

# When Does Communication Improve Coordination?

By TORE ELLINGSEN AND ROBERT ÖSTLING\*

This version: April 16, 2009

*We study costless pre-play communication of intentions among inexperienced players. Using the level- $k$  model of strategic thinking to describe players' beliefs, we fully characterize the effects of pre-play communication in symmetric  $2 \times 2$  games. One-way communication weakly increases coordination on Nash equilibrium outcomes, although average payoffs sometimes decrease. Two-way communication further improves payoffs in some games, but is detrimental in others. Moving beyond the class of symmetric  $2 \times 2$  games, we find that communication facilitates coordination in common interest games with positive spillovers and strategic complementarities, but there are also games in which any type of communication hampers coordination.*

*JEL: C72.*

*Keywords: Pre-play communication, coordination games, Stag Hunt, level- $k$ , bounded rationality.*

Some people find themselves in a new strategic situation. How can they best coordinate their actions? Since they cannot rely on precedence, maybe they should start talking? If so, what are the exact reasons why communication helps? These fundamental questions crop up in many disciplines, including evolutionary biology, psychology, political science, and economics.<sup>1</sup>

Farrell (1987, 1988) and Matthew Rabin (1990, 1994) provide formal analyses of costless communication, or cheap talk, as a means to convey intentions and thereby improve coordination among rational players in games with complete information.<sup>2</sup> While the models are insightful,

\* Ellingsen: Department of Economics, Stockholm School of Economics, Box 6501, SE-113 83 Stockholm, Sweden (e-mail: [tore.ellingsen@hhs.se](mailto:tore.ellingsen@hhs.se)). Östling: Institute for International Economic Studies, Stockholm University, SE-106 91 Stockholm, Sweden (e-mail: [robert.ostling@iies.su.se](mailto:robert.ostling@iies.su.se)). We have benefited greatly from the detailed comments from Vincent P. Crawford and several anonymous referees. We are also grateful for helpful discussions with Colin F. Camerer, Drew Fudenberg, Botond Köszegi, Joseph Tao-yi Wang and many seminar participants. Financial support from the Torsten and Ragnar Söderberg Foundation and the Jan Wallander and Tom Hedelius Foundation is gratefully acknowledged.

<sup>1</sup>For evolutionary biology, see e.g., Steven Pinker and Paul Bloom (1990) and Martin A. Nowak (2006, Chapter 13); for psychology, see e.g., Norbert L. Kerr and Cynthia M. Kaufmann-Gilliland (1994); for political science, see e.g., Thomas C. Schelling (1966, Chapter 7), and for economics, see e.g., Joseph Farrell and Garth Saloner (1988) or David Genesove and Wallace P. Mullin (2001).

<sup>2</sup>For a non-technical introduction to the literature on cheap talk about intentions, see Farrell and Rabin (1996), especially pages 110–116. An early precursor is Robert J. Aumann (1974). See also Roger Myerson (1989), who emphasizes that cheap talk can communicate both own intended actions (“promises”) and desires about others’ actions (“requests”). Like most of the literature, we focus on the former. Note that we ignore the communication of private information; Vincent P. Crawford and Joel Sobel (1982) and Jerry R. Green and Nancy L. Stokey (2007) (originally written in 1981) are seminal contributions to the study of strategic information transmission.

Throughout, we take for granted that players have access to a common language. A substantial fraction of the literature on cheap talk starts from the opposite presumption that messages are not inherently meaningful; instead, messages may or may not acquire meaning in equilibrium—where equilibrium is a steady state of an evolutionary process of random matches between pre-programmed players; see, for example, Akihiko Matsui (1991), Karl Wärneryd (1991), Yong-Gwan Kim and Sobel (1995), Luca Anderlini (1999) and Abhijit Banerjee and Jörgen W. Weibull (2000). While the evolutionary approach can explain how language emerges in “old” games, it is less appropriate for our question of how an existing language will be used in “new” games. For

	<i>H</i>	<i>L</i>
<i>H</i>	0, 0	3, 1
<i>L</i>	1, 3	0, 0

FIGURE 1: BATTLE OF THE SEXES

	<i>H</i>	<i>L</i>
<i>H</i>	9, 9	0, 8
<i>L</i>	8, 0	7, 7

FIGURE 2: STAG HUNT

we argue that they make problematic assumptions concerning players' beliefs, and that recent models of strategic thinking offer alternative assumptions that better fit our intuitions and available experimental evidence.

To put our arguments into perspective, let us briefly review some of the literature. Farrell (1987) studies communication in a Battle of the Sexes game (Figure 1). Farrell assumes that behavior will correspond to the symmetric mixed strategy Nash equilibrium if players cannot communicate. He also assumes that message pairs (“*H*”, “*L*”) and (“*L*”, “*H*”) that are consistent with a pure strategy equilibrium will induce play of that equilibrium. Based on these assumptions, he shows that with two-way communication there are better symmetric mixed strategy equilibria than the no-communication equilibrium. Payoffs improve with the number of communication rounds, but full efficiency is unattainable because players have conflicting interests over the two efficient equilibria.

Although the symmetric Nash equilibrium assumption makes some sense in the Battle of the Sexes, it is not generally a convincing assumption about the behavior of inexperienced players in the absence of communication; rationality and reasoning alone is insufficient to support equilibrium. In their subsequent work, Farrell and Rabin instead make the weaker assumption that the outcome without communication will be rationalizable. One drawback of this approach is that rationalizability provides no prediction about behavior in many games, including Battle of the Sexes. Without a prediction for the game without communication, a specific prediction for the game with communication does not suffice to say whether communication improves coordination or not.

Another objection to Farrell and Rabin's approach is that one-way communication sometimes works excessively well, especially in coordination games like Stag Hunt (Figure 2).<sup>3</sup> In Stag Hunt, all outcomes are again rationalizable without communication. Farrell (1988) suggests that one-way communication suffices to attain coordination on the efficient outcome (*H*, *H*), because the message “*H*” is *self-committing*. That is, if sending the message “*H*” convinces the receiver that the sender intends to play *H*, the best response is for the receiver to play *H*, and thus the sender has an incentive to play according to the own message. Aumann (1990) objects that even a sender who has decided to play *L* has an incentive to induce the opponent to play *H*. That is, the message “*H*” is not *self-signaling*. Relatedly, Pei-yu Lo (2007) demonstrates that the message “*L*” is weakly dominated under the two assumptions that players have common knowledge about the meaning of the language and believe their opponent to behave rationally given this knowledge. As a consequence, both action *H* and *L* survive iterated elimination of

evolutionary models in which language has some pre-existing meaning, see Andreas Blume (1998) and Stefano Demichelis and Weibull (2008).

<sup>3</sup>Since Stag Hunt is the prototype representation of coordinated hunting situations, it is an apt touchstone for theories of communication. Indeed the benefits from coordinated hunting of large animals has been proposed as an explanation for why language have emerged (see, e.g., Pinker and Bloom 1990, Section 5.3).

	<i>H</i>	<i>L</i>
<i>H</i>	11, 11	0, 10
<i>L</i>	10, 10	10, 9

FIGURE 3: VULNERABILITY GAME

weakly dominated strategies.

Addressing Aumann’s critique, Farrell and Rabin (1996, page 114) acknowledge that their theory is not entirely satisfactory, but think that it has the right implications: “[A]lthough we see the force of Aumann’s argument, we suspect that cheap talk will do a good deal to bring [the players] to the stag hunt.” Are Farrell and Rabin right? The experimental evidence on behavior in Stag Hunt games is somewhat conflicting, but it consistently shows that communication improves coordination. For example, in an experiment by Gary Charness (2000) one-way communication induces substantial coordination on the efficient equilibrium. In the prior experiment by Russell Cooper et al. (1992) one-way communication improves players willingness to play *H*, but two-way communication does so to a greater extent; see Section III for a more detailed discussion of the evidence.

To summarize, Farrell and Rabin’s approach leaves open at least three questions. First and foremost, when does communication improve coordination? Second, why does communication of intentions matter even in situations when messages apparently fail to be self-signaling? Third, might two-way communication sometimes generate more coordination than one-way communication, and if so why?

In a nutshell, we argue that rationalizability is an inappropriate assumption about inexperienced players’ beliefs, and that more realistic assumptions help to answer all the above three questions. One shortcoming of rationalizability we have already mentioned: Too often, it imposes no restriction on beliefs. Another shortcoming is that it sometimes imposes unrealistically strong restrictions on beliefs. To fix ideas, consider the Vulnerability Game in Figure 3.<sup>4</sup>

In the Vulnerability Game, playing *H* is strategically risky for the row player, whereas *L* is safe. We intuitively believe that many row players would be unwilling to risk losing 10 in order to gain 1.<sup>5</sup> We also believe that communication by the column player can increase the row player’s willingness to play *H*. By sending the pre-play message “*H*”, the column player provides some reassurance to the row player, who as a result comes to regard action *H* less risky. We think it is this sort of intuition that explains why communication has an effect in Stag Hunt. However, the intuition is inadmissible in Farrell and Rabin’s framework. Since *L* is a dominated action for the column player, (*H*, *H*) is the unique rationalizable outcome, and thus ought to obtain whether players communicate or not. Rationalizability assumes not only that players are rational, but also that players believe with probability 1 that their opponents are rational.<sup>6</sup> In the Vulnerability Game, the row player’s firm belief that the column player is rational eliminates the need for communication. We propose instead that it is exactly the doubt about the column player’s rationality that induces the row player to pay attention to the column player’s message.

If some players doubt that their opponent is rational, what do they believe? Data from Beauty Contest games led Rosemarie Nagel (1995) to suggest that people’s implicit beliefs about others’ behavior can often be characterized as follows: Some people believe that their opponents randomize uniformly. Other more advanced people believe that their opponents

<sup>4</sup>The Vulnerability Game is inspired by the game in Figure 1.4 of Drew Fudenberg and Jean Tirole (1991).

<sup>5</sup>For related arguments and evidence, see Robert Rosenthal (1981), and T. Randolph Beard and Richard O. Beil Jr. (1994).

<sup>6</sup>Indeed, all players are assumed to believe that all players believe that all players believe that...etc...all players are rational.

believe that opponents randomize uniformly. Others again believe that their opponents believe that opponents believe that opponents randomize uniformly. The corresponding formal model is known as the level- $k$  model. Level-0 players, who may or may not be assumed to exist in reality, randomize uniformly. Level-1 players believe that their opponent is level-0. Level-2 players believe that their opponents are level-1, and so on. The level- $k$  model was first studied by Dale O. Stahl and Paul W. Wilson (1994, 1995) and Nagel (1995) and has the virtue of offering a structural non-equilibrium approach to the analysis of people's initial behavior in unfamiliar games. A natural extension of the level- $k$  model is to assume that a level- $k$  player believes that the opponent is drawn from a distribution of more primitive player types; see Colin F. Camerer, Teck-Hua Ho and Juin-Kuan Chong (2004) for an analysis of the ensuing cognitive hierarchy model. These models successfully organize data on the behavior of inexperienced players in a wide variety of settings.<sup>7</sup>

The level- $k$  model straightforwardly explains why some row players would play  $L$  in the Vulnerability Game. A level-1 row player thinks it is equally likely that the column player plays  $L$  and  $H$ , and since  $10 > 5.5$  it is better for the row player to play  $L$ . A drawback with the level- $k$  model is that higher level row players all believe with probability 1 that their opponent understands the game well enough to not play a dominant strategy. Thus, these players all prefer  $H$ , and would continue to do so independently of the difference between their  $(H, H)$ -payoff and their  $(L, H)$ -payoff, as long as it is positive. The cognitive hierarchy model, on the other hand, predicts that the behavior of higher level player types is sensitive to this payoff difference, since no player completely rules out the possibility that the opponent is of level-0. These considerations notwithstanding, we focus attention on the error-free level- $k$  model due to its greater simplicity. See Appendix 3 for a detailed analysis applying the cognitive hierarchy model.

Observe that both the level- $k$  model and the cognitive hierarchy model assume that all players at level 1 or higher behave rationally given their beliefs. Thus, players' messages will also be chosen to maximize expected payoffs. However, to fully pin down behavior we need to specify the choice of message when players are indifferent between several messages. Following Demichelis and Weibull (2008), we assume that whenever the truthful message is in the indifference set, players are truthful. That is, they have a weak (lexicographic) preference for being honest.<sup>8</sup> This weak preference for honesty is key to several of our results. For example, in the Vulnerability Game, a level-1 row player is affected by the column player's message precisely because it is believed to be honest.

The lexicographic truthfulness assumption is strong enough to determine the messages of level-0. Since level-0 players are indifferent between all actions, they are also truthful. Observe that an alternative and more direct justification of level-0 truthfulness is to assume credulity on the part of level-1; this is the essentially the approach taken by Crawford (2003). If level-0 does not exist, except in the minds of level-1 players, the two assumptions are behaviorally similar.

In symmetric two-player games with one-way communication, we show that the truthfulness of

<sup>7</sup>See Stahl and Wilson (1994, 1995), Nagel (1995), Ho, Camerer and Keith Weigelt (1998), Miguel A. Costa-Gomes, Crawford and Bruno Broseta (2001), Camerer, Ho and Chong (2004), Costa-Gomes and Crawford (2006), Crawford and Nagore Iriberri (2007*a*), Costa-Gomes, Crawford and Iriberri (2009) for various normal form game applications of the level- $k$  and cognitive hierarchy models based on laboratory data. Crawford and Iriberri (2007*b*) use the level- $k$  model to explain the winners' curse and overbidding in private-value auctions and Crawford et al. (2009) use it to study optimal auction design. Toshiji Kawagoe and Hirokazu Takizawa (2008*b*) apply the level- $k$  model to an extensive form game, the centipede game. Östling et al. (2008) and Alexander L. Brown, Camerer and Dan Lovallo (2009) estimate cognitive hierarchy and level- $k$  models using field data. See also footnote 9 for references to level- $k$  analyses of communication.

<sup>8</sup>There is considerable experimental evidence that many people assign strictly positive utility to behaving honestly (e.g., Ellingsen and Magnus Johannesson, 2004*b* and the references therein), and our results would be largely the same with positive utility from honesty. However, the analysis is simpler if the preference is lexicographically small.

level-0 is contagious: A level-1 receiver plays a best response to the received message. Since level-1 behavior constitutes level-2 players' model of the world, and the game is symmetric, a level-2 sender will send a truthful message that corresponds to the sender's favorite Nash equilibrium. Analogous reasoning proves that, in this class of games, all player types communicate their intentions honestly under one-way communication. However, contagious honesty does not imply that one-way communication suffices to induce an efficient outcome. For example, in the Stag Hunt game above, level-1 players would send and play  $L$ .

Let us now briefly describe our main results. For parameter choices that are typical in the level- $k$  literature, the following is true for symmetric  $2 \times 2$  games: (i) One-way communication improves average payoffs in Stag Hunt games with a conflict between efficiency and strategic risk, such as that in Figure 2, and in some but not all mixed motive (Chicken) games. (ii) Two-way communication may yield higher average payoffs than one-way communication, but only in Stag Hunt games with a conflict between efficiency and strategic risk and in mixed motive games with high miscoordination payoffs. (iii) In mixed motive games with high miscoordination payoffs, average payoffs can be lower with communication than without. An additional finding is that if players are sufficiently sophisticated, both one-way and two-way communication suffices to attain the efficient outcome in Stag Hunt. This conclusion holds not only in the limit as sophistication goes to infinity; it suffices that both players perform at least two thinking steps.

Extending our analysis to larger games and/or relaxing the symmetry assumption, we find that both one-way and two-way communication facilitates coordination in all two-player common interest games: When both players make at least two thinking steps, there is always coordination on (the best) Nash equilibrium in these games. If there are more than two players, a similar result holds under the additional assumption of positive spillovers and strategic complementarities.

On the other hand, it is easy to identify games in which communication erodes coordination. The reason is that players have an incentive to deceive the opponent by misrepresenting their intentions. Even if the game has a unique pure strategy equilibrium, players can obtain large non-equilibrium payoffs if they successfully fool their opponent. When players are not too sophisticated, they may end up playing non-equilibrium strategies that are either more or less profitable than equilibrium.

Crawford (2003) is the seminal study of communication of intentions with level- $k$  beliefs. Crawford studies a special class of zero-sum games, namely Hide and Seek games, with one-way communication. Our work adapts Crawford's approach in order to study a different (and larger) class of games, while considering both one-way and two-way communication.<sup>9</sup> The resulting sets of applications are quite different. Where Crawford's paper studies deception, ours predominantly studies mutually beneficial coordination.

## I. Model

Let  $G = \langle N, A, u \rangle$  denote some normal form game between  $|N|$  players where  $N$  denotes the set of players,  $A_i$  denotes the finite set of actions for player  $i$ ,  $A = \times_{j \in N} A_j$  denotes the set of feasible action profiles,  $u_i : A \rightarrow \mathbb{R}$  is player  $i$ 's von Neumann-Morgenstern utility function, and  $u$  is the vector of all players' utility functions. We refer to  $G$  as an *action game*. Let  $\Gamma_{N^*}(G)$  denote the game  $G$  preceded by one round of pre-play communication, where the subset  $N^* \subseteq N$  of the players are allowed (by nature) to send a message. Let  $M_i = A_i$  be the set of feasible messages for a communicating player  $i$  and let  $\mathcal{M}_i = A_i \cup \emptyset$ . Let  $M = \times_{j \in N} M_j$  and  $\mathcal{M} = \times_{j \in N} \mathcal{M}_j$  denote the corresponding sets of feasible message profiles. By convention, a non-communicating player sends an empty message  $\emptyset$ . The nonempty messages are assumed to

<sup>9</sup>Recently, Erik Wengström (2008) has applied the level- $k$  model to study communication in a price competition game. Previously, Hongbin Cai and Joseph Tao-Yi Wang (2006) and Kawagoe and Takizawa (2008a) have adapted Crawford's model to study one-sided cheap talk in sender-receiver games with private information.

articulate a statement about the sender's intention (rather than for example a statement about which action the sender desires from the receiver). Let player  $i$ 's message and action be denoted  $m_i$  and  $a_i$ , respectively. A strategy of the game  $\Gamma_{N^*}(G)$  for player  $i$  is a message  $m_i \in \mathcal{M}_i$  and a mapping  $f_i : \mathcal{M} \rightarrow A_i$  defining the action for any message profile.

To begin with, we focus attention on symmetric and generic  $2 \times 2$  games.<sup>10</sup> The two actions are labeled  $H$  and  $L$ . The utilities associated with each outcome are denoted  $u_{HH}$ ,  $u_{HL}$ ,  $u_{LH}$  and  $u_{LL}$ . In the game  $G$  preceded by one-way communication,  $\Gamma_I(G)$ , one of the players is allowed to send one of two messages,  $h$  and  $l$ , before the action game  $G$  is played. Although from player  $i$ 's point of view, the game in which he is a receiver is quite separate from the game in which he is a sender, it is convenient to write the two strategies jointly. From now on we thus say that a strategy  $s_i$  for player  $i$  of the full game  $\Gamma_I(G)$  prescribes what message  $m_i$  to send and action  $a_i$  to take as sender, and a mapping  $f_i : \{h, l\} \rightarrow \{H, L\}$  from received messages to actions as receiver. We write a pure strategy of player  $i$  (given the received message  $m_j$ ) as

$$s_i = \langle m_i, a_i, f_i(m_j = h), f_i(m_j = l) \rangle.$$

For example,  $s_1 = \langle h, H, L, L \rangle$  means that player 1 sends the message  $h$  and takes the action  $H$  if he is the sender, while playing  $L$  whenever acting as receiver.

Observe that we neglect unused strategy components by restricting attention to the reduced normal form. In other words, we do not specify what action a player would take in the counterfactual case when he sends another message than the message specified by his strategy.

In the game with two-way communication,  $\Gamma_{II}(G)$ , both players simultaneously send a message  $m_i \in \{h, l\}$  before  $G$  is played.<sup>11</sup> A strategy  $s_i$  for player  $i$  of the full game is therefore given by a message  $m_i$  and a mapping  $f_i : \{h, l\} \rightarrow \{H, L\}$  from the opponent's message to actions. A pure strategy of player  $i$  (given the message  $m_j$  sent by player  $j$ ) can thus be written

$$s_i = \langle m_i, f_i(m_j = h), f_i(m_j = l) \rangle.$$

For example,  $s_1 = \langle h, H, L \rangle$  means that player 1 sends the message  $h$ , but plays according to the received message (i.e., plays  $H$  if player 2 sends message  $h$  and  $L$  if player 2 plays message  $l$ ).

Players' behavior depends on their degree of sophistication. A player of type 0 (or level-0), henceforth called a  $T_0$  player, is assumed to understand only the set of strategies, and not how these strategies map into payoffs. Thus,  $T_0$  makes a uniformly random action plan, sticking to this plan independently of any message from the opponent. (Hearing the opponent's intended action is of little help to a player who does not understand which game is being played.) Importantly, since  $T_0$  players do not understand how their own or their opponent's actions map into payoffs, or how their messages may affect their opponent's action, they are indifferent concerning their own messages.

For positive integers  $k$ , a  $T_k$  player chooses a best response to (the behavior that the  $T_k$  player expects from) a  $T_{k-1}$  opponent. In particular,  $T_1$  plays a best response to  $T_0$ . When  $k \geq 2$ ,  $T_k$  players will sometimes observe unexpected messages. In this case  $T_k$  assumes that the message comes from a  $T_{k-l}$  player, where  $l \leq k$  is the smallest integer that makes  $T_k$ 's inference consistent. (As we shall see,  $T_0$  sends all messages with positive probability, so  $l \in \{1, \dots, k\}$ )

<sup>10</sup>There is a tension between genericity and symmetry, but none of our results are knife-edge with respect to symmetry. For the purpose of this paper, we consider a game to be generic if no player obtains exactly the same payoff for two different pure strategy profiles. We restrict attention to symmetric and generic games merely in order to keep down the number of cases under consideration. In section II.B, however, we discuss an asymmetric  $2 \times 2$  game.

<sup>11</sup>Simultaneous messages may appear to be an artificial assumption. However, besides preserving symmetry, the case of simultaneous messages may capture the notion from models with sequential communication that the first and the last speaker may both have an impact.

always exists.) Let  $p_k$  denote the proportion of type  $k$  in the player population. As we shall see, players who perform more than one thinking step often, but not always, behave alike. Therefore, it is convenient to let  $T_{k+}$  denote player types that perform at least  $k$  thinking steps.

When a player is indifferent about actions in  $G$ , we assume that the player randomizes uniformly. However, when the player is indifferent about what pre-play message to send, we assume that there is randomization only in case the player is unable to predict the own action—which can only happen under two-way communication. Otherwise, indifferent players send truthful messages (or more precisely, a message that conveys the action that the player expects to be playing). The assumption reflects the notion that people are somewhat averse to lying, but it does so without incurring the notational burden of introducing explicit lying costs into the model. While such lexicographic preference for truthfulness is an apparently weak assumption, one of its immediate implications is that the message by  $T_0$  reveals the intended action. Or to put it even more starkly,  $T_1$  believes in received messages. (In Section I.C we explore alternative assumptions regarding how  $T_0$  treat messages.)

In Appendix 1 we explicitly characterize the strategies of all player types. However, it is common to argue that  $T_0$  does not accurately describe the behavior of any significant portion of real adult people and that actual players are best described by a distribution with support only on  $T_1, T_2$  and  $T_3$  (e.g. Costa-Gomes, Crawford and Broseta 2001 and Costa-Gomes and Crawford 2006). For some of our results we thus refer to type distributions consisting exclusively of players of these three types. Accordingly, we say that  $p = (p_0, p_1, \dots)$  is a *standard type distribution* if  $p_k > 0$  for all  $k \in \{1, 2, 3\}$  and  $p_k = 0$  for all  $k \notin \{1, 2, 3\}$ .

### A. Examples

Consider the Stag Hunt game in Figure 2. Absent communication,  $T_1$  best responds to the uniformly randomizing  $T_0$  by playing the risk dominant action  $L$ . Understanding this, the best response of  $T_2$  is to play  $L$  as well. Indeed, by induction any player  $T_{1+}$  plays  $L$ . For any type distributions with  $p_0 = 0$ , the unique outcome is the risk dominant equilibrium  $(L, L)$ .<sup>12</sup> The level- $k$  model hence provides a rationale for why players play the risk dominant equilibrium in coordination games without communication.

If players can communicate, one-way communication suffices to induce play of  $H$  by all types  $T_{2+}$ . The analysis starts by considering the behavior of  $T_0$  (as imagined by  $T_1$ ). By assumption, a  $T_0$  sender randomizes uniformly over  $L$  and  $H$ , while sending the corresponding truthful message. A  $T_0$  receiver randomizes uniformly over  $L$  and  $H$ . As a sender,  $T_1$  best responds by playing the risk dominant action  $L$ , and due to the lexicographic preference for truthfulness sends the honest message  $l$ . As a receiver,  $T_1$  believes that messages are honest and thus plays  $L$  following the message  $l$  and  $H$  following the message  $h$ . Consider now  $T_2$ . A  $T_2$  sender believes to be facing a  $T_1$  receiver who best responds to the message, so  $T_2$  sends  $h$  and plays  $H$ . A  $T_2$  receiver, expects to receive an  $l$  message and therefore play  $L$ . If receiving a counterfactual  $h$  message,  $T_2$  thinks it is sent by a truthful  $T_0$  sender and therefore plays  $H$ . It is easily checked that all  $T_{2+}$  behave like  $T_2$ , implying that there will be coordination on the payoff dominant equilibrium whenever two  $T_{2+}$  players meet and communicate. In other words, the level- $k$  model not only shows that it is feasible for advanced players to coordinate on the payoff dominant equilibrium, but that the *unique* outcome is that they will do so. Note in particular how reassurance plays a crucial role in the example. When a receiver gets a message  $h$ , the receiver is reassured that the sender will play  $H$ , and is therefore also willing to play  $H$ . Even if the message  $h$  is actually only self-signaling for (the non-existing) level-0 senders, it is self-committing for all other types, and this suffices to attain efficient coordination as long as both

<sup>12</sup>Note that this is not about equilibrium selection in the ordinary sense. Players do not select among the set of equilibria, but best-respond to the behavior of lower-step thinkers. Their behavior ultimately results from the uniform randomization of  $T_0$ , which explains the parallel to risk dominance.

TABLE 1: ACTION PROFILES PLAYED IN STAG HUNT WITH COMMUNICATION

$\Gamma_I(G)$ ( <i>one-way communication</i> )				$\Gamma_{II}(G)$ ( <i>two-way communication</i> )			
	$0R$	$1R$	$\geq 2R$	$0$	$1$	$\geq 2$	
$0S$	Uniform	$\frac{1}{2}HH, \frac{1}{2}LL$	$\frac{1}{2}HH, \frac{1}{2}LL$	$0$	Uniform	$\frac{1}{2}HH, \frac{1}{2}LL$	$\frac{1}{2}HH, \frac{1}{2}LH$
$1S$	$\frac{1}{2}LL, \frac{1}{2}LH$	$LL$	$LL$	$1$	$\frac{1}{2}HH, \frac{1}{2}LL$	Uniform	$HH$
$\geq 2S$	$\frac{1}{2}HH, \frac{1}{2}HL$	$HH$	$HH$	$\geq 2$	$\frac{1}{2}HH, \frac{1}{2}HL$	$HH$	$HH$

	$H$	$L$
$H$	$0, 0$	$3, 1$
$L$	$1, 3$	$a, a$

FIGURE 4: MIXED MOTIVE GAME

parties perform at least two thinking steps.

In Stag Hunt, the reassurance role of communication is strengthened even more when both players send messages. Under such two-way communication,  $T_1$  trusts the received message and responds optimally to it. Expecting to play either action with equal probability,  $T_1$  sends both messages with equal probability.  $T_2$  believes that the opponent listens to messages, and therefore sends  $h$  and plays  $H$  irrespective of the received message.  $T_{3+}$  players believe that the opponent will play  $H$  and they therefore play  $H$  and send an  $h$  message. If they receive an unexpected  $l$  message, they believe it comes from  $T_1$  and therefore play  $H$  anyway (as  $T_1$  will respond to the received  $h$  message by playing  $H$ ). Note that under two-way communication,  $T_{2+}$  players are so certain that the opponent will play  $H$  that they play  $H$  irrespective of the received message. This is an important case in which the cognitive hierarchy model predicts a different strategy. Because  $T_2$  players in the cognitive hierarchy model find it likely that the opponent is a truthful  $T_0$  player, they respond to messages under reasonable parameter assumptions (see Appendix 3 for details).

Table 1 summarizes the action profiles that will result in the Stag Hunt under one-way and two-way communication. The notation  $1S$  indicates a player of type 1 in the role of sender, and so on. “Uniform” indicates that all four outcomes are equally likely.

Communication entails perfect coordination on the payoff dominant equilibrium whenever  $T_{2+}$  players meet. However, one-way and two-way communication differ in two respects whenever  $T_1$  players are involved. With one-way communication,  $T_1$  senders play  $L$  and the risk dominant equilibrium therefore results whenever  $T_1$  senders play (since  $T_0$  does not exist). Under two-way communication, however, there is miscoordination in half of the cases when two  $T_1$  players meet. Thus, there is a trade-off when choosing the optimal communication structure between coordination on either equilibria and achieving the payoff dominant equilibrium more often. For standard type distributions, two-way communication entails higher expected payoffs than one-way communication as long as  $p_1 \in (0, 2/3)$ .

In Stag Hunt, communication increases players’ payoff because it brings sufficiently much reassurance for players to coordinate on the risky but payoff dominant equilibrium. In mixed motive games such as Battle of the Sexes and Chicken, communication instead serves the role of conflict resolution. To see this, consider the mixed motive game depicted in Figure 4, where  $a < 3$  and  $a \neq 2$ . If  $a = 0$ , then this is a Battle of the Sexes, whereas it is a Chicken game if  $a > 0$ . The outcome for this game depends on whether  $L$  or  $H$  is the risk dominant action, i.e., whether  $a \geq 2$ . For simplicity, we disregard the possibility that  $a = 2$ , but allow the “Battle of the Sexes” possibility that  $a = 0$  (although this makes the game non-generic).

First consider the case of no communication.  $T_1$  then plays the risk dominant action, i.e.,  $L$  if  $a > 2$  and  $H$  if  $a < 2$ .  $T_2$  responds optimally by playing  $H$  if  $a > 2$  and  $L$  if  $a < 2$ . The



TABLE 2: ACTION PROFILES PLAYED IN MIXED MOTIVE GAMES ( $a > 2$ )

$G$ (no communication)				$\Gamma_I(G)$ (one-way communication)			
	0	Odd	Even	0R	0R	1R	$\geq 2R$
0	Uniform	$\frac{1}{2}HL, \frac{1}{2}LL$	$\frac{1}{2}LH, \frac{1}{2}HH$	0S	Uniform	$\frac{1}{2}HL, \frac{1}{2}LH$	$\frac{1}{2}HL, \frac{1}{2}LH$
Odd	$\frac{1}{2}LH, \frac{1}{2}LL$	LL	LH	1S	$\frac{1}{2}LH, \frac{1}{2}LL$	LH	LH
Even	$\frac{1}{2}HL, \frac{1}{2}HH$	HL	HH	$\geq 2S$	$\frac{1}{2}HL, \frac{1}{2}HH$	HL	HL

behavior of more advanced players continues to alternate, odd types playing  $L$  if  $a > 2$  and  $H$  otherwise, whereas even types play  $H$  if  $a > 2$  and  $L$  otherwise. The outcome therefore depends on the type distribution, but there will generally be many instances of miscoordination.<sup>13</sup>

One-way communication powerfully resolves the conflict inherent in such games with two pure asymmetric equilibria. If  $H$  is the risk dominant action, then  $T_{1+}$  senders send  $h$  and play  $H$ , whereas  $T_{1+}$  receivers optimally respond to messages. If instead  $L$  is risk dominant, a  $T_1$  sender sends  $l$  and plays  $L$ , whereas  $T_{2+}$  senders continue to send  $h$  and play  $H$ . One-way communication therefore implies that  $T_{1+}$  players always coordinate on an equilibrium. Except in the case when  $L$  is risk dominant and the sender is of type  $T_1$ , coordination is on the sender's preferred equilibrium.

It is unsurprising that one-way communication can resolve the conflict and achieve coordination in games with two asymmetric equilibria. However, our analysis also reveals the novel possibility that in some versions of Chicken some players propose and play their least favorite equilibrium.  $T_1$  senders play their risk-dominant action which may not correspond to their preferred equilibrium, whereas  $T_2$  senders are confident in reaching their preferred equilibrium. Table 2 shows the outcomes that will result without communication and with one-way communication, demonstrating the improved coordination on equilibrium outcomes.

Although one-way communication entails more equilibrium coordination than no communication, more coordination need not raise players' average payoffs. If  $a > 2$ , then players prefer the  $(L, L)$  outcome to ending up in either equilibrium with equal probability. If the type distribution is such that the  $(L, L)$  outcome results sufficiently often without communication, average payoffs are thus higher without communication. For example, when  $a = 5/2$  and there is a standard type distribution with  $p_2 < 1/3$ , then average payoffs are lower under one-way communication than under no communication.

Suppose players could choose whether to engage in communication or not, and that the allocation of roles is random. Each player type  $k$  would then consider the own expected payoff in each regime conditional on meeting a player of type  $k - 1$ . To illustrate that players may prefer not to communicate, we consider the case when  $a = 0$ , i.e., the Battle of the Sexes. Absent communication,  $T_3$  believes that the opponent will play  $L$  and thus obtains the preferred equilibrium payoff. With one-way communication and a random allocation of roles, however,  $T_3$  expects to end up in either equilibrium with equal probability. That is,  $T_3$  expects to be better off if communication is impossible.

## B. Results

In this section we generalize the findings from the previous section to all symmetric and generic  $2 \times 2$  games, disregarding (the measure zero class of) games in which neither action is risk dominant. There are three broad classes of such games. The first class of games are the dominance solvable ones, like Prisoners' Dilemma. We use the convention of labelling the dominant action of these games  $H$ (igh). The second class are coordination games, where we

<sup>13</sup>The outcome without communication does generally not resemble the symmetric mixed strategy equilibrium, but may happen to do so for certain combinations of payoff configurations and type distributions.

follow the example above and label the actions corresponding to the payoff dominant equilibrium  $H$ (igh). The third class of games are mixed motive games like the one in Figure 4. For this class of games, we label the action corresponding to a player's preferred equilibrium  $H$ (igh). In Appendix 1, we completely characterize behavior of all player types  $k \in \mathbb{N}$  for these three classes of games. These characterizations provide the foundation for the results in this section, where we focus on average outcomes under standard type distributions.

Our first result states the conditions under which one-way communication serves to increase players' average payoffs relative to no communication.

**PROPOSITION 1:** *Given a standard type distribution, the average payoff associated with  $\Gamma_I(G)$  exceeds the average payoff associated with  $G$  if and only if (i)  $G$  is a coordination game with a conflict between risk and payoff dominance, or (ii)  $G$  is a mixed motive game that satisfies either*

*a.  $L$  is risk dominant and*

$$\left(\frac{1}{2} - p_2(1 - p_2)\right)(u_{HL} + u_{LH}) > p_2^2 u_{HH} + (1 - p_2)^2 u_{LL},$$

*or*

*b.  $H$  is risk dominant and*

$$\left(\frac{1}{2} - p_2(1 - p_2)\right)(u_{HL} + u_{LH}) > (1 - p_2)^2 u_{HH} + p_2^2 u_{LL}.$$

**PROOF:**

In Appendix 2.

If we replace  $p_2$  by  $p_E$ , the probability that players think an even number of steps, Proposition 1 generalizes straightforwardly to all type distributions in which  $p_0 = 0$ . In our examples, we have already explained why one-way communication improves average payoffs in Stag Hunt, and indicated why it sometimes fails to improve payoffs in mixed motive games. A straightforward implication of Proposition 1 is that one-way communication raises the average payoff in the Battle of the Sexes.<sup>14</sup> (To see this, recall that in Battle of the Sexes  $0 = u_{HH} = u_{LL} < u_{LH} < u_{HL}$ , which implies that  $H$  is risk dominant and that condition (b) in Proposition 1 is satisfied.) Proposition 1 also implies that communication does not improve average payoffs in dominance solvable games. For Chicken, the impact of communication hinges more delicately on parameters, and communication may even serve to reduce payoffs.

**COROLLARY 2:** *Given a standard type distribution, the average payoff associated with  $\Gamma_I(G)$  is smaller than the average payoff of  $G$  if and only if  $G$  is a game of Chicken that satisfies either*

*a.  $L$  is risk dominant and*

$$\left(\frac{1}{2} - p_2(1 - p_2)\right)(u_{HL} + u_{LH}) < p_2^2 u_{HH} + (1 - p_2)^2 u_{LL},$$

*or*

*b.  $H$  is risk dominant and*

$$\left(\frac{1}{2} - p_2(1 - p_2)\right)(u_{HL} + u_{LH}) < (1 - p_2)^2 u_{HH} + p_2^2 u_{LL}.$$

<sup>14</sup>Note that this does not contradict the statement at the end of Section I.A that  $T_3$  prefers not to communicate in the Battle of the Sexes. Proposition 1 refers to payoffs averaged across player types, while the earlier remark referred only to  $T_3$ 's payoff given that he is certain that he faces a  $T_2$  opponent.

PROOF:

In Appendix 2.

Since  $H$  is risk dominant in Battle of the Sexes, one-way communication suffices to attain perfect coordination on the speaker's preferred equilibrium outcome. Thus, we here have a case in which the prediction from the level- $k$  model coincides with the prediction from Farrell (1988). The ineffectiveness of cheap talk in dominance solvable games is also analogous. More generally, the two approaches share the property that communication, if anything, pulls players towards Nash equilibria in symmetric  $2 \times 2$  games.

PROPOSITION 3: *For any distribution of types, the frequency of coordination on pure strategy Nash equilibrium profiles is weakly greater in  $\Gamma_I(G)$  than in  $G$ .*

PROOF:

In Appendix 2.

The pull towards Nash equilibria is so strong that one-way communication results in equilibrium play whenever  $T_{1+}$  meet. Moreover,  $T_{2+}$  always play the action corresponding to the sender's preferred equilibrium.

COROLLARY 4: *For type distributions with  $p_0 = 0$ , players in  $\Gamma_I(G)$  always coordinate on pure strategy Nash equilibrium profiles. If in addition  $p_1 = 0$ , players in  $\Gamma_I(G)$  always coordinate on the sender's preferred equilibrium.*

PROOF:

Follows directly from Tables A1 to A4 in the proof of Proposition 3.

Unlike one-way communication, two-way communication may destroy not only average payoffs but also coordination on equilibrium outcomes. For example, suppose there are only  $T_1$  players and let  $G$  be a coordination game in which payoff and risk dominance coincide. Then  $\Gamma_{II}(G)$  entails miscoordination in half of the cases, because  $T_1$  sends random messages while listening to received messages. By contrast, in  $G$  and in  $\Gamma_I(G)$  two  $T_1$  players always play the (payoff and risk) dominant equilibrium. Our model therefore captures the intuition that two-way communication can bring noise in the form contradictory messages.

Nevertheless, there are important classes of games in which two-way communication outperforms one-way communication.

PROPOSITION 5: *Given a standard type distribution, the average payoff associated with  $\Gamma_{II}(G)$  exceeds the average payoff associated with  $\Gamma_I(G)$  if and only if (i)  $G$  is a coordination game in which  $L$  is the risk dominant action and  $(4 - 3p_1)u_{HH} + p_1(u_{LH} + u_{HL}) > (4 - p_1)u_{LL}$ , or (ii)  $G$  is a mixed motive game with a type distribution satisfying the following condition:*

$$1 + \frac{2(p_1 - 1)(p_1 - 1 + 2p_3)}{p_1^2 + 4p_3^2} < \frac{u_{LL} - u_{HH}}{u_{LH} + u_{HL} - 2u_{HH}}.$$

PROOF:

In Appendix 2.

The Stag Hunt game in Figure 2 belongs to the first class of games identified by Proposition 5. For that particular game, two-way communication yields higher expected payoff than one-way communication whenever  $p_1 \in (0, 2/3)$ . The second class of games identified in Proposition 5 is harder to specify because of the cycling patterns of behavior under two-way communication

in mixed motive games. However, for two-way communication to be beneficial, the payoff when both players play  $L$  must be sufficiently high (at least  $(u_{HL} + u_{LH})/2$ ) and in addition the type distribution has to be such that the miscoordination outcome  $(L, L)$  happens sufficiently often with two-way communication. For example, with only  $T_3$  players, the outcome is  $(L, L)$  under two-way communication, whereas such players coordinate on an asymmetric equilibrium with one-way communication.

### C. Robustness

How robust are our results to the assumptions that we have made about players' behavior?

The largest difference in comparison with other level- $k$  applications is that we assume that players have a weak preference for truthfulness. If players have no preference for truthfulness, communication ceases to have any effect whatsoever in our model: behavior is the same in  $\Gamma_I(G)$ ,  $\Gamma_{II}(G)$  and  $G$ . This specification is strongly at odds with the evidence that communication matters in many game experiments.

Another alternative hypothesis is that all players prefer to be truthful, but that the most primitive types also respond systematically to received messages. The idea is that (if the actions of both players have the same label), the receiver could imitate or differentiate based on the sender's message. The most natural way to account for such imitation is to allow heterogeneous  $T_1$  players, some believing that receivers randomize, others believing that receivers imitate.<sup>15</sup> With one-way communication, this implies that some  $T_1$  players believe  $T_0$  receivers randomize, whereas others believe that they imitate. With two-way communication, some  $T_1$  players believe that opponents are truthful, whereas others believe they imitate. Let us now consider the consequences of this specification.

First consider the Stag Hunt in Figure 2. Under one-way communication,  $T_1$  senders who believe that receivers imitate send the message  $h$  and play  $H$ . This in turn implies that  $T_2$  receivers respond to messages as if they were truthful irrespective of which kind of  $T_1$  sender they think they face. Under one-way communication, the only difference compared to our original assumption is that there will be somewhat more coordination on  $(H, H)$  since some  $T_1$  senders now play  $H$ . Under two-way communication,  $T_1$  players who believe that opponents imitate send  $h$  and play  $H$  instead of responding to received messages.  $T_2$  players therefore optimally send  $h$  and play  $H$  irrespective of which type of  $T_1$  player they meet. Since miscoordination only occurs whenever two  $T_1$  players that send random messages meet, there will now be more equilibrium coordination compared to the standard case.

Second, consider one-way communication in the Battle of the Sexes. While  $T_1$  receivers, and hence  $T_2$  senders, behave as before,  $T_1$  senders that believe they face imitators now send  $l$  and play  $H$ . In the previous footnote, we have already argued that this behavior is implausible and that the fraction of such  $T_1$  players must therefore be small. However, irrespective of how small a proportion they constitute,  $T_2$  receivers now play  $L$  irrespective of what message they receive. This implies that  $T_3$  senders send  $h$  and play  $H$ . Under a standard type distribution, the outcome in terms of observed action profiles is thus the same as before.

Although some details of the analysis change with the introduction of heterogeneous  $T_1$  players, we conclude that the main mechanisms are robust to this modification.

Another cause for concern is our assumption about how unexpected messages are treated. An alternative assumption to ours is that  $T_k$  believes that unexpected messages are truthful.

<sup>15</sup>An alternative is to let  $T_1$  assume that some fraction of  $T_0$  imitates rather than randomizes. In this case,  $T_1$  is sophisticated enough to consider heterogeneity among  $T_0$ . We do not think this is plausible, and the consequences are counterfactual too: Consider one-way communication in the Battle of the Sexes. If there is heterogeneity among  $T_0$ ,  $T_1$  will send  $l$  and play  $H$ —believing that some opponents ignore their message, whereas others imitate their message and play  $L$ . Since  $p_1$  is typically estimated to be quite high, the implication is that sending  $l$  and playing  $H$  would be a relatively common practice. Cooper et al. (1989) studies one-way communication in Battle of the Sexes. They find that only 2 percent of all senders even sent an  $l$  message.

	X	Y	Z
U	5*, 5*	-50, -50	2, 4
V	-50, -50	2, 4*	4*, 3
W	4, 4*	3*, 3	3, 3

FIGURE 5: HIGH RISK GAME

This would not change any of our results for one-way communication, but it would imply that  $T_{3+}$  responds to messages (rather than always playing  $H$ ) in coordination games as well as slightly different behavior of  $T_{2+}$  in mixed motive games. (See also the discussion following Observation 7 about the sensitivity to the assumption about unexpected messages for the behavior of  $T_{2+}$  in mixed motive games.) In Appendix 3, we consider the cognitive hierarchy model in which no messages are unexpected (because all players take into account that the opponent might be  $T_0$ ) and show that, with the exception of  $T_{2+}$  in mixed motive games, our results are robust.

## II. Extensions

So far, we have confined attention to symmetric  $2 \times 2$  games. In principle it is straightforward to extend the analysis to games with more players and strategies. In this section, we show that the reassurance property of communication extends to two-player games in which players' interests are sufficiently well aligned. When attractive non-equilibrium outcomes are present, however, senders might try to obtain these by deceiving the opponent. The possibility of deception implies that one-way communication may hamper coordination on Nash equilibria. In addition, we show that multilateral communication in  $N$ -player games facilitates coordination in a class of common interest games.

### A. Two-player common interest games

The Stag Hunt example illustrates that pre-play communication facilitates the play of a risky payoff dominant equilibrium. Since our model does not assume equilibrium play, it is also applicable to situations in which players realistically fail to play a unique and efficient Nash equilibrium—such as the High Risk game, devised by Margaret Gilbert (1990) and reproduced in Figure 5 (in which best replies are marked with asterisks).<sup>16</sup> Absent communication, the level- $k$  model predicts that two  $T_{5+}$  players coordinate on the unique pure strategy equilibrium  $(U, X)$ , whereas all less sophisticated players fail to do so.<sup>17</sup> In contrast, one-way and two-way communication implies that  $T_{2+}$  coordinate on equilibrium. That is, much less sophistication is required to reach equilibrium with communication than without.<sup>18</sup>

<sup>16</sup>Experimental results of Anthony Burton and Martin Sefton (2004) confirm the prevalence of coordination failure in one-shot play of the High Risk game, but demonstrate that players learn to play the equilibrium after having played a number of practice rounds with the same opponent.

<sup>17</sup>To see this, note that  $T_1$  plays  $W$  and  $Z$  since these are the risk dominant actions. Using the best responses indicated in Figure 5 it follows that  $T_2$  plays  $V$  and  $X$ ,  $T_3$  plays  $U$  and  $Y$ ,  $T_4$  plays  $W$  and  $X$ , and finally that  $T_{5+}$  plays  $U$  and  $X$ .

A referee makes the following additional observation: “Replace 50 by  $x$ . If anyone is playing  $U$  or  $V$  in a first encounter with the game [...], the number of such players should decline as  $x$  increases.” We agree. This is another case in which the cognitive hierarchy model offers a richer and more realistic prediction than level- $k$ .

<sup>18</sup>To see this, first consider one-way communication. A  $T_1$  row sender sends  $w$  and plays  $W$ , while a column sender sends  $z$  and plays  $Z$ . A  $T_1$  receiver best responds to messages. A  $T_{2+}$  row sender therefore sends  $u$  and plays  $U$ , while a column sender sends  $x$  and plays  $X$ , while a  $T_{2+}$  receiver best responds to messages. Now consider two-way communication.  $T_1$  believes the opponent is truthful and therefore best responds to messages and randomize what message to send. A  $T_{2+}$  row player therefore sends  $u$  and plays  $U$  while a column player sends  $x$  and plays  $X$ .

	Y	Z
W	3*, 2*	4*, 0
X	0, 0	3, 3*

FIGURE 6: ASYMMETRIC 2 x 2 GAME

The positive effect of communication in the High Risk game extends to all finite and normal form two-player games which has a payoff dominant equilibrium that gives strictly higher payoffs to both players than all other outcomes of a game, i.e., to all *common interest games*. For this class of games it is straightforward to show that  $T_{2+}$  coordinate on the payoff dominant equilibrium. The underlying mechanism is that since  $T_1$  listens and best responds to messages,  $T_2$  can achieve the best possible outcome by sending and playing the payoff dominant equilibrium.

**PROPOSITION 6:** *Let  $G$  be a two-player common interest game. For type distributions with  $p_0 = p_1 = 0$ , players in  $\Gamma_I(G)$  and  $\Gamma_{II}(G)$  always coordinate on the payoff dominant Nash equilibrium.*

**PROOF:**

See Appendix 2.

### B. Other two-player games

In common interest games and in symmetric  $2 \times 2$  games with one-way communication, players always represent their intentions truthfully. In other classes of games, however, this is not necessarily the case. Crawford (2003) already shows how deception arises naturally in a level- $k$  model of communication in Hide-and-Seek games. Deception can also arise in an asymmetric dominance solvable  $2 \times 2$  game with a unique pure strategy equilibrium. Consider the game in Figure 6.

The game's unique pure strategy equilibrium is  $(W, Y)$ . Since  $W$  and  $Y$  are the risk dominant actions,  $T_{1+}$  players coordinate on the  $(W, Y)$  equilibrium if they are not allowed to communicate. Now consider one-way communication. Suppose that the row player acts as sender and the column player acts as receiver. The  $T_1$  sender sends  $w$  and plays  $W$ , while a  $T_1$  receiver best responds to received messages. A  $T_2$  sender therefore sends  $x$ , but plays  $W$ , while a  $T_2$  receiver best responds to messages.  $T_3$  sends  $x$  but plays  $W$ , while a  $T_3$  receiver ignores messages and always plays  $Y$ . Whenever  $T_{3+}$  players meet, the resulting outcome is equilibrium play, but not when less sophisticated players play. In contrast to Proposition 3, one-way communication leads to less equilibrium coordination than no communication unless all players carry out three or more thinking steps.

Proposition 3 does not generalize to symmetric two-player games with more than two actions either. To see this, consider the game in Figure 7.<sup>19</sup> This symmetric  $3 \times 3$  game has a unique pure strategy equilibrium,  $(H, H)$ , for all  $q > 1$ , but the game also has the asymmetric outcomes  $(H, L)$  and  $(L, H)$  that are attractive either to the row or column player. Since there is a third strategy,  $D$ , which has  $L$  as its best response, some senders will try to use this strategy to deceive the other player into playing  $L$ .

Specifically, consider the case when  $q = 1$  and pre-play communication is not possible. In that case  $T_1$  would play  $H$  since it is the best action to take if the opponent randomizes uniformly, and  $T_{2+}$  would best respond by playing  $H$ . One-way communication, however, makes it more difficult to reach equilibrium. A  $T_1$  sender sends  $h$  and plays  $H$ , while a  $T_1$  receiver best responds (as indicated by the asterisks in Figure 7) to the received message. A  $T_2$  sender sends  $d$ , but

<sup>19</sup>This game is non-generic, but the analysis is analogous in the generic case.

	$H$	$L$	$D$
$H$	$4/q^*, 4/q^*$	$(4 + 1/q)^*, 0$	$0, 0$
$L$	$0, (4 + 1/q)^*$	$0, 0$	$1^*, 1$
$D$	$0, 0$	$1, 1^*$	$0, 0$

FIGURE 7: SYMMETRIC  $3 \times 3$  GAME

	$H_1$	$L_1$	$D_1$	$H_2$	$L_2$	$D_2$	$\dots$	$H_Q$	$L_Q$	$D_Q$
$H_1$	4, 4	5, 0	0, 0	0, 0	0, 0	0, 0	$\dots$	0, 0	0, 0	0, 0
$L_1$	0, 5	0, 0	1, 1	0, 0	0, 0	0, 0	$\dots$	0, 0	0, 0	0, 0
$D_1$	0, 0	1, 1	0, 0	0, 0	0, 0	0, 0	$\dots$	0, 0	0, 0	0, 0
$H_2$	0, 0	0, 0	0, 0	2, 2	4.5, 0	0, 0	$\dots$	0, 0	0, 0	0, 0
$L_2$	0, 0	0, 0	0, 0	0, 4.5	0, 0	1, 1	$\dots$	0, 0	0, 0	0, 0
$D_2$	0, 0	0, 0	0, 0	0, 0	1, 1	0, 0	$\dots$	0, 0	0, 0	0, 0
$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\vdots$	$\ddots$	0, 0	0, 0	0, 0
$H_Q$	0, 0	0, 0	0, 0	0, 0	0, 0	0, 0	0, 0	$\frac{4}{Q}, \frac{4}{Q}$	$4 + \frac{1}{Q}, 0$	0, 0
$L_Q$	0, 0	0, 0	0, 0	0, 0	0, 0	0, 0	0, 0	$0, 4 + \frac{1}{Q}$	0, 0	1, 1
$D_Q$	0, 0	0, 0	0, 0	0, 0	0, 0	0, 0	0, 0	0, 0	1, 1	0, 0

FIGURE 8: SYMMETRIC  $3Q \times 3Q$  GAME

plays  $H$ , while a  $T_2$  receiver best responds to received messages. A  $T_3$  sender sends  $d$  and plays  $H$ , while a  $T_3$  receiver plays  $H$  irrespective of the received message. A  $T_{4+}$  sender is indifferent about what message to send and is thus truthful, sending  $h$  and playing  $H$ ; a  $T_{4+}$  receiver ignores messages and plays  $H$ . We conclude that  $T_{3+}$  coordinate on  $(H, H)$  and that one-way communication consequently lowers equilibrium coordination unless all players make three or more thinking steps.

A modification of the game illustrates how the number of thinking-steps required to reach equilibrium may increase linearly with the size of the game. Consider the  $3Q \times 3Q$  game shown in Figure 8. It has the game in Figure 7 on the main diagonal and zero payoffs elsewhere.

Let messages be denoted  $m_q$ , with  $m \in \{h, l, d\}$  and  $q \in \{1, 2, \dots, Q\}$ . Without communication,  $T_{1+}$  plays  $H_1$  as in the  $3 \times 3$  game. However, when one-way communication is allowed, all players must make at least  $2Q + 2$  thinking steps in order to coordinate on the unique equilibrium  $(H_1, H_1)$ . To see why, note first that  $T_1$  through  $T_3$  will behave as in the  $3 \times 3$  game, but that receivers will best-respond to all messages  $m_q$  with  $q \in \{2, 3, \dots, Q\}$ , believing those messages to come from  $T_0$ . A  $T_4$  sender therefore sends  $d_2$  and plays  $H_2$  in order to get the outcome  $(H_2, D_2)$  which is preferred over  $(H_1, H_1)$ .  $T_5$  receivers do not believe in  $d_2$  messages and therefore play  $H_2$  if either  $h_2, l_2$  or  $d_2$  is played. In turn,  $T_6$  senders send  $d_3$  and play  $H_3$  in order to induce the  $(H_3, L_3)$  outcome. The inductive argument continues like this up until  $T_{2Q+1}$  sends  $d_Q$  and plays  $H_Q$ . A  $T_{2Q+2}$  sender cannot hope to get anything better than  $(H_1, H_1)$  and therefore sends  $h_1$  and plays  $H_1$ , whereas a  $T_{2Q+2}$  receiver plays  $H_q$  whenever  $h_q, d_q$  or  $l_q$  is played (for all  $q$ ).

This example illustrates that the degree of sophistication required to play equilibrium increases with the size of the game. Since the degree of sophistication required is unrealistically high, in these games players coordinate better if they are unable to communicate.

### C. Other communication protocols

Like much of the cheap talk literature, we have here considered communication of intentions. Messages are of the form “I plan to play...”. What would happen if players communicated requests instead, that is if messages were of the form “I want you to play...”? While the model still admits a notion of truthfulness, the analysis would be quite different. For example, it is no longer clear that  $T_1$  players should care about the messages that they receive, since  $T_0$  players’ requests may reveal nothing about their intentions. We thus expect that credulity will play a more important role than truthfulness in this case. Specifically, communication might now affect behavior if  $T_1$  senders believe that receivers are credulous in the sense that they fulfill requests. Preliminary investigations suggest that the ensuing analysis offers a perspective on how cheap talk may be used to understand cheating in games, but we leave a fuller analysis for a separate paper.

The paper only considers the communication of intentions by interested parties. A natural avenue for future research is to study the communication of desires or recommendations, by players themselves as well as by more or less interested third parties such as managers.<sup>20</sup>

Another natural extension is to consider multiple rounds of communication. Crawford (2007) has already used the level- $k$  model to analyze longer conversations in the Battle of Sexes. He demonstrates that longer bilateral conversations improve coordination rates in a way that is qualitatively similar to, but quantitatively and intuitively different from, the equilibrium analysis of Farrell (1987).

### D. Multilateral communication

Communication may also facilitate play of a potentially risky payoff dominant equilibrium in games with more than two players. In this section we show that for all finite common interest games with strategic complementarities and positive spillovers, multilateral pre-play communication facilitates play of the payoff dominant equilibrium whenever  $T_{2+}$  play the game.

We restrict attention to games with unique best responses. The actions of each player are assigned integers  $\{1, 2, \dots, \bar{a}_i\}$ . Such a game has *strategic complementarities* if best responses are non-decreasing in the opponents’ actions, i.e., if  $a_{-i} \geq a'_{-i}$  implies  $BR_i(a_{-i}) \geq BR_i(a'_{-i})$ . Finally, a game has *positive spillovers* if the own payoff increases in the opponents’ actions, i.e., if  $a_{-i} \geq a'_{-i}$  implies  $\pi_i(a_i, a_{-i}) \geq \pi_i(a_i, a'_{-i})$ . Note that the payoff dominant equilibrium of a common interest game involves all players choosing their highest actions,  $\bar{a} = (\bar{a}_1, \bar{a}_2, \dots, \bar{a}_n)$ .<sup>21</sup>

**PROPOSITION 7:** *Let  $G$  be a finite common interest game with unique best responses, strategic complementarities and positive spillovers. For type distributions with  $p_0 = p_1 = 0$ , players in  $\Gamma_N(G)$  always coordinate on the payoff dominant Nash equilibrium of  $G$ .*

**PROOF:**

See Appendix 2.

To better understand the intuition behind Proposition 7, consider the Weak-link game. In a Weak-link game, each player picks an integer from 1 to  $M$ . Payoffs are such that all players want to play the minimum of what the opponents play, but all players are better off if everybody

<sup>20</sup>For experimental evidence bearing on these issues, see for example John B. van Huyck, Ann B. Gillette and Raymond C. Battalio (1992), Roberto Weber et al. (2001) and Jordi Brandts and David J. Cooper (2007).

<sup>21</sup>To see this, suppose the payoff dominant equilibrium is some profile  $a^* \neq \bar{a}$ . Then at least one player has an action  $a_i > a_i^*$  available that by positive spillovers gives the opponents the same or higher payoffs, contradicting the assumption that  $a^*$  is the payoff dominant equilibrium that yields strictly higher payoffs to all players than all other outcomes of the game.



chooses higher numbers. Any strategy profile in which all players choose the same number constitutes a Nash equilibrium, and the payoff dominant equilibrium involves all players playing  $M$ . Note that the Weak-link game is essentially a Stag Hunt game with more than two strategies. (For a more detailed exposition of the Weak-link game see for example Camerer 2003, Chapter 7.) If the Weak-link game is preceded by multilateral communication,  $T_1$  sends a random message and best responds by playing the minimum of the received messages. A  $T_2$  player faces  $N - 1$  opponents that play the minimum of the received messages, so  $T_{2+}$  best responds by also playing the minimum of all received messages, but sends the message  $m$  (so that the best outcome occurs if all other players sent  $m$ ). Hence, as long as there are no  $T_0$  and  $T_1$  players, there will be perfect coordination on the payoff dominant equilibrium.

Note that the logic of this argument breaks down if only a subset of the players is allowed to send a message. To see this, suppose that all but one player is allowed to send a message. Then it is generally no longer optimal for  $T_1$  to play the minimum of the received message profile since one opponent's action is unpredictable, which in turn implies that  $T_2$  does not play the minimum of the received messages, which would be required to guarantee play of the payoff dominant equilibrium.

### III. Evidence

The level- $k$  model of pre-play communication is primarily a model to explain initial responses, i.e., the behavior of players that play a game for the first time. If players gain experience of the game and the population of players, they are likely to change their model of opponents' behavior or perhaps think further and proceed to higher levels of reasoning. In experimental work on pre-play communication, players typically play the same game in several rounds. Strictly speaking, most of the available evidence is thus inadequate for our purposes.

Another difficulty is that experimenters rarely elicit subjects' von Neumann-Morgenstern utilities. Instead, payoffs are typically monetary. In order to interpret the behavior as evidence of beliefs, experimenters thus have to assume a particular relationship between monetary allocations and utility. For example, they may assume that subjects maximize their own expected monetary payoff. However, subjects frequently have other goals; for example, it would be ludicrous to interpret Dictator game giving as evidence that subjects are confused about the game's payoffs. In principle, we should always distinguish *games* (involving utilities) from *game forms* (involving monetary payoffs).

With these caveats in mind, and continuing to conflate games and game forms, let us briefly discuss some of the most relevant communication experiments.

Two papers contrast one-way and two-way communication in Stag Hunt games. Cooper et al. (1992) report that average coordination on the payoff dominant equilibrium is 0 percent without communication, 53 percent with one-way communication and 91 percent with two-way communication. This study therefore suggests that communication plays a reassurance role, as emphasized by Crawford (1998).<sup>22</sup> By contrast, in a Stag Hunt game with somewhat different relative payoffs, Burton, Graham Loomes and Sefton (2005) find that one-way communication results in 52 percent coordination on the payoff dominant equilibrium, whereas two-way communication entailed average coordination on the payoff dominant equilibrium of only 34 percent. Both papers find that behavior varies substantially across sessions, indicating that heterogeneity in early rounds of the game affect players choices in later rounds. Burton, Loomes and Sefton (2005) also collect data on some of their individual subjects' complete strategies (plans). By far the most common strategy, in our notation, is  $\langle h, H, L \rangle$ . According to the level- $k$  model, this strategy should only be used by half of the  $T_1$  players. On the other hand, the strategy is used by all  $T_{2+}$  in the cognitive hierarchy model (under the weak assumption that the average of the

<sup>22</sup>Relatedly, Ellingsen and Johannesson (2004a) identifies a reassurance role of communication in hold-up games with multiple equilibria.

type distribution is below 7). As usual when the two models yield conflicting predictions, the cognitive hierarchy model's prediction is preferable.

In addition to the two studies comparing one-way and two-way communication, there are also a few studies of the Stag Hunt game that investigate either one-way or two-way communication. John Duffy and Nick Feltovich (2002) find that one-way communication entails coordination on the payoff dominant equilibrium in 84 percent of the cases with one-way communication and in 61 percent of the cases without communication. Charness (2000) studies the effect of one-way communication in three versions of the Stag Hunt and finds 86 percent coordination on the payoff dominant equilibrium with one-way communication. Kenneth Clark, Stephen Kay and Sefton (2001) study two-way communication in two different Stag Hunt games. In the first game, based on Cooper et al. (1992), playing  $L$  yields the same payoff irrespective of the opponents behavior. In this game, coordination on the payoff dominant equilibrium is 2 percent without communication and 70 with two-way communication. In a more standard Stag Hunt game, they find that coordination on the payoff dominant equilibrium occurs in only 19 percent of the cases with two-way communication. Hence, it appears possible that the beneficial effects of two-way communication in Cooper et al. (1992) is sensitive to their choice of payoff matrix.

On the other hand, the hypothesis that multilateral communication plays a major role in creating reassurance is more consistently supported by evidence from Weak-link games (see the previous subsection for a definition); see in particular Blume and Andreas Ortmann (2007). Camerer and Weber (2007) summarize the existing evidence as follows: "Taken together, the above results suggest that communication can help solve even the most difficult coordination problems, with relatively large numbers of players and where the minimum effort determines the entire group's output. However, the communication required to get large groups to efficiency is extreme—players must all send messages and have public knowledge of messages." This is precisely what Proposition 7 predicts, under the additional prerequisite that all players are sufficiently sophisticated.

For mixed motive games the picture also seems clear, although the consistency in this case may be due to the low number of studies. Cooper et al. (1989) find that one-way communication results in a high degree of coordination in Battle of the Sexes. Averaged over several rounds of play, Cooper et al. (1989) report that one-way communication increases coordination from 48 percent without communication to 95 percent with one-way communication. With one round of two-way communication, coordination is 55 percent.<sup>23</sup> For a comparison of this evidence with the prediction of Rabin's (1994) cheap talk model, see Costa-Gomes (2002).

To summarize, we believe that more experimental work is needed in order to test the theory laid out in this paper. Such a test should focus on players' initial responses to several different games, which would allow a clearer separation of types. Costa-Gomes and Crawford (2006) illustrates how this can be done. It would also be useful to directly test the assumption about  $T_0$  players. Since  $T_0$  players mainly exist in the minds of other players, we need data on players' beliefs. Such data can be generated not only through belief elicitation (e.g., Costa-Gomes and Georg Weizsäcker 2008), but also by response time measurement (e.g., Camerer et al. 1993 and Ariel Rubinstein 2007), information search (e.g., Camerer et al. 1993, Costa-Gomes, Crawford and Broseta 2001 and Costa-Gomes and Crawford 2006) and through neuroimaging (e.g., Meghana Bhatt and Camerer 2005).

<sup>23</sup>It should be noted, however, that Cooper et al. (1989) allow the players to be silent and that 27 percent of the players in the two-way treatment, and 5 percent in the one-way treatment, choose to do so. We have not allowed silence in our analysis. It is of course possible to extend the message space to allow for voluntary silence, but we have chosen not to do so. Since players are assumed to have a slight preference for truthfulness, they might want to be silent when they don't know what action they are going to take in the action game (as  $T_1$  under two-way communication in coordination games).

#### IV. Conclusion

Coordination of behavior in new strategic situations is facilitated by communication. Since communication seeks to affect the beliefs of others, assumptions about initial beliefs are central to the analysis. Our starting point is that prevailing assumptions about initial beliefs in the strategic communication literature, as captured by the rationalizability assumption, are problematic. Thus, we consider the role of communication within the two other general models of initial beliefs that have won widespread acceptance, namely the level- $k$  and cognitive hierarchy models. Our analysis demonstrates that these models generate sharp predictions that are often, if not always, intuitively plausible.

We see two immediate avenues for future research. First, there is a need for evidence that systematically distinguishes the effects of preferences, beliefs and rationality. In particular, identification of beliefs is only possible when preferences and rationality are controlled for. The task is not easy, since direct elicitation of beliefs tends to yield quite different measures than revealed measures of beliefs, possibly because direct elicitation affect the depth of subjects' strategic thinking; see Costa-Gomes and Weizsäcker (2008). Cleaner evidence would help us to evaluate existing models of beliefs and to suggest new ones, and it would clarify the status of previously puzzling experimental evidence concerning the effect of communication in games like Stag Hunt. Secondly, we entirely lack evidence concerning the effects of communication in many of the other games studied in this paper. Experimental evidence on such games would allow a true out-of-sample test of the theory.

#### REFERENCES

- Anderlini, Luca.** (1999). 'Communication, Computability, and Common Interest Games', *Games and Economic Behavior* 27, 1–37.
- Aumann, Robert J.** (1974). 'Subjectivity and Correlation in Randomized Strategies', *Journal of Mathematical Economics* 1(1), 67–96.
- Aumann, Robert J.** (1990), Nash Equilibria are not Self-Enforcing, in **Jean J. Gabszewicz, Jean-François Richard and Laurence A. Wolsey.**, eds, 'Economic Decision-Making: Games, Econometrics and Optimization', Elsevier, Amsterdam, chapter 10, pp. 201–206.
- Banerjee, Abhijit and Weibull, Jörgen W.** (2000). 'Neutrally Stable Outcomes in Cheap-Talk Coordination Games', *Games and Economic Behavior* 32(1), 1–24.
- Beard, T. Randolph and Beil Jr., Richard O.** (1994). 'Do People Rely on the Self-interested Maximization of Others? An Experimental Test', *Management Science* 40(2), 252–262.
- Bhatt, Meghana and Camerer, Colin F.** (2005). 'Self-Referential Thinking and Equilibrium as States of Mind in Games: fMRI Evidence', *Games and Economic Behavior* 52(2), 424–459.
- Blume, Andreas.** (1998). 'Communication, Risk, and Efficiency in Games', *Games and Economic Behavior* 22(2), 171–202.
- Blume, Andreas and Ortmann, Andreas.** (2007). 'The Effects of Costless Pre-play Communication: Experimental Evidence from Games with Pareto-ranked Equilibria', *Journal of Economic Theory* 132(1), 274–290.
- Brandts, Jordi and Cooper, David J.** (2007). 'It's What You Say, Not What You Pay: An Experimental Study of Manager-Employee Relationships in Overcoming Coordination Failure', *Journal of the European Economic Association* 5(6), 1223–1268.
- Brown, Alexander L., Camerer, Colin F. and Lovo, Dan.** (2009), 'To Review or Not to Review? Limited Strategic Thinking at the Movie Box Office', Mimeo.  
**URL:** <http://ssrn.com/abstract=1281006>

- Burton, Anthony, Loomes, Graham and Sefton, Martin.** (2005), Communication and Efficiency in Coordination Game Experiments, in **John Morgan.**, ed., 'Experimental and Behavioral Economics', Vol. 13 of *Advances in Applied Microeconomics*, JAI Press, pp. 63–85.
- Burton, Anthony and Sefton, Martin.** (2004). 'Risk, Pre-play Communication and Equilibrium', *Games and Economic Behavior* 46(1), 23–40.
- Cai, Hongbin and Wang, Joseph Tao-Yi.** (2006). 'Overcommunication in Strategic Information Transmission Games', *Games and Economic Behavior* 56(1), 7–36.
- Camerer, Colin F.** (2003), *Behavioral Game Theory*, Princeton University Press, Princeton.
- Camerer, Colin F., Ho, Teck-Hua and Chong, Juin-Kuan.** (2004). 'A Cognitive Hierarchy Model of Games', *Quarterly Journal of Economics* 119(3), 861–898.
- Camerer, Colin F., Johnson, Eric J., Rymon, Talia and Sen, Sankar.** (1993), Cognition and Framing in Sequential Bargaining for Gains and Losses, in **Ken Binmore, Alan Kirman and Piero Tani.**, eds, 'Frontiers of Game Theory', MIT Press, Boston, pp. 27–48.
- Camerer, Colin F. and Weber, Roberto.** (2007), 'Experimental Organizational Economics', Forthcoming in Handbook of Organizational Economics.
- Charness, Gary.** (2000). 'Self-Serving Cheap Talk: A Test of Aumann's Conjecture', *Games and Economic Behavior* 33(2), 177–194.
- Clark, Kenneth, Kay, Stephen and Sefton, Martin.** (2001). 'When are Nash Equilibria Self-Enforcing? An Experimental Analysis', *International Journal of Game Theory* 29(4), 495–515.
- Cooper, Russell, DeJong, Douglas V., Forsythe, Robert and Ross, Thomas W.** (1989). 'Communication in the Battle of the Sexes Game: Some Experimental Results', *RAND Journal of Economics* 20(4), 568–587.
- Cooper, Russell, DeJong, Douglas V., Forsythe, Robert and Ross, Thomas W.** (1992). 'Communication in Coordination Games', *Quarterly Journal of Economics* 107(2), 739–771.
- Costa-Gomes, Miguel A.** (2002). 'A Suggested Interpretation of Some Experimental Results on Pre-play Communication', *Journal of Economic Theory* 104, 104–136.
- Costa-Gomes, Miguel A. and Crawford, Vincent P.** (2006). 'Cognition and Behavior in Two-Person Guessing games: An Experimental Study', *American Economic Review* 96, 1737–1768.
- Costa-Gomes, Miguel A., Crawford, Vincent P. and Broseta, Bruno.** (2001). 'Cognition and Behavior in Normal-Form Games: An Experimental Study', *Econometrica* 69(5), 1193–1235.
- Costa-Gomes, Miguel A., Crawford, Vincent P. and Iriberry, Nagore.** (2009). 'Comparing Models of Strategic Thinking in Van Huyck, Battalio, and Beil's Coordination Games', *Journal of the European Economic Association* 7, forthcoming.
- Costa-Gomes, Miguel A. and Weizsäcker, Georg.** (2008). 'Stated Beliefs and Play in Normal-form Games', *Review of Economic Studies* 75(3), 729–762.
- Crawford, Vincent P.** (1998). 'A Survey of Experiments on Communication via Cheap Talk', *Journal of Economic Theory* 78(2), 286–298.
- Crawford, Vincent P.** (2003). 'Lying for Strategic Advantage: Rational and Boundedly Rational Misrepresentation of Intentions', *American Economic Review* 93(1), 133–149.
- Crawford, Vincent P.** (2007), 'Let's Talk It Over: Coordination via Preplay Communication with Level-k Thinking', Mimeo.
- Crawford, Vincent P. and Iriberry, Nagore.** (2007a). 'Fatal Attraction: Salience, Naivete, and Sophistication in Experimental Hide-and-Seek Games', *American Economic Review* 97(5), 1731–1750.

- Crawford, Vincent P. and Iriberry, Nagore.** (2007*b*). ‘Level-k Auctions: Can a Non-Equilibrium Model of Strategic Thinking Explain the Winner’s Curse and Overbidding in Private-Value Auctions?’, *Econometrica* 75(6), 1721–1770.
- Crawford, Vincent P., Kugler, Tamar, Neeman, Zvika and Pauzner, Ady.** (2009). ‘Behaviorally Optimal Auction Design: An Example and Some Observations’, *Journal of the European Economic Association* 7, forthcoming.
- Crawford, Vincent P. and Sobel, Joel.** (1982). ‘Strategic Information Transmission’, *Econometrica* 50, 1141–1152.
- Demichelis, Stefano and Weibull, Jörgen W.** (2008). ‘Language, Meaning and Games - A Model of Communication, Coordination and Evolution’, *American Economic Review* 98(4), 1292–1311.
- Duffy, John and Feltovich, Nick.** (2002). ‘Do Actions Speak Louder Than Words? An Experimental Comparison of Observation and Cheap Talk’, *Games and Economic Behavior* 39(1), 1–27.
- Ellingsen, Tore and Johannesson, Magnus.** (2004*a*). ‘Is There a Hold-Up Problem?’, *Scandinavian Journal of Economics* 106, 475–494.
- Ellingsen, Tore and Johannesson, Magnus.** (2004*b*). ‘Promises, Threats, and Fairness’, *Economic Journal* 114, 397–420.
- Farrell, Joseph.** (1987). ‘Cheap Talk, Coordination, and Entry’, *RAND Journal of Economics* 18(1), 34–39.
- Farrell, Joseph.** (1988). ‘Communication, Coordination and Nash Equilibrium’, *Economic Letters* 27(3), 209–214.
- Farrell, Joseph and Rabin, Matthew.** (1996). ‘Cheap Talk’, *Journal of Economic Perspectives* 10(3), 103–118.
- Farrell, Joseph and Saloner, Garth.** (1988). ‘Coordination Through Committees and Markets’, *RAND Journal of Economics* 19(2), 235–252.
- Fudenberg, Drew and Tirole, Jean.** (1991), *Game Theory*, MIT Press, Cambridge.
- Genesove, David and Mullin, Wallace P.** (2001). ‘Rules, Communication, and Collusion: Narrative Evidence from the Sugar Institute Case’, *American Economic Review* 91(3), 379–398.
- Gilbert, Margaret.** (1990). ‘Rationality, Coordination, and Convention’, *Synthese* 84(1), 1–21.
- Green, Jerry R. and Stokey, Nancy L.** (2007). ‘A Two-person Game of Information Transmission’, *Journal of Economic Theory* 135(1), 90–104.
- Ho, Teck-Hua, Camerer, Colin F. and Weigelt, Keith.** (1998). ‘Iterated Dominance and Iterated Best Response in Experimental “p-Beauty Contests”’, *American Economic Review* 88(4), 947–969.
- Kawagoe, Toshiji and Takizawa, Hirokazu.** (2008*a*). ‘Equilibrium Refinement vs. Level-k Analysis: An Experimental Study of Cheap-talk Games with Private Information’, *Games and Economic Behavior* p. forthcoming.
- Kawagoe, Toshiji and Takizawa, Hirokazu.** (2008*b*), ‘Level-k Analysis of Experimental Centipede Games’, Mimeo.  
URL: <http://ssrn.com/abstract=1289514>
- Kerr, Norbert L. and Kaufman-Gilliland, Cynthia M.** (1994). ‘Communication, Commitment and Cooperation in Social Dilemmas’, *Journal of Personality and Social Psychology* 66(3), 513–529.
- Kim, Yong-Gwan and Sobel, Joel.** (1995). ‘An Evolutionary Approach to Pre-play Communication’, *Econometrica* 63(5), 1181–1193.
- Lo, Pei-yu.** (2007), ‘Language and Coordination Games’, Mimeo.

- Matsui, Akihiko.** (1991). ‘Cheap Talk and Cooperation in Society’, *Journal of Economic Theory* 54(2), 245–258.
- Myerson, Roger.** (1989). ‘Credible Negotiation Statements and Coherent Plans’, *Journal of Economic Theory* 48(1), 264–303.
- Nagel, Rosemarie.** (1995). ‘Unraveling in Guessing Games: An Experimental Study’, *American Economic Review* 85(5), 1313–1326.
- Nowak, Martin A.** (2006), *Evolutionary Dynamics*, Belknap Press of Harvard University Press, Cambridge.
- Östling, Robert, Wang, Joseph Tao-yi, Chou, Eileen and Camerer, Colin F.** (2008), ‘Strategic Thinking and Learning in the Field and Lab: Evidence from Poisson LUPI Lottery Games’, SSE/EFI Working Paper Series in Economics and Finance No. 671, Stockholm School of Economics.  
URL: <http://ssrn.com/abstract=1007181>
- Pinker, Steven and Bloom, Paul.** (1990). ‘Natural Language and Natural Selection’, *Behavioral and Brain Sciences* 13(4), 707–784.
- Rabin, Matthew.** (1990). ‘Communication Between Rational Agents’, *Journal of Economic Theory* 51(1), 144–170.
- Rabin, Matthew.** (1994). ‘A Model of Pre-Game Communication’, *Journal of Economic Theory* 63(2), 370–391.
- Rosenthal, Robert.** (1981). ‘Games of Perfect Information, Predatory Pricing and the Chain Store Paradox’, *Journal of Economic Theory* 25(1), 92–100.
- Rubinstein, Ariel.** (2007). ‘Instinctive and Cognitive Reasoning: A Study of Response Times’, *Economic Journal* 117, 1243–1259.
- Schelling, Thomas C.** (1966), *Arms and Influence*, Yale University Press, New Haven.
- Stahl, Dale O. and Wilson, Paul W.** (1994). ‘Experimental Evidence on Players’ Models of Other Players’, *Journal of Economic Behavior and Organization* 25(3), 309–327.
- Stahl, Dale O. and Wilson, Paul W.** (1995). ‘On Players’ Models of Other Players: Theory and Experimental Evidence’, *Games and Economic Behavior* 10(1), 33–51.
- van Huyck, John B., Gillette, Ann B. and Battalio, Raymond C.** (1992). ‘Credible Assignments in Coordination Games’, *Games and Economic Behavior* 4(1), 606–626.
- Wärneryd, Karl.** (1991). ‘Evolutionary Stability in Unanimity Games with Cheap Talk’, *Economic Letters* 36(4), 375–378.
- Weber, Roberto, Camerer, Colin F., Rottenstreich, Yuval and Knez, Marc.** (2001). ‘The Illusion of Leadership: Misattribution of Cause in Coordination Games’, *Organization Science* 12(5), 582–598.
- Weibull, Jörgen W.** (1995), *Evolutionary Game Theory*, MIT Press.
- Wengström, Erik.** (2008). ‘Price Competition, Level-k Theory and Communication’, *Economics Bulletin* 3(66), 1–15.

### Appendix 1: Characterization of Behavior

We here characterize behavior in all symmetric and generic  $2 \times 2$  games using the level- $k$  model. Consider the symmetric  $2 \times 2$  game in Figure A1.

We assume that this game is generic in the sense that none of the four different payoffs ( $u_{HH}, u_{HL}, u_{LH}$  and  $u_{LL}$ ) are identical. Depending on the relations  $u_{HH} \leq u_{LH}$  and  $u_{LL} \leq u_{HL}$ , we can divide the class of generic  $2 \times 2$  games into three familiar types of games as shown in Figure A2.<sup>24</sup>

<sup>24</sup>The classification of symmetric games follows Weibull (1995) closely. To understand how this classification

	$H$	$L$
$H$	$u_{HH}, u_{HH}$	$u_{HL}, u_{LH}$
$L$	$u_{LH}, u_{HL}$	$u_{LL}, u_{LL}$

FIGURE A1: SYMMETRIC 2 x 2 GAME

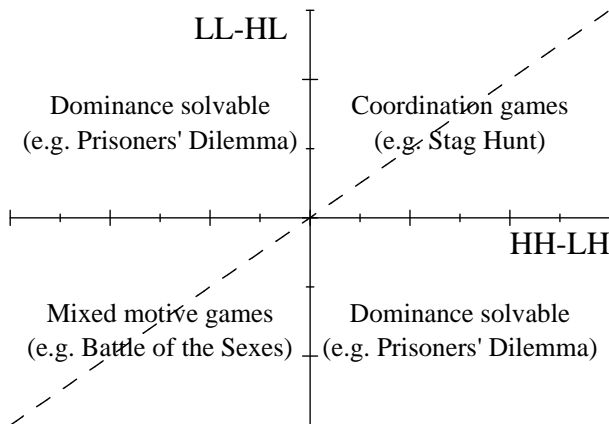


FIGURE A2: THE FOUR TYPES OF GENERIC AND SYMMETRIC 2 x 2 GAMES

If we were only interested in Nash equilibria, there would be only one prediction for each of these games. For the level- $k$  model, however, these games will be divided into subclasses with different predictions. The most important distinction is indicated by the dashed line in Figure A2. This condition corresponds to whether  $u_{LL} - u_{HL} \leq u_{HH} - u_{LH}$ , i.e., whether  $u_{LH} + u_{LL} \leq u_{HH} + u_{HL}$ . This means that action  $H$  is risk dominant above the dashed line in Figure A2, whereas action  $L$  is risk dominant below it. For tractability, we disregard the cases when neither action is risk dominant throughout the paper.

#### *Dominance solvable games*

Dominance solvable games are easiest to analyze, but also least interesting. In a dominance solvable game, players always have an incentive to play the dominant action, and neither one-way or two-way communication affect the actions players take.

arises, note that if we were only interested in Nash equilibria of  $2 \times 2$  games, we could have subtracted  $u_{LH}$  from both action  $H$  and  $L$  when the other player plays  $H$  and  $u_{HL}$  from both actions when the other player plays  $L$ . This would leave the equilibria of the game unchanged, whereas it affects the prediction for level- $k$  models. The main reason is that in a level- $k$  model, strategic uncertainty plays a role due to the randomization of level-0 players and we can therefore not use the sure-thing principle to transform the game. After the transformation, the game is the following.

	$H$	$L$
$H$	$u_{HH} - u_{LH}, u_{HH} - u_{LH}$	$0, 0$
$L$	$0, 0$	$u_{LL} - u_{HL}, u_{LL} - u_{HL}$

>From this game it is clear why the class of symmetric games can be classified by two real numbers,  $u_{HH} - u_{LH}$  and  $u_{LL} - u_{HL}$ .

We assume  $u_{HL} > u_{LL}$  and  $u_{HH} > u_{LH}$  so that  $H(\text{igh})$  is the dominant action. The case when  $L$  is the dominant action is symmetric.

**OBSERVATION 1:** *If players cannot communicate,  $T_{1+}$  plays the dominant action  $H$ . If players can communicate, then both one-way and two-way communication implies that  $T_{1+}$  sends  $h$  and plays  $H$  irrespective of any received messages.*

**PROOF:**

Since  $H$  is a dominant action,  $T_{1+}$  players play  $H$  irrespective of the believed behavior of the opponent. With the possibility to communicate, this also implies that there are no players that respond to messages, and  $T_{1+}$  players are therefore indifferent about sending  $h$  or  $l$ . (Sending  $l$  would have been beneficial if some players responded to messages and  $u_{HL} > u_{HH}$  as in the Prisoners' Dilemma.) However, since players have a lexicographic preference for truthfulness, they send  $h$ .

For dominance solvable  $2 \times 2$  games, communication plays no role. Except for some miscoordination due to  $T_0$  playing the dominated action, all players play the dominant action. Since the proof only relies on the fact that each player has a strictly dominant strategy, the result extends to all normal form two-player games in which both players have a strictly dominant action.

### *Coordination games*

Behavior in coordination games depends crucially on payoff and risk dominance. Since we restrict attention to generic games, one of the equilibria has to be payoff dominant. Let us without loss of generality assume that  $H(\text{igh})$  is the payoff dominant equilibrium, i.e.,  $u_{HH} > u_{LL}$ .

**OBSERVATION 2:** *(No communication)  $T_{1+}$  plays the risk dominant action.*

**PROOF:**

$T_1$  players believe that the opponent randomizes uniformly and therefore plays the risk dominant action.  $T_2$  players best respond and play the same risk dominant action, and so on.

Absent communication,  $T_1$  plays the best response to a uniformly randomizing  $T_0$  opponent, which is the risk dominant action. Since this is a coordination game, more advanced players best respond by playing the same action.

**OBSERVATION 3:** *(One-way communication) If  $H$  is the risk dominant action,  $T_{1+}$  sends  $h$  and plays  $H$  as sender and responds to messages as receiver. If  $L$  is the risk dominant action,  $T_1$  sends  $l$  and plays  $L$  as sender and responds to messages as receiver.  $T_{2+}$  sends  $h$  and plays  $H$  as sender and responds to messages as receiver.*

**PROOF:**

First consider the case when  $H$  is risk dominant.  $T_1$  plays  $\langle h, H, H, L \rangle$  (facing randomizing  $T_0$  receivers and truthful  $T_0$  senders). A  $T_2$  sender believes that the receiver best-responds to the sent message and therefore sends  $h$  and plays  $H$ . A  $T_2$  receiver believes that the sender will send  $h$  and play  $H$ , but if  $T_2$  receives message  $l$ , he believes it comes from a truthful  $T_0$  sender.  $T_{2+}$  therefore plays  $\langle h, H, H, L \rangle$ .

Now consider the case when  $L$  is risk dominant. Then,  $T_1$  plays  $\langle l, L, H, L \rangle$ .  $T_{2+}$  believes that the opponent responds to messages and that all messages are truthful and therefore play  $\langle h, H, H, L \rangle$ .



When risk and payoff-dominance coincide, one-way communication is sufficient to achieve coordination among  $T_{1+}$  players. When there is a conflict between risk and payoff dominance, there is still perfect coordination among  $T_{1+}$  players, but there is more play of the risk dominant equilibrium (since a  $T_1$  sender plays the action corresponding to that equilibrium).

**OBSERVATION 4:** *(Two-way communication)  $T_1$  randomizes messages and responds to received messages, whereas  $T_{2+}$  sends  $h$  and plays  $H$ .*

**PROOF:**

$T_1$  believes that the opponent is truthful and therefore best responds to the received message, while sending random messages (not knowing what action will be taken).  $T_2$  believes that the opponent responds to messages and therefore sends and plays  $H$  irrespective of the message received (since  $T_1$  sends a random message).  $T_3$  therefore sends  $h$  and plays  $H$ . Receiving an unexpected  $L$  message,  $T_3$  also plays  $H$ , believing the opponent to be  $T_1$ . More advanced players reason in the same way and thus also play  $\langle h, H, H \rangle$ .

### *Mixed motive games*

Two common examples of  $2 \times 2$  mixed motive games are Chicken or Hawk-Dove and Battle of the Sexes. In order for the game to have mixed motive, we assume  $u_{HL} > u_{LL}$  and  $u_{LH} > u_{HH}$ . Without loss of generality, we further assume that  $u_{HL} > u_{LH}$  so that each player prefer the equilibrium where he is the one to play  $H$ (igh). If  $u_{LL} = u_{HH} = 0$ , then this game is the Battle of the Sexes, whereas it is a Chicken game if  $u_{LL} > u_{HH}$ . Battle of the Sexes is a non-generic game, but the results in this section hold also for the Battle of the Sexes.

**OBSERVATION 5:** *(No communication) If  $H$  is the risk dominant action, then  $T_k$  plays  $H$  if  $k$  is odd and  $L$  if  $k$  is even. If  $L$  is the risk dominant action, then  $T_k$  plays  $L$  if  $k$  is odd and  $H$  if  $k$  is even.*

**PROOF:**

$T_1$  plays the risk dominant action and  $T_k$  best-responds to the behavior of  $T_{k-1}$ , which generates the alternating behavior.

With no possibility to communicate, there is little players can do to coordinate on either of the asymmetric equilibria and behavior therefore alternates over thinking steps. One-way communication, on the other hand, provides a way to break the symmetry inherent in the game.

**OBSERVATION 6:** *(One-way communication) If  $H$  is the risk dominant action, then  $T_{1+}$  sends  $h$  and plays  $H$  as sender and responds to messages as receiver. If  $L$  is the risk dominant action, then  $T_1$  sends  $l$  and plays  $L$  as sender and responds to messages as receiver.  $T_{2+}$  sends  $h$  and plays  $H$  as sender and responds to messages as receiver.*

**PROOF:**

First let  $H$  be the risk dominant action. A  $T_1$  sender faces a randomizing receiver and therefore plays  $H$  and sends  $h$ . A  $T_1$  receiver, on the other hand, responds to the sent message, believing it comes from a truthful  $T_0$  opponent.  $T_{2+}$  can get the preferred equilibrium as sender and therefore sends  $h$  and plays  $H$ , while responding to messages as receiver. If instead  $L$  is the risk dominant action, a  $T_1$  sender instead sends and plays  $L$ , but otherwise behavior is unchanged.

In general, senders play their preferred equilibrium and receivers yield and play their least preferred equilibrium. However, if the preferred equilibrium does not coincide with the risk dominant action,  $T_1$  senders send and play their least preferred equilibrium.<sup>25</sup>

OBSERVATION 7: (*Two-way communication*)  $T_1$  sends  $h$  and  $l$  with equal probabilities and responds to messages. The behavior of  $T_{2+}$  players cycles in thinking steps of six as follows:  $\langle h, H, H \rangle, \langle l, L, L \rangle, \langle h, L, H \rangle, \langle h, H, H \rangle, \langle l, L, H \rangle, \langle h, L, H \rangle$ .

PROOF:

$T_1$  believes that the opponent is truthful and therefore sends random messages, but responds to the message sent.  $T_2$  believes that the opponent responds to messages and therefore plays  $\langle h, H, H \rangle$ .  $T_3$  expects to receive a truthful  $h$  message, and thus sends  $l$  and plays  $L$ . If receiving an  $l$  message,  $T_3$  believes it comes from a  $T_1$  opponent and therefore plays  $L$  (believing the opponent will play  $H$ ).  $T_4$  expects to play  $H$  and therefore sends  $h$ . If receiving the message  $h$ ,  $T_4$  believes it comes from a  $T_2$  opponent and therefore responds by playing  $L$ .  $T_5$  thinks the opponent responds to messages and therefore plays  $H$  and sends  $h$ . Believe an  $l$  message comes from a  $T_2$  opponent,  $T_5$  subsequently plays  $H$ .  $T_6$  expects to play  $L$  and therefore sends  $l$ , but plays  $H$  upon receiving an  $l$  message (believing it comes from a  $T_2$  opponent).  $T_7$  expects to play  $H$  and sends an  $h$  message, playing  $L$  if receiving an  $h$  message.  $T_8$  sends  $h$  and plays  $H$ , playing  $H$  if he receives an  $l$  message, just like  $T_2$ .  $T_9$  plays  $\langle l, L, L \rangle$  just like  $T_3$ . Since the behavior of eight and nine-level players is just like two- and three-level players, and the rationale for  $T_{4+}$  did not depend on the behavior of  $T_0$  or  $T_1$ , behavior continues to cycle like this.

Note that the behavior of  $T_0, T_1, T_2$ , and  $T_3$  is identical to Crawford (2007). However,  $T_4$  responds to received messages in our model, but always plays  $H$  in Crawford (2007). The difference stems from the fact that we assume that whenever  $T_4$  receives the message  $h$ , the inference is that it comes from a  $T_2$  player that will actually play  $H$ , whereas Crawford (2007) assumes that  $T_4$  believes an  $h$  message is a mistake by a  $T_3$  opponent who will play  $H$  anyway.<sup>26</sup>

Comparing one-way and two-way communication, it is clear that two-way communication will lead to several instances of miscoordination. However, as pointed out by Crawford (2007), the degree of coordination may still be higher than predicted by Farrell (1987) and Rabin (1994).

Finally, note the parallel to coordination games that risk-dominance only plays a role with one-way communication. The underlying reason is the strategic uncertainty resulting from randomizing  $T_0$  receivers.

## Appendix 2: Proofs

### *Proof of Proposition 1*

>From Observation 1 we know that communication has no effect in dominance solvable games. Similarly, for coordination games when  $H$  is risk dominant, Observation 2 and 3 show that communication has no effect. In coordination games when  $L$  is risk dominant, however, Observation 2 and 3 show that one-way communication results in either  $(L, L)$  or  $(H, H)$ ,

<sup>25</sup>The result when  $L$  is risk dominant is sensitive to the assumption that  $T_{1+}$  players have lexicographic preferences for truthfulness. Without that preference, level-1 senders would send random messages. Then, the behavior of more advanced players would alternate and entail many instances of miscoordination.

<sup>26</sup>Also note that although our  $T_3$  behaves as in Crawford (2007), the rationale for their behavior is slightly different.  $T_3$  in our framework believes an  $l$  message comes from a  $T_1$  opponent that sends random messages. Since  $T_3$  sent the message  $l$ , the player believes that the opponent will play  $H$  and they therefore play  $L$ . In Crawford (2007), a  $T_3$  player that receives the counterfactual message  $l$  believes that it was a mistake by the  $T_2$  opponent and therefore plays  $L$  anyway.

whereas no communication results in  $(L, L)$ . As long as there is a positive fraction of  $T_{2+}$  players, one-way communication therefore results in higher expected payoffs.

For mixed motive games, first suppose  $L$  is risk dominant. From Observation 6 we know that one-way communication always induces coordination when  $T_{1+}$  play, so the expected payoff for a player playing the game is  $(u_{HL} + u_{LH})/2$ . However, as noted in Observation 5, no communication results in miscoordination when two odd-level players meet as well as when two even-level players meet. Under the standard type distribution, a player's average payoff is

$$p_2^2 u_{HH} + p_2(1-p_2)u_{HL} + (1-p_2)p_2 u_{LH} + (1-p_2)^2 u_{LL}.$$

One-way communication results in higher expected payoff whenever

$$\left(\frac{1}{2} - p_2(1-p_2)\right)(u_{HL} + u_{LH}) > p_2^2 u_{HH} + (1-p_2)^2 u_{LL}.$$

A sufficient condition is that  $u_{LL} < u_{HL}$  (we already know that  $u_{HH} < u_{LH}$ ), but the necessary condition depends on  $p_2$ . Now let  $H$  be the risk dominant outcome. The expected payoff for communicating players is unchanged, whereas the condition for one-way communication to result in higher expected payoff is

$$\left(\frac{1}{2} - p_2(1-p_2)\right)(u_{HL} + u_{LH}) > (1-p_2)^2 u_{HH} + p_2^2 u_{LL}.$$

*Proof of Corollary 2*

>From the proof of Proposition 1 it follows directly that one-way communication only decreases average payoffs if one of the conditions hold with opposite inequality. To see why the corresponding game is a Chicken, suppose first that  $L$  is risk dominant. The first condition in Proposition 1 for one-sided communication to decrease expected payoffs is

$$(1) \quad \left(\frac{1}{2} - p_2(1-p_2)\right)(u_{HL} + u_{LH}) < p_2^2 u_{HH} + (1-p_2)^2 u_{LL}.$$

We know that  $u_{HL} > u_{LL}$ ,  $u_{LH} > u_{HH}$  and  $u_{HL} > u_{LH}$ . This implies that  $u_{HH} < (u_{LH} + u_{HL})/2$ . Suppose that  $u_{LL} \leq (u_{LH} + u_{HL})/2$ . Then the right hand side of (1) satisfies

$$\begin{aligned} p_2^2 u_{HH} + (1-p_2)^2 u_{LL} &< p_2^2 \frac{1}{2}(u_{LH} + u_{HL}) + \frac{1}{2}(1-p_2)^2(u_{LH} + u_{HL}) \\ &= \left(\frac{1}{2} - p_2(1-p_2)\right)(u_{LH} + u_{HL}). \end{aligned}$$

This implies that (1) cannot hold, and therefore the condition must fail unless  $u_{LL} > \frac{1}{2}(u_{LH} + u_{HL})$ . This implies that  $u_{LL} > u_{HH}$ , which implies that it is a Chicken. An analogous argument can be made when  $H$  is risk dominant.

*Proof of Proposition 3*

>From Observation 1 we know that communication has no effect in dominance solvable games.

>From Observation 2 and 3, we know that the outcomes of coordination games in which  $L$  is the risk dominant action. These are given in Table A1. Pairwise comparison of the cells in

TABLE A1: ACTION PROFILES PLAYED IN COORDINATION GAMES (L RISK DOMINANT)

$G$ (no communication)			$\Gamma_I(G)$ (one-way communication)			
	0	$\geq 1$	0S	0R	1R	$\geq 2R$
0	Uniform	$\frac{1}{2}LL, \frac{1}{2}LH$	0S	Uniform	$\frac{1}{2}HH, \frac{1}{2}LL$	$\frac{1}{2}HH, \frac{1}{2}LL$
$\geq 1$	$\frac{1}{2}LL, \frac{1}{2}LH$	LL	1S	$\frac{1}{2}LL, \frac{1}{2}LH$	LL	LL
			$\geq 2S$	$\frac{1}{2}HH, \frac{1}{2}HL$	HH	HH

TABLE A2: ACTION PROFILES PLAYED IN COORDINATION GAMES (H RISK DOMINANT)

$G$ (no communication)			$\Gamma_I(G)$ (one-way communication)		
	0	$\geq 1$	0S	0R	$\geq 1R$
0	Uniform	$\frac{1}{2}HH, \frac{1}{2}HL$	0S	Uniform	$\frac{1}{2}HH, \frac{1}{2}LL$
$\geq 1$	$\frac{1}{2}HH, \frac{1}{2}HL$	HH	$\geq 1S$	$\frac{1}{2}HH, \frac{1}{2}HL$	HH

Table A1 reveals that one-way communication entails weakly more coordination.

If instead  $H$  is risk dominant, the outcomes are given in Table A2. The degree of coordination is again the same or higher with one-way communication than without communication.

Now consider mixed motive games. Observations 5 and 6 yield the outcomes reported in Table A3 when  $L$  is risk dominant. Pairwise comparisons of cells reveal that the degree of coordination is higher with one-way communication.

Finally, when  $H$  is risk dominant, the outcomes are given in Table A4. Again the degree of coordination is the same or higher for one-way communication for all combinations of types.

### Proof of Proposition 5

As Observation 1 shows, communication plays no role in dominance solvable games, so two-way communication cannot increase expected payoffs. In coordination games in which  $H$  is risk dominant, Observation 3 and 4 imply that  $\Gamma_I(G)$  and  $\Gamma_{II}(G)$  yield identical outcomes unless two  $T_1$  players meet. In  $\Gamma_I(G)$ , players then coordinate on  $(H, H)$ , whereas there is miscoordination in  $\Gamma_{II}(G)$ . Thus  $\Gamma_I(G)$  is weakly better than  $\Gamma_{II}(G)$  in this case. When instead  $L$  is the risk dominant action,  $T_1$  senders always play  $L$ . The average payoff associated with  $\Gamma_I(G)$  is thus

$$p_1(1-p_1)u_{LL} + p_1(1-p_1)u_{HH} + (1-p_1)(1-p_1)u_{HH} + p_1^2u_{LL}.$$

The average payoff associated with  $\Gamma_{II}(G)$  is

$$2p_1(1-p_1)u_{HH} + (1-p_1)(1-p_1)u_{HH} + \frac{1}{4}p_1^2(u_{LL} + u_{HH} + u_{LH} + u_{HL}).$$

TABLE A3: ACTION PROFILES PLAYED IN MIXED MOTIVE GAMES (L RISK DOMINANT)

$G$ (no communication)				$\Gamma_I(G)$ (one-way communication)			
	0	Odd	Even	0S	0R	1R	$\geq 2R$
0	Uniform	$\frac{1}{2}HL, \frac{1}{2}LL$	$\frac{1}{2}LH, \frac{1}{2}HH$	0S	Uniform	$\frac{1}{2}HL, \frac{1}{2}LH$	$\frac{1}{2}HL, \frac{1}{2}LH$
Odd	$\frac{1}{2}LH, \frac{1}{2}LL$	LL	LH	1S	$\frac{1}{2}LH, \frac{1}{2}LL$	LH	LH
Even	$\frac{1}{2}HL, \frac{1}{2}HH$	HL	HH	$\geq 2S$	$\frac{1}{2}HL, \frac{1}{2}HH$	HL	HL

TABLE A4: ACTION PROFILES PLAYED IN MIXED MOTIVE GAMES (H RISK DOMINANT)

$G$ (no communication)				$\Gamma_I(G)$ (one-way communication)		
	0	Odd	Even	0R	Uniform	$\geq 1R$
0	Uniform	$\frac{1}{2}LH, \frac{1}{2}HH$	$\frac{1}{2}HL, \frac{1}{2}LL$	0S	$\frac{1}{2}HL, \frac{1}{2}LH$	
Odd	$\frac{1}{2}HL, \frac{1}{2}HH$	HH	HL	$\geq 1S$	$\frac{1}{2}HL, \frac{1}{2}HH$	HL
Even	$\frac{1}{2}LH, \frac{1}{2}LL$	LH	LL			

TABLE A5: ACTION PROFILES PLAYED IN MIXED MOTIVE GAMES

$\Gamma_{II}(G)$ (two-way communication)			
	1	2	3
1	Uniform	LH	HL
2	HL	HH	HL
3	LH	LH	LL

Two-way communication thus yields higher payoff whenever

$$(4 - 3p_1) u_{HH} + p_1 (u_{LH} + u_{HL}) > (4 - p_1) u_{LL}.$$

Now consider mixed motive games. Observation 6 shows that for  $T_{1+}$  players,  $\Gamma_I(G)$  entails perfect coordination, implying an average payoff of  $(u_{LH} + u_{HL})/2$ . As shown in Observation 7, matters are generally more complicated for  $\Gamma_{II}(G)$  since behavior cycles over six thinking steps. Table A5 provides the resulting outcomes when confining attention to standard type distributions.

We know that  $(u_{LH} + u_{HL})/2 > u_{HH}$ . However, if  $u_{LL} > (u_{LH} + u_{HL})/2$  then two-way communication might be preferable. Two-way communication is preferable to one-way communication whenever

$$\begin{aligned} \left( p_2 p_1 + p_1 p_3 + p_2 p_3 + \frac{1}{4} p_1^2 \right) (u_{HL} + u_{LH}) + \left( p_3^2 + \frac{1}{4} p_1^2 \right) u_{LL} + \left( p_2^2 + \frac{1}{4} p_1^2 \right) u_{HH} \\ > \frac{1}{2} (u_{LH} + u_{HL}). \end{aligned}$$

Letting  $p_2 = (1 - p_1 - p_3)$  we can rewrite this as

$$\frac{u_{LL} - u_{HH}}{u_{LH} + u_{HL} - 2u_{HH}} > 1 + \frac{2(p_1 - 1)(p_1 - 1 + 2p_3)}{p_1^2 + 4p_3^2}.$$

A necessary condition for this inequality to hold is that  $u_{LL} > (u_{LH} + u_{HL})/2$ . This follows from the fact that the minimum of the right hand side is  $1/2$ , whereas the left hand side can only be larger than  $1/2$  if  $u_{LL} > (u_{LH} + u_{HL})/2$ .

#### Proof of Proposition 6

First consider  $\Gamma_I(G)$ . A  $T_1$  sender sends and plays the action that is optimal given that the opponent randomizes uniformly over actions. If there are several optimal actions,  $T_1$  plays each of them with equal probability and sends a truthful message. As a receiver,  $T_1$  best responds to messages. Since the payoff dominant equilibrium gives the highest possible payoff,  $T_2$  sends

and plays the corresponding action as sender, while best responding to messages as receiver. It follows that  $T_{3+}$  behaves as  $T_2$ . Now consider  $\Gamma_{II}(G)$ .  $T_1$  believes the opponent is truthful and therefore best responds to messages, but sends a random message.  $T_{2+}$  believes the opponent best responds and therefore sends and plays the payoff dominant equilibrium irrespective of the received message.

*Proof of Proposition 7*

First suppose that  $\bar{a}_i$  is not strictly dominant for any player. A  $T_1$  player faces truthful  $T_0$  senders, so  $T_1$  sends a random message and best responds to the message profile received. In particular, if  $m_{-i} = \bar{a}_{-i}$ , then  $T_1$  plays  $\bar{a}_i$  since  $\bar{a}$  is the payoff dominant equilibrium.  $T_2$  consequently believes that the opponents best-respond to messages. Since  $G$  has strategic complementarities,  $BR_i(m_{-i})$  is non-decreasing in  $m_{-i}$ , and since there are positive spillovers, it is weakly dominant for a  $T_2$  player to send the message  $m_i = \bar{a}_i$ . However, since  $\bar{a}$  gives strictly higher payoff than all other outcomes of the game, it is strictly dominant to send the message  $m_i = \bar{a}_i$  and play  $\bar{a}_i$  if  $m_{-i} = \bar{a}_{-i}$ .  $T_{3+}$  similarly achieves the highest payoff by sending  $\bar{a}_i$  and playing  $\bar{a}_i$  if  $m_{-i} = \bar{a}_{-i}$ .

If the action  $\bar{a}_i$  is strictly dominant for some player  $i$ , a  $T_1$  player  $i$  truthfully sends the message  $m_i = \bar{a}_i$  and plays  $\bar{a}_i$ . If  $\bar{a}_j$  is strictly dominant for all other players  $j \neq i$ , then a  $T_2$  player  $i$  is indifferent about what message to send, but since a  $T_2$  player  $i$  expects to play  $\bar{a}_i$ , he sends  $\bar{a}_i$ . (If only some other players have strictly dominant strategies, the same argument for the behavior of  $T_2$ 's message as in the previous paragraph hold.)  $T_{2+}$  consequently sends  $\bar{a}_i$  and plays  $\bar{a}_i$  if  $m_{-i} = \bar{a}_{-i}$  also in the presence of strictly dominant actions.

### Appendix 3: Cognitive Hierarchy

As a robustness check, we conduct our analysis with the cognitive hierarchy model of Camerer, Ho and Chong (2004). There, the distribution of types is Poisson distributed, i.e., the proportion of  $T_k$  is given by

$$p_k = \frac{e^{-\tau} \tau^k}{k!}.$$

$T_k$  best responds given the belief that the others players are  $T_0$  up to  $T_{k-1}$ .  $T_k$ 's belief about the proportion of  $T_{l < k}$  is

$$g_k(l) = \frac{p_l}{\sum_{h=0}^{k-1} p_h}.$$

The cognitive hierarchy model is developed for normal form games only. In order to adapt the model to games with pre-play communication we must specify how beliefs are updated after messages have been received. For reasons of familiarity, we assume Bayesian updating. For the games preceded by one round of communication, let  $q_{ki}(m_i)$  denote the probability a  $T_k$  player  $i$  sends the message  $m_i$  (and is allowed to send a message).  $T_k$ 's belief that the sender  $i$  is a  $T_{l < k}$  player conditional upon receiving the message  $m_i$  is

$$g_{ki}(l|m_i) = \frac{g_k(l) q_{li}(m_i)}{\sum_{h=0}^{k-1} g_k(h) q_{hi}(m_i)} = \frac{p_l q_{li}(m_i)}{\sum_{h=0}^{k-1} p_h q_{hi}(m_i)},$$

where the latter equality follows from the definition of  $g_k(l)$ .

We retain the assumption that players randomize uniformly when indifferent, but that they prefer to be honest if it does not affect expected payoffs. This implies that the behavior of  $T_0$  and  $T_1$  is the same in the cognitive hierarchy model and in the level- $k$  model.

A feature of the cognitive hierarchy model is that if  $T_k$  plays a strategy that is a best response to  $T_k$  opponent in a two-player game with one round of pre-play communication, then  $T_{m>k}$  will play the same strategy. Since this result will be used repeatedly we state it separately in Lemma 1.

LEMMA 1: *Let  $G$  be a symmetric two-player normal form game. If  $T_k$  plays a strategy profile that is a best response to a  $T_k$  opponent in  $G$ ,  $\Gamma_I(G)$  or  $\Gamma_{II}(G)$ , then  $T_{m>k}$  play this strategy too.*

PROOF:

Consider the case of one-way communication (the proof for two-way communication and without communication is analogous). Let the strategy played by  $T_k$  be denoted  $s^* = \langle m^*, a^*, f^*(m) \rangle$ . Consider a  $T_k$  player that received the message  $m$ . We know that  $f^*(m)$  is the action that maximizes expected payoff conditional on receiving  $m$  given the belief that the opponent is  $T_{l<k}$  with probability

$$g_k(l|m) = \frac{p_l q_l(m)}{\sum_{h=0}^{k-1} p_h q_h(m)}.$$

Similarly, a  $T_{k+1}$  player that receives the same message  $m$  best responds given the belief that the opponent is a  $T_{l<k+1}$  player with probability

$$g_{(k+1)}(l|m) = \frac{p_l q_l(m)}{\sum_{h=0}^k p_h q_h(m)}.$$

Since  $f^*(m)$  maximizes the expected payoff of  $T_k$  and is a best response to a  $T_k$  sender, by linearity of expected payoffs it must be a best response also to the mixture of types  $T_{k+1}$  believes to be facing (note that this argument does not extend to more than two players).

Now consider the communication stage of the game. The message  $m^*$  followed by the action  $a^*$  is a best response given the belief that the opponent is  $T_{l<k}$  with probability

$$g_k(l) = p_l / \sum_{h=0}^{k-1} p_h.$$

Similarly a  $T_{k+1}$  player believes that the opponent is  $T_{l<k+1}$  with probability

$$g_{k+1}(l) = p_l / \sum_{h=0}^k p_h.$$

Since  $m^*$  and  $a^*$  maximizes the payoff of  $T_k$  and is a best response against another  $T_k$  player, it must be a best response also to the mixture of types  $T_{k+1}$  believes to be facing.

By induction this reasoning holds for all  $T_{m>k}$  players.

In the cognitive hierarchy model, predicted behavior depends both on the payoff configuration and the average of the type distribution,  $\tau$ . A complete characterization of behavior is therefore intractable, and the remainder of this appendix focuses on  $T_2$  and  $T_3$  in the class of symmetric and generic  $2 \times 2$  games. However, a general characterization of the behavior of  $T_3$  in mixed motive games with two-way communication is also intractable, so in this case we focus on  $T_2$  only. For simplicity, we finally disregard cases in which the combination of  $\tau$  and the payoff structure of the game implies that  $T_{2+}$  is indifferent between strategies as well as games in which neither action is risk dominant.

Two general findings emerge from the analysis. First, when  $\tau$  is close to zero,  $T_2$  and  $T_3$  players are practically certain that the opponent is  $T_0$  and consequently play the same action as  $T_1$ . However,  $T_2$  and  $T_3$  may send another message since they take into account the possibility that the opponent is (a responsive)  $T_1$ . Second, for sufficiently large  $\tau$ ,  $T_2$  and  $T_3$  play as in the level- $k$  model in all games except in mixed motive games with two-way communication. In this part of the parameter space, the level- $k$  model is robust to the assumption about lexicographic beliefs. For intermediate levels of  $\tau$ ,  $T_2$  and  $T_3$  best-respond to the mixture of lower-level types they believe they are facing.

An interesting new finding is that the behavior in the Stag Hunt hypothesized by Aumann (1990) emerges endogenously in the model. With the payoffs in Aumann's original example, depicted in Figure 2, a  $T_{3+}$  player sends the message  $h$  and plays  $L$  as sender, and plays  $L$  irrespective of the received message as receiver, whenever  $\tau$  is between 0.547 and 1.646. As a sender,  $T_{3+}$  does so in order to induce  $T_1$  and  $T_2$  to play  $H$ , but believes that there such a high probability of meeting a randomizing  $T_0$  that it is better to play  $L$ .  $T_{3+}$  ignores the received message because of the likelihood of meeting a  $T_2$  opponent, who sends  $h$  messages that are not self-signalling.

In dominance solvable games,  $T_1$  sends and plays the dominant strategy, so by Lemma 1,  $T_{1+}$  does so too (irrespective of whether communication is possible). We now proceed to characterize the behavior in the two remaining classes of games.

### *Coordination games*

As before, we assume that  $H$ (igh) is the payoff dominant equilibrium, i.e.,  $u_{HH} > u_{LL}$ .

**OBSERVATION 8:** *(No communication)  $T_{1+}$  plays the risk dominant action.*

**PROOF:**

$T_1$  plays the risk dominant action. Consequently, by Lemma 1 all higher level types do the same.

Absent communication, behavior is the same in the level- $k$  and cognitive hierarchy models.

**OBSERVATION 9:** *(One-way communication) If  $H$  is the risk dominant action,  $T_{1+}$  sends  $h$  and plays  $H$  as sender and responds to messages as receiver. If  $L$  is the risk dominant action,  $T_1$  sends  $l$  and plays  $L$  as sender and respond to messages as receiver, but the behavior of  $T_{2+}$  depends on the payoff structure of the game:*

**Case 1** ( $u_{LH} > u_{LL}$ ): Let  $\alpha \equiv (u_{LL} - u_{HL}) / (u_{HH} - u_{LH})$ . If  $\tau < (\alpha - 1) / 2$ , then  $T_2$  plays  $\langle h, L, H, L \rangle$  and  $T_3$  plays as follows:

$T_3$  plays  $\langle h, L, H, L \rangle$  if  $\tau < \sqrt{\alpha} - 1$  and  $\tau < (\sqrt{\alpha + 1} + 1) / \alpha$ ,

$T_{3+}$  plays  $\langle h, L, L, L \rangle$  if  $(\sqrt{\alpha + 1} + 1) / \alpha < \tau < \sqrt{\alpha} - 1$ ,

$T_{3+}$  plays  $\langle h, H, H, L \rangle$  if  $\sqrt{\alpha} - 1 < \tau < (\sqrt{\alpha + 1} + 1) / \alpha$ ,

$T_3$  plays  $\langle h, H, L, L \rangle$  if  $\tau > \sqrt{\alpha} - 1$  and  $\tau > (\sqrt{\alpha + 1} + 1) / \alpha$ .

If  $\tau > (\alpha - 1) / 2$ , then  $T_{2+}$  play  $\langle h, H, H, L \rangle$ .

**Case 2** ( $u_{LH} < u_{LL}$ ): Let  $\beta \equiv (u_{LH} - u_{HL}) / (u_{HH} - u_{LL})$ . If  $\tau < (\beta - 1) / 2$ , then  $T_2$  plays  $\langle l, L, H, L \rangle$  and  $T_3$  plays  $\langle l, L, H, L \rangle$  if  $\tau < \sqrt{\beta} - 1$  and  $\langle h, H, H, L \rangle$  if  $\tau > \sqrt{\beta} - 1$ . If  $\tau > (\beta - 1) / 2$ , then  $T_{2+}$  plays  $\langle h, H, H, L \rangle$ .

**PROOF:**

First consider the case when  $H$  is risk dominant. As in the level- $k$  model,  $T_1$  plays  $\langle h, H, H, L \rangle$  (facing randomizing  $T_0$  receivers and truthful  $T_0$  senders). Since this strategy is a best-response to itself,  $T_{2+}$  plays the same strategy.



Now consider the case when  $L$  is risk dominant so that  $T_1$  plays  $\langle l, L, H, L \rangle$ . For  $T_2$  senders, the strategy  $\langle l, H \rangle$  is dominated by  $\langle h, H \rangle$ , so we need not consider that strategy. The expected payoff for the remaining three sender strategies are

$$\begin{aligned}\pi(\langle l, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LL}, \\ \pi(\langle h, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LH}, \\ \pi(\langle h, H \rangle) &= g_2(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_2(1) u_{HH}.\end{aligned}$$

If  $u_{LH} > u_{LL}$ , then it is clear that  $T_2$  senders play either  $\langle h, L \rangle$  or  $\langle h, H \rangle$ . The payoff from playing  $\langle h, L \rangle$  is higher whenever  $\tau$  is sufficiently low,

$$\tau < \frac{(u_{LL} + u_{LH}) - (u_{HL} + u_{HH})}{2(u_{HH} - u_{LH})} = (\alpha - 1)/2.$$

Similarly, if  $u_{LH} < u_{LL}$ , then  $T_2$  senders prefer  $\langle l, L \rangle$  over  $\langle h, H \rangle$  whenever

$$\tau < \frac{(u_{LL} + u_{LH}) - (u_{HL} + u_{HH})}{2(u_{HH} - u_{LL})} = (\beta - 1)/2.$$

$T_2$  receivers face truthful  $T_1$  and  $T_2$  senders, so they respond to messages. It is clear that for sufficiently high  $\tau$ ,  $T_{2+}$  plays  $\langle h, H, H, L \rangle$ .

We now go on to consider the behavior of  $T_3$  when  $\tau$  is below the thresholds above. First suppose that  $u_{LH} > u_{LL}$  and  $\tau < (\alpha - 1)/2$ . Then  $T_3$  senders prefer  $\langle h, L \rangle$  over  $\langle h, H \rangle$  whenever

$$\begin{aligned}g_3(0) \frac{1}{2} (u_{LL} + u_{LH}) + (g_3(1) + g_3(2)) u_{LH} \\ > g_3(0) \frac{1}{2} (u_{HL} + u_{HH}) + (g_3(1) + g_3(2)) u_{HH},\end{aligned}$$

which simplifies to  $(1 + \tau/2)\tau < (\alpha - 1)/2$ . Since the left hand side is larger than  $\tau$ , this condition may or may not hold. Both sides of the inequality are positive, so the condition is equivalent to  $\tau < \sqrt{\alpha} - 1$ . Suppose now that  $u_{LH} < u_{LL}$ . Then  $T_3$  senders prefer  $\langle l, L \rangle$  over  $\langle h, H \rangle$  whenever

$$(1 + \tau/2)\tau < \frac{(u_{LL} + u_{LH}) - (u_{HL} + u_{HH})}{2(u_{HH} - u_{LL})} = (\beta - 1)/2$$

Both sides of the inequality are positive, so this condition is equivalent to  $\tau < \sqrt{\beta} - 1$ .

Finally,  $T_3$ 's behavior as receivers depend on the  $T_2$  senders. It is only when  $T_2$  senders send  $h$ , but play  $L$  that  $T_3$  may not respond to messages. If  $T_3$  receives a  $l$  message, it comes from a  $T_0$  player and  $T_3$  best responds by playing  $L$ . The payoff from each action upon receiving  $h$  is

$$\begin{aligned}\pi(H|h) &= g_3(0|h) u_{HH} + g_3(1|h) u_{HH} + g_3(2|h) u_{HL}, \\ \pi(L|h) &= g_3(0|h) u_{LH} + g_3(1|h) u_{LH} + g_3(2|h) u_{LL}.\end{aligned}$$

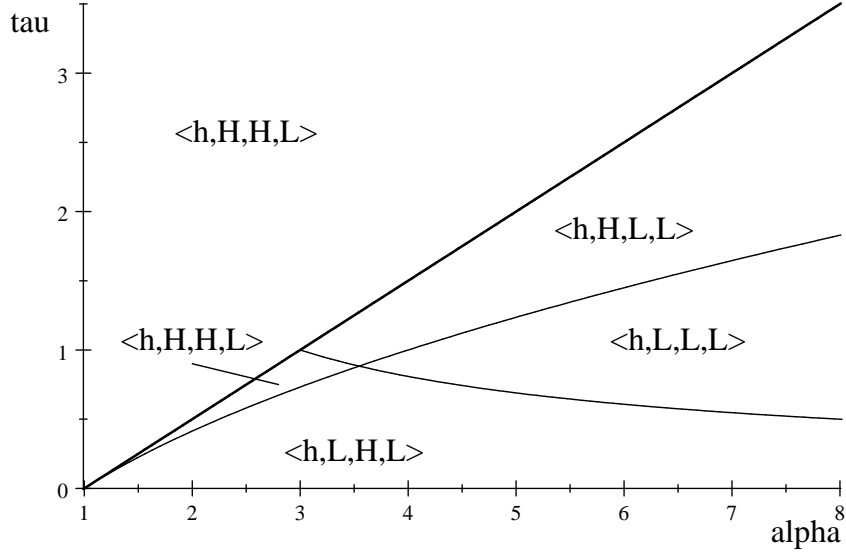


FIGURE A3: LEVEL-3 IN COORDINATION GAMES

Playing  $L$  is preferable whenever

$$\frac{\tau^2}{1+2\tau} > \frac{u_{HH} - u_{LH}}{u_{LL} - u_{HL}} = 1/\alpha.$$

To illustrate the first case when  $L$  is risk dominant, Figure A3 displays the behavior of  $T_3$  as a function of  $\tau$  and the payoffs of the game. First note that  $\alpha$  has to be larger than 1 because  $L$  is risk dominant. Above the thick line in Figure A3,  $T_{2+}$  plays  $\langle h, H, H, L \rangle$  and below it  $T_2$  plays  $\langle h, L, H, L \rangle$ . Figure A3 shows the four different cases for the behavior of  $T_3$  in the latter case. For example, for the Stag Hunt depicted in Figure 2,  $\alpha = 7$ , implying that  $T_{3+}$  plays  $\langle h, L, L, L \rangle$  whenever  $0.547 < \tau < 1.646$ .

**OBSERVATION 10:** (*Two-way communication*)  $T_1$  randomizes messages and responds to received messages. Let  $\lambda \equiv (u_{HH} - u_{LH}) / (2u_{LL} - u_{LH} - u_{HH})$ . If  $L$  is the risk dominant action and  $0 < \lambda < (\beta - 1) / 2$ , then  $T_{2+}$  plays  $\langle h, H, L \rangle$  if  $\tau < \lambda$ ,  $\langle l, L, L \rangle$  if  $\lambda < \tau < (\beta - 1) / 2$ , and  $\langle h, H, H \rangle$  if  $\tau > (\beta - 1) / 2$ . If both inequalities are violated,  $T_{2+}$  plays  $\langle h, H, L \rangle$  if  $\tau < \alpha$  and  $\langle h, H, H \rangle$  if  $\tau > \alpha$ . If  $H$  is the risk dominant action, the behavior of  $T_2$  depends on the payoff structure of the game:

**Case 1** ( $u_{LH} + u_{HH} > u_{LL} + u_{HL}$ ):  $T_{2+}$  plays  $\langle h, H, L \rangle$  if  $\tau < \alpha$  and  $\langle h, H, H \rangle$  if  $\tau > \alpha$ .

**Case 2** ( $u_{LH} + u_{HH} < u_{LL} + u_{HL}$ ): Let  $\gamma \equiv (u_{LL} - u_{HL}) / (2u_{HH} - u_{LL} - u_{HL})$  and  $\delta \equiv (u_{HH} - u_{LL}) / (u_{LL} - u_{HL})$ . If  $\tau < \gamma$ , then  $T_2$  plays  $\langle l, H, L \rangle$ ;  $T_3$  plays  $\langle l, H, L \rangle$  if in addition  $\tau < (\sqrt{4\delta\gamma^2 + 1} - 1) / 2\gamma\delta$ , but plays  $\langle h, H, H \rangle$  if  $\tau > (\sqrt{4\delta\gamma^2 + 1} - 1) / 2\gamma\delta$ . If  $\tau > \gamma$ , then  $T_{2+}$  plays  $\langle h, H, H \rangle$ .

**PROOF:**

$T_1$  believes that the opponent is truthful and therefore best responds to the received message, while sending random messages (not knowing what action will be taken).

$T_2$  faces truthful  $T_0$  and responding  $T_1$ . Since zero-step and one-step thinkers send both messages with equal probabilities,  $g_2(l|m) = g_2(l)$ . The strategy  $\langle l, H, H \rangle$  is clearly dominated by  $\langle h, H, H \rangle$  and  $\langle h, L, L \rangle$  is dominated by  $\langle h, H, L \rangle$ . The remaining strategies gives the following expected payoff:

$$\begin{aligned}\pi(\langle h, H, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{HH}) + g_2(1) \frac{1}{2} (u_{LH} + u_{HH}), \\ \pi(\langle h, H, H \rangle) &= g_2(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_2(1) u_{HH}, \\ \pi(\langle l, L, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LL}, \\ \pi(\langle l, H, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{HH}) + g_2(1) \frac{1}{2} (u_{LL} + u_{HL}).\end{aligned}$$

First suppose that  $L$  is risk dominant. This implies that  $u_{LL} + u_{LH} > u_{HL} + u_{HH}$  and consequently that  $u_{LH} > u_{HL}$  so that  $\langle h, H, L \rangle$  dominates  $\langle l, H, L \rangle$ .  $T_2$  prefers  $\langle h, H, H \rangle$  over  $\langle h, H, L \rangle$  whenever

$$\tau > (u_{LL} - u_{HL}) / (u_{HH} - u_{LH}) = \alpha,$$

and  $\langle h, H, H \rangle$  over  $\langle l, L, L \rangle$  whenever  $\tau > (\beta - 1) / 2$ . Finally,  $T_2$  prefers  $\langle l, L, L \rangle$  over  $\langle h, H, L \rangle$  whenever

$$\tau > \frac{u_{HH} - u_{LH}}{2u_{LL} - u_{LH} - u_{HH}},$$

given that the right hand side is positive.

Second, suppose that  $H$  is risk dominant and  $u_{LH} + u_{HH} > u_{LL} + u_{HL}$  so that  $\langle h, H, L \rangle$  dominates  $\langle l, H, L \rangle$  and  $\langle h, H, H \rangle$  dominates  $\langle l, L, L \rangle$ .  $T_2$  plays  $\langle h, H, H \rangle$  rather than  $\langle h, H, L \rangle$  if  $\tau > \alpha$ .

Finally, suppose that  $H$  is risk dominant and  $u_{LH} + u_{HH} < u_{LL} + u_{HL}$ . Now  $\langle l, H, L \rangle$  dominates  $\langle h, H, L \rangle$  and  $\langle h, H, H \rangle$  dominates  $\langle l, L, L \rangle$ . Therefore,  $T_2$  plays  $\langle h, H, H \rangle$  if  $\tau > (u_{LL} - u_{HL}) / (2u_{HH} - u_{LL} - u_{HL})$ .

Since  $\langle l, L, L \rangle$ ,  $\langle h, H, L \rangle$  and  $\langle h, H, H \rangle$  are best responses if the opponent plays the same strategies, by Lemma 1,  $T_{3+}$  play like  $T_2$  in these cases. Finally, consider  $T_3$  when  $T_2$  plays  $\langle l, H, L \rangle$ . In this case, whenever  $T_3$  receives an  $h$  message, he believes that it comes from a  $T_0$  or  $T_1$  player. Suppose first  $T_3$  receives the message  $h$ . If  $T_3$  sent  $h$ , then it is optimal to play  $H$ . If  $T_3$  sent the message  $l$ , the payoffs from playing  $L$  and  $H$  are

$$\begin{aligned}\pi(\langle l, L, \cdot \rangle | h) &= g_3(0) u_{LH} + g_3(1) u_{LL}, \\ \pi(\langle l, H, \cdot \rangle | h) &= g_3(0) u_{HH} + g_3(1) u_{HL}.\end{aligned}$$

Playing  $H$  is preferred whenever  $\tau < (u_{HH} - u_{LH}) / (u_{LL} - u_{HL}) = 1/\alpha$ . Since  $H$  is risk dominant,  $\alpha < 1$  and since  $\gamma < 1$ , this condition always hold. Now consider the case when  $T_3$  receives the message  $l$ . If  $T_3$  sent  $l$ , then it is optimal to play  $L$  (since  $T_0$  is truthful and  $T_1$  and  $T_2$  best-responds). Suppose that  $T_3$  sent  $h$ . Then expected payoffs are:

$$\begin{aligned}\pi(\langle h, \cdot, L \rangle | l) &= g_3(0|l) u_{LL} + g_3(1|l) u_{LH} + g_3(2|l) u_{LH} \\ \pi(\langle h, \cdot, H \rangle | l) &= g_3(0|l) u_{HL} + g_3(1|l) u_{HH} + g_3(2|l) u_{HH}\end{aligned}$$

Playing  $L$  is preferred whenever  $\tau(1 + \tau) < \alpha$ .

Which message will  $T_3$  send? Suppose first that  $\tau(1 + \tau) < \alpha$  so that  $T_3$  plays  $\langle h, H, L \rangle$  or  $\langle l, H, L \rangle$ . These strategies give the following ex ante payoffs

$$\begin{aligned}\pi(\langle h, H, L \rangle) &= g_3(0) \frac{1}{2} (u_{LL} + u_{HH}) + g_3(1) \frac{1}{2} (u_{LH} + u_{HH}) + g_3(2) u_{LH}, \\ \pi(\langle l, H, L \rangle) &= g_3(0) \frac{1}{2} (u_{LL} + u_{HH}) + g_3(1) \frac{1}{2} (u_{LL} + u_{HL}) + g_3(2) u_{LL}.\end{aligned}$$

It follows from the condition  $u_{LH} + u_{HH} < u_{LL} + u_{HL}$  that  $\langle l, H, L \rangle$  dominates  $\langle h, H, L \rangle$ . Now consider the case when  $\tau(1 + \tau) > \alpha$  so that  $T_3$  play either  $\langle h, H, H \rangle$  or  $\langle l, H, L \rangle$ . The payoff from each strategy is

$$\begin{aligned}\pi(\langle h, H, H \rangle) &= g_3(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_3(1) u_{HH} + g_3(2) u_{HH}, \\ \pi(\langle l, H, L \rangle) &= g_3(0) \frac{1}{2} (u_{LL} + u_{HH}) + g_3(1) \frac{1}{2} (u_{LL} + u_{HL}) + g_3(2) u_{LL}.\end{aligned}$$

Sending  $h$  is preferred whenever  $\tau + \tau^2\delta\gamma > \gamma$ , i.e. when  $\tau > \left(\sqrt{4\gamma^2\delta + 1} - 1\right)/2\gamma\delta$  (since  $\gamma > 0$  and  $\delta > 1$ ).

Note that two-way communication always entails play of the payoff dominant equilibrium in the Stag Hunt game depicted in Figure 2. For that particular game,  $\alpha = 7$  so that  $T_{2+}$  plays  $\langle h, H, L \rangle$  if  $\tau < 7$  and  $\langle h, H, H \rangle$  otherwise.

#### *Mixed motive games*

As before, we assume without loss of generality that  $u_{HL} > u_{LH}$  so that each player prefers the equilibrium where he is the one to play  $H$  (igh).

**OBSERVATION 11:** (*No communication*) Let  $\theta = (u_{LH} - u_{HH}) / (u_{HL} - u_{LL})$ . If  $H$  is the risk dominant action,  $T_1$  plays  $H$ . If  $\tau < (1/\theta - 1)/2$ ,  $T_2$  plays  $H$ ;  $T_3$  plays  $H$  if in addition  $\tau + \tau^2/2 < (1/\theta - 1)/2$ , but plays  $L$  if  $\tau + \tau^2/2 > (1/\theta - 1)/2$ . If  $\tau > (1/\theta - 1)/2$ ,  $T_2$  plays  $L$ ;  $T_3$  plays  $H$  if in addition  $(2 - \tau/\theta)\tau < 1/\theta - 1$ , but plays  $L$  if  $(2 - \tau/\theta)\tau > 1/\theta - 1$ . If  $L$  is the risk dominant action,  $T_1$  plays  $L$ . If  $\tau < (\theta - 1)/2$ ,  $T_2$  plays  $L$ ;  $T_3$  plays  $L$  if in addition  $\tau + \tau^2/2 < (\theta - 1)/2$ , but plays  $H$  if  $\tau + \tau^2/2 > (\theta - 1)/2$ . If  $\tau > (\theta - 1)/2$ ,  $T_2$  plays  $H$ ;  $T_3$  plays  $L$  if in addition  $(2 - \tau\theta)\tau < \theta - 1$ , but plays  $H$  if  $(2 - \tau\theta)\tau > \theta - 1$ .

**PROOF:**

First suppose  $H$  is risk dominant (which implies that  $\theta < 1$ ).  $T_2$  plays  $H$  rather than  $L$  if

$$g_2(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_2(1) u_{HH} > g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LH},$$

which is equivalent to  $1 + 2\tau < 1/\theta$ . Suppose this holds so that  $T_2$  plays  $H$ . Then  $T_3$  prefers  $H$  over  $L$  whenever

$$\begin{aligned}g_3(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_3(1) u_{HH} + g_3(2) u_{HH} \\ > g_3(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_3(1) u_{LH} + g_3(2) u_{LH},\end{aligned}$$

which simplifies to  $1 + 2\tau + \tau^2 < 1/\theta$ . Suppose instead  $T_2$  plays  $L$ . Then  $T_3$  prefers  $H$  over  $L$  whenever

$$\begin{aligned} g_3(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_3(1) u_{HH} + g_3(2) u_{HL} \\ > g_3(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_3(1) u_{LH} + g_3(2) u_{LL}, \end{aligned}$$

which is equivalent to  $(2 - \tau/\theta)\tau < 1/\theta - 1$ .

Now suppose  $L$  is risk dominant. Then  $T_2$  plays  $H$  rather than  $L$  if

$$g_2(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_2(1) u_{HL} > g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LL},$$

which is equivalent to  $1 + 2\tau > \theta$ . Suppose that this holds so that  $T_2$  plays  $H$ . Then  $T_3$  prefers  $H$  over  $L$  whenever

$$\begin{aligned} g_3(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_3(1) u_{HL} + g_3(2) u_{HH} \\ > g_3(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_3(1) u_{LL} + g_3(2) u_{LH}, \end{aligned}$$

which simplifies to  $(2 - \tau\theta)\tau > \theta - 1$ . If  $T_2$  instead plays  $L$ , then  $T_3$  prefers  $H$  over  $L$  whenever  $\tau + \tau^2/2 > (\theta - 1)/2$ .

Note that some of the conditions above are quadratic, implying that they may be satisfied both for low and high values of  $\tau$ .

**OBSERVATION 12:** (*One-way communication*) *If  $H$  is the risk dominant action, then  $T_{1+}$  sends  $h$  and plays  $H$  as sender and responds to messages as receiver. If  $L$  is the risk dominant action, then  $T_1$  sends  $l$  and plays  $L$  as sender and responds to messages as receiver. The behavior of  $T_{2+}$  depends on the payoff structure of the game:*

**Case 1** ( $u_{LH} > u_{LL}$ ): *Let  $\eta \equiv (u_{LL} + u_{LH} - u_{HL} - u_{HH})/2(u_{HL} - u_{LH})$ . If  $\tau < \eta$ , then  $T_2$  plays  $\langle l, L, L, H \rangle$  and  $T_3$  plays  $\langle l, L, L, H \rangle$  if  $(1 + \tau/2)\tau < \eta$  whereas  $T_{3+}$  plays  $\langle h, H, L, H \rangle$  if  $(1 + \tau/2)\tau > \eta$ . If  $\tau > \eta$ , then  $T_{2+}$  plays  $\langle h, H, L, H \rangle$ .*

**Case 2** ( $u_{LH} < u_{LL}$ ): *If  $\tau < (\theta - 1)/2$ , then  $T_2$  plays  $\langle h, L, L, H \rangle$  and  $T_3$  plays as follows:*

*$T_3$  plays  $\langle h, L, L, H \rangle$  if  $(1 + \tau/2)\tau < (\theta - 1)/2$  and  $\tau < \sqrt{\theta}$ ,*

*$T_{3+}$  plays  $\langle h, L, L, L \rangle$  if  $(1 + \tau/2)\tau < (\theta - 1)/2$  and  $\tau > \sqrt{\theta}$ ,*

*$T_{3+}$  plays  $\langle h, H, L, H \rangle$  if  $(1 + \tau/2)\tau > (\theta - 1)/2$  and  $\tau < \sqrt{\theta}$ ,*

*$T_{3+}$  plays  $\langle h, H, L, L \rangle$  if  $(1 + \tau/2)\tau > (\theta - 1)/2$  and  $\tau > \sqrt{\theta}$ .*

*If  $\tau > (\theta - 1)/2$ , then  $T_{2+}$  plays  $\langle h, H, L, H \rangle$ .*

**PROOF:**

First let  $H$  be the risk dominant action. A  $T_1$  sender faces a randomizing receiver and therefore plays  $H$  and sends  $h$ . A  $T_1$  receiver, on the other hand, responds to the sent message, believing it comes from a truthful  $T_0$  opponent. By Lemma 1,  $T_{2+}$  plays the same strategy as  $T_1$ .

If instead  $L$  is the risk dominant action, a  $T_1$  sender instead sends and plays  $L$ , but responds to messages as receiver. A  $T_2$  sender faces a tradeoff between playing  $L$  (the best response against  $T_0$ ) and sending  $h$  and playing  $H$  (the best response against  $T_1$ ). The expected payoffs

from the three relevant sender strategies are:

$$\begin{aligned}\pi(\langle l, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LH}, \\ \pi(\langle h, H \rangle) &= g_2(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_2(1) u_{HL}, \\ \pi(\langle h, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LL}.\end{aligned}$$

Suppose  $u_{LH} > u_{LL}$  so that  $\langle l, L \rangle$  dominates  $\langle h, L \rangle$ . Then a  $T_2$  sender plays  $\langle l, L \rangle$  if

$$\tau < \frac{(u_{LL} + u_{LH}) - (u_{HL} + u_{HH})}{2(u_{HL} - u_{LH})} = \eta,$$

but plays  $\langle h, H \rangle$  otherwise.  $T_2$  receivers face truthful  $T_0$  and  $T_1$  senders, so they respond to messages. If  $\tau$  is above the threshold above,  $T_{2+}$  play  $\langle h, H, L, H \rangle$ . However, if  $\tau$  is below the threshold,  $T_3$  senders trade off truthfully playing  $L$  or  $H$ . They play  $L$  if  $(1 + \tau/2)\tau < \eta$  and otherwise play  $H$ .

Suppose now that  $u_{LL} > u_{LH}$  so that  $T_2$  senders prefer sending  $h$  when they intend to play  $L$ . They prefer doing so over  $\langle h, H \rangle$  whenever

$$\tau < \frac{(u_{LL} + u_{LH}) - (u_{HL} + u_{HH})}{2(u_{HL} - u_{LL})} = (\theta - 1)/2.$$

$T_2$  receivers face truthful senders, so they respond to messages. If  $\tau > (\theta - 1)/2$ ,  $T_{2+}$  plays  $\langle h, H, L, H \rangle$ . A  $T_3$  sender plays  $\langle h, L \rangle$  rather than  $\langle h, H \rangle$  if  $(1 + \tau/2)\tau < (\theta - 1)/2$ .

A  $T_3$  receiver believes that an  $l$  message is truthful, so they play  $H$  in that case. An  $h$  message comes either from a  $T_0$  or  $T_3$ . When receiving a  $h$  message, the payoff from each action is:

$$\begin{aligned}\pi(H|h) &= g_3(0|h) u_{HH} + g_3(2|h) u_{HL}, \\ \pi(L|h) &= g_3(0|h) u_{LH} + g_3(2|h) u_{LL}.\end{aligned}$$

So,  $T_3$  play  $\langle L, H \rangle$  if

$$\tau < \sqrt{\frac{u_{LH} - u_{HH}}{u_{HL} - u_{LL}}} = \sqrt{\theta},$$

and play  $\langle H, H \rangle$  otherwise.

Two-way communication in mixed motive games is particularly cumbersome to characterize generally. The following observation therefore focuses on the behavior of  $T_2$ . (For a particular payoff configuration, however, it is straightforward to derive the behavior of  $T_{3+}$  players.)

**OBSERVATION 13:** (*Two-way communication*)  $T_1$  sends  $h$  and  $l$  with equal probabilities and responds to messages. Let  $\nu \equiv (u_{HL} - u_{LL}) / (2u_{LH} - u_{LL} - u_{HL})$ . If  $L$  is risk dominant and  $\eta > \nu > 0$ , then  $T_2$  plays  $\langle h, L, H \rangle$  if  $\tau < \nu$ ,  $\langle l, L, L \rangle$  if  $\nu < \tau < \eta$ , and  $\langle h, H, H \rangle$  if  $\tau > \eta$ . If both inequalities are violated, then  $T_2$  plays  $\langle h, L, H \rangle$  if  $\tau < \theta$  and  $\langle h, H, H \rangle$  if  $\tau > \theta$ . If  $H$  is risk dominant and  $u_{LH} + u_{HH} > u_{LL} + u_{HL}$ , then  $T_2$  plays  $\langle l, L, H \rangle$  if  $\tau < (u_{LH} - u_{HH}) / (2u_{HL} - u_{LH} - u_{HH})$  and  $\langle h, H, H \rangle$  if  $\tau > (u_{LH} - u_{HH}) / (2u_{HL} - u_{LH} - u_{HH})$ . If instead  $u_{LH} + u_{HH} < u_{LL} + u_{HL}$ ,  $T_2$  plays  $\langle h, L, H \rangle$  if  $\tau < \theta$  and  $\langle h, H, H \rangle$  if  $\tau > \theta$ .

**PROOF:**

$T_1$  believes that the opponent is truthful and therefore sends random messages, but responds to the received message. The strategy  $\langle l, H, H \rangle$  is dominated by  $\langle h, H, H \rangle$  and the expected payoff for  $T_2$ 's other strategies are:

$$\begin{aligned}\pi(\langle h, L, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LL}, \\ \pi(\langle h, L, H \rangle) &= g_2(0) \frac{1}{2} (u_{HL} + u_{LH}) + g_2(1) \frac{1}{2} (u_{LL} + u_{HL}), \\ \pi(\langle h, H, H \rangle) &= g_2(0) \frac{1}{2} (u_{HL} + u_{HH}) + g_2(1) u_{HL}, \\ \pi(\langle l, L, H \rangle) &= g_2(0) \frac{1}{2} (u_{HL} + u_{LH}) + g_2(1) \frac{1}{2} (u_{LH} + u_{HH}), \\ \pi(\langle l, L, L \rangle) &= g_2(0) \frac{1}{2} (u_{LL} + u_{LH}) + g_2(1) u_{LH}.\end{aligned}$$

Suppose  $H$  is risk dominant. Then  $\langle h, H, H \rangle$  dominates  $\langle l, L, L \rangle$  and  $\langle h, L, L \rangle$ . First suppose that  $u_{LH} + u_{HH} > u_{LL} + u_{HL}$  so that  $\langle l, L, H \rangle$  dominates  $\langle h, L, H \rangle$ .  $T_2$  plays  $\langle h, H, H \rangle$  rather than  $\langle l, L, H \rangle$  if

$$\tau > \frac{u_{LH} - u_{HH}}{2u_{HL} - u_{LH} - u_{HH}}.$$

If instead  $u_{LL} + u_{HL} > u_{LH} + u_{HH}$ , then  $T_2$  plays  $\langle h, H, H \rangle$  rather than  $\langle h, L, H \rangle$  if

$$\tau > \frac{u_{LH} - u_{HH}}{u_{HL} - u_{LL}} = \theta.$$

Now consider the case when  $L$  is risk dominant. This implies that  $u_{LL} > u_{HH}$ , so  $\langle h, L, H \rangle$  dominates  $\langle l, L, H \rangle$  and  $\langle h, L, H \rangle$  dominates  $\langle h, L, L \rangle$ . There are three remaining strategies to consider.  $T_2$  plays  $\langle h, H, H \rangle$  rather than  $\langle h, L, H \rangle$  if  $\tau > \theta$ .  $T_2$  may prefer to play  $\langle l, L, L \rangle$ .  $\langle l, L, L \rangle$  preferred over  $\langle h, L, H \rangle$  whenever

$$\tau > \frac{u_{HL} - u_{LL}}{2u_{LH} - u_{LL} - u_{HL}} = \nu,$$

given that the right hand side is positive (otherwise the condition cannot hold).  $\langle l, L, L \rangle$  is preferred over  $\langle h, H, H \rangle$  whenever

$$\tau < \frac{(u_{LL} + u_{LH}) - (u_{HL} + u_{HH})}{2(u_{HL} - u_{LH})} = \eta.$$

Hence, in order for  $\langle l, L, L \rangle$  to be optimal,  $\tau$  must be between  $\nu$  and  $\eta$  and the payoffs must satisfy  $\eta > \nu > 0$ .