal* strategy.

In the third round, the Sender knows that if she says “club,” the Receiver will go to the *Club* and thus it’s best for her to go to the *Club*, and if she says “opera”, the Receiver will go to the *Opera* and thus it’s best for her to go to the *Opera* as well. Since she likes (*Opera, Opera*) better than (*Club, Club*), the optimal strategy for her is to say “opera” and go the the *Opera*. Thus, we obtain the unique outcome that they coordinate on the Sender’s preferred equilibrium.

## 2.2 Coordination with positive spillovers

To understand the role of the self-signalling criterion, let’s look at the Investment game in figure 2.2. As in the Battle-of-the-Sexes game, there are two Nash equilibria in this game: (*Invest, Invest*) and (*Not, Not*). The promise “I’m going to invest” is self-committing because if the Sender believes that the Receiver is going to believe the statement and play his best response, it is optimal for the Sender to carry out the promise and play the strategy *Invest*. In Farrell’s point of view, this message is thus credible and should be believed. Aumann argues that this promise is not self-signalling and hence is not credible. Even if the Sender intends to play *Not*, possibly due to lack of confidence that Receiver is really going to *Invest*, she still prefers the Receiver to use the strategy *Invest*. Therefore, she would like the Receiver to believe her promise regardless of her intended action. If she is pessimistic about the effect of communication and believes that, with high probability, the Receiver is going to *Not Invest* regardless of what she says, then she would prefer to *Not Invest*. However, if the probability that the Receiver uses the strategy *Invest* is higher after hearing the promise “I’m going to invest”, the Sender would like to make that promise even though she does not intend to carry it out.

Let’s look at the cheap talk extension game in detail. Suppose  $M =$

		Receiver	
		<i>Invest</i>	<i>NotInvest</i>
Sender's actions	<i>Invest</i>	2, 2	-1, 1
	<i>NotInvest</i>	1, -1	0, 0

Table 3: Investment Game

	<i>"invest"</i>	<i>"not"</i>
<i>Always Invest</i>	Invest	Invest
<i>Never Invest</i>	Not	Not
<i>Literal</i>	Invest	Not
<i>Opposite</i>	Not	Invest

Table 4: Receiver's Strategies in Investment Game

$\{“invest”, “not”\}$ . Then the set of strategies for the Sender is

$$S^S := \{ (“invest”, *Invest*), (“not”, *Not*), (“invest”, *Not*), (“not”, *Invest*) \},$$

while the set of strategies for the Receiver is listed in table 2.2. Suppose it is common knowledge that the Receiver follows the language convention and never uses the strategy *Opposite*. In the transformed game  $G_L$ , the set of strategies for the Receiver is thus

$$S_L^R := \{ *Always Invest*, *Never Invest*, *Literal* \}.$$

We will now show that every outcome remains after one round of deletion of weakly dominated strategies, when the iterative process stops. Sending the message “not” and using the strategy *Invest* is weakly dominated by sending the message “invest” and using the strategy *Invest*, because when the Sender invests, she prefers the Receiver to invest, and whenever talking affects the outcome, she gets her preferred action only by saying “invest.” Since the Sender has the same preference over the Receiver's actions regardless of the action she takes, the same argument shows that  $(“not”, *Not*)$  is weakly dominated by  $(“invest”, *Not*)$ . Thus, after the first round of deletion, only the message “invest” survives. The process of iterative deletion of weakly dominated strategies stops after the first round, because the Receiver, after receiving the message, still does not know what the Sender is going to play, and thus might play *Never Invest* if he is pessimistic about the Sender's intention, and either *Literal* or *Always Invest* if he is optimistic. After the first round of deletion, the two Receiver strategies, *Literal* and *Always Invest*, are payoff-equivalent for the Receiver because they differ only in the action taken after the message “not”, which is reached with probability zero.

Unlike in the Battle-of-the-Sexes game, when there are positive spillovers, pre-game communication does not eliminate strategic uncertainties. These two examples illustrate the role of the self-signalling criterion.

### 2.3 Partial Common Interest

The Fighting-Couple game in table 5 shows that communication can help players avoid bad equilibria, even though their preferences are not fully aligned. The game has one pure strategy equilibrium:  $(Home, Home)$ . In one of the two mixed strategy equilibria, both go to the *Opera* with probability  $\frac{1}{2}$  and go to the *Club* with probability  $\frac{1}{2}$ . In the other mixed strategy equilibrium, both go to the *Opera* with probability  $\frac{1}{8}$ , go to the *Club* with probability  $\frac{1}{8}$  and stay *Home* with probability  $\frac{3}{4}$ . The mixed-strategy equilibrium where both staying *Home* with probability 0 is the efficient one. Both going to the *Opera* and going to the *Club* is consistent with going out, as opposed to staying home. One prefers to stay *Home* if and only if the other stays *Home*. Moreover, avoiding staying *Home* and restricting themselves to the submatrix  $\{Opera, Club\} \times \{Opera, Club\}$  is mutually beneficial for both players.

Now suppose the Sender has an opportunity to leave a voice message before they play the one-shot game in table 5. She cannot possibly persuade the Receiver to go to the *Opera*, nor can she persuade the Receiver to go to the *Club*, because they have conflict of interest regarding the two actions. However, it is self-committing for her to say “you should go out,” in the sense that if the Receiver is persuaded and goes out, the Sender will go out, i.e., she will choose to either go to the *Opera* or go to the *Club*, in which case, the Receiver prefers to go out. In addition, the suggestion “you should go out” is also self-signalling in the sense that the Sender prefers the Receiver to go out only if she plans to go out herself.

Consider the suggestion “Definitely go out tonight, dear. Regarding where to go, you should go to the opera.” We can write this suggestion as a 2-sequence of decreasing subset:  $\{Opera, Club\} \{Opera\}$ . The suggestion “Definitely go out tonight, dear. Regarding where to go, you should go to the club” is slightly different from the previous one. Another possible suggestion, “You should stay home,” on the other hand, is drastically different from the previous two. We can write this message as “ $\{Home\}$ ”. If the Receiver plays the same action after receiving both suggestions “ $\{Opera, Club\} \{Opera\}$ ” and “ $\{Home\}$ ”, then the Receiver ignores the first layer of literal distinction between going out and staying home. Intuitively, the fine literal difference between the two messages “ $\{Opera, Club\} \{Opera\}$ ” and “ $\{Opera, Club\} \{Club\}$ ” should also



be ignored by the Receiver. That is, the Receiver should play exactly the same action after receiving the message “ $\{Opera, Club\} \{Club\}$ ”, the messages “ $\{Opera, Club\} \{Opera\}$ ” and “ $\{Home\}$ ”. Suppose the set of messages is

$$M = \{“\{Opera, Club\} \{Opera\}”, “\{Opera, Club\} \{Club\}”, “\{Fight\}”\}.$$

Then the preceding discussion suggests the type of language assumption that restricts the Receiver’s strategies to those in table 6.

We will now show that, in the language game  $G_L$ , no player uses the action *Home* after three rounds of deletion of weakly dominated strategies. In the first round of deletion, the strategy (“ $\{Opera, Club\} \{Opera\}$ ”, *Home*) is weakly dominated by (“ $\{Home\}$ ”, *Home*) for the Sender. In addition, Both

$$(“\{Opera, Club\} \{Opera\}”, *Opera*)$$

and (“ $\{Home\}$ ”, *Opera*) are weakly dominated by the Sender strategy

$$(“\{Opera, Club\} \{Club\}”, *Opera*).$$

Therefore, after the first round of deletion of weakly dominated strategies, if the Sender suggests any non-violent action (either *Opera* or *Club*), she definitely plans to not fight; if the Sender suggests to fight, she definitely plans to fight. Thus, it is weakly dominated for the Receiver to fight after a non-violent suggestion; it is also weakly dominated for the Receiver to play a non-violent action after the suggestion to fight. However, it is also weakly dominated in the second round of deletion for the Receiver to play the *Completely Literal* strategy, because they have opposing interests when restricting the game to  $\{Opera, Club\} \times \{Opera, Club\}$ . It can be easily checked that the set of strategies that survive the second round of deletion of weakly dominated strategies is thus  $\{Opera\ Home, Club\ Home\}$ . Therefore, the Sender can be guaranteed a non-violent response if she says either “ $\{Opera, Club\} \{Opera\}$ ” or “ $\{Opera, club\} \{Club\}$ ”. Since she prefer any outcome in the submatrix

$$\{Opera, Club\} \times \{Opera, Club\}$$

to anyone outside of that matrix, the strategy (“ $\{Fight\}$ ”, *Fight*) is strictly dominated in the third round.

The strategy set for the Sender that survives iterative admissibility is

$$\{(“\{Opera, Club\} \{Opera\}”, *Club*), (“\{Opera, Club\} \{Club\}”, *Opera*)\},$$

while the strategy set for the Receiver that survives iterative admissibility is

$$\{Opera\ Home, Club\ Home\}.$$

Pregame communication guarantees both players a payoff of at least 2.

		Receiver's actions		
		<i>Opera</i>	<i>Club</i>	<i>Home</i>
Sender's actions	<i>Opera</i>	2,4	4,2	0,0
	<i>Club</i>	4,2	2,4	0,0
	<i>Home</i>	0,0	0,0	1,1

Table 5: Fighting-Couple Game

	“{opera,club} {opera}”	“{opera,club} {club}”	“{home}”
<i>Always Opera</i>	Opera	Opera	Opera
<i>Always Club</i>	Club	Club	Club
<i>Always Home</i>	Home	Home	Home
<i>Opera &amp; Home</i>	Opera	Opera	Home
<i>Club &amp; Home</i>	Club	Club	Home
<i>Completely Literal</i>	Opera	Club	Home

Table 6: Receiver's Strategies in the Fighting Couple Game

### 3 The Model

In this paper, we focus on one-sided communication extension to finite two-player games with complete information. The Sender (S) and the Receiver (R) simultaneously choose an action  $a^S, a^R$  from a finite set  $A^S$  and  $A^R$  respectively. Their payoffs are given by  $g^S : A^S \times A^R \rightarrow R$  and  $g^R : A^S \times A^R \rightarrow R$ . Write  $g = (g^S, g^R)$ . We will abuse the notation and denote the stage game also by  $g$ . In the one-sided cheap talk extension game  $G$ , the Sender sends a message from a finite set  $M$  before they play the stage game  $g$ . A strategy for the Sender in the reduced-form cheap talk extension game  $G$ , denoted by  $s^S$ , is a message  $m \in M$  and an action  $a^S \in A^S$ . A strategy for the Receiver in  $G$ , denoted by  $s^R$ , is a mapping from  $M$  to  $A^R$ . To characterize the set of communication outcomes with an existing language, we first transform the cheap talk game  $G$  into the language game  $G_L$  by directly restricting the set of strategies for the Receiver to  $S_L^R$ . We motivate and describe the restricted set of Receiver strategies in section 3.2. We then characterize the set of iteratively admissible outcomes in the language game  $G_L$ .

We focus on games where changing exactly one player's action in the action profile changes the payoff. This assumption will be carried throughout the paper.

**Assumption** The stage game payoff  $g^i : A^S \times A^R \rightarrow R$  is such that

$$\begin{aligned} g^i(a^S, a^R) &\neq g^i(a'^S, a^R); \\ g^i(a^S, a^R) &\neq g^i(a^S, a'^R) \end{aligned}$$

for any  $a^S \neq a'^S$ ,  $a^R \neq a'^R$ , and  $i = S, R$ .

This assumption implies that, in particular, the stage-game best response correspondences for both players are well-defined functions. This condition is weaker than genericity, which is a common assumption, and does not exclude any of the motivating games in section 2.

We denote the stage-game best response functions by  $b^i : A^i \rightarrow A^j$  where  $i, j \in \{S, R\}$  and  $i \neq j$ . When we mention “best response” later in the paper, it refers to best response in the language game. Let  $X^S \times X^R$  be a subset of  $S^S \times S^R_L$ , and  $i, j \in \{S, R\}$  where  $i \neq j$ . Two strategies  $s_1^i$  and  $s_2^i$  are equivalent w.r.t.  $X^j$  if they give the same payoff to  $i$  under any of  $j$ 's strategies that belong to  $X^j$ . An observation that will be useful later is that, if every Receiver strategy in  $X^R$  is constant on a message subset  $E$ , then no two Sender strategies that both use messages in  $E$  but take different actions can be equivalent w.r.t.  $X^R$ . A strategy  $s^i$  is said to be a strict best response to  $\sigma^j$  w.r.t.  $X^S \times X^R$  if

$$U^i(s^i, \sigma^j) \geq U^i(\tilde{s}^i, \sigma^j)$$

for every  $\tilde{s}^i \in X^i$ , and strict inequality holds for every  $\tilde{s}^i$  which is not equivalent to  $s^i$  w.r.t.  $X^j$ .

### 3.1 Solution Concept

For ease of exposition, we rewrite here the definition of the solution concept of iterative admissibility taken from Brandenburger et al (2004).

**Definition 1** Fix  $(X^j)_{j \in I} \subseteq (S^j)_{j \in I}$ . A strategy  $s^i$  is weakly dominated with respect to  $X^{-i}$  if there exists  $\hat{\sigma}^i \in \Delta X^i$  such that  $U^i(\hat{\sigma}^i, s^{-i}) \geq U^i(s^i, s^{-i})$  for every  $s^{-i} \in X^{-i}$  and that  $U^i(\hat{\sigma}^i, \hat{s}^{-i}) > U^i(s^i, \hat{s}^{-i})$  for some  $\hat{s}^{-i} \in X^{-i}$ . Otherwise, say that  $s^i$  is admissible with respect to  $(X^j)_{j \in I}$ . If  $s^i$  is admissible w.r.t.  $(S^j)_{j \in I}$ , simply say that  $s^i$  is admissible.

**Definition 2** Set  $S^i(0) = S^i$  for  $i \in I$  and iteratively define

$$S^i(k+1) = \left\{ s^i \in S^i(k) : \begin{array}{l} s^i \text{ is not weakly dominated with respect to } (S^i(k))_{i \in I} \end{array} \right\}.$$

Write  $\cap_{k=0}^{\infty} S^i(k) = S^i(\infty)$  and  $\cap_{k=0}^{\infty} S(k) = S(\infty)$ . A strategy  $s^i \in S^i(\infty)$  is called iteratively admissible.

Denote by  $\Delta X$  the set of probability distribution on  $X$ , and by  $\Delta^+ X$  the set of probability distribution which puts positive weight on every element of  $X$ .

Brandenburger et al (2004) show that if there are only two players, say player  $S$  and player  $R$ , a strategy is weakly dominated if and only if it is never a best response to a totally mixed strategy. For completeness of arguments, this equivalence result is restated as Lemma 1 below.

**Lemma 1 (Brandenburger et al (2004))** *A strategy  $\hat{s}^i \in X^R$  where  $i \in \{S, R\}$  is admissible with respect to  $X^S \times X^R$  if and only if there exists  $\hat{\sigma}^j \in \Delta^+ S^j$  where  $j \neq i$  such that  $U^R(\hat{\sigma}^S, \hat{s}^R) \geq U^R(\hat{\sigma}^S, s^R)$  for every  $s^R \in X^R$ .*

For our purpose we strengthen the characterization of admissible strategies as follows.

**Corollary 1** *A strategy  $\hat{s}^i \in X^R$  where  $i \in \{S, R\}$  is admissible with respect to  $X^S \times X^R$  if only if there exists  $\hat{\sigma}^j \in \Delta^+ S^j$  where  $j \neq i$  to which  $\hat{s}^i$  is a strict best response w.r.t.  $X^S \times X^R$ .*

### 3.2 Incorporating Language

Consider a language  $L$ . Suppose  $L$  contains an expression for a subset of Receiver actions  $B$ . Denote this expression by  $\xi_0$ . If  $L$  also contains an expression for logical negation “not,” then  $L$  contains an expression for the idea “do not do  $B$ .” Denote this expression by  $\xi_1$ . In another language  $L'$ , the expression  $\xi_0$  may mean “please do  $B$ ,” while the expression  $\xi_1$  may mean “please do not do  $B$ .” Since messages are costless and are only means to convey information, it does not matter which language the Sender and the Receiver are speaking, as long as it is common knowledge that they speak the same language. Suppose the common language that the Sender and the Receiver speak is  $L$ . If the Receiver decides to ignore the Sender’s messages, whatever the Sender says does not matter and the Receiver takes the same action regardless. If the Receiver decides to respond to  $\xi_0$  and  $\xi_1$  differently because he thinks the Sender conveys information through her messages, he refers to his own knowledge  $L$  and responds to message  $\xi_0$  with an action in  $B$  while to message  $\xi_1$  with an action not in  $B$ .<sup>1</sup>

---

<sup>1</sup>Some may argue that the Receiver would want to take the *Opposite* strategy in a matching penny game. However, if the Sender knows the payoff structure of a matching penny game, and if she knows that the Receiver uses language  $L'$  in a matching penny game, she will give recommendations according to  $L'$ , thereby destroying the incentive for the Receiver to use language  $L'$ . In this case, one may argue that the Receiver randomizes his actions after

In the Battle-of-the-Sexes game, if we let  $B$  refer to going to the opera, then “not  $B$ ,” i.e. “not go to the opera,” is equivalent to “go to the club,” since this is the only choice other than going to the opera. In the Fighting-Couple game, the expression, “go out”, is saying exactly the same thing as, “go to the opera or go to the club”, and the expression, “do not go out”, says the same thing as, “go home”. If the Receiver responds differently to the two recommendations, “go out”, and, “go home”, then we see from previous discussion that he responds to, “go out”, by going out. However, the Receiver still has to decide whether to go to the opera or the club. Carrying this idea forward, lets suppose that the subset of Receiver actions  $B$  contains a strict subset  $B_2$ , and the language  $L$  contains an expression for  $B_2$ . Suppose further that  $L$  contains an expression  $\xi_{00}$  that is simply a concatenation of  $\xi_0$  and the expression for  $B_2$ . Then with the expression for logical negation,  $L$  contains an expression  $\xi_{01}$  which is the concatenation of  $\xi_0$  meaning “do not do  $B$ ”, and the expression for, “within  $B_1$ , do not do  $B_2$ ”. That the receiver may decide that the messages, “go out”, and, “go home”, convey separate information, but decide to ignore the finer differences between, “go out; furthermore, go to the opera”, and “go out; and then go to the club’. Then the Receiver takes the same action after receiving both message  $\xi_{00}$  and  $\xi_{01}$ . However, if the Receiver decides not to ignore the finer difference between the recommendations  $\xi_{00}$  and  $\xi_{01}$ , he refers to his knowledge of  $L$  and responds to message  $\xi_{00}$  with action  $B_2$  and to message  $\xi_{01}$  with an action in  $B_1$  but not in  $B_2$ .

Let  $M$  denote every message that the Sender could possibly utter. We also assume that the language  $L$  the players commonly speak contains an expression for every subset of Receiver actions, an expression for logical negation, and an expression for concatenation. Then, the language  $L$  contains an expression for every strictly decreasing sequence of subsets of Receiver actions  $A_1 A_2 \dots A_n$ . As a convention, let  $A_0 = A^R$  and  $A_{n+1} = \emptyset$ . Each sequence can be seen as a sequence of instructions with finer and finer details. The set of all such sequences where the last subset has only one element is called the set of hierarchical receiving a message. This could be achieved by randomizing between the *Always B* strategy, and the *Never B* strategy.

Some also argue that the Receiver plays the *Literal* and *Opposite* strategies at the same time. This argument is supported by observing the game being played many times. Throughout these observations, there are incidents where the Receiver takes action  $B$  after message  $\xi_0$  and after message  $\xi_1$ . There are also incidents where the Receiver takes action not in  $B$  after message  $\xi_0$  and after message  $\xi_1$ . These observations do not refute the hypothesis that the Receiver does not play the *Opposite* strategy because all of the aforementioned outcomes may be realizations of a Receiver strategy that randomizes between *Always B* and *Never B*. Finally, in a matching penny game, the Receiver actually has no incentive to respond differently with the Sender’s messages, because he knows that the Sender will not convey any information about her intention.

recommendations, denoted by  $M^h$ . Given a hierarchical recommendation  $m = A_1 \dots A_n$ , we call  $A_j$  the  $j^{th}$  level of instruction. Define  $M(A_1 \dots A_j)$  to be the set of all messages that start with the strictly decreasing sequence  $A_1 \dots A_j$ . Every message  $m$  in  $M(A_1 \dots A_j)$  express the same idea of “Do  $A_1$ . Further more, take an action in  $A_2$ . ....To be even more precise, do  $A_j$ ”.

Let  $s^R$  be a mapping from the set of messages  $M$  to the set of Receiver actions  $A^R$ , and  $m = A_1 \dots A_n$  a hierarchical recommendation. Let  $a^R = s^R(m)$ . Let  $\gamma$  be the highest level of instruction the action  $a^R$  is consistent with according to  $m$ . That is,  $a^R \in A_\gamma \setminus A_{\gamma+1}$ . Therefore, within the subset of  $A_\gamma$ , the action  $a^R$  is “opposite to” the instruction of  $A_{\gamma+1}$ . Our previous discussion suggests that, if  $s^R$  is a language-based Receiver strategy, then either  $s^R$  takes the same action after both expressions for  $A_1 \dots A_\gamma (A_{\gamma+1})$  and expressions for  $A_1 \dots A_\gamma (A_\gamma \setminus A_{\gamma+1})$ , or  $s^R$  responds to expressions for  $A_1 \dots A_\gamma (A_{\gamma+1})$  with actions in  $A_{\gamma+1}$  and expressions for  $A_1 \dots A_\gamma (A_\gamma \setminus A_{\gamma+1})$  with actions in  $A_\gamma \setminus A_{\gamma+1}$ . Therefore, if  $s^R$  is language-based and  $s^R(m) \in A_\gamma \setminus A_{\gamma+1}$ ,  $s^R$  must ignore differences between  $A_{\gamma+1}$  and  $A_\gamma \setminus A_{\gamma+1}$ , and takes the same action after both expressions for  $A_1 \dots A_\gamma (A_{\gamma+1})$  and expressions for  $A_1 \dots A_\gamma (A_\gamma \setminus A_{\gamma+1})$ . That is,  $s^R$  takes the “opposite” action to the instruction  $A_{\gamma+1}$ . The preceding discussion suggests that, if  $s^R$  is a language-based Receiver strategy, then

$$s^R(m') = s^R(m)$$

for every message  $m' \in M(A_1 \dots A_\gamma A_{\gamma+1}) \cup M(A_1 \dots A_\gamma (A_\gamma \setminus A_{\gamma+1}))$ . Call the set

$$M(A_1 \dots A_\gamma A_{\gamma+1}) \cup M(A_1 \dots A_\gamma (A_\gamma \setminus A_{\gamma+1}))$$

the constrained message set given message  $m$  and action  $a^R$ . Formally, given  $m \in M$  and  $a^R \in A^R$ , define

$$M^{cstr}(m, a^R) \equiv \begin{cases} M(A_1 \dots A_\gamma A_{\gamma+1}) & \text{if } m \in M^h \text{ and } \gamma \text{ such that} \\ \cup M(A_1 \dots A_\gamma (A_\gamma \setminus A_{\gamma+1})) & a^R \in A_\gamma \setminus A_{\gamma+1} \\ m & \text{otherwise} \end{cases} .$$

Now we formally define our language assumptions.

**Definition 3**  $s^R : M \rightarrow A$  is a language-based Receiver strategy, denoted by  $s^R \in S_L^R$ , if and only if  $s^R$  is constant on  $M^{cstr}(m, s^R(m))$ , for every  $m \in M$ .

This definition is best illustrated with graphs. Suppose  $A^R = \{A, B, C, D\}$ . Figure 2 shows some hierarchical recommendations in this game. There are many different ways to group  $A^R$ . The first level of instruction can be about taking action  $D$  or not taking action  $D$ , as shown by the two branches  $\{D\}$

and  $\{A, B, C\}$  that diverge from each other. The first layer of instruction can also tell you whether to take actions in  $\{A, C\}$  or not, as shown by the two branches  $\{A, C\}$  and  $\{B, D\}$ . Given a first-layer instruction  $\{A, B, C\}$ , the second layer of instruction could be about whether to take action  $A$  or not, as shown by the two branches  $\{A\}$  and  $\{B, C\}$  that diverge one node on the branch of  $\{A, B, C\}$ . In general, expressions that are “opposite to” each other at some level of instruction are drawn to diverge from the same node. We call all the messages that diverge from the same node a message bundle. A message branch  $A_1 \dots A_j$ , on the other hand, consists of all messages that start with instructions  $A_1 \dots A_j$ . For example, all the messages in the circle in figure 2 constitutes message branch  $\{A, B, C\}$ , while the message “ $\{D\}$ ” combined with all messages in the red circle constitute a message bundle. There can be several parallel message bundles on a branch, which represent different ways to subdivide the set of Receiver actions relevant for the branch.

Suppose we choose the branch of  $\{A, B, C\}$  and then choose  $\{B\}$ , we end up with the message “ $\{A, B, C\} \{B\}$ .” The set  $M(\{A, B, C\} \{A, C\})$  consists of two messages: “ $\{A, B, C\} \{A, C\} \{A\}$ ” and “ $\{A, B, C\} \{A, C\} \{C\}$ .” Within the broad instruction  $\{A, B, C\}$ , these two messages are both “opposite to” message “ $\{A, B, C\} \{B\}$ .”

Given message  $\{A, B, C\} \{B\}$  and Receiver action  $C$ , we first find that action  $C$  belongs to the subset  $\{A, B, C\}$  but not to the subset  $\{B\}$ . According to the definition, the constrained message set given  $m$  and  $a^R$  is thus the following set of messages:

$$\{ \text{“}\{A, B, C\} \{B\}\text{”}, \text{“}\{A, B, C\} \{A, C\} \{A\}\text{”}, \text{“}\{A, B, C\} \{A, C\} \{C\}\text{”} \}.$$

If a language-based Receiver strategy  $s^R$  responds to message “ $\{A, B, C\} \{B\}$ ” with action  $C$ , then by definition,  $s^R$  takes action  $C$  after receiving message  $\{A, B, C\} \{B\}$ , “ $\{A, B, C\} \{A, C\} \{A\}$ ” and “ $\{A, B, C\} \{A, C\} \{C\}$ .”

We call  $M^{cstr}(m, b^R(a^S))$  the constrained message set of the pure Sender strategy  $(m, a^S)$ . Every message in  $M^{cstr}(m, b^R(a^S))$  recommends actions in  $A_j$  which includes the Receiver’s best response to  $a^S$ . Let  $(m_1, a_1^S)$  and  $(m_2, a_2^S)$  be two Sender strategies with overlapping constrained message sets where  $b^R(a_1^S) \neq b^R(a_2^S)$ . It is easy to see that the constrained message set of one Sender strategy must contain that of the other. Let  $A_j$  be the last common level of instruction in the larger of the two constrained message sets. In a way, both message  $m_1$  and  $m_2$  share the recommendation of  $A_j$ . If we see all messages in the larger constrained message set as one compound message, then these two Sender strategies essentially use the same message for different

intentions. It is immediate from the definition of constrained message set that no Receiver strategy consistent with language can be a best response to both  $(m_1, a_1^S)$  and  $(m_2, a_2^S)$ .

On the other hand, if  $(m_1, a_1^S)$  and  $(m_2, a_2^S)$  has disjoint constrained message set and  $b^R(a_1^S) \neq b^R(a_2^S)$ , given any Receiver strategy in a certain class, we can modify it in a minimal way so that it is a best response to both  $(m_1, a_1^S)$  and  $(m_2, a_2^S)$ .

**Lemma 2** *Let  $(m_1, a_1^S)$  and  $(m_2, a_2^S)$  be two Sender strategies in  $S^S(j)$  with disjoint constrained message set where  $b^R(a_1^S) \neq b^R(a_2^S)$ . Let  $E$  and  $B$  be the smallest message branch and message bundle respectively, containing the constrained message sets of both Sender strategies. Then*

1. *there exists a Receiver strategy in  $S^R(j+1)$  non-constant on  $B$ , and*
2. *there exists a mapping  $\psi_B : S^R(j) \rightarrow S^R(j)$  such that  $\psi_B(s^R) = s^R$  for every  $s^R \in S^R(j)$  constant on  $B$ , while for  $s^R \in S^R(j)$  non-constant on  $B$ ,  $\psi_B(s^R)$  is equal to  $s^R$  outside of  $E$  and  $\psi_B(s^R)$  responds to message  $m_i$  with  $b^R(a_i^S)$ ,  $i = 1, 2$ .*

## 4 Results

In this section we generalize the intuition gained from the contrast between the Battle-of-the-Sex game and the Investment game. Section 4.1 gives sufficient conditions for one-sided pre-game communication to guarantee coordinated play in a coordination game. Section 4.2 shows that, when the Sender's preference over the Receiver's actions is independent of the Sender's own action, every rationalizable outcome in the stage game is an iteratively admissible outcome of the language game.

### 4.1 A Sufficient Condition to Guarantee Stackelberg Pay-off for the Sender

In Farrell's definition, messages are about intended actions. In this chapter, we focus on messages that serve as recommendations of actions to the Receiver. We can easily translate a message about the speaker's intended action into a recommendation for the Receiver, since the payoff matrix of the stage game is common knowledge, and thus the Receiver can infer from the speaker's claim about her intended action what the speaker wants the Receiver to do. For



example, the message, “I will take action  $a^S$ ”, is equivalent to a recommendation for the Receiver to take his best response to  $a^S$ .

Let  $b^i$  denote the best reply correspondence for player  $i$  in the stage game,  $i = S, R$ . Since we focus on games where changing only one player’s action changes both players’ payoff, the aforementioned best response correspondence  $b^i$  is in fact a function.

For ease of comparison, we re-write the formal definition of the condition of self-committing by Baliga and Morris (2002) in the following. We then give our version of the definition.

**Definition 4 (Baliga and Morris (2002))** *Claim about intended action  $a^S$  is self-committing if  $b^S(b^R(a^S)) = a^S$ .*

**Definition 5** *Recommendation  $a^R \in A^R$  is self-committing if  $b^R(b^S(a^R)) = a^R$ .*

**Definition 6** *The stage game  $g$  is self-committing if every recommendation  $a^R \in A^R$  is self-committing.*

It is straightforward to see that the recommendation  $a^R$  is self-committing if and only if the claim about intended action  $b^S(a^R)$  is self-committing because  $b^S(b^R(b^S(a^R))) = b^S(a^R)$ .

The definition Aumann (1990) gives for self-signalling criterion is as follows. A statement is self-signalling if the speaker would want it to be believed only if it is true. We can thus say that a recommendation is self-signalling if the speaker would want it to be followed only if she plans to take the action which makes the recommendation optimal for the Receiver. This definition implies that the speaker would NOT want her recommendation  $b^R(a^S)$  to be followed if her planned action would not make this recommendation optimal for the Receiver, that is, if she planned to take an action different from  $a^S$ . This suggests that the self-signalling condition is a property on the stage game as whole, not one about individual actions.

Baliga and Morris formalizes the definition as follows.

**Definition 7 (Baliga and Morris (2002))** *The game  $g$  is self-signalling (for the Sender) if  $g^S(a^S, b^R(a^S)) > g^S(a^S, a^R)$  for every  $a^S \in A^S$ , and  $a^R \in A^R$  where  $a^R \neq b^R(a^S)$ .*

The following proposition gives a sufficient condition for the Sender to be guaranteed her Stackelberg payoff.

		Receiver's actions		
		A	B	C
Sender's actions	a	2, 3	1, 2	-1, -9
	b	0, 0	4, 3	-1, 2
	c	1, -9	2, 2	3, 3

Table 7: A Stage Game with Three Receiver Actions

**Proposition 1** *If the stage game  $g$  is self-signalling and self-committing, then any strategy profile  $((m, a^S), s^R)$  that survives iterative deletion of weakly dominated strategies in the language game gives the Sender her Stackelberg payoff, that is,*

$$g^S(a^S, a^R) = \max_{a^S} u^S(a^S, b^R(a^S))$$

for every  $((m, a^S), s^R) \in S_L(\infty)$ .

#### 4.1.1 An Example

To see the main idea behind the proof, it is easy to start with a simple example. The game is shown in table 7. This game is generic and has three pure strategy Nash-equilibria:  $(a, A)$ ,  $(b, B)$  and  $(c, C)$ . It is obvious that every recommendation is self-committing, and this game is self-signaling.

For ease of exposition, we will assume that the Sender can only give hierarchical recommendations that start with either “ $\{B\}$ ” or “ $\{A, C\}$ .” The very top of figure 1 lists every such message. The bottom table in figure 1 lists every Sender strategy that survives the first, third, and fifth round of deletion of weakly dominated strategies. Action  $a$  is listed in the cell at the intersection of row  $S^S(1)$  and column “ $\{A, C\}\{A\}$ ”, while action  $b$  and  $c$  is not listed in that cell. This indicates that taking action  $a$  after sending the recommendation “ $\{A, C\}\{A\}$ ” survives the first round of deletion of weakly dominated strategies, while taking action  $b$  or action  $c$  after sending the recommendation “ $\{A, C\}\{A\}$ ” does not. The table in the middle of figure 1 lists Receiver strategies that survive the  $0^{th}$ , the second and the fourth round of deletion of weakly dominated strategies. For example, the Receiver strategy, *First Layer and A*, shown in the fourth row in the right panel of figure 1, responds to message “ $\{B\}$ ” with action  $B$ , and to both message “ $\{A, C\}\{A\}$ ” and message “ $\{A, C\}\{C\}$ ” with action  $A$ . By definition, every language-based Receiver strategy survives the  $0^{th}$  round of deletion of weakly dominated strategies in the language game.

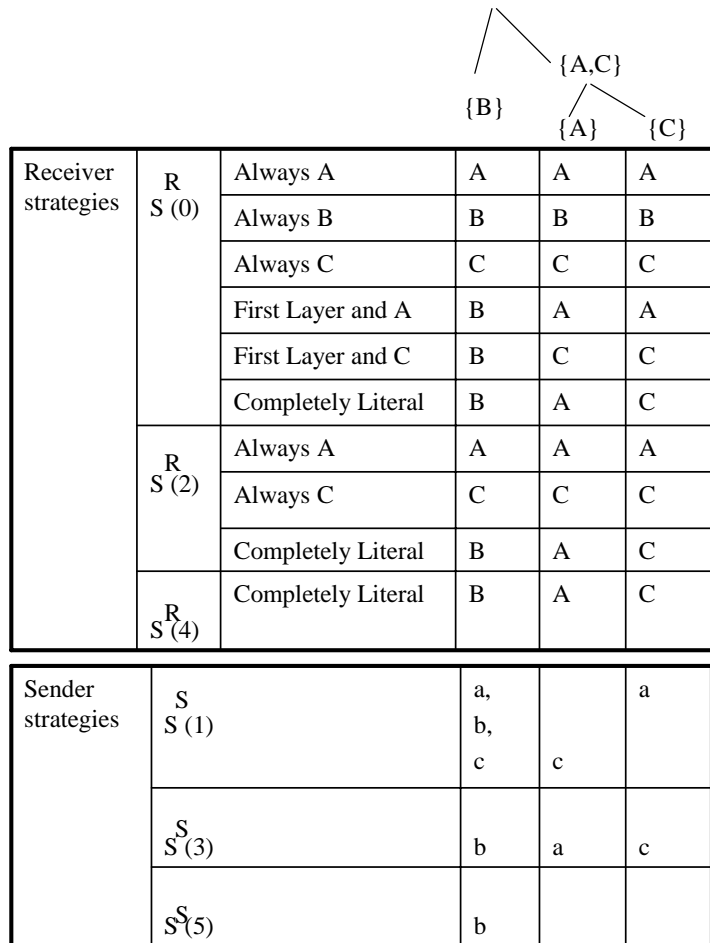


Figure 1: The Iterative Process for a Game with Three Receiver Actions

The Sender strategy that takes action  $c$  after sending the recommendation “ $\{A, C\} \{A\}$ ” does not survive the first round of elimination for the Sender because by the self-signalling condition, the Sender prefers the Receiver action  $C$  to any other Receiver action if the Sender is going to take action  $c$ , and thus taking action  $c$  after sending the recommendation “ $\{A, C\} \{A\}$ ” is weakly dominated by taking the same action  $c$  while sending the recommendation “ $\{A, C\} \{C\}$ .” It can be shown in a similar way that taking action  $b$  after sending the recommendation “ $\{A, C\} \{A\}$ ” is weakly dominated by the Sender strategy that takes action  $b$  but sends the recommendation “ $\{B\}$ .”

However, it is not the case that every Sender strategy that takes an action which makes the recommendation suboptimal for the Receiver is weakly dominated in the first round of deletion. For example, taking action  $c$  after giving the recommendation “ $\{B\}$ ” is not weakly dominated in the first round because it is the best response to the belief

$$(1 - \varepsilon) \textit{Always } C + \varepsilon (1 - \varepsilon) \textit{First Layer and } A + \varepsilon^2 \sigma^R$$

for  $\varepsilon$  sufficiently small and any totally mixed Receiver strategy  $\sigma^R$ .

Proceeding to the second round of deletion, *First Layer and B*, is weakly dominated by the Receiver strategy, *Literal*, in the second round, because 1) these two strategies are not equivalent since they differ only in their response to message “ $\{A, C\} \{A\}$ ” and message “ $\{A, C\} \{A\}$ ” is used by a Sender strategy in  $S^S(1)$ , and 2) every Sender strategy in  $S^S(1)$  that uses message “ $\{A, C\} \{A\}$ ” involves taking action  $a$ , and Receiver action  $A$  does strictly better than Receiver action  $C$  given that the Sender plays action  $a$ .

In the third round of deletion, the Sender strategy (“ $\{B\}$ ”,  $c$ ) is weakly dominated by the Sender strategy (“ $\{A, C\} \{C\}$ ”,  $c$ ) because any Receiver strategy surviving the second round of deletion that takes different action after the two messages responds to message “ $\{A, C\} \{C\}$ ” with action  $C$ , which is most preferred by the Sender given her action  $c$ .

The Sender strategy (“ $\{B\}$ ”,  $b$ ) survives the third round of deletion because it is a best response to the Receiver strategy

$$(1 - \varepsilon) \textit{Completely Literal} + \varepsilon^2 \sigma^R$$

where  $\varepsilon$  is very small and  $\sigma^R$  is any totally mixed Receiver strategy in  $S^R(2)$ . This in turn leads to the deletion of Receiver strategies *Always C* and *Always A* in the fourth round by *Completely Literal*.

It follows that, after four rounds of deletion of weakly dominated strategies, the message “ $\{B\}$ ” will certainly induce action  $B$ . Since the strategy profile  $(b, B)$  gives the Sender her highest payoff, the Sender strategy that sends

message “ $\{B\}$ ” and takes action  $b$  strictly dominates any other Sender strategy remaining after four rounds of deletion of weakly dominated strategies. Therefore, the unique outcome surviving iterative deletion of weakly dominated strategies gives the Sender her Stackelberg payoff.

This example illustrates two points. First, Sender strategies which use a message whose ultimate recommended action is not optimal for the Receiver given the Sender’s intention may still survive the first round of deletion, even though the stage game is a pure coordination game. One such Sender strategy in this particular example is (“ $\{B\}$ ”,  $c$ ). This is because the Sender is afraid that the Receiver may follow those layered recommendations only halfway. Once those Receiver strategies that follow recommendations like “ $\{A, C\} \{C\}$ ” halfway are eliminated, those Sender strategies that do not recommend the Receiver’s best response to the Sender’s intention may subsequently have a chance to be eliminated. Second, to continue the iterative process and eliminate those Sender strategies, we need to show that a Sender strategy that serves as a dominator remains when those Receiver strategies that follow only halfway are eliminated.

#### 4.1.2 The Proof

Denote the size of the set of Receiver strategies  $A^R$  by  $N$ . We can arbitrarily order Receiver actions and write

$$A^R = \{a_1^R, \dots, a_N^R\}.$$

Define  $a_i^S$  to be  $b^S(a_i^R)$ . Let  $\phi$  denote a permutation of  $\{1, 2, \dots, N\}$  and  $\Phi$  the set of all permutations of  $\{1, \dots, N\}$ , with  $id$  being the identity permutation. Define

$$m_{\phi(N-k)} = A_1 \dots A_{N-k-1} \left\{ a_{\phi(N-k)}^R \right\}$$

where  $A_j = A_{j-1} \setminus \left\{ a_{\phi(j)}^R \right\}$  for  $j = 1, \dots, N - k - 1$ . Define  $M_{\phi(N-k)}$  to be the set of hierarchical messages that share in common the first  $N - k - 1$  levels of instruction which eliminates one action at a time from the previous level according to  $\phi$ . Formally, define

$$M_{\phi(N-k)} := \left\{ \begin{array}{l} m = A_1 \dots A_{N-k} A_{N-k+1} \dots A_n | \\ A_j \supseteq A_{j+1}, \forall j = 1, \dots, n-1; \\ A_j = A_{j-1} \setminus \left\{ a_{\phi(j)}^R \right\} \forall j = 1, \dots, N-k \\ n \geq N-k \end{array} \right\}.$$

Figure 2 shows a partial set of hierarchical recommendations. Let  $\phi$  be such that  $a_{\phi(1)}^R = D$ ,  $a_{\phi(2)}^R = B$ ,  $a_{\phi(3)}^R = C$  and  $a_{\phi(4)}^R = D$ . Then the messages in

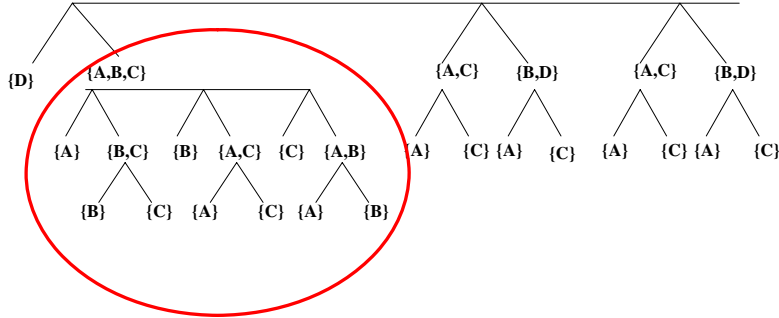


Figure 2: Partial Set of Hierarchical Recommendations,  $A^R = \{A, B, C, D\}$ .

the circle of figure 2 constitute the set  $M_{\phi(1)}$ , while the message “ $\{D\}$ ” is the message  $m_{\phi(1)}$ .

The following observation basically says that, if the Sender intends to take some action  $\hat{a}^S$ , then it is weakly dominated for the Sender to use a message that has, through a process of one-by-one exclusion, recommended the Receiver NOT to take his best response to  $\hat{a}^S$ . For example, it implies that taking action  $d$ , to which the Receiver’s best response is  $D$ , and sending a message in the circle in figure 2 is weakly dominated in the first round by the Sender strategy that takes action  $d$  but uses the message “ $\{D\}$ .”

**Observation** Given a Sender action  $a^S$  and a permutation  $\phi$  and any message  $m \in M_{\phi(q)}$  where  $q$  is such that  $b^R(a^S) = a_{\phi(q)}^R$ , then the Sender strategy  $(m, a^S)$  is weakly dominated in the first round by the Sender strategy  $(m_{\phi(q)}, a^S)$ .

This follows immediately by noticing that every language-based Receiver strategy that takes different actions after receiving message  $m_{\phi(q)}$  and  $m \in M_{\phi(q)}$  responds to message  $m_{\phi(q)}$  with action  $a_{\phi(q)}^R$ , which is most preferred by the Sender given that she takes action  $a^S$  due to the self-signalling condition, and the two messages are not equivalent in the first round. It implies that for every Sender strategy  $(m, a^S)$  surviving the first iteration where  $m \in M_{\phi(q)}$ , its constrained message set is contained in  $M_{\phi(q)}$ .

Communication may fail to achieve coordinated play only if the Receiver makes wrong inference about the Sender’s intention from the message he re-

ceived. This may happen only if the message in use is associated with at least two different Sender actions. Proving that every iteratively admissible outcome gives the Sender her Stackelberg payoff involves ruling out situations like this.

**Definition 8** *A subset of messages  $E$  is intention-clear given a collection of Sender strategies  $\Xi$  if*

1. *the constrained message set of every Sender strategy  $(m, a^S)$  in  $\Xi$  where  $m \in E$  is contained in  $E$ , and*
2. *every pair of Sender strategies in this collection  $\Xi$  that use messages in  $E$  either share the same intention or have disjoint constrained message sets. That is,  $\forall (m_1, a_1^S), (m_2, a_2^S) \in \Xi$ , either  $a_1^S = a_2^S$  or*

$$M^{cstr}(m_1, b^R(a_1^S)) \cap M^{cstr}(m_2, b^R(a_2^S)) = \emptyset.$$

Let  $E$  be either a message bundle or a message branch and  $B$  be the smallest message bundle strictly containing  $E$ . An example is  $B$  as the message set consisting of the circle and the message “ $\{D\}$ ” in figure 2 and  $E$  as the message set consisting of “ $\{B\}$ ” and messages starting with “ $\{AC\}$ ,” or  $E$  is the set of messages in the circle. Suppose  $E$  is intention-clear given  $S^S(k)$ . Lemma 2 implies that, for every Receiver strategy  $s^R$  in  $S^R(k+1)$  that is non-constant on  $B$ ,  $S^R(k+1)$  contains a Receiver strategy  $\psi_E(s^R)$  equal to  $s^R$  outside of  $E$  but is a best response to every Sender strategy in  $S^S(k)$  that uses a message in  $E$ . It follows from intention-clearness of  $E$  given  $S^S(k)$  that if  $s^R$  is not a best response to some Sender strategy in  $S^S(k)$  that uses a message in  $E$ ,  $s^R$  would be weakly dominated by  $\psi_E(s^R)$ , a contradiction to the construction of  $s^R$ . This is stated formally in lemma 3.

**Lemma 3** *Let  $E$  be either a message bundle or a message branch and  $F$  be the smallest message bundle strictly containing  $E$ . If  $E$  is intention-clear given  $S^S(k)$ , then given any Sender strategy  $(m, a^S)$  in  $S^S(k)$  that uses a message in  $E$  and any Receiver strategy  $s^R \in S^R(k+1)$  that is non-constant on  $F$ , we have*

$$s^R(m) = b^R(a^S).$$

We can deduce that  $M_{\phi(N-3)}$  is intention clear given  $S^S(3)$  by lemma 3, observation 4.1.2 and the observation that the Sender strategy  $(m_{\phi(N-l)}, a_{\phi(N-l)}^S)$  belongs to  $S^S(2)$  for every  $l = 0, \dots, N-l$  where  $a_{\phi(N-l)}^S = b^S(a_{\phi(N-l)}^R)$ .

Lemma 4 is the key step in establishing proposition 3.

**Lemma 4**  $\forall k = 1, \dots, N-1$ , the message set  $M_{\phi(N-k-1)}$  is intention-clear given  $S^S(4k-1)$ .

Since  $M \supset M_{\phi(0)}$ , after sufficiently many rounds of deletion, the Sender can convey her intention without fail, and thus every iteratively admissible outcome achieves coordinated play. In particular, this implies that if there exists an admissible Sender strategy that uses a Stackelberg action, then this Sender strategy yields the Stackelberg payoff under any admissible belief. Therefore, if there exists an admissible Sender strategy that uses a Stackelberg action, then every admissible Sender strategy uses a Stackelberg action. Proposition 3 immediately follows because there must exist an admissible Sender strategy that uses a Stackelberg action, since no such Sender strategy can be weakly dominated by a Sender strategy that does not use a Stackelberg action.

We prove lemma 4 by induction. To visualize the inductive process, think of  $M_{\phi(N-k-1)}$  as the messages in the circle in figure 2. Assuming that the circle is intention-clear given  $S^S(4k-1)$ , we first show that the circle plus the message “ $\{D\}$ ” is intention-clear given  $S^S(4k+1)$ , and then show that the set of all messages shown in figure 2 ( $M_{\phi(N-k-2)}$ ) is intention clear given  $S^S(4k+3)$ . We need to show that, given a pair of Sender strategies in  $S^S(4k-1)$  which use messages in  $M_{\phi(N-k-2)}$  and have overlapping constrained message sets, at least one of them is weakly dominated w.r.t.  $S(4k+3)$ . To show that  $(m, a^S)$  is weakly dominated at or before the  $(4k+3)^{th}$  round of deletion, we need to show the existence of a Sender strategy  $(\hat{m}, \hat{a}^S)$  that may weakly dominate  $(m, a^S)$  w.r.t.  $S(4k+3)$ . In the case where  $m = m_{\phi(N-l)}$  for some  $l = 0, \dots, N-1$ , lemma 5 gives conditions for existence of dominators of  $(m, a^S)$  that use messages in  $M_{\phi(N-k-1)}$ . In all other cases, the constrained message set of  $(m, a^S)$  belongs to some message bundle  $E$ . Lemma 6 gives conditions for existence of dominators of  $(m, a^S)$  that use messages in a message set parallel to  $E$ .

To establish existence of dominators of  $(m, a^S)$ , we need to show that the desired potential dominator “outlives”  $(m, a^S)$ . We do so by minimally modifying an anti-best response of  $(m, a^S)$  and show that its Sender best response gives us the desired properties. Given a message bundle  $E$ , let  $B$  denote the smallest message bundle strictly containing  $E$ . We define  $\Psi_B$  to be a mapping from  $S_L^R$  to a subset in  $S_L^R$  such that  $\Psi_B(s^R) = \{s^R\}$  for  $s^R$  constant on  $B$ , while for  $s^R$  non-constant on  $B$ ,  $\Psi_B(s^R)$  is the set of all Receiver strategies in  $S_L^R$  that is equal to  $s^R$  outside of  $B$ .  $\Psi_B$  is a collection of minimal modifications.



Let

$$\Psi_B(\sigma^R) = \left\{ \begin{array}{l} \sigma'^R \in \Delta S_L^R : \\ \forall s^R \in S_L^R, \exists \psi_B(s^R) \in \Psi_B(s^R) \text{ such that} \\ \sigma'^R(\psi_B(s^R)) = \sigma^R(s^R) \end{array} \right\}.$$

Given any  $\sigma^R \in \Delta S_L^R$ ,  $\Psi_B(\sigma^R)$  is the set of strategies “minimally modified” from  $\sigma^R$  that will give us potential dominators that use messages in  $B$ .

Suppose we want a potential dominator of  $(m, a^S)$  to use the same action and send a message in the same message bundle  $E$  as  $m$ , then the highest payoff to the potential dominator given that the Receiver uses strategies in  $\Psi_B(\sigma^R)$  for some  $\sigma^R \in \Delta S_L^R$  is

$$\begin{aligned} \chi(\sigma^R, B, a^S) &\equiv \sum_{\substack{s^R \text{ constant} \\ \text{on } B}} \sigma^R(s^R) g^S(a^S, s^R(m)) \\ &+ \sum_{\substack{s^R \text{ non-constant} \\ \text{on } B}} \sigma^R(s^R) g^S(a^S, b^R(a^S)). \end{aligned}$$

This level is obtained if the Sender uses message  $m$  in  $B$  and the Receiver uses strategy  $\psi(\sigma^R) \in \Psi_B(\sigma^R)$  such that  $\psi(s^R)(m) = b^R(a^S)$  for every  $s^R$  non-constant on  $B$  which receives positive weight from  $\sigma^R$ . In general, the potential dominator uses a message in  $E$  and an action  $\hat{a}^S$  that maximizes  $\chi(\sigma^R, B, a'^S)$  over a subset of Receiver actions. In the iteration where we aim to show weak dominance of  $(m, a^S)$ , typically the payoff to  $(m, a^S)$  given  $\sigma^R$  is strictly lower than  $\chi(\sigma^R, B, a^S)$ , and the potential dominator gives a payoff at least as high as  $\chi(\sigma^R, B, a^S)$ .

To understand lemma 5, we need to define an action-strict best response.

**Definition 9**  $(\hat{m}, \hat{a}^S)$  is an action strict best response to  $\hat{\sigma}^R$  w.r.t.  $S^R(k)$  for some iteration  $k$  if it gives the Sender a strictly higher payoff against  $\hat{\sigma}^R$  than any other Sender strategy that uses some action different from  $\hat{a}^S$ .

**Lemma 5** Suppose a Sender strategy  $(m_{\phi(N-l)}, a^S)$  survives the  $k^{\text{th}}$  round of deletion of weakly dominated strategies, where  $b^R(a^S) = a_{\phi(q)}^R$  and  $q \neq N-l$ . Write  $a^S = a_{\phi(q)}^S$ . If

1. **(Unclear Intention)**  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$  is not intention-clear given  $S^S(k)$ , and
2. **(No Sincere Recommendation)** no Sender strategy that takes action  $a_{\phi(q)}^S$  while using a message in  $M_{\phi(N-l)}$  survives the  $k^{\text{th}}$  round of deletion,

then for every Receiver strategy  $\sigma^R \in \Delta S^R(k-1)$  to which the Sender strategy

$$\left(m_{\phi(N-l)}, a_{\phi(q)}^S\right)$$

is a best response, there exists a Sender strategy  $(\hat{m}, \hat{a}^S)$  in  $S^S(k)$  and a mixed Receiver strategy  $\psi(\sigma^R) \in \Delta S^R(k-1)$  that belongs to

$$\Psi_{m_{\phi(N-l)} \cup M_{\phi(N-l)}}(\sigma^R)$$

such that

A  $\hat{m} \in M_{\phi(N-l)}$  and  $\hat{a}^S \neq a_{\phi(q)}^S$ ;

B  $(\hat{m}, \hat{a}^S)$  is a best response to  $\psi(\sigma^R)$ , and

C

$$\begin{aligned} & \chi\left(\sigma^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_{\phi(q)}^S\right) \\ & \leq \chi\left(\sigma^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \hat{a}^S\right). \end{aligned}$$

To visually grasp the idea in lemma 6, think of message sets  $E$ ,  $F_1$  and  $F_2$  as those denoted in figure 2.

**Lemma 6** *Let  $E$  be a message bundle in  $M_{\phi(N-l)}$ , and  $F_1$  and  $F_2$  be two other parallel message bundles in  $M_{\phi(N-l)}$ . Let  $(\hat{m}, \hat{a}^S)$  be a Sender strategy in  $S^S(k)$  for some iteration  $k$  where  $\hat{m} \in E$ , and  $\hat{\sigma}^R$  be a totally mixed strategy in  $S^R(k-1)$  to which  $(\hat{m}, \hat{a}^S)$  is a best response. If there exists  $s^R \in S^R(k)$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$  such that  $s^R(\hat{m}) \neq b^R(\hat{a}^S)$ , then there exists two Sender strategies  $(m_1, a_1^S)$ ,  $(m_2, a_2^S)$  in  $S^S(k)$  where*

1. message  $m_1$  belongs to  $F_1$ , message  $m_2$  belongs to  $F_2$ , and
2. Sender actions  $a_1^S$  and  $a_2^S$  both maximize

$$\chi\left(\hat{\sigma}^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a^S\right)$$

$$\text{over } \left\{a^S : b^R(a^S) = a_{\phi(i)}^S, \text{ for some } i > N-l\right\}.$$

**Proof of lemma 4.** It is trivially true for  $k=1$  because  $M_{\phi(N-2)} = \{m_{\phi(N-1)}, m_{\phi(N)}\}$ , and if the Sender strategy  $(m_{\phi(N-1)}, a^S)$  belongs to  $S^S(1)$ , then  $a^S = a_{\phi(N-1)}^S$ .

Suppose it is true for  $k=1, \dots, \bar{k}$ .

**Claim 1**  $m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$  is intention-clear given  $S^S(4\bar{k}+1)$ .

**Proof.** Suppose to the contrary that there exists two Sender strategies  $(m_1, a_1^S)$  and  $(m_2, a_2^S)$  in  $S^S(4\bar{k}+1)$  with overlapping constrained message sets where  $m_1, m_2$  belong to  $m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$  and  $b^R(a_1^S) \neq b^R(a_2^S)$ . We can assume w.l.o.g. that  $M^{cstr}(m_1, b^R(a_1^S))$  is larger. By the assumption that  $M_{\phi(N-\bar{k}-1)}$  is intention-clear given  $S^S(4\bar{k}-1)$ , it cannot be the case that both  $m_1$  and  $m_2$  belong to  $M_{\phi(N-\bar{k}-1)}$ . It follows that the constrained message set of  $(m_2, a_2^S)$  must be  $m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$ . That implies  $m_2 = m_{\phi(N-\bar{k}-1)}$  and  $b^R(a_2^S) = a_{\phi(q)}^R$  where  $q \neq N - \bar{k} - 1$ . Observation 4.1.2 further implies that  $q > N - \bar{k} - 1$ . Either there exists  $(m'_2, a_2^S)$  in  $S^S(4\bar{k}+1)$  where  $m'_2 \in M_{\phi(N-\bar{k}-1)}$ , or the assumptions of lemma 5 holds. It follows that, given any  $\sigma_2^R \in \Delta^+ S^R(4\bar{k})$  to which  $(m_2, a_2^S)$  is a strict best response w.r.t.  $S(4\bar{k})$ , there exists a Sender strategy  $(\hat{m}_2, \hat{a}_2^S) \in S^S(4\bar{k}+1)$  where  $\hat{m}_2 \in M_{\phi(N-\bar{k}-1)}$  such that

$$\begin{aligned} & \chi(\sigma_2^R, m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}, \hat{a}_2^S) \\ & \geq \chi(\sigma_2^R, m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}, a_2^S). \end{aligned}$$

The assumption that  $M_{\phi(N-\bar{k}-1)}$  is intention-clear given  $S^S(4\bar{k}-1)$  combined with lemma 3 implies that  $s^R(\hat{m}_2) = b^R(\hat{a}_2^S)$  for every  $s^R \in S^R(4\bar{k})$  non-constant on  $m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$ . It follows that

$$\begin{aligned} & u^S((\hat{m}_2, \hat{a}_2^S), \sigma_2^R) \\ & = \chi(\sigma_2^R, m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}, \hat{a}_2^S) \\ & \geq \chi(\sigma_2^R, m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}, a_{\phi(q)}^S) \tag{1} \\ & \geq u^R((m_{\phi(N-\bar{k}-1)}, a_2^S), \sigma_2^R), \tag{2} \end{aligned}$$

Strict inequality holds in line 2 if  $S^R(4\bar{k})$  contains a Receiver strategy non-constant on  $m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$ . In that case, we obtain a contradiction to the construction of  $\sigma_2^R$ . If  $m_1 = m_{\phi(N-\bar{k}-1)}$  and  $b^R(a_1^S) = a_{\phi(N-\bar{k}-1)}^R$ , then  $(m_1, a_1^S)$  and  $(\hat{m}_2, \hat{a}_2^S)$  have disjoint constrained message sets, and by lemma 2,  $S^R(4\bar{k})$  contains a Receiver strategy non-constant on  $m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$ . If  $b^R(a_1^S) \neq a_{\phi(N-\bar{k}-1)}^R$ , then there exists  $(\hat{m}_1, \hat{a}_1^S) \in S^S(4\bar{k}+1)$  with  $\hat{m}_1 \in M_{\phi(N-\bar{k}-1)}$  that is related to  $(m_1, a_1^S)$  in the same way  $(\hat{m}_2, \hat{a}_2^S)$  is to  $(m_2, a_2^S)$ . The assumption of unclear intention implies that either  $b^R(\hat{a}_1^S) \neq b^R(\hat{a}_2^S)$ , in which case we apply lemma 2 and we are done, or  $b^R(\hat{a}_1^S) = b^R(\hat{a}_2^S)$  for some

$i = 1, 2$ . Assume w.l.o.g. that  $i = 2$ . But then  $\hat{a}_2^S \neq a_2^S$ , and the construction that  $(m_2, a_2^S)$  is a strict best response to  $\sigma_2^R$  implies that at least one strict inequality holds in line 1 and 2. ■

**Claim 2**  $M_{\phi(N-\bar{k}-2)}$  is intention-clear given  $S^S(4\bar{k}+3)$ .

**Proof.** Suppose to the contrary that there exists two Sender strategies  $(m_1, a_1^S)$  and  $(m_2, a_2^S)$  in  $S^S(4\bar{k}+1)$  with overlapping constrained message sets where  $m_1, m_2$  belong to  $M_{\phi(N-\bar{k}-2)}$  and  $b^R(a_1^S) \neq b^R(a_2^S)$ . Claim 1 implies that it cannot be the case that  $m_1$  and  $m_2$  both belong to  $m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$ . Observation 4.1.2 and the definition of *constrained message sets* imply that  $m_1$  and  $m_2$  must belong to the same message bundle in  $M_{\phi(N-\bar{k}-2)}$ . Denote this message bundle by  $E$ . Then from lemma 6, for  $i = 1, 2$ , given  $\sigma_i^R \in \Delta^+ S^R(4\bar{k}+2)$  to which  $(m_i, a_i^S)$  is a strict best response w.r.t.  $S(4\bar{k}+2)$ , there exists  $(\hat{m}_i, \hat{a}_i^S) \in S^S(4\bar{k}+3)$  where  $\hat{m}_i \in m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$  and

$$\begin{aligned} & \chi\left(\sigma_i^R, m_{\phi(N-\bar{k}-2)} \cup M_{\phi(N-\bar{k}-2)}, \hat{a}_i^S\right) \\ & \geq \chi\left(\sigma_i^R, m_{\phi(N-\bar{k}-2)} \cup M_{\phi(N-\bar{k}-2)}, a_i^S\right). \end{aligned}$$

But from claim 1,  $m_{\phi(N-\bar{k}-1)} \cup M_{\phi(N-\bar{k}-1)}$  is intention-clear given  $S^S(4\bar{k}+1)$ . Lemma 3 then implies that  $s^R(\hat{m}_i) = b^R(\hat{a}_i^S)$  for every  $s^R \in S^R(4\bar{k}+2)$  non-constant on  $M_{\phi(N-\bar{k}-2)}$ . Therefore,

$$\begin{aligned} & u^S((\hat{m}_i, \hat{a}_i^S), \sigma_i^R) \\ & = \chi\left(\sigma_i^R, m_{\phi(N-\bar{k}-2)} \cup M_{\phi(N-\bar{k}-2)}, \hat{a}_i^S\right) \\ & \geq \chi\left(\sigma_i^R, m_{\phi(N-\bar{k}-2)} \cup M_{\phi(N-\bar{k}-2)}, a_i^S\right) \quad (3) \\ & \geq u^R((m_i, a_i^S), \sigma_i^R). \quad (4) \end{aligned}$$

There exists  $j \in \{1, 2\}$  such that  $b^R(\hat{a}_1^S) \neq b^R(a_j^S)$  because  $b^R(a_1^S) \neq b^R(a_2^S)$ . Since  $\hat{m}_1$  belongs to a message bundle parallel to  $E$ , the constrained message set of  $(\hat{m}_1, \hat{a}_1^S)$  is disjoint from that of  $(m_j, a_j^S)$ . Lemma 2 implies that  $S^R(4\bar{k}+2)$  contains Receiver strategies non-constant on  $m_{\phi(N-\bar{k}-2)} \cup M_{\phi(N-\bar{k}-2)}$ . Hence inequality 4 holds strictly for both  $i = 1, 2$ , contradiction to the construction of  $\sigma_i^R$ . ■ ■

## 4.2 Games with Positive Spillovers

The self-signalling criterion implies that the Sender's preference over the Receiver's actions differ with her own intention. Our language assumption com-

bined with iterative admissibility connects different messages with different preferences. This then separates one intention from the other and guarantees the Sender her Stackleberg payoff.

It seems natural then that the Sender cannot convey any information about her intention through cheap talk if the Sender's preference over the Receiver's actions is invariant with her own intention.

If the stage game is self-committing, then for every  $a^R \in A^R$ ,  $b^R(b^S(a^R)) = a^R$ . Therefore,  $A^R(1) = A^R$ . It follows that  $A^R(\infty) = A^R$  and  $A^S(\infty) = A^S(1)$  that contains  $b^S(A^R)$ .

**Theorem 1** *If the stage game is self-committing and the Sender's preference over the Receiver's actions is independent of her own action, then for every  $(a^S, a^R) \in A(\infty)$ , there exists  $(m, a^S) \in S^S(\infty)$  and  $s^R \in S^R(\infty)$  such that  $s^R(m) = a^R$ .*

**Proof.** Define  $M(k)$  to be the projection of  $S^S(k)$  onto  $M$ . So  $M(k)$  is the set of messages that are used in  $S^S(k)$ .

**Claim 3**  $s^R(1) = S_L^R$ .

**Proof.** Given  $s^R \in S_L^R$ , let  $\sigma^S$  be a totally mixed Sender strategy in

$$\Delta \{ (m, b^S(s^R(m))) : m \in M \}.$$

For every  $m \in M$ ,  $(m, b^S(s^R(m)))$  is the only Sender strategy given positive weight under  $\sigma^S$ . Therefore,  $s^R$  is a best response to every pure strategy given strictly positive weight by  $\sigma^S$ , and thus is a strict best response to  $\sigma^S$ . It follows that  $s^R \in S^R(1)$ . ■

**Claim 4**  $S^S(1) = M(1) \times A^S(1)$ .

**Proof.** Suppose  $(\hat{m}, \hat{a}^S) \in S^S(1)$ . Write  $\hat{m} = A_1 \dots A_n$ , and define  $A_0 = A^R$ . Then for every  $j = 1, \dots, n$ , there exists  $\bar{a}_j \in A_j$  and  $\underline{a}_j \in A_{j-1} \setminus A_j$  where  $g^S(\hat{a}^S, \bar{a}_j) > g^S(\hat{a}^S, \underline{a}_j)$ . Otherwise,  $(\hat{m}, \hat{a}^S)$  would be weakly dominated by  $(m', \hat{a}^S)$  for any  $m' \in M(A_1 \dots A_{j-1}(A_{j-1} \setminus A_j))$ .

Define a partial relation  $>$  on  $A^R$  by the preference order of the Sender. That is,  $a_2^R > a_1^R$  iff  $g^S(a^S, a_2^R) > g^S(a^S, a_1^R)$ . Define

$$s_1^R(m) = \begin{cases} \bar{a}_1 & m \in M(A_1) \\ \underline{a}_1 & m \in M(A_1^c) \\ \min A^R & \text{otherwise} \end{cases}$$

and

$$s_j^R(m) = \begin{cases} \bar{a}_j & m \in M(A_1 \dots A_{j-1} A_j) \\ \underline{a}_j & m \in M(A_1 \dots A_{j-1} (A_{j-1} \setminus A_j)) \\ s_{j-1}^R(m) & m \notin M(A_1 \dots A_{j-1}) \\ \min A_{j-1} & \text{otherwise} \end{cases}$$

for  $j = 2, \dots, n$ . It follows that, for  $j = 1, \dots, n$ , for every  $a^S$ ,

$$\begin{aligned} u^S((m, a^S), s_j^R) &= g^S(a^S, \bar{a}_j) \\ &> g^S(a^S, \underline{a}_j) \\ &= u^S((m, a^S), s_j^R) \end{aligned}$$

for every  $m \in M(A_1 \dots A_{j-1} (A_{j-1} \setminus A_j))$  and

$$\begin{aligned} u^S((\hat{m}, a^S), s_j^R) &= g^S(a^S, \bar{a}_j) \\ &> g^S(a^S, \min A^R) \\ &= u^S((m, a^S), s_j^R) \end{aligned}$$

for every  $m \notin M(A_1) \cup M(A_1^c)$ . Define

$$\hat{\sigma}_\varepsilon^R := \sum_{j=1}^{n-1} \varepsilon^{j-1} (1 - \varepsilon) s_j^R + \varepsilon^{n-1} s_n^R.$$

Therefore,

$$u^S((\hat{m}, a^S), \hat{\sigma}_\varepsilon^R) > u^S((m, a^S), \hat{\sigma}_\varepsilon^R)$$

for every  $m \neq \hat{m}$ , every  $a^S$  and every  $\varepsilon$  sufficiently small. Let  $a^R$  denote also the constant Receiver strategy that takes the action  $a^R$  upon receiving every message. There exists  $\alpha \in \Delta A^R$  to which  $a^S$  is a best response in the stage game, since  $a^S \in A^S(1)$ . Let  $\alpha$  also denote the Receiver strategy that puts weights  $\alpha(a^R)$  on the constant strategy  $a^R$ . Then for  $\varepsilon$  sufficiently small,

$$\begin{aligned} &u^S((\hat{m}, \tilde{a}^S), (1 - \varepsilon)\alpha + \varepsilon\hat{\sigma}_\varepsilon^R) \\ &> u^S(m, a^S) \end{aligned}$$

for every  $(m, a^S) \neq (\hat{m}, a^S)$ . We have thus established that  $(\hat{m}, a^S) \in S^S(1)$  for every  $a^S \in A^S(1)$ . ■

Since  $S^R(1) = S_L^R$ , it is immediate that  $S^S(\infty) = S^S(1)$  and  $S^R(\infty) = S_L^R$ . We are done. ■

		Receiver's Action	
		Invest	Not
Sender's	Invest	$10+x, 10+x$	$-90, x$
Action	Not	$x, -90$	$0, 0$

Table 8: leading example in Baliga Morris (2002)

		Receiver's Action	
		Invest	Not
Sender's	Invest	$10+x, 10+x$	$-90, x$
Action	Not	$x, -90$	$0, 0$
<i>Low Cost</i>			
		Receiver's Action	
		Invest	Not
Sender's	Invest	$-10+x, 10+x$	$-110, x$
Action	Not	$x, -90$	$0, 0$
<i>High Cost</i>			

Table 9: Incomplete Information Investment Game

### 4.3 Comparison with Baliga and Morris

To formally formulate the role of the self-signalling criterion, Baliga and Morris (2002) transforms the complete information game into a coordination game with incomplete information, and use the solution concept of perfect Bayesian equilibrium. The counterfactual “what would the Sender have said had she intended to play action  $a'$  instead of  $a$ ” does not really have a role in the solution concept of Nash equilibrium in complete information games. However, the solution concept of perfect Bayesian equilibrium addresses the question “what would the Sender have said were she of type  $t'$ ?”

The easiest way to see the comparison is to look at the leading example in Baliga and Morris (2002). The game is shown in table 8.

In this stage game, both action *Invest* and action *Not* are self-committing. If  $x < 0$ , the stage-game is self-signalling, while the game exhibits positive spillovers if  $x > 0$ . To formally study the role of self-signalling, Baliga and Morris (2002) study the following incomplete information game where with probability  $1 - p$  the Sender is of *Low Cost* and with probability  $p > 0$  the Sender is of *High Cost*. The *Low Cost* type has the same payoff matrix as in the complete information game of table 8. However, the *High Cost* Sender has a dominant strategy to not invest. The Receiver's payoff depends only on the action taken by the Sender, not on the Sender's type. Therefore, the Receiver

cares about the type of the Sender only insofar as it conveys information about the action the Sender would take. For example, if the Receiver knew that the Sender is of *High Cost*, the Receiver would infer that the Sender would not invest, and thus his best response would be to not invest. Hence, the hypothetical Sender who intends to not invest is equated with the *High Cost* Sender who has a dominant strategy to not invest. Since the prior puts strictly positive weight on the *High Cost* type, the strategy of the *High Cost* type, or equivalently, the strategy of the hypothetical Sender who intends to not invest, has to be taken into account by the Receiver.

They show that when  $x < 0$ , there exists a perfect Bayesian equilibrium where the *Low Cost* Sender sends a different message from the *High Cost* Sender and both the Sender and the Receiver invest when the Sender is of *Low Cost*, while neither of them invest when the Sender is of *High Cost*. However, when  $x > 0$ , there can be no perfect Bayesian equilibrium where the outcomes are type-dependent. Conditional on the Sender being *Low Cost* type, when  $x < 0$ , there exists an equilibrium where the outcome is *(Invest, Invest)*. On the other hand, when  $x > 0$ , conditional on the Sender being *Low Cost* type, the unique equilibrium outcome is *(Not, Not)* if the probability of the *High Cost* type ( $p$ ) is greater than  $\frac{1}{10}$ , while the equilibrium outcome could be either *(Invest, Invest)* or *(Not, Not)* if  $p < \frac{1}{10}$ . Since the stage-game is self-signalling only when  $x < 0$ , this illustrates the role of the self-signalling criterion.

When  $x < 0$ , our approach predicts that the unique outcome is for both to invest, which coincides with the prediction of Baliga and Morris (2002). When  $x > 0$ , our approach predicts that every action profile is possible. This is natural because there is not a fixed probability attached to the pessimistic Sender who is going to not invest, and players may have incorrect belief about each other.

The formal model in Baliga and Morris (2002) is as follows. The Sender is one of a finite set of possible types  $T$ . The Sender's utility function is  $\tilde{g}^S : A^S \times A^R \times T \rightarrow R$ ; the Receiver's utility function is  $g^R : A^S \times A^R \rightarrow R$ . For ease of comparison, I rewrite the positive result in Baliga and Morris (2002) here.

**Proposition 2** *If (1) for each  $a^S \in A^S$ , there exists a type  $\tau(a^S) \in T$  such that  $a^S$  the dominant strategy for the Sender in the game  $g^S(., t)$ ; and (2) for each action  $a^R \in A^R$ , there exists  $a^S \in A^S$  such that  $a^R = b^R(a^S)$ , then there exists a full revelation perfect Bayesian equilibrium in the one-sided cheap talk game if and only if*



1.  $a^S$  is a self-committing action for the Sender in the game  $g^S(\cdot, \tau(a^S))$ ;
2.  $a^S$  is the Stackelberg action for the Sender in the game  $g^S(\cdot, \tau(a^S))$ ;
3.  $a^S$  is self-signalling for the Sender in the game  $g^S(\cdot, \tau(a^S))$ .

Let's take the complete information stage game  $g = (A^S, A^R, g^S, g^R)$  where  $A^R = \{b^R(a^S) : a^S \in A^S\}$ . Let  $T = A^S$ . For clarity, let  $\tau$  be the bijective function from  $A^S$  to  $T$ . Let  $a_*^S = \arg \max_{a^S} g^S(a^S, b^R(a^S))$ . Then  $a_*^S$  is the Stackelberg action for the Sender in the game  $G$ . Define

$$\bar{d} := \max_{(a^S, a^R) \neq (a'^S, a'^R)} |g^S(a^S, a^R) - g^S(a'^S, a'^R)|.$$

$\bar{d}$  is thus the maximum payoff difference for the Sender. Expand the utility function  $g^S : A^S \times A^R \rightarrow R$  into  $\tilde{g}^S : A^S \times A^R \times T \rightarrow R$  as follows.  $\tilde{g}^S(\cdot, \tau(a_*^S)) = g(\cdot)$ , and for every  $a^S \neq a_*^S$ ,  $\tilde{g}^S(a^S, a'^R, \tau(a^S)) = g^S(a^S, a'^R)$  for every  $a'^R \in A^R$ , and  $\tilde{g}^S(a'^S, a'^R, \tau(a^S)) = g^S(a'^S, a'^R) - 2\bar{d}$ , for every  $a'^S \neq a^S$  and every  $a'^R \in A^R$ . Denote the one-sided cheap talk extension game of  $g$  with language by  $G_L$ , and the one-sided cheap talk extension game of  $\tilde{g}$  by  $\tilde{G}$ . Then proposition 2 implies that  $\tilde{G}$  has a full revelation perfect Bayesian equilibrium if and only if the complete information stage game  $g$  is self-signalling and every action  $a^S$  for the Sender is self-committing. In this equilibrium, the type  $\tau(a_*^S)$ , whose payoff matrix is  $g$ , gets her Stackelberg payoff. Our positive result equivalently states that if every  $a^S$  in  $g$  is self-committing and  $g$  is self-signalling, the unique iterative admissible outcome in  $G_L$  gives the Sender her Stackelberg payoff. The negative result in Baliga and Morris (2002) says that there is no communication in any equilibrium of  $\tilde{G}$  if  $g$  exhibits binary action positive spillovers. Equilibrium outcomes of  $\tilde{G}$  in such games depend on the common prior over  $T$ . If there is no common prior, and we allow any prior over  $T$ , we can span every rationalizable outcome. Our negative result relaxes the condition to any finite games with positive spillovers, and states that every rationalizable outcome is consistent with iterative admissibility in  $G_L$ .

## 5 Conclusion

By modeling the idea that there exists common knowledge of language, we are able to formally distinguish coordination games with and without the self-signaling condition first used by Aumann (1990) within the framework of complete information games. We show that, if the stage game is self-committing and self-signalling, every iterative admissible outcome in the language game achieves

coordination and gives the Sender her Stackelberg payoff. On the other hand, if the stage game is self-committing but the Sender's preference for the Receiver's actions does not depend on her intended action, every rationalizable stage game outcome is also an iteratively admissible outcome in the language game.

## 6 Appendix

### 6.1 Proof for lemma 2

Every Receiver best response to

$$\frac{1}{2} (m_1, a_1^S) + \frac{1}{2} (m_2, a_1^S)$$

responds to message  $m_i$  with  $b^R(a_i^S)$  for  $i = 1, 2$ , and thus is non-constant on  $B$  since  $b^R(a_1^S) \neq b^R(a_2^S)$ . At least one of the best responses survive the  $(j + 1)^{th}$  iteration and statement 1 follows.

Given  $s^R \in S^R(j)$  non-constant on  $B$ , there exists  $\sigma^S \in \Delta^+ S^S(j - 1)$  to which  $s^R$  is a best response. Then every best response to

$$\frac{1}{2} (1 - \varepsilon) \sum_{i=1,2} (m_i, a_i^S) + \varepsilon \sigma^S|_{M \setminus E}$$

gives us the desired property in statement 2.

### 6.2 Proof for lemma 6

We need the following lemma for the proof.

**Lemma 7** *Let  $F_1$ ,  $F_2$ , and  $F_3$  be three parallel message bundles, and  $B$  be the smallest message bundle that strictly contains  $F_1$ . If*

1.  $(m_1, a_1^S) \in S^S(j - 1)$  where

$$M^{cstr}(m_1, b^R(a_1^S)) \subset F_1,$$

2.  $S^S(j - 1)$  contains Sender strategies that use messages in  $F_2$  and  $F_3$  respectively, and
3.  $S^R(j)$  contains Receiver strategies  $s_2^R$  and  $s_3^R$  non-constant on  $B$  such that either  $s_2^R|_{F_2}$  or  $s_3^R|_{F_3}$  is not equivalent w.r.t.  $S^S(j - 1)$  for the Receiver to a constant of  $b^R(a_1^S)$ ,

then there exists a mapping  $\psi : S^R(j) \rightarrow S^R(j)$  that belongs to  $\Psi_B$  such that for every  $s^R \in S^R(j)$  non-constant on  $B$ ,  $\psi(s^R)$  is equal to  $s_i^R$  on  $F_i$  for  $i = 2, 3$ , is equal to  $s^R$  outside of  $E \cup F_1 \cup F_2$ , and is a best response to the Sender strategy  $(m_1, a_1^S)$ .

**Proof.** Let  $\sigma_i^S$  be a totally mixed Sender strategy in  $S^S(j-1)$  to which  $s_i^R$  is a best response, for  $i = 2, 3$ . Let  $\sigma^S|_{\tilde{F}}$  denote the mixed Sender strategy that is the probability distribution of  $\sigma^S$  conditional on sending messages in  $\tilde{F}$ . Given  $s^R \in S^R(j)$  non-constant on  $B$ , let  $\sigma^S$  be a totally mixed Sender strategy in  $S^S(j-1)$  to which  $s^R$  is a best response. Then for  $\varepsilon$  sufficiently small, every Receiver best response to

$$(1 - \varepsilon)(m_1, a_1^S) + \frac{\varepsilon(1 - \varepsilon)}{2}\sigma_2^S|_{F_2} + \frac{\varepsilon(1 - \varepsilon)}{2}\sigma_3^S|_{F_3} + \varepsilon^2\sigma^S|_{M \setminus (F_2 \cup F_3)} \quad (5)$$

responds to message  $m_1$  with action  $b^R(a_1^S)$ , is equivalent to  $s_i^R$  on  $F_i$  w.r.t.  $S^R(j-1)$  for the Receiver, for  $i = 2, 3$ , and is equal to  $s^R$  outside of  $F_1 \cup F_2 \cup F_3$ . Condition 3 implies that such a best response must also be non-constant on  $B$ . Since the Sender strategy in expression 5 is totally mixed on  $S^S(j-1)$ , at least one best response to expression 5 belongs to  $S^R(j)$ . Define  $\psi(s^R)$  to be a best response to expression 5 in  $S^R(j)$  and we are done. ■

If  $l \leq 2$ , then  $M_{\phi(N-l)}$  consists of only one parallel message bundle, and lemma 6 holds automatically. Hence we are concerned only with  $l \geq 3$ . It is easy to see that the statement holds for every  $l = 3, \dots, N$  for  $k = 1$ . Suppose the statement is true for  $l = 3, \dots, N$  and  $k = 1, \dots, \bar{k}$ . Suppose  $S^S(\bar{k} + 1)$  contains a Sender strategy  $(\hat{m}, \hat{a}^S)$  where  $\hat{m}$  belongs to a message bundle  $E$  in  $M_{\phi(N-l)}$  and there exists  $s'^R \in S^R(\bar{k})$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$  such that  $s_{w,(\hat{m}, \hat{a}^S)}^R(\hat{m}) \neq b^R(\hat{a}^S)$ . Let  $\hat{\sigma}^R$  be a totally mixed strategy in  $\Delta S^R(\bar{k})$  to which  $(\hat{m}, \hat{a}^S)$  is a best response. Let  $A^{\max}$  denote the set of Sender actions that maximize

$$\chi(\hat{\sigma}^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a^S)$$

over  $\{a_{\phi(i)}^S : i > N - l\}$ . Given any two message bundles  $F_1$  and  $F_2$  parallel to  $E$ , by assumption, there exist two Sender strategies  $(m_1, a_1^S)$ ,  $(m_2, a_2^S)$  in  $S^S(\bar{k})$  where message  $m_1$  belongs to  $F_1$ , message  $m_2$  belongs to  $F_2$ , and Sender actions  $a_1^S$  and  $a_2^S$  both belong to  $A^{\max}$ . It suffices to show that  $S^S(\bar{k} + 1)$  contains a Sender strategy  $(m_{1*}, a_{1*}^S)$  where  $m_{1*} \in F_1$  and  $a_{1*}^S \in A^{\max}$ . We would be done if  $(m_1, a_1^S) \in S^S(\bar{k} + 1)$ . Suppose  $(m_1, a_1^S) \notin S^S(\bar{k} + 1)$ .

**Claim 5** Given any  $(\tilde{m}, \tilde{a}^S) \in S^S(\bar{k})$  where  $\tilde{a}^S \in A^{\max}$  and  $\tilde{m} \in E \cup F_2$ , there exists  $\psi_{d,(\tilde{m}, \tilde{a}^S)} : S^R(\bar{k}) \rightarrow S^R(\bar{k})$  that belongs to  $\Psi_{m_{\phi(N-l)} \cup M_{\phi(N-l)}}$

such that for every  $s^R$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ ,  $\psi_{d,(\tilde{m},\tilde{a}^S)}(s^R)$  is equal to  $s^R$  outside of  $E \cup F_1 \cup F_2$ , while  $\psi_{d,(\tilde{m},\tilde{a}^S)}(s^R)(m_1) = b^R(a_1^S)$  and  $\psi_{d,(\tilde{m},\tilde{a}^S)}(s^R)(\tilde{m}) \neq b^R(\tilde{a}^S)$ .

**Proof.** Suppose that  $s^R(\tilde{m}) = b^R(\tilde{a}^S)$  for every Receiver strategy in  $S^R(\bar{k})$  which is non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ , then

$$\begin{aligned} & u^S\left((\tilde{m}, \tilde{a}^S), \hat{\sigma}^R\right) \\ &= \chi\left(\hat{\sigma}^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_1^S\right) \\ &\geq \chi\left(\hat{\sigma}^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \tilde{a}^S\right) \\ &> u^S\left((\hat{m}, \hat{a}^S), \hat{\sigma}^R\right), \end{aligned}$$

which contradicts the assumption that  $(\hat{m}, \hat{a}^S)$  is a best response to  $\hat{\sigma}^R$ . Therefore, there exists  $s_{w,(\tilde{m},\tilde{a}^S)}^R \in S^R(\bar{k})$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$  where  $s_{w,(\tilde{m},\tilde{a}^S)}^R(m_j) \neq b^R(\tilde{a}^S)$ . Analogously, there exists  $s_{w,i}^R \in S^R(\bar{k})$  with the same properties w.r.t.  $(m_i, a_i^S)$ .

Assume w.l.o.g. that  $\tilde{m} \in E$ . If either  $s_{w,(\tilde{m},\tilde{a}^S)}^R|_E$  or  $s_{w,(m_2,a_2^S)}^R|_{F_2}$  is not equivalent to a constant Receiver strategy of  $b^R(a_1^S)$  w.r.t.  $S^S(\bar{k}-1)$ , then we can apply lemma 7 and we are done. Otherwise,  $b^R(a_2^S) \neq b^R(a_1^S)$ . Since  $F_1$  and  $F_2$  are parallel in  $M_{\phi(N-l)}$ , constrained message sets of  $(m_i, a_i^S)$  are disjoint, and lemma 2 implies that there exists  $s_2^R \in S^R(\bar{k})$  which is a best response to  $(m_i, a_i^S)$ , for  $i = 1, 2$ . It follows that  $s_2^R$  is non-constant on  $M_{\phi(N-l)}$  and  $s_2^R|_{F_2}$  is not equivalent to a constant of  $b^R(a_1^S)$ . We can again apply lemma 7 and we are done. ■

Following claim 5, we can then define  $T_d : S^R(\bar{k}) \rightarrow S^R(\bar{k})$  where  $T_d(s^R)$  puts strictly positive probability on every strategy in the set

$$\left\{ \begin{array}{l} \psi_{d,(\tilde{m},\tilde{a}^S)} : (\tilde{m}, \tilde{a}^S) \in S^S(\bar{k}+1) \\ \tilde{m} \in E \cup F_2 \text{ and } \tilde{a}^S \in A^{\max} \end{array} \right\}.$$

Not that for every  $s^R$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ ,  $T_d(s^R)(m_1) = b^R(a_1^S)$ ,  $T_d(s^R)$  is equal to  $s^R$  outside if  $E \cup F_1 \cup F_2$ , and

$$T_d(s^R)(\tilde{m}) \neq b^R(\tilde{a}^S) \text{ if } (\tilde{m}, \tilde{a}^S) \in S^S(\bar{k}) \text{ and } \tilde{m} \in E \cup F_2. \quad (6)$$

Therefore, for every Sender strategy  $(m, a^S)$  where  $m \notin E \cup F_1 \cup F_2$ ,

$$\begin{aligned}
u^S \left( (m, a^S), T_d \left( \hat{\sigma}^R \right) \right) &= u^S \left( (m, a^S), \hat{\sigma}^R \right) \\
&\leq u^S \left( (\hat{m}, \hat{a}^S), \hat{\sigma}^R \right) \\
&< \chi \left( \hat{\sigma}^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \hat{a}^S \right) \\
&\leq \chi \left( \hat{\sigma}^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_1^S \right) \\
&= u^S \left( (m_1, a_1^S), T_d \left( \hat{\sigma}^R \right) \right).
\end{aligned}$$

If  $m \in E \cup F_1 \cup F_2$ , then

$$\begin{aligned}
&u^S \left( (m, a^S), T_d \left( \hat{\sigma}^R \right) \right) \\
&\leq \chi \left( \hat{\sigma}^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a^S \right) \tag{7} \\
&\leq \chi \left( \hat{\sigma}^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_1^S \right) \tag{8} \\
&= u^S \left( (m_1, a_1^S), T_d \left( \hat{\sigma}^R \right) \right).
\end{aligned}$$

If  $a^S \notin A^{\max}$ , then strictly inequality holds in line 8. If  $a^S \in A^{\max}$ ,  $m \in E \cup F_2$ , and  $(m, a^S) \in S^S(\bar{k})$ , then strict inequality holds in line 7 because of inequality 6 and the construction of  $\hat{\sigma}^R$  and  $T_d$  such that  $T_d \left( \hat{\sigma}^R \right)$  puts positive probability on Receiver strategies non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ . If  $m \in F_1$  but  $a^S \notin A^{\max}$ , then it follows that

$$u^S \left( (m, a^S), T_d \left( \hat{\sigma}^R \right) \right) < u^S \left( (m_1, a_1^S), T_d \left( \hat{\sigma}^R \right) \right).$$

Therefore, if  $(m_*, a_*^S)$  is a best response in  $S^S(\bar{k})$  to  $T_d \left( \hat{\sigma}^R \right) \in \Delta S^R(\bar{k})$ , it has to be the case that message  $m_* \in F_1$  and Sender action  $a_*^S \in A^{\max}$ . Since at least one Sender best response to  $T_d \left( \hat{\sigma}^R \right)$  survives the  $(\bar{k} + 1)^{th}$  iteration, we are done.

### 6.3 Proof for lemma 5

Both statements are true for  $k = 1$ , for all  $N - l \neq q$ . Since  $M_{\phi(N-3)}$  is intention clear given  $S^S(3)$ , we are concerned only if  $l \geq 3$ . Suppose both statements are true for every  $N - l \neq q$ , for  $k = 1, \dots, \bar{k}$ . Let  $(m_{\phi(N-l)}, a_{\phi(q)}^S)$  be a Sender strategy that survives the  $(\bar{k} + 1)^{th}$  round of deletion, where  $b^R \left( a_{\phi(q)}^S \right) = a_{\phi(q)}^R$ . Observation 4.1.2 implies that  $q \geq N - l$ . Suppose  $(m_{\phi(N-l)}, a_{\phi(q)}^S)$  satisfies both the *Unclear Intention* condition and the *No Sincere Recommendation* condition, . This implies that  $q \neq N - l$ .

Let  $\sigma_q^R$  be a Receiver strategy in  $\Delta S^R(\bar{k})$  to which  $(m_{\phi(N-l)}, a_{\phi(q)}^S)$  is a best response w.r.t.  $S(\bar{k})$ . By assumption, there exists  $(\hat{m}, \hat{a}^S) \in S^S(\bar{k})$  where  $\hat{m} \in M_{\phi(N-l)}$  and

$$\begin{aligned} & \chi(\sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \hat{a}^S) \\ & \geq \chi(\sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_{\phi(q)}^S), \end{aligned}$$

and there exists  $\hat{\sigma}_q^R \in \Psi_{m_{\phi(N-l)} \cup M_{\phi(N-l)}}(\sigma_q^R)$  to which  $(\hat{m}, \hat{a}^S)$  is a best response where  $s^R(\hat{m}) = b^R(\hat{a}^S)$  for every  $s^R$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$  that receives positive weights under  $\hat{\sigma}_q^R$ . If  $S^R(\bar{k})$  contains only Receiver strategies constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ , then for  $\sigma_q^R$  totally mixed in  $\Delta S^R(\bar{k})$ ,

$$\begin{aligned} & U^S((\hat{m}, \hat{a}^S), \sigma_q^R) \\ & = \chi(\sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \hat{a}^S) \\ & \geq \chi(\sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_{\phi(q)}^S) \\ & = U^S((m_{\phi(N-l)}, a_{\phi(q)}^S), \sigma_q^R). \end{aligned}$$

Thus  $(\hat{m}, \hat{a}^S)$  belongs to  $S^S(\bar{k} + 1)$  and we are done.

Suppose  $S^R(\bar{k})$  contains Receiver strategies non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ . Then there exists  $s_w^R \in S^R(\bar{k})$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$  such that  $s_w^R(\hat{m}) \neq b^R(\hat{a}^S)$ . Otherwise, for  $\sigma_q^R$  totally mixed,

$$\begin{aligned} & U^S((\hat{m}, \hat{a}^S), \sigma_q^R) \\ & = \chi(\sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \hat{a}^S) \\ & \geq \chi(\sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_{\phi(q)}^S) \\ & > U^S((m_{\phi(N-l)}, a_{\phi(q)}^S), \sigma_q^R), \end{aligned}$$

because  $s^R(m_{\phi(N-l)}) \neq b^R(a_{\phi(q)}^S)$  for every  $s^R$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ . This contradicts the construction of  $\sigma_q^R$ .

Let  $E$  be the message bundle which  $\hat{m}$  belongs to. By lemma 6, there exists  $(m_1, a_1^S), (m_2, a_2^S) \in S^S(\bar{k})$  where each  $m_i$  belongs to a different message bundle  $F_i$  parallel to  $E$ , and  $a_i^S$  belongs to

$$\arg \max_{a^S \in (b^R)^{-1}(a_{\phi(j)}^R)} \chi(\hat{\sigma}_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a^S),$$

for  $i = 1, 2$ . Therefore, for  $i = 1, 2$ ,

$$\begin{aligned}
& \chi \left( \hat{\sigma}_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_i^S \right) \\
& \geq \chi \left( \hat{\sigma}_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \hat{a}^S \right) \\
& = \chi \left( \sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \hat{a}^S \right) \\
& \geq \chi \left( \sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_{\phi(q)}^S \right).
\end{aligned}$$

Similarly, there must exist  $s_{wi}^R \in S^R(\bar{k})$  such that  $s_{wi}^R(m_i) \neq b^R(a_i^S)$ . If  $b^R(a_i^S) \neq b^R(\hat{a}^S)$  for some  $i \in \{1, 2\}$ , then we use lemma 2, otherwise, we use lemma 7 to establish existence of a mapping  $\psi : S^R(k) \rightarrow S^R(k)$  that belongs to  $\Psi_{m_{\phi(N-l)} \cup M_{\phi(N-l)}}$  such that  $\psi(s^R)(\hat{m}) = b^R(\hat{a}^S)$ , for every  $s^R \in S^R(k)$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ . It follows that

$$\begin{aligned}
& u^S \left( (\hat{m}, \hat{a}^S), \psi(\sigma_q^R) \right) \\
& = \chi \left( \sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, \hat{a}^S \right) \\
& \geq \chi \left( \sigma_q^R, m_{\phi(N-l)} \cup M_{\phi(N-l)}, a_{\phi(q)}^S \right) \\
& > u^S \left( \left( m_{\phi(N-l)}, a_{\phi(q)}^S \right), \psi(\hat{\sigma}^R) \right) \\
& = u^S \left( \left( m_{\phi(N-l)}, a_{\phi(q)}^S \right), \hat{\sigma}^R \right) \tag{9} \\
& \geq u^S \left( (m, a^S), \hat{\sigma}^R \right) \\
& = u^S \left( (m, a^S), \psi_d(\hat{\sigma}^R) \right)
\end{aligned}$$

for every  $(m, a^S)$  where  $m \notin M_{\phi(N-l)}$ . Equality holds in line 9 because  $s^R(m_{\phi(N-l)}) = a_{\phi(N-l)}^R$  for every  $s^R$  non-constant on  $m_{\phi(N-l)} \cup M_{\phi(N-l)}$ . Therefore, any best response to  $\psi(\sigma_q^R)$  uses a message in  $M_{\phi(N-l)}$ . Since at least one of the Sender best responses to  $\psi(\sigma_q^R)$  survives the  $(\bar{k} + 1)^{th}$  iteration, we are done.

## References

- [1] R. Aumann, "Nash Equilibria are not Self-Enforcing," *Economics Decision-Making: Games, Econometrics and Optimization* (J.J. Gabszewicz, J.-F. Richard, and L. A. Wolsey, Eds.), 1990.
- [2] Brandenburger, A., Friedenberg, A., and H.J. Keisler., "Admissibility in Games," mimeo, 2004. Available at [www.stern.nyu.edu/~abranden](http://www.stern.nyu.edu/~abranden).
- [3] S. Baliga and S. Morris, "Coordination, Spillover, and Cheap Talk," *Journal of Economic Theory*, 2002, 450-468.
- [4] J. Farrell, "Communication, Coordination and Nash Equilibrium," *Economic Letter*, 1988, 209-214.