# Almost-Truthful Interim-Biased Mediation Enables Information Exchange between Agents with Misaligned Interests

## Dmitry Sedov<sup>\*</sup>

\*Northwestern University, Evanston, USA

March 3, 2022

#### Abstract

Information necessary for decision-making is often distributed among agents with misaligned interests. In such settings receiving information is desirable for the agents, while revealing it may be privately harmful. This paper constructs a class of almost-truthful interim-biased mediation protocols that incentivize information exchange in a succinct model capturing such conflicts. The protocols in this class receive signal reports from the agents and send private messages back, almost always transmitting the received signal reports without any distortions. Each rare distorted message is deliberately designed to prevent deviations from truth-telling. Specifically, each mediator's distortion aims at implicitly encouraging a truthful agent to take the action that is interim-optimal given her private signal report only. A deviating agent, however, receives an encouragement based on an untruthful report and thus shifts her action away from the truly interim-optimal one when facing such a distortion. As a result, the deviating agent is put to a disadvantage when the mediator distorts the signals, which is enough to ensure truthful communication when the misalignment of interests between the agents is sufficiently small.

#### JEL classification: C72; D82; D83.

*Keywords:* Communication; Misaligned interests; Information; Mediation; Cheap talk; Mechanism design.

*E-mail address:* dsedov@u.northwestern.edu; *URL*: dsedov.io

Address for correspondence: Department of Economics, Northwestern University, 2211 Campus Drive, Evanston, IL 60208-2600, USA.

Acknowledgements The author thanks Timur Abbiasov, Yingni Guo, Gaston Illanes, Sergei Izmalkov, Riccardo Marchingiglio, Alessandro Pavan, Robert Porter, Dmitry Sorokin, Xavier Vives, Asher Wolinsky, Andrey Zhukov and seminar participants at Northwestern University for helpful comments and discussions. Special thanks to Anna Algina, Francisco Poggi and Quitzé Valenzuela-Stookey. The author is also grateful to the Editor, Moritz Meyer-ter-Vehn, and to the anonymous reviewers for useful suggestions.

## **1** Introduction

Information relevant to decision-making is often distributed among agents with misaligned interests. For illustration, consider the following examples. In knowledge-intensive organizations complementary information is fragmented and held by employees who may be directly or indirectly affected by their colleagues' choices. As another example, intelligence agencies gather incomplete material relevant for investigations, but are also involved in competition for influence and authority. In these cases, *receiving* information is privately desirable for the agents, while *revealing* it may be privately harmful.

When such a conflict is present, how can mutually advantageous information exchange be organized? The present paper answers this question in a succinct model of distributed information and misaligned interests.

The model features a one-stage game without monetary transfers, in which agents obtain private signals from a finite set, then take actions and receive payoffs that depend on the combination of signals and actions. Agents lack commitment power and may have misaligned interests regarding each other's actions, while the payoffs are assumed to be separable in actions. Depending on the exact preferences that each agents has regarding her counterpart's actions, direct communication can be hard in this model. In fact, welfare-improving direct communication cannot be sustained in equilibrium at least in case when the payoffs are such that an action change benefiting one agent necessarily harms the other agent. However, this paper shows that even for such preferences (but also for other preferences regarding the counterpart's actions), a special class of almost-truthful interim-biased mediation protocols can facilitate communication provided that each agent's payoff depends substantially less on her counterpart's action relative to her own.

Agent interactions mediated by the protocols in the almost-truthful interim-biased class are structured as follows. First, agents observe private signals and make cheap-talk signal reports to the mediator. Next, the mediator sends a private message back to each agent. Then, each agent takes an action based on her private signal and the mediator's message. Finally, the combination of actions and actual signals determines the payoffs.

What does the mediator's private message to an agent contain? Almost always, it contains the actual signal report submitted by the other agent (thus the mediation is labelled as almost-truthful) and, with a very small probability, it contains a distorted signal report. Distorted messages aim at implicitly encouraging a truthful agent to take the action that is interim-optimal given her private signal report only (thus the mediation is labelled as interim-biased). Since the two types of messages (with and without a distortion) take value in the same set, the agents cannot distinguish them with certainty.

Why do such messages ensure the existence of a truth-telling equilibrium? To begin with, since the distortion probability is small, each agent is almost certain that the mediator's message coincides with the other agent's actual signal report and selects an action based on this belief. The mediator then exploits this belief when using distorted messages by shifting agents' actions and affecting payoffs in a way that prevents deviations from truth-telling. Specifically, a distorted message harms an agent more when she reports untruthfully rather than truthfully. This occurs because a truthful agent is implicitly encouraged by the distorted message to take the interim-optimal action conditional on her private signal, while a deviating agent combines her private signal with the encouragement based on an untruthful report and shifts her action away from the truly interimoptimal one. The possibility of such an undesirable shift makes revealing the private signal truthfully to the mediator optimal from the perspective of selecting one's own action. While a deviating agent

may still benefit from the change in the counterpart's action caused by the deviation, the incentives for truthful communication dominate provided that the misalignment of interests between the agents is small enough.

Beyond showing that the almost-truthful interim-biased mediation allows for a truthful information exchange under certain condition, the paper includes several auxiliary results. These include developing an optimal mediation protocol for the illustrative example used in the paper and showing that the main results of the paper hold when the assumption of agents' payoffs being separable in actions is slightly relaxed.

Overall, this paper contributes to the existing literature on information exchange by developing a novel tool that facilitates communication between agents with a sufficiently small misalignment of interests. Reiterating on the discussion above, two notable features of almost-truthful interimbiased mediation are worth highlighting. First, an almost-truthful interim-biased mediator is able to punish deception by sometimes distorting the transmitted signals in a way that makes a truthful agent choose the action that's optimal given her private information only, while shifting a deceitful agent's action away from such an interim-optimal action. Second, the fact that an almost-truthful interim-biased mediator almost always transmits information without any distortions, ensures that it can actually shift agents' actions in the desired direction when using the distortions.

The rest of the paper is organized as follows. Section 2 briefly reviews the relevant literature. Section 3 provides an illustrative example revealing the intuition behind the main result. Section 4 presents the baseline model, shows that welfare-improving direct communication is not possible at least in the special case of agents benefiting from their counterpart mistakes, characterizes the class of almost-truthful interim-biased mediation protocols and demonstrates that such protocols enable information transmission. Section 5 concludes by summarizing the results and discussing the modest considerations for communication in organizations and for information exchange between intelligence agencies.

## 2 Literature review

This section discusses two major branches of research on transmission of unverifiable information connected with the present paper. Contributing to that research, this paper adds a novel tool into the communication facilitation toolbox, the class of almost-truthful interim-biased mediation protocols. To put the model and the mediation protocol developed in the present paper in context, the review below lists several key information transmission papers and highlights the major differences in assumptions and incentives for truthful communication between the almost-truthful interim-biased mediation and the constructions in those papers.

First, there is a vast literature on communication with informed agents not being able to influence decisions directly. Unmediated cheap talk has been shown to allow information transmission in the seminal paper by Crawford and Sobel (1982) (henceforth, CS). Moreover, important spin-offs have been explored, including, but not limited to, multiple senders by Austen-Smith (1993) and Krishna and Morgan (2001), multiple receivers by Farrell and Gibbons (1989) and Goltsman and Pavlov (2011), multiple rounds of communication by Krishna and Morgan (2004) and Goltsman et al. (2009), unbounded multidimensional state space by Battaglini (2002), bounded multidimensional state space by Ambrus and Takahashi (2008) and communication error by Blume et al. (2007). Communication through a mediator has been explored as well. Among others, Goltsman et al. (2009) characterize the optimal mediated communication in the canonical CS setting. Ivanov (2010) and Ambrus et al. (2013) explore communication via strategic mediators. In the settings

listed above truthful information transmission is incentivized by potential actions of the decisionmaker. That is, informed agents prefer to tell the truth, since the corresponding action is preferred. To a large extent, these results are based on the existence of some "common interest" shared by informed agents and decision-makers, i.e. there exist actions over which the preferences of the sender and the receiver coincide. In the present paper the possibility of communication does not rely on the existence of "common interests"<sup>0</sup>. The only incentive for communication is the higher benefit of information received through the mediation protocol in case of truth-telling (i.e. the lower downside of distorted messages).

Second, communication between partially informed agents, who also take actions, is studied in a number of papers. An early paper proposing the communication equilibrium solution concept is Forges (1986) (see also Myerson (1997, Chapter 6)). On the more applied side, Galeotti et al. (2013) extend the CS model to the case of multiple decision-makers with private information regarding the state of nature. In their model telling lies is again precluded by the corresponding unfavorable change in other agents' actions. Alonso et al. (2008) model communication between partially informed managers who care about the profits of own divisions and action coordination. Communication via cheap talk is possible in that model, since the managers prefer their actions to be close to each other and thus have an incentive to reveal some private information. In the industrial organization context Goltsman and Pavlov (2014) look into the case of communication between Cournot oligopolists<sup>1</sup>, who may share unverifiable private information about costs, and show that no information transmission occurs in the cheap-talk game, but information can be transmitted through a neutral third party. The question explored by Goltsman and Pavlov (2014) is similar to that of the present paper, but in their case the competitor's information is relevant, because it affects her action and actions are strategic substitutes. The mediator is able to exploit the coordination motives to achieve communication: some types of agents report truthfully in order to make the opponent less aggressive. The incentives for truth-telling provided by the mediation protocol in the present paper are purely informational: reporting truthfully leads to a higher benefit of information received back. The "secret sharing" game presented as an example in Vida and Forges (2013) has a similar structure to the illustrative example in this paper. However, the example setting in Vida and Forges (2013) features (i) information that is independent across players; (ii) communication between the players relying on the availability of verifiable "signatures" that allow to detect deception of each individual agent; (iii) the possibility of full information transmission. In the present paper, instead, (i) players' types can be correlated; (ii) the mediation protocol solely relies on the information structure and is still able to provide incentives for truth-telling; (iii) only partial information transmission is possible. Kolotilin et al. (2017) explore persuasion of a privately informed receiver. Similarly, in the present paper, after observing a report from one of the agents, the mediation protocol essentially tries to persuade the other one with a caveat that the other agent is informed herself. Kolotilin et al. (2017) show that private persuasion by the sender (asking the receiver for a type report and returning a private message) is equivalent to the sender broadcasting information without asking the receiver for a type report. This is not the case in the present paper:

<sup>&</sup>lt;sup>0</sup>In fact, the preferences may be completely "opposite": the payoff structure, such that an increase in one agent's utility may always lead to the decrease of the other agent's utility, is allowed.

<sup>&</sup>lt;sup>1</sup>Further examples of the literature on communication in oligopoly include Novshek and Sonnenschein (1982), Vives (1984), Gal-Or (1985), Li (1985), Shapiro (1986), Vives (1990) and Raith (1996). See Kühn and Vives (1995) and Vives (2001) for extensive reviews. This strand of research typically assumes commitment power or verifiable private information. A notable exception is Ziv (1993), who shows that conveying credible information in the oligopoly setting is also possible if "money-burning" or transfers are allowed. In the present paper agents lack commitment power, information is non-verifiable, and both "money-burning" and transfers are assumed out.

the mediation protocol *does* need to condition her recommendation on agents' reports in order to sustain informative communication.

To sum up, on the one hand, the present paper differs from the literature discussed above in terms of underlying assumptions, results and intuition. On the other hand, it contributes to that literature by offering a novel tool that facilitates communication between agents with a sufficiently small misalignment of interests.

## **3** Illustrative example

This section introduces an example capturing the main intuition of the results in the present paper. In the example players receive binary signals and need to guess the mean of the two signals. While direct communication is not possible, the class of almost-truthful interim-biased mediation protocols is shown to enable information exchange when (i) signals are correlated, and (ii) the misalignment of interests is sufficiently small.

#### 3.1 Setup

Consider the following game  $\Gamma_E$ . Each of the two agents  $k \in \{1, 2\}$  obtains a binary signal  $s_k \in S = \{0, 1\}$  with the following joint distribution  $\pi$  over  $S^2 = S \times S$  parametrized by  $r \in (1/2, 1)$ :

$$\frac{\pi}{s_2 = 0 \quad s_2 = 1}$$

$$s_1 = 0 \quad \frac{r}{2} \quad \frac{1 - r}{2}$$

$$s_1 = 1 \quad \frac{1 - r}{2} \quad \frac{r}{2}$$

Together, these signals determine the correct action  $s^* = \frac{1}{2} \sum_k s_k$ . Both agents would like to guess  $s^*$  by choosing an action in the set  $\mathcal{A}_k = \{0, 1/2, 1\}$ . The agents have misaligned interests and prefer the opponent not to be able to guess the correct action. The payoffs representing such preferences are given by

$$u_k(a,s) = \mathbb{1} \{ a_k = s^* \} - \alpha \times \mathbb{1} \{ a_{3-k} = s^* \}$$
(1)

where  $\alpha \in (0, 1)$  parametrizes the degree of interest misalignment between the agents: the higher  $\alpha$ , the more each player is hurt by the competitor's correct guess.

Notice that under no communication the optimal strategy of each agent is choosing an action that coincides with the observed  $a_k = s_k$ . Also note that the lack of communication is suboptimal: if agents were able to disclose signals to each other, the expected payoffs in the game would go from  $r(1 - \alpha)$  to  $(1 - \alpha)$ . Proposition A.1 and Proposition A.2 formally establish these two results in Appendix A.

However, the agents are not able to communicate in a game with simultaneous message exchange. If there was a message one agent could send and shift the action of the counterpart, such message would be used to deceive the counterpart. The deceitful agent could benefit from the counterpart's mistake, and would suffer no losses due to the lack of coordination motives and the unchanged counterpart's messaging strategy. See Proposition A.3 of Appendix A for a formal treatment.

#### 3.2 Almost-truthful interim-biased mediation

While direct communication is impossible, this subsection introduces the notion of almost-truthful interim-biased mediation, which facilitates information exchange for low enough misalignment of interests.

**Mediation setup** The almost-truthful interim-biased mediation protocol receives signal reports from the agents and sends private messages back to them. It almost always sends agent (3 - k)'s signal report to agent k without any distortions, thus the almost-truthful label. However, with a positive probability  $\varepsilon$  the mediation protocol returns to agent k a message that coincides with (3 - k)'s most likely signal given k's own report. If agent k's report is truthful, such a message leads to k choosing the interim-optimal action, thus the interim-biased label.

Formally, let  $m_k(\hat{s}_k, \hat{s}_{3-k})$  be a collection of binary distributions over mediator's messages to agent k when the mediator receives reports  $\hat{s}_k, \hat{s}_{3-k}$  from the agents:

$$m_k(\hat{s}_k, \hat{s}_{3-k}) = \begin{cases} \hat{s}_{3-k} & \text{with probability} \quad 1-\varepsilon \\ \hat{s}_k & \text{with probability} \quad \varepsilon, \end{cases}$$

where

$$\varepsilon \in \left(0, \frac{1-r}{r}\right)$$

Notice that the two types of messages (with and without a distortion) take value in the same set, and thus the agents cannot distinguish them with certainty.

**Truthful equilibrium** Such a mediation protocol ensures that there exists an equilibrium in which both agents transmit their information truthfully. Two observations are necessary for this result.

First, if agent (3 - k) reports truthfully, agent k chooses the average between the private signal and the mediator's message as her action irrespective of whether k herself reports truthfully. To see this, note that the posterior probability on the event  $s_{3-k} = m$  where m is the observed mediator's message is greater than 1/2. Indeed, if the mediator's message does not coincide with k's report,  $m \neq \hat{s}_k$ , then agent k for sure knows that  $m = s_{3-k}$  and  $\mathbb{P}_{\hat{s}_i} \left[ s_{3-k} = m | m_k(\hat{s}_k, s_{3-k}) = m, s_k \right] = 1$ . If the mediator's message does coincide with k's report,  $m = \hat{s}_k$ , then

$$\mathbb{P}_{\hat{s}_{i}}\left[s_{3-k} = m | m_{k}(\hat{s}_{k}, s_{3-k}) = m, s_{k}\right] = \frac{\mathbb{P}\left[m_{k}(\hat{s}_{k}, s_{3-k}) = m | s_{3-k} = m, s_{k}\right] \mathbb{P}\left[s_{3-k} = m | s_{k}\right]}{\mathbb{P}\left[m_{k}(\hat{s}_{k}, s_{3-k}) = m | s_{k}\right]} \\ = \begin{cases} \frac{r}{r + \varepsilon(1-r)} & \text{if } \hat{s}_{k} = s_{k} \\ \frac{1-r}{(1-r) + \varepsilon r} & \text{if } \hat{s}_{k} \neq s_{k} \end{cases} \\ > \frac{1}{2},$$

since

$$\varepsilon < \frac{1-r}{r}$$

As a result, for low values of  $\varepsilon$  the agent optimally chooses action  $a_k = \frac{1}{2} (s_k + m_k(\hat{s}_k, s_{3-k}))$ , the average of the private signal and mediator's message.

Second, since k's optimal actions conditional on deviating and not are known, the consequences of deviating and not are simple to predict. If the mediator does not distort the information, agent k chooses action  $s^*$ . If the mediator distorts the information agent k chooses the interim-optimal

action  $s_k$  in case she doesn't deviate and chooses action 1/2 in case she does deviate. That is, if k doesn't deviate, she makes a mistake when the mediator distorts and her interim-optimal action is not ex-post optimal, which happens with probability  $(1 - r)\varepsilon$ . If k deviates, she makes a mistake when the mediator distorts and her interim-optimal action was actually ex-post optimal, which happens with probability  $r\varepsilon$ . Thus the expected *partial payoff from k's own actions* in case of a truthful report equals  $V^T = 1 - (1 - r)\varepsilon$  and equals  $V^D = 1 - r\varepsilon$  in case of a deviation.

It remains to notice that the expected partial payoff in case of reporting truthfully is higher:  $\Delta V = V^T - V^D = \varepsilon(2r - 1) > 0$ . Provided that  $\alpha$  is low enough, the positive partial difference  $\Delta V$  dominates the total payoff difference, and reporting truthfully is optimal for each agent. The existence of a truthful equilibrium is thus shown.

Also, the expected payoff of each agent is higher in case of communication enabled by the mediation protocol above relative to the no-communication case. Specifically, without communication the expected payoff of each player is equal to  $r(1 - \alpha)$  (a player doesn't make a mistake only when the signals coincide). With almost-truthful mediation, each player also stops making a mistake when signals don't coincide and there is no distortion from the mediator, resulting in the expected payoff of  $r(1 - \alpha) + (1 - r)(1 - \varepsilon)(1 - \alpha) > r(1 - \alpha)$ . Thus, almost-trustful mediation leads to a welfare improvement in the illustrative example case.

**Highlighting incentives and assumptions** First, the mediator is able to control agent's beliefs by being sufficiently truthful. Such control provides the mediator with the opportunity to deliberately shift agents' actions when the mediator distorts. Second, when the mediator does distort the information, the action of the deviating player is shifted away from the interim-optimal action, while the action of the non-deviating player is not. Note that one can also interpret the distorted message as an implicit encouragement for a truthful agent to take the action that is interim-optimal given her private signal report only (as the distorted message coincides with the report). At the same time, a deviating agent receives an encouragement based on an untruthful report and shifts her action away from the truly interim-optimal one when receiving a distorted message. As a result, the mediator hurts the deviating player more, when distorting the information: effectively, the deviating player imposes an inefficient action upon herself and is put to a disadvantage.

This last point depends on the assumption of r > 1/2, which imposes positive correlation between the agents' signals. If r = 1/2 (signals are uncorrelated), then the mediator cannot shift the action away from the interim-optimal action<sup>2</sup>. Thus the mediator's distortion is equally harmful for the agent irrespective of whether she reports truthfully or not, and communication breaks down. In fact r = 1/2 is problematic in two ways: (i) non-uniqueness of interim-optimal action, (ii) different types have same beliefs about the counterpart's signals. While there is no distinction between the two in the illustrative example, the main result of this paper will separately impose the assumptions of uniqueness (see Assumption 4.1 and Assumption 4.2) and sensitivity of beliefs to private information (see Assumption 4.4). Jointly these assumptions guarantee that the mediator's messages can shift agents' actions.

Additionally, the payoff structure in (1) ensures that the agent's optimal action is sensitive to the counterpart's information. If this was not the case, agents would have little incentive to report their signals truthfully as getting additional information from the counterpart would be worthless.

<sup>&</sup>lt;sup>2</sup>Which is not unique: both the private signal  $s_k$  and 1/2 are optimal actions for every agent.

## 4 Main result

This section generalizes the illustrative example by constructing a model that allows for arbitrary joint distributions over states of nature. When agents' preferences are separable in actions, the almost-truthful interim-biased mediation enables communication, while, in general, it is not true that welfare-improving direct communication is possible under the maintained assumptions. Sufficient conditions for the almost-truthful interim-biased mediation protocol allowing truthful information exchange are established. These conditions are: (i) action sensitivity to counterpart's information, (ii) sufficient variation in interim beliefs across agent types, and (iii) sufficiently weak misalignment of interests.

### 4.1 Model

The model consists of two agents who are endowed with private information regarding the state of nature and have to take an action. Each action affects the payoffs of both agents. A form of additive separability is assumed: the action's effect on the other party's payoff only depends on the action itself and the state of nature, but not on that other party's action; this assumption is slightly relaxed in Appendix C.

**Baseline** Formally, consider a 2-agents setup with a finite state space  $S = S_1 \times S_2$ . Each agent k learns the realization of  $s_k \in S_k$  (the signal), but does not learn  $s_{3-k}$ . Let  $\pi(\cdot)$  be the common prior over S. Also, let  $\pi_k(\cdot|s_k)$  and  $\mathbb{E}_k [\cdot |s_k]$  be agent k's posterior and expectation operator upon learning  $s_k$  respectively. Each agent k chooses an action  $a_k$  from a finite action space  $\mathcal{A}_k$ . It is assumed that the agents' preferences over action profiles are represented by a state-dependent utility function

$$u_k(s,a) = v_k(a_k,s) - \alpha \times c_k(a_{3-k},s)$$
<sup>(2)</sup>

 $v_k(a_k, s)$  captures agent k's value from taking action  $a_k$  in a given state  $s \in S$ . Without loss of generality assume  $v_k(a_k, s) \ge 0$  for every k,  $a_k$  and s. The preferences of agent k with respect to (3-k)'s actions in a given state are captured by the cost function  $\alpha \times c_k(a_{3-k}, s)$ . The cost function  $c_k$  captures the idea of agents' misaligned interests: as in the illustrative example above, a change in agent (3-k)'s action that is beneficial for agent (3-k) is allowed to increase the cost function  $c_k$  (and thus to be harmful) for agent k. The cost component of the utility function is further parametrized by the parameter  $\alpha > 0$  that captures the degree of the misalignment of interests<sup>3</sup>. It is also worth noting that the separability of agents' utility functions in each other's actions implies that actions are neither strategic complements nor strategic substitutes. This eliminates the option to reveal information about the counterpart's action as a potential leverage that the mediator can use to elicit the truth (this leverage is used, for example, in Goltsman and Pavlov (2014)).

**Definitions** The definitions that simplify the notation in the remaining part of the paper are now introduced. First, in a fixed state *s* in S each agent *k* can maximize her payoff by taking the same action for all actions of (3 - k). That is, due to the assumption of separability of agents'

<sup>&</sup>lt;sup>3</sup>This paper is primarily motivated by the situations where agents have misaligned interests, but the model formally allows for aligned agents' preferences as well (i.e. a change in one player's action being beneficial for both players). Exploiting the specifics of such an alignment of interests can, in principle, be used to facilitate communication, but constructing such setting-specific schemes is beyond the scope of the present paper. The mediation protocol class introduced further in the paper can enable information exchange in the situations of aligned interests (under a set of assumptions, also discussed below), but does not depend on such an alignment.

preferences with respect to each other's actions, agent k can choose the action that maximizes the  $v_k$ -component of her utility. Definition 4.1 below introduces formal notation for such state-specific *correct actions*:

DEFINITION 4.1. For agent k in state  $s \in S$  let  $a_k^*(s)$  be the set of state-specific *correct actions*. That is,  $a_k^*(s) = \arg \max_{a_k} v_k(a_k, s)$ .

Similarly, for every privately observed state  $s_k \in S_k$  the *interim-correct actions* are defined:

DEFINITION 4.2. For each agent k and signal  $s_k \in S_k$  let  $\tilde{a}_k(s_k)$  be the set of *interim-correct actions*. That is,  $\tilde{a}_k(s_k) = \arg \max_{a_k} \mathbb{E}_k \left[ v(a_k, s) | s_k \right]$ .

Definition 4.2 captures the notion of the best possible actions in autarky. Such actions would be taken by each agent in the absence of any information exchange. Definition 4.3 below introduces the notion of *interim-correct counterpart signals* that links correct actions and interim-correct actions.

DEFINITION 4.3. For each agent k and signal  $s_k \in S_k$ , let  $\tilde{\sigma}_k(s_k)$  be the set of *interim-correct* counterpart signals. That is,  $\tilde{\sigma}_k(s_k) = \{s_{3-k} | \tilde{a}_k(s_k) = a_k^*(s_k, s_{3-k})\}$ .

Specifically, if agent (3-k)'s signal was revealed to belong to the set  $\tilde{\sigma}_k(s_k)$ , then agent k endowed with signal  $s_k$  would have no incentive to take an action other than the interim-correct one. At this point,  $\tilde{\sigma}_k(s_k)$  can be an empty set, Assumption 4.2 below ensures it is non-empty.

**A preliminary result** The following lemma establishes a natural result that will be useful throughout the rest of the paper. It states that an action that is correct for a given signal of the counterpart will be chosen for high enough belief on this signal.

LEMMA 4.1. Let  $\tilde{\pi}_k$  be agent k's belief over  $S_{3-k}$ . There exists a  $\bar{\delta}_k < 1$  such that for all  $s_{3-k}$  if  $\tilde{\pi}_k(s_{3-k}) \ge \bar{\delta}_k$ , then  $\arg \max_{a_k} \mathbb{E}_{\tilde{\pi}_k} \left[ v_k(a_k, s) \right] = a_k^*(s_k, s_{3-k})$ .

*Proof.* Take an arbitrary  $s_{3-k}$  and notice that for every  $a^* \in a_k^*(s_k, s_{3-k})$ , when  $\tilde{\pi}_k$  is such that  $\tilde{\pi}_k(s_{3-k}) = 1$ , then

$$\mathbb{E}_{\tilde{\pi}_{k}} \left[ v_{k}(a^{*}, s) \right] = \sum_{t_{3-k}} \tilde{\pi}_{k}(t_{3-k}) v_{k}(a^{*}, s)$$
  
=  $v_{k}(a^{*}, s)$   
>  $v_{k}(a', s)$   
=  $\sum_{t_{3-k}} \tilde{\pi}_{k}(t_{3-k}) v_{k}(a', s)$   
=  $\mathbb{E}_{\tilde{\pi}_{k}} \left[ v_{k}(a', s) \right]$ 

for every  $a' \notin a_k^*(s_k, s_{3-k})$  by Definition 4.1. Thus by continuity of  $\sum_{t_{3-k}} \tilde{\pi}_k(t_{3-k})v_k(a, s)$  with respect to  $\tilde{\pi}_k(t_{3-k})$ , there exists a  $\bar{\delta}_k(s_{3-k}) < 1$  such that the same strict inequality holds for all  $\tilde{\pi}_k$  such that  $\tilde{\pi}_k(s_{3-k}) \ge \bar{\delta}_k(s_{3-k})$ . The proof of the lemma is completed by defining  $\bar{\delta}_k = \max_{s_{3-k}} \{\bar{\delta}_k(s_{3-k})\}$ .

#### 4.2 Assumptions

The additional assumptions stated below limit the scope of the results of the present paper to settings in which the prior has full support and the action space is of intermediate coarseness. As will be clear from the main result of the paper, intermediate coarseness guarantees two things. First, actions can be shifted by additional information. Second, concealment of information can

appear like provision of additional information. It should be noted that while the assumptions below impose restrictions on endogenous objects, such a presentation leads to a succinct description of the setup's features that are necessary for the main result of the current paper.

First, an assumption regarding the structure of the correct action set is made.

Assumption 4.1. (i)  $a_k^*(s)$  is a singleton for every k and  $s \in S$ . (ii) For every agent k and signal  $s_k$ , if  $s'_{3-k} \neq s_{3-k}$ , then  $a_k^*(s_k, s'_{3-k}) \neq a_k^*(s_k, s_{3-k})$ .

Part (i) precludes the existence of actions that lead to the same consequences in a given state and can be interpreted as a "no redundant actions" requirement. Part (ii) is a sensitivity assumption which ensures that the action space is rich enough so that the choice of action can be adjusted for alternative states of nature.

Then, an assumption regarding the structure of the interim-correct set is made.

Assumption 4.2. (i)  $\tilde{a}_k(s_k)$  is a singleton for every k and  $s_k \in S_k$ . (ii) For every agent k and signal  $s_k$  there exists  $s_{3-k}$  such that  $\tilde{a}_k(s_k) = a_k^*(s_k, s_{3-k})$ .

Part (i) of Assumption 4.2 is a joint assumption on the action space and the information structure. It ensures no indifferences on the interim stage. Under Assumption 4.1 this assumption is guaranteed to hold when the posteriors  $\pi_k$  are close enough to degenerate ones:  $\pi_k(\cdot|s_k) > 0^+$  for the true signal  $s_{3-k}$  only<sup>4</sup>. The results of the present paper are thus guaranteed to hold when agents' signals exhibit a sufficient degree of dependence.

Part (ii) of Assumption 4.2 is a coarseness assumption which guarantees that the interim action cannot be perfectly adjusted to the non-degenerate interim information of the agents. It is guaranteed to hold if there are no "redundant" actions in the action set  $(\mathcal{A}_k = \{a_k^*(s)|s \in S\})$  for every k, and if the  $v_k$ -component is monotone with respect to own-signal (for every agent k, signals  $s_k, s_{3-k}$  and  $s'_{3-k}, v_k(a_k^*(s_k, s'_{3-k}), (s_k, s_{3-k})) \ge v_k(a_k^*(s'_k, s'_{3-k}), (s_k, s_{3-k}))$ . Note that Assumption 4.1 and Assumption 4.2 jointly ensure that the set of interim-correct counterpart signals  $\tilde{\sigma}_k(s_k)$  is a singleton for all agents k and private states  $s_k$ .

Finally, the assumption regarding the information structure is made.

Assumption 4.3. (i)  $|S_1| = |S_2|$ . (ii)  $\pi(s) > 0$  for all  $s \in S$ .

Jointly, these assumptions can be interpreted as similarity of agents' private signal quality. Part (i) of Assumption 4.3 limits the scope of the paper to settings where each agent can receive the same number of different signals. It will help ensuring that different types of each agent hold sufficiently different beliefs about the counterpart's signals<sup>5</sup>. Part (ii) of Assumption 4.3 prevents complete elimination of uncertainty and thus each type of every agent is *not* entirely informed about the state of nature upon realization of the private signal.

### 4.3 On direct communication

Can the two agents help each other directly by simultaneously sending messages that are at least partially informative? If one relies on the intuition from the illustrative example case (where agents benefit from each other's mistakes), the answer should be negative, see Proposition A.3 in Appendix A. In that case, conditional on receiving a message from the counterpart and the

<sup>&</sup>lt;sup>4</sup>To see this formally, one can utilize Lemma 4.1.

<sup>&</sup>lt;sup>5</sup>The setup under part (i) of Assumption 4.3 can be interpreted as follows. Agents agree on the possible states of the world:  $S_1 = S_2 = S$ , and receive noisy signals about the true state with the signal space coinciding with the states of the world set.

corresponding posterior over the counterpart's signal (which governs the action choice), each agent has no incentives to send truthful information as it benefits the counterpart and, correspondingly, hurts the agent herself. This section establishes that, more generally, for cases of "opposite" preferences, such that, state by state, one agent is better off when the other one is worse off, direct communication cannot increase agents' payoffs. That is, it is not true, in general, that welfare-improving direct communication is possible under the maintained assumptions.

**Proposition 4.1.** Consider direct communication extensions of the baseline game, i.e., games from Section 4.1 extended with finite sets  $W_k$  (k = 1, 2) of messages that each agent can send to her counterpart upon observing the private signal and before taking action. There exist cost functions  $c_k(\cdot, \cdot)$  such that, for any  $\alpha > 0$  and for any such direct communication extension of the baseline game, both agents' expected payoffs are not higher in equilibrium relative to the baseline game.

*Proof.* Consider the case of "opposite" preferences:  $c_k(a_{3-k}, s) = v_{3-k}(a_{3-k}, s)$ . To show that payoff-improving direct communication is not possible in this case, the following approach is used. Under the assumption that there exists an equilibrium with direct communication increasing at least one agent's expected payoff, it is shown that some type of one of the agents necessarily has a profitable deviation, in which she always sends a fixed message and shifts the counterpart's action away from the interim-correct one.

Formally, consider an extended game in which agents can simultaneously exchange messages from predetermined sets  $W_k$  after observing the private signals and before taking actions. Suppose that there is an equilibrium of such an extended game in which the expected payoff of at least one agent is strictly higher relative to the baseline game. Then for at least one agent (3 - k), the ex-ante expected value of the own-choice component of the utility  $v_{3-k}$  has to be strictly greater than under no communication<sup>6</sup>. The strictly higher expected value of the  $v_{3-k}$ -component relative to no communication means that some type of agent (3 - k) necessarily shifts her action away from the interim-correct one when receiving some message  $w_k$  from agent k.

Then, in turn, agent k can exploit (3 - k)'s reaction to message  $w_k$  and profitably deviate from the equilibrium under consideration. k can do so by (1) sending this same message  $w_k$  irrespective of her own signal and (2) following the same own-action plan (possibly contingent on the received messages) as in the equilibrium under consideration. Note that under such a deviation, the ex-ante expected value of the  $v_{3-k}$ -component of agent (3 - k)'s payoff becomes lower than without any communication at all: since  $w_k$  is sent regardless of k's actual signal, at least one type of agent (3 - k) shifts her action away from the interim-correct action for every k's signal. This deviation is thus profitable for agent k. Indeed, it ensures that the expected value of the  $v_{3-k}$ -component is strictly lower than under no communication (by Assumption 4.2 the interim-correct action is unique and thus a shift to a different action means a reduction in the  $v_{3-k}$ -component), and thus the expected value of k's  $c_k$ -component of the utility  $-\alpha c_k = -\alpha v_{3-k}$  is higher. Given that k's own action plan is unchanged and the expected value of the  $v_k$ -component is fixed, the proposed deviation is profitable for agent k.

Thus there does not exist an equilibrium of the direct communication extension increasing at least

<sup>&</sup>lt;sup>6</sup>If the expected value of  $v_k$  is weakly lower than under no communication for both k = 1 and k = 2, there are two possible cases. In the first case, the expected value of the  $v_k$ -component is the same as without communication for both agents, which means no welfare improvement coming from direct communication. In the second case, the expected value of the  $v_k$ -component is strictly lower relative to the no-communication situation for some k. The latter, however, is impossible in equilibrium, since returning to the interim-correct actions would be profitable for such agent k. Thus, there has to exist k such the expected value of the  $v_k$ -component is strictly greater in the welfare-improving equilibrium of the direct communication extension, provided that such an equilibrium exists.

one agent's expected payoffs relative to the baseline game, and direct communication does not allow for welfare improvements in case of "opposite" preferences.

Note that Proposition 4.1 demonstrates that welfare-improving direct communication is impossible for the case of each agent benefiting from the losses of the counterpart:  $c_k(a_{3-k}, s) = v_{3-k}(a_{3-k}, s)$ . While this implies that, in general, it is not true that welfare-improving direct communication is possible under the maintained assumptions, for some cost functions, direct communication may actually be possible (consider, for instance, the case of  $c_k(a_{3-k}, s) = 0 \forall a_{3-k}, s$ ). Almost-truthful interim-biased mediation protocols presented below, in turn, enable information exchange for any cost function, provided that  $\alpha$  is low enough.

#### 4.4 Almost-truthful interim-biased mediation

This subsection introduces interim-biased and almost-truthful interim-biased protocol classes. Such protocols take signal reports from both players and send private messages to each agent. The set of messages  $M_k$  to agent k coincides with the set  $S_{3-k}$ . This captures the notion that the mediator may share some of (3 - k)'s private information with agent k.

Interim-biased mediation protocols introduced in Definition 4.4 below are randomized for each pair of signal reports. Such protocols transmit agent (3 - k)'s report to agent k with probability  $1 - \varepsilon_k$  and send a message with k's interim-correct counterpart signal with the complementary probability  $\varepsilon_k$ .

DEFINITION 4.4. Let an *interim-biased mediation protocol*  $m^b$  be a collection of random variables  $\{m_k^b\}_{k=1,2}$  with

$$m_{k}^{b}(\hat{s}_{k},\hat{s}_{3-k}) = \begin{cases} \hat{s}_{3-k} & \text{with probability} \quad 1 - \varepsilon_{k} \\ \tilde{\sigma}_{k}(\hat{s}_{k}) & \text{with probability} \quad \varepsilon_{k}, \end{cases}$$
(3)

for some  $\varepsilon_k \in [0, 1]$ .

The following lemma establishes that there exist positive probabilities  $\varepsilon_k$  such that, conditional on agent (3 - k) reporting truthfully, agent k believes that the mediator's message coincides with (3 - k)'s report irrespective of k's own report. Formally,

LEMMA 4.2. For every agent k and  $1/2 < \delta < 1$ , there exists an  $\bar{\varepsilon}_k(\delta) > 0$  such that if an interimbiased mediation protocol satisfies  $\varepsilon_k \leq \bar{\varepsilon}_k(\delta)$ , then  $\mathbb{P}\left[s_{3-k} = m|s_k, m_k^b(\hat{s}_k, s_{3-k}) = m\right] \geq \delta$  for all  $s_k, \hat{s}_k \in S_k$  and  $m \in S_{3-k}$ .

*Proof.* Consider the beliefs of agent k upon reporting signal  $\hat{s}_k$  and receiving message m from the mediator.

1. If  $m \neq \tilde{\sigma}_k(\hat{s}_k)$ , then the posterior probability that agent (3 - k)'s signal report is equal to *m* can be computed as follows.

$$\mathbb{P}\left[s_{3-k} = m | s_k, m_k^b(\hat{s}_k, s_{3-k}) = m\right] = \frac{\mathbb{P}\left[m_k^b(\hat{s}_k, s_{3-k}) = m | s_k, s_{3-k} = m\right] \times \mathbb{P}\left[s_{3-k} = m | s_k\right]}{\mathbb{P}\left[m_k^b(\hat{s}_k, s_{3-k}) = m | s_k\right]}$$
$$= \frac{(1 - \varepsilon_k) \times \mathbb{P}\left[s_{3-k} = m | s_k\right]}{\mathbb{P}\left[m_k^b(\hat{s}_k, s_{3-k}) = m | s_k\right]}$$

Notice that given Definition 4.4,  $\mathbb{P}\left[m_k^b(\hat{s}_k, s_{3-k}) = m|s_k\right] = (1 - \varepsilon_k \times \mathbb{P}\left[s_{3-k} = m|s_k\right]$  and thus  $\mathbb{P}\left[s_{3-k} = m|s_k, m_k^b(\hat{s}_k, s_{3-k}) = m\right] = 1 \ge \delta$ . Consequently, for the lemma to be true only the case of  $m = \tilde{\sigma}_k(\hat{s}_k)$  needs to be considered.

2. If  $m = \tilde{\sigma}_k(\hat{s}_k)$ , then

$$\mathbb{P}\left[s_{3-k} = m | s_k, m_k^b(\hat{s}_k, s_{3-k}) = m\right] = \frac{1}{1 + \sum_{t_{3-k} \neq m} \varepsilon_k R_\pi(t_{3-k}, m | s_k)},$$

where

$$R_{\pi}(t_{3-k}, m|s_k) = \frac{\pi_k(s_{3-k} = t_{3-k}|s_k)}{\pi_k(s_{3-k} = m|s_k)}$$

Note that  $R_{\pi}(t_{3-k}, m|s_k)$  is a positive finite number for every *m* and  $t_{3-k}$  under Assumption 4.3. Thus if  $\varepsilon_k$  satisfies

$$\varepsilon_k \leqslant \bar{\varepsilon}_k = \min_{s_k, \hat{s}_k} \left[ \frac{1-\delta}{\delta} \times \frac{1}{\sum\limits_{t_{3-k} \neq m} R_\pi(t_{3-k}, m|s_k)} \right] = \frac{1-\delta}{\delta} \times \min_{s_k, \hat{s}_k} \left[ \frac{\pi_k(s_{3-k} = m|s_k)}{1-\pi_k(s_{3-k} = m|s_k)} \right],$$

then  $\mathbb{P}\left[s_{3-k} = m | s_k, m_k^b(\hat{s}_k, s_{3-k}) = m\right] \ge \delta$  for every  $s_k, \hat{s}_k$ .

The proof is completed by observing that  $\bar{\varepsilon}_k \in (0, 1)$ .

Now the class of *almost-truthful interim-biased mediation protocols* is introduced. Such protocols belong to the interim-biased mediation class and share a defining common feature: the probability  $1 - \varepsilon_k$  of truthful transmission of each agent's report is close to 1. Formally,

DEFINITION 4.5. Let an *almost-truthful interim-biased mediation protocol*  $m^a$  be a collection of random variables  $\{m_k^a\}_{k=1,2}$  with  $m_k^a(\hat{s}_k, \hat{s}_{3-k}) = m_k^b(\hat{s}_k, \hat{s}_{3-k})$  for some  $\varepsilon_k \in (0, \overline{\varepsilon}_k(\overline{\delta}_k)]$ .

Notice that due to Lemma 4.1 and Lemma 4.2, each almost-truthful interim-biased mediation ensures agents k's posterior belief has a high enough weight on m when the mediator's message is m. Thus agents take action  $a_k^*(s_k, m)$  upon receiving m from the mediator. Importantly, this fact only depends on agent (3 - k) (but not agent k) reporting truthfully. Lemma 4.3 below states this result formally.

LEMMA 4.3. For each agent k, signal  $s_k$ , signal report  $\hat{s}_k$  and realization m of mediator's message  $m_k^a(\hat{s}_k, s_{3-k})$ , agent k's optimal action coincides with the correct action in state  $(s_k, m)$ : arg  $\max_{a_k} \mathbb{E}_k \left[ v_k(a_k, s) | s_k, m_k^a(\hat{s}_k, s_{3-k}) = m \right] = a_k^*(s_k, m)$ .

*Proof.* By Definition 4.5 agent k's posterior belief over agent (3 - k)'s signal places a higher than  $\bar{\delta}_k$  weight on m. By Lemma 4.1,  $\arg \max_{a_k} \mathbb{E}_k \left[ v_k(a_k, s) | s_k, m_k^a(\hat{s}_k, s_{3-k}) = m \right] = a_k^*(s_k, m).$ 

#### 4.5 Equilibrium with information exchange

This subsection establishes that protocols from the almost-truthful interim-biased class allow truthful information exchange between the agents. An additional assumption required for this result ensures enough variation in intermediate beliefs across different types of agents:

Assumption 4.4. For each agent k and pair of signals  $s'_k \neq s_k$ ,  $\tilde{\sigma}_k(s'_k) \neq \tilde{\sigma}_k(s_k)$ .

This assumption is interpreted as follows: the interim-correct actions are rationalized by different counterpart's signals. Notice that Assumption 4.4 implies  $|S_1| = |S_2|$ , a feature of the setting that was assumed earlier under Assumption 4.3. Assumption 4.4 is guaranteed to hold if receiving different signals leads to sufficiently different beliefs about the opponents' signals. For example, if the posterior weight on the most likely opponent's signal is sufficiently strong and under different

observed signals different counterpart's signals are most likely, then the result of Lemma 4.1 and Assumption 4.4 hold<sup>7</sup>.

To demonstrate the existence of truth-telling equilibrium, first, the  $v_k$ -component of the utility is shown to be strictly maximized by truthful reporting. Let  $V_{m_k^a}(s_k, \hat{s}_k)$  be the expectation of agent k's  $v_k$ -component of the utility conditional on reporting  $\hat{s}_k$  to an almost-truthful mediation protocol  $m_k^a$  when the true signal is  $s_k$  and agent (3 - k) reports truthfully. Notice that by Lemma 4.3,

$$V_{m_k^a}(s_k, \hat{s}_k) = \sum_{t_{3-k}} \pi_k(t_{3-k}|s_k) \times (1 - \varepsilon_k) \times v_k(a_k^*(s_k, t_{3-k}), (s_k, t_{3-k})) + \sum_{t_{3-k}} \pi_k(t_{3-k}|s_k) \times \varepsilon_k \times v_k(a_k^*(s_k, \tilde{\sigma}_k(\hat{s}_k)), (s_k, t_{3-k}))$$

The following lemma establishes that truthful reporting maximizes the  $v_k$ -component of the utility:

LEMMA 4.4. Suppose that Assumption 4.4 holds. For each agent k, signals  $s_k \neq \hat{s}_k$  and mediation protocol  $m^a$ ,  $V_{m_k^a}(s_k, s_k) > V_{m_k^a}(s_k, \hat{s}_k)$ .

*Proof.* Define  $\Delta V_{m_k^a}(s_k, \hat{s}_k) = V_{m_k^a}(s_k, s_k) - V_{m_k^a}(s_k, \hat{s}_k)$ . Notice that

$$\begin{aligned} \Delta V_{m_k^a}(s_k, \hat{s}_k) &= \\ &= \sum_{t_{3-k}} \pi_k(t_{3-k}|s_k) \times v_k(a_k^*(s_k, t_{3-k}), (s_k, t_{3-k})) \times [\varepsilon_k - \varepsilon_k] \\ &+ \sum_{t_{3-k}} \pi_k(t_{3-k}|s_k) \times \varepsilon_k \times \left[ v_k(a_k^*(s_k, \tilde{\sigma}_k(s_k)), (s_k, t_{3-k})) - v_k(a_k^*(s_k, \tilde{\sigma}_k(\hat{s}_k)), (s_k, t_{3-k})) \right] \\ &= \varepsilon_k \sum_{t_{3-k}} \pi_k(t_{3-k}|s_k) \times \left[ v_k(a_k^*(s_k, \tilde{\sigma}_k(s_k)), (s_k, t_{3-k})) - v_k(a_k^*(s_k, \tilde{\sigma}_k(\hat{s}_k)), (s_k, t_{3-k})) \right] \end{aligned}$$

Rearranging,

$$\Delta V_{m_k^a}(s_k, \hat{s}_k) = \varepsilon_k \sum_{t_{3-k}} \pi_k(t_{3-k}|s_k) \times v_k(a_k^*(s_k, \tilde{\sigma}_k(s_k)), (s_k, t_{3-k})) - \varepsilon_k \sum_{t_{3-k}} \pi_k(t_{3-k}|s_k) \times v_k(a_k^*(s_k, \tilde{\sigma}_k(\hat{s}_k)), (s_k, t_{3-k}))$$

First notice that  $a_k^*(s_k, \tilde{\sigma}_k(s_k)) = \tilde{a}_k(s_k)$  by Definition 4.2 and part (i) of Assumption 4.2. Next, by Assumption 4.4 different intermediate actions are rationalized by different counterpart's signals and thus  $\tilde{\sigma}_k(s_k) \neq \tilde{\sigma}_k(\hat{s}_k)$ . Also, by Assumption 4.1 the optimal action varies for different signals of the counterpart and thus  $a_k^*(s_k, \tilde{\sigma}_k(s_k)) \neq a_k^*(s_k, \tilde{\sigma}_k(\hat{s}_k))$ . Since by Assumption 4.2 there is a unique interim-correct action,  $a_k^*(s_k, \tilde{\sigma}_k(\hat{s}_k)) \neq \tilde{a}_k(s_k) = a_k^*(s_k, \tilde{\sigma}_k(s_k))$  Combining, for  $a_k' = a_k^*(s_k, \tilde{\sigma}_k(\hat{s}_k))$ ,

$$\Delta V_{m_k^a}(s_k, \hat{s}_k) = \underbrace{\varepsilon_k}_{>0} \underbrace{\left( \mathbb{E}_{\pi_k} \left[ v_k(\tilde{a}_k(s_k), s) \right] - \mathbb{E}_{\pi_k} \left[ v_k(a'_k, s) \right] \right)}_{>0} > 0,$$

where the first inequality  $\varepsilon_k > 0$  is ensured by Definition 4.5 of almost-truthful mediation protocols, and the second inequality  $\mathbb{E}_{\pi_k} \left[ v_k(\tilde{a}_k(s_k), s) \right] - \mathbb{E}_{\pi_k} \left[ v_k(a'_k, s) \right] > 0$  is ensured by Definition 4.2 of the interim-optimal action.

Consider now an extended game in which agents simultaneously send reports to an almost-truthful mediation protocol, observing the private signal, receive the mediator's messages back and then

<sup>&</sup>lt;sup>7</sup>Notice that Assumption 4.4 does not add restrictions on the cost functions  $c_k$ , so the result on the impossibility of direct communication with "opposite" preferences still holds with this additional assumption.

simultaneously choose actions. Let  $U_{m_k^a}(s_k, \hat{s}_k)$  be the expectation of agent k's utility conditional on reporting  $\hat{s}_k$  to an almost-truthful mediation protocol  $m_k^a$  when the true signal is  $s_k$ . Let  $C_{m_k^a}(s_k, \hat{s}_k)$  be the expectation of agent k's  $c_k$ -component of the utility conditional on reporting  $\hat{s}_k$  to an almost-truthful mediation protocol  $m_k^a$  when the true signal is  $s_k$ .

For low enough conflict of interest, as parametrized by  $\alpha$ , reporting truthfully is strictly optimal in such an extended game:

**Theorem 4.1.** Suppose that Assumption 4.4 holds. There exists an  $\bar{\alpha}$  such that for all  $\alpha \in (0, \bar{\alpha})$ , for each agent k, signals  $s_k \neq \hat{s}_k$  and almost-truthful mediation protocol  $m_k^a$ ,  $U_{m_k^a}(s_k, s_k) > U_{m_k^a}(s_k, \hat{s}_k)$ .

*Proof.* Define  $\Delta U_{m_{k}^{a}}(s_{k}, \hat{s}_{k})$  and  $\Delta C_{m_{k}^{a}}(s_{k}, \hat{s}_{k})$  analogously to  $\Delta V_{m_{k}^{a}}(s_{k}, \hat{s}_{k})$ . Notice that

$$\Delta U_{m_k^a}(s_k, \hat{s}_k) = \Delta V_{m_k^a}(s_k, \hat{s}_k) - \alpha \times \Delta C_{m_k^a}(s_k, \hat{s}_k)$$

Let

$$\bar{\Delta}V_{m_k^a} = \min_{s_k, \hat{s}_k} \left[ \Delta V_{m_k^a}(s_k, \hat{s}_k) \right]$$

and

$$\bar{\Delta}C_{m_k^a} = \max\left\{0, \max_{s_k, \hat{s}_k} \left[\Delta C_{m_k^a}(s_k, \hat{s}_k)\right]\right\}$$

Notice that  $\overline{\Delta}V_{m_k^a} > 0$  by Lemma 4.4 and  $\overline{\Delta}C_{m_k^a} \ge 0$  by construction. Define

$$\bar{\alpha} = \begin{cases} 1 & \text{if } \forall k \,\bar{\Delta} C_{m_k^a} = 0\\ \min_k \left[ \frac{\bar{\Delta} V_{m_k^a}}{\bar{\Delta} C_{m_k^a}} \right] & \text{if } \exists k \,\bar{\Delta} C_{m_k^a} \neq 0 \end{cases}$$

It remains to notice that for  $\alpha \in (0, \bar{\alpha})$ ,  $\Delta U_{m_k^a}(s_k, \hat{s}_k) > 0$  for every k and  $s_k \neq \hat{s}_k$ , which completes the proof.

Lemma 4.4 and Theorem 4.1 share the main intuition with the illustrative example in Section 3. Lemma 4.4 establishes each agent can optimize her own action by reporting truthfully. The reasons for that are: (i) the agents trust the mediator's message since it is almost always truthful, (ii) the mediators distortions harm a deviating agent more by shifting her action away from the interimoptimal action. Note that, similarly to the illustrative example case, a distorted message  $\tilde{\sigma}_k(\hat{s}_k)$  implicitly encourages an agent who reports truthfully ( $\hat{s}_k = s_k$ ) to select the action that is interimoptimal, since  $a_k^*(s_k, \tilde{\sigma}_k(s_k)) = \tilde{a}_k(s_k)$ . An agent, who reports untruthfully ( $\hat{s}_k \neq s_k$ ), receives an implicit encouragement  $\tilde{\sigma}_k(\hat{s}_k)$  based on an untruthful report, which results in an action other than the interim-optimal one by Assumption 4.4 and Assumption 4.1.

Theorem 4.1 establishes that when the misalignment of interests is small enough, the incentives to report truthfully dominate as optimizing own action is relatively more important than benefiting from a shift in the counterpart's actions. Reiterating on the assumptions required for the results above, the ability of the mediation protocols to shift actions, so that the deviating agent is put to a disadvantage, relies on three details. First, the almost-truthful design of the protocols guarantees that agents' posterior beliefs regarding the counterpart's signal place most of the weight on the mediator's message (Definition 4.5). Second, the deviating and non-deviating behavior need to result in different distorted messages by the mediator, which is guaranteed, if agents' interim beliefs are sufficiently sensitive to their private information (Assumption 4.4). Third, the mediator's messages actually shift agents' decisions, which relies on the assumption that each agent's optimal action is sensitive to the counterpart's information (Assumption 4.1). These three details (a design

feature and two assumptions) are crucial for the mediated communication scheme developed in this paper.

A note of caution regarding the welfare implications of almost-truthful interim-biased mediation is worth pointing out. On the one hand, there clearly exist cases of strict welfare improvements created by such mediation (see the illustrative example case in Section 3.2). On the other hand, given relatively few assumptions on functions v and c, there are also cases when the proposed mechanism does not lead to a welfare improvement.

**Remark 4.1.** Appendix C demonstrates that the above results can be slightly generalized to the settings where agents' payoffs are not fully separable in actions. More concretely, it is shown in Appendix C that almost-truthful interim-biased mediation protocols can facilitate communication when the additional action interaction payoff component  $z_k(a_k, a_{3-k}, s)$  varies little (compared to the main payoff component  $v_k$ ) in response to a change of agent k's action  $a_k$ .

### 4.6 Generalized almost-truthful interim-biased mediation

The almost-truthful mediation can be generalized, allowing the probability  $\varepsilon_k$  of information distortion to depend on the report profile submitted by the agents:

DEFINITION 4.6. Let a *generalized almost-truthful interim-biased mediation protocol* be a collection of random variables

 $m_k^g(\hat{s}_k, \hat{s}_{3-k}) = \begin{cases} \hat{s}_{3-k} & \text{with probability} \quad 1 - \varepsilon_k(\hat{s}_k, \hat{s}_{3-k}) \\ \tilde{\sigma}_k(\hat{s}_k) & \text{with probability} \quad \varepsilon_k(\hat{s}_k, \hat{s}_{3-k}), \end{cases}$ 

such that  $\varepsilon_k(\hat{s}_k, \hat{s}_{3-k}) \in (0, \bar{\varepsilon}_k(\bar{\delta}_k)]$  for all  $\hat{s}_k, \hat{s}_{3-k}$ .

**Remark 4.2.** Appendix *B* finds the optimal mediation protocol for the illustrative example of Section 3 and demonstrates that it belongs to the class of generalized almost-truthful interimbiased mediation protocols. The optimal mediation protocol sends distorted messages relatively more often when the signals reported by the agents are jointly unlikely. This feature leads to stricter informational punishments for deviations, while permitting more accurate information transmission when the agents report truthfully.

### 4.7 On the case of public mediation

By assumption, messages sent by the protocols in the almost-truthful interim-biased class introduced in Definition 4.5 are private: agent k observes the realization of  $m_a^k(\hat{s}_k, \hat{s}_{3-k})$ , while agent (3 - k)does not. An interesting question suggested by a reviewer is whether the mediator's messages can be public. While public announcements can be useful in applications, they also change each agent's information structure and, in fact, do not allow the key result on the existence of an equilibrium with information exchange for sufficiently small misalignment of interests to hold in its current form.

Specifically, if the announcements are public, an almost-truthful interim-biased mediation protocol from Definition 4.5 may, in general, fail to exist. That is, irrespective of how small the probabilities of the mediator's distortions are, agents do not place arbitrarily high weight on the message received from the mediator "directly" and rely on the public message sent to the counterpart instead, at least for some combinations of agents' preferences and exogenous information structures satisfying the earlier assumptions of Section 4.2. Intuitively, it can be impossible to convince an agent that the mediator's message is an undistorted transmission of the counterpart's signal when the mediator's

message to the counterpart is observable and reveals additional information on the counterpart's private signal.

Formally, with public mediation (and other assumptions being the same as in the private mediation case) Lemma 4.2 is no longer true in general, while Definition 4.5 and Lemma 4.3 are ill-posed. Below, these claims are established with the help of the illustrative example case, for which almost-truthful interim-biased protocols do not exist under public mediation.

To see that Lemma 4.2 is no longer true when the mediator's announcements are public, consider the illustrative example. Under private mediation with a sufficiently low distortion probability, agents place sufficiently high beliefs on the mediator's message (see Section 3.2 and notice that weights placed on the message get arbitrarily close to 1 as  $\varepsilon$  converges to 0). If messages are public, however, no matter how small the probability of distortion is, some combinations of agents' signals and mediator's message pairs lead to a situation in which agent k's posterior is fully determined by the message intended for agent (3 - k), while zero weight is placed on the message intended for agent k herself ("direct" message). Specifically, consider agent k and the case when agents' true signals are different: k observes  $s_k$  and (3 - k) observes  $1 - s_k$ . Suppose that both agents report truthfully and that the mediator publicly sends distorted signals to both agents ( $s_k$  is sent to agent k and  $1 - s_k$  is sent to (3 - k)). This event has a non-zero probability as the distortion probability is bounded away from 0. Note that in this situation agent k can actually recover (3 - k)'s true signal and thus disregard the distorted message received. First, k knows that the message sent to (3 - k) was distorted (since it does not coincide with k's own report). Second, by the structure of the mediation protocol, k also knows that the distorted message sent by the mediator to agent (3 - k) perfectly reveals (3 - k)'s true signal, which is  $1 - s_k$ . Thus agent k will place 0 weight on the message  $s_k$  received from the mediator, which shows that Lemma 4.2 no longer holds under public mediation<sup>8</sup>.

Since Lemma 4.2 breaks downs under public announcement of the mediator's messages and probabilities  $\bar{\varepsilon}_k(\delta)$  leading to beliefs arbitrarily close to 1 on the mediator's message may fail to exist, Definition 4.5 (that relies on such probabilities) is ill-posed. Thus public almost-truthful interim-biased mediation protocols do not necessarily exist: agents no longer place high enough beliefs on the "direct" message received from the mediator irrespective of how low the distortion probabilities are. While formally Lemma 4.3 is also ill-posed under public mediation as it relies on the existence of an interim-biased almost-truthful protocol  $m_k^a$ , one could also describe it as being false, since there may no longer exist such small distortion probabilities that lead to the agents acting on the mediator's "direct" messages. To see this, consider the illustrative example case again. Once k knows for sure that (3 - k)'s signal is  $(1 - s_k)$  in the situation described above, k also knows that the correct action is  $\frac{1}{2}(s_k + (1 - s_k)) = \frac{1}{2}$ . It is thus in k's best interest to select action  $\frac{1}{2}$  rather than to choose action  $s_k$  based on the mediator's message and to receive a payoff of 0. Summing up, under public mediation, the agents can no longer simply act on the mediator's "direct" message even when the distortion probabilities are very small.

Since the main result of the paper, Theorem 4.1, relies on Lemma 4.3 (through Lemma 4.4), the existence of an equilibrium in which the agents submit truthful reports is no longer guaranteed under public mediation. This observation implies that the almost-truthful interim-biased mediation protocol class introduced in the present paper *does* rely on the announcements being private. While one could explore whether interim-biased protocols can help information transmission in case of

<sup>&</sup>lt;sup>8</sup>For completeness, note that  $\tilde{\sigma}(s_k) = s_k$  in the illustrative example case while assumptions Assumption 4.1, Assumption 4.2, Assumption 4.3 and Assumption 4.4 are trivially satisfied. Also note that the described argument similarly holds in case when agent k misreports to the mediator.

public announcements without inducing arbitrarily high beliefs on the mediator's messages, such an exploration would require an updated strategy of showing that an equilibrium with information transmission exists. A full exploration of that case is out of scope of the present paper: agents placing high weights on mediator's messages helps to establish the existence of an equilibrium with communication, but the ability to act on the mediator's message in a straightforward manner is arguably a desirable mediation property from the agents' perspective. The arguments above demonstrate that one needs to give up on such a property if one is to explore public announcements in the context of interim-biased mediation protocols with rare distortions.

## 5 Conclusion

This section concludes the paper by discussing the considerations for practical information exchange, providing a summary of the results and outlining directions for future research.

#### 5.1 Some considerations for practical information exchange

The considerations raised by this paper's results may potentially be of relevance when arranging information exchange in practice, specifically, in settings where the misalignment of interests between different parties is impossible (or too costly) to avoid by changing the organizational structure altogether. Two suggestive cases are discussed below, although one should be aware that the real-world complications most likely preclude the direct application of almost-truthful mediation. The first case suggests using its elements at least as a starting point in the context of employees sharing information relevant for project choice decisions. The second case illustrates the difficulty of direct communication between parties with misaligned interests as well as the use of mediation (broadly defined) to relax the corresponding incentive issues.

As mentioned in the introductory example, many organizations try to encourage knowledge sharing among their employees, see, for example, the Harvard Business Review discussion of knowledge sharing by Myers (2015) and the Wall Street Journal article on knowledge hiding by Deal (2018)). Appropriately structured compensation schemes, along with trust, corporate culture and management support have been recognized as important forces that drive information exchange in knowledge-intensive enterprises, see Wang and Noe (2010) for an extensive review. The mediation protocol designed in the present paper suggests another way to facilitate employee knowledge sharing with the help of an organizational leader acting as a communication intermediary for her subordinates. Consider the following simplified real-life setup and the corresponding entities in the model above. A team is comprised of two employees (agents k = 1, 2 in the model) each choosing a project to work on, and a manager (the mediator). The employees are partially informed about the circumstances that affect the success chances of all potential projects (which is captured by the  $s_k$  signals in the model); additionally, this knowledge is often "soft" in the fast-moving work environment and can arguably be modeled as unverifiable information. Each employee needs to choose her project (action  $a_k$ ), which ultimately determines the employee's personal success (the  $v_k$  component of the agent's payoff). Two considerations may affect the employees' willingness to share this information regarding the projects' prospects with the counterpart. First, employees may expect to bear the cost of envy, if their teammate succeeds, see the Harvard Business Review article by Menon and Thompson (2010) who define envy in the workplace context as "the distress people feel when others get what they want". Second, the project choice of the first employee may also result in alternative levels of help / consulting needed from the second employee, resulting in a misalignment of interests: depending on the first employee's project, the second employee will spend different amounts of resources without a direct personal payoff. These two considerations can be captured by the separable  $c_k$  component of the agents' payoffs, since such an adjustment to the main payoff component  $v_k$  can arguably pick up both the envy towards the counterpart as well as the cost of help, both of which largely do not directly interact with employee's own project choice<sup>9</sup>. The almost-truthful interim-biased mediation protocols defined above hint at a strategy of information dissemination that can be used by the manager in this example. In particular, the manager can mediate the information exchange by requesting opinion reports  $(\hat{s}_k)$  on the circumstances affecting the potential projects from both employees for review and privately communicating the review results back. When doing so, the manager can in most cases simply transmit the opinion reports without modifications. In a small share of cases chosen at random, the manager can instead announce to an employee such counterpart's opinion that encourages the optimal employee's project choice implied by her own report (that is, the interim-correct counterpart signal  $\tilde{\sigma}(\hat{s}_k)$  in the model). Theorem 4.1 indicates that (under suitable assumptions discussed in Section 4.2) such behavior may work as the proper mediation tactic creating the incentives for the employees to share truthful opinion reports. A truthful employee actually selects the project that is best given her private information only when facing a (distorted) encouragement. If an employee lies, she acts on the combination of private information and on the encouragement based on her deceitful report, thus her project choice is shifted away from the optimal one given private information only. The shift of the project choice is ensured by the trust towards the manager, who only rarely distorts the transmitted opinion reports. Summing up, while information is most often transferred without any distortions, the rare (implicit) encouragements to select projects that would be optimal conditional on the reported information only, can ensure that the employees who engage in knowledge hiding are put to a disadvantage, thus preventing improper communication. While, given the real-world complications, this procedure likely can't be applied as is, it can at least be used as a starting point to develop more tailored tactics to promote information exchange.

Another consideration raised in the present paper is the difficulty of direct information exchange between the agents with "opposite" preferences such that each one benefits from the counterpart's mistakes (formally demonstrated by Proposition 4.1). To see how similar frictions become relevant for practical information exchange, consider the case of competing intelligence agencies. An example of such a situation is described by a think-tank Council on Foreign Relations (2006) which claims that "fundamental cultural differences and turf wars have long hindered cooperation between the two agencies [FBI and CIA]". This conflict may have contributed to the fact that "agencies did not adequately share relevant counter-terrorism information, prior to September 11", see Finding 9 of the Joint Inquiry into Intelligence Community Activities (2002) (also see the relevant discussion by Garicano and Posner (2005, p. 161)). As a response to this concern, the Intelligence Reform and Terrorism Prevention Act of 2004 established the position of Director of National Intelligence, the responsibilities of whose Office include overseeing the Information Sharing Environment that facilitates exchange of intelligence across various governmental agencies. This decision is, perhaps, a real-world illustration of using mediation as a partial remedy for direct communication difficulties and as such, parallels the use of mediation in this paper (and in related research, e.g., Goltsman et al. (2009)). Of course, it should be noted that, given the model specification, the present paper can

<sup>&</sup>lt;sup>9</sup>Note also that the results in Appendix C ensure that this discussion also applies when the interaction between project decisions can slightly affect employees' payoffs. Specifically, Appendix C demonstrates that almost-truthful interimbiased mediation protocols can also facilitate communication in cases when agents' payoffs are not fully separable in actions, provided that the the payoff changes due to action interaction are relatively small compared to payoff changes induced by the own-action component  $v_k$ .

capture only a small subset of the interactions between intelligence agencies: those that (i) happen in informal contexts such that cheap talk is actually a reasonable approximation to the real-world communication and (ii) are dominated by the direct action effects rather than by the interplay of the actions. Such interactions can occur when (i) agencies are working on parallel cases, (ii) are only communicating informally (e.g. when the security protocols prevent sharing hard evidence, at least in the short term) and (iii) each agency may be subjectively hurt by the rival's success, similarly to the specification of "opposite" preferences in the model,  $c_k(a_{3-k}, s) = v_{3-k}(a_{3-k}, s)$ and paralleling the case of employee envy discussed above. This type of situation is not unheard of: the former head of the U.S. State Department's Bureau of Counterterrorism Nathan A. Sales mentions that "intelligence agencies worry that, if a competitor uses shared data to enhance its analytical products, the credit for any intelligence breakthroughs will go to the recipient rather than the originator", see Sales (2010, p. 309). Again, the present paper formally shows why direct communication in such cases is problematic and establishes that mediation can be of help, mirroring some of the solutions strategies used by the actual intelligence organizations.

## 5.2 Concluding remarks

This paper studies communication between partially informed agents with misaligned interests. For such agents receiving information is desirable, while revealing it may be privately harmful. The paper (i) offers a simple model that captures the main attributes of such a tradeoff; (ii) characterizes the class of almost-truthful interim-biased mediation protocols that enable communication provided that the misalignment of interests between the agents is sufficiently small, while beliefs and actions are sufficiently sensitive to information; (iii) highlights the main leverage that allows communication: almost-truthful information transmission provides the mediator with the opportunity to distort actions in a deliberate manner so that a deviating agent is put to a disadvantage; (iv) discusses considerations raised by the results that may be of relevance in the settings of organizational knowledge hiding and sharing of intelligence information.

Two generalizations of this paper's results can be interesting. First, increasing the number of agents in the model would lead to relaxed truth-telling incentives as individual deviations are easier to detect (under correlated signals). However, it is interesting, whether any adjustments to the mediation design of the present paper (that does not rely on detecting deviations) are needed to extend the truth-telling equilibrium existence result result to multi-player settings. Second, generalizing the class of almost-truthful mediation protocol to a continuous state and action space with an appropriate payoff structure is another potential direction for further research. This would require ensuring that (i) the agents cannot distinguish between accurate and distorted messages of the mediator (e.g. the distorted message cannot be a non-random function of the agent's report); (ii) the deviating agent is shifted away from the interim-optimal action *more*<sup>10</sup> than the truthful agent when distorted messages are sent. While an exploration of these directions is interesting, it is outside of the scope of the present paper, which introduces almost-truthful interim-biased mediation and highlights the corresponding intuition behind the incentives for truthful communication.

<sup>&</sup>lt;sup>10</sup>In the continuous action space both agents' optimal actions conditional on the available information will reflect the possibility that the mediator's message was distorted.

## References

- Alonso, R., W. Dessein, and N. Matouschek (2008). When Does Coordination Require Centralization? *American Economic Review* 98(1), 145–79.
- Ambrus, A., E. M. Azevedo, and Y. Kamada (2013). Hierarchical Cheap Talk. *Theoretical Economics* 8(1), 233 261.
- Ambrus, A. and S. Takahashi (2008). Multi-Sender Cheap Talk with Restricted State Spaces. *Theoretical Economics* 3(1), 1–27.
- Austen-Smith, D. (1993). Interested Experts and Policy Advice: Multiple Referrals under Open Rule. *Games* and Economic Behavior 5(1), 3 43.
- Battaglini, M. (2002). Multiple Referrals and Multidimensional Cheap Talk. *Econometrica* 70(4), 1379 1401.
- Bergemann, D. and S. Morris (2018). Information Design: A Unified Perspective. *Journal of Economic Literature*. Forthcoming.
- Blume, A., O. J. Board, and K. Kawamura (2007). Noisy Talk. Theoretical Economics 2(4), 395 440.
- Council on Foreign Relations (2006, January). FBI and Law Enforcement. Accessed: 2017-11-19. Archived at https://web.archive.org/web/20171119172108/https://www.cfr.org/ backgrounder/fbi-and-law-enforcement.
- Crawford, V. P. and J. Sobel (1982). Strategic Information Transmission. Econometrica 50(6), 1431 1451.
- Deal, J. (2018, August). How Leaders Can Stop Employees from Deliberately Hiding Information. *The Wall Street Journal*. Accessed: 2018-10-14. Archived at https://web.archive.org/web/ 20181015152315/https://blogs.wsj.com/experts/2018/08/14/how-leaders-can-stopemployees-from-deliberately-hiding-information/.
- Farrell, J. and R. Gibbons (1989). Cheap Talk with Two Audiences. *The American Economic Review* 79(5), 1214 1223.
- Forges, F. (1986). An Approach to Communication Equilibria. *Econometrica* 54(6), 1375–1385.
- Gal-Or, E. (1985). Information Sharing in Oligopoly. *Econometrica* 53(2), 329–343.
- Galeotti, A., C. Ghiglino, and F. Squintani (2013). Strategic information transmission networks. *Journal of Economic Theory* 148(5), 1751 1769.
- Garicano, L. and R. A. Posner (2005). Intelligence Failures: An Organizational Economics Perspective. Journal of Economic Perspectives 19(4), 151–170.
- Goltsman, M., J. Hörner, G. Pavlov, and F. Squintani (2009). Mediation, Arbitration and Negotiation. *Journal of Economic Theory* 144(4), 1397 – 1420.
- Goltsman, M. and G. Pavlov (2011). How to Talk to Multiple Audiences. *Games and Economic Behavior* 72(1), 100 122.
- Goltsman, M. and G. Pavlov (2014). Communication in Cournot Oligopoly. *Journal of Economic Theory* 153, 152 176.
- Ivanov, M. (2010). Communication Via a Strategic Mediator. Journal of Economic Theory 145(2), 869 884.
- Kolotilin, A., T. Mylovanov, A. Zapechelnyuk, and M. Li (2017). Persuasion of a Privately Informed Receiver. *Econometrica* 85(6), 1949–1964.
- Krishna, V. and J. Morgan (2001). A Model of Expertise. The Quarterly Journal of Economics 116(2), 747 – 775.
- Krishna, V. and J. Morgan (2004). The Art of Conversation: Eliciting Information from Experts through Multi-Stage Communication. *Journal of Economic Theory* 117(2), 147 – 179.
- Kühn, K.-U. and X. Vives (1995). *Information Exchange Among Firms and Their Impact On Competition*. Luxembourg: Office for Official Publications of the European Communities.
- Li, L. (1985). Cournot Oligopoly with Information Sharing. *The RAND Journal of Economics* 16(4), 521–536.

- Menon, T. and L. Thompson (2010, April). Envy at Work. *Harvard Business Review*. Accessed: 2020-11-30. Archived at https://web.archive.org/web/20211128132300/https://hbr.org/2010/04/ envy-at-work.
- Myers, C. G. (2015, November). Is Your Company Encouraging Employees to Share What They Know? Harvard Business Review. Accessed: 2018-10-14. Archived at https://web.archive.org/web/20181015154019/https://hbr.org/2015/11/is-yourcompany-encouraging-employees-to-share-what-they-know.
- Myerson, R. B. (1982). Optimal Coordination mechanisms in Generalized Principal-Agent Problems. *Journal of Mathematical Economics* 10(1), 67 – 81.
- Myerson, R. B. (1986). Multistage Games with Communication. Econometrica 54(2), 323–358.
- Myerson, R. B. (1997). Game theory: analysis of conflict. Harvard University Press.
- Novshek, W. and H. Sonnenschein (1982). Fulfilled Expectations Cournot Duopoly with Information Acquisition and Release. *The Bell Journal of Economics* 13(1), 214–218.
- Raith, M. (1996). A General Model of Information Sharing in Oligopoly. *Journal of Economic Theory* 71(1), 260–288.
- Sales, N. A. (2010). Share and Share Alike: Intelligence Agencies and Information Sharing. Geo. Wash. L. Rev. 78(2), 279–352.
- Shapiro, C. (1986). Exchange of Cost Information in Oligopoly. *The Review of Economic Studies* 53(3), 433–446.
- U.S. Congress (2004). Intelligence Reform and Terrorism Prevention Act of 2004. 108th Cong., 2d sess., Public Law 108-458. Washington, D.C.: U.S. G.P.O.
- U.S. Congress, Senate, Select Committee on Intelligence. U.S. Congress, House, Permanent Select Committee on Intelligence (2002). Joint Inquiry into Intelligence Community Activities before and after the Terrorist Attacks of September 11, 2001: Report of the U.S. Senate Select Committee on Intelligence and U.S. House Permanent Select Committee on Intelligence together with Additional Views. 107th Cong., 2d sess., S. Rep. 107-351, H. Rep. 107-792. Washington, D.C.: U.S. G.P.O.
- Vida, P. and F. Forges (2013). Implementation of Communication Equilibria by Correlated Cheap Talk: The Two-Player Case. *Theoretical Economics* 8(1), 95–123.
- Vives, X. (1984). Duopoly Information Equilibrium: Cournot and Bertrand. *Journal of Economic Theory* 34(1), 71–94.
- Vives, X. (1990). Trade Association Disclosure Rules, Incentives to Share Information, and Welfare. *The RAND Journal of Economics* 21(3), 409–430.
- Vives, X. (2001). Oligopoly Pricing: Old Ideas and New Tools. Cambridge, Mass. ; London, England: MIT Press, 2001.
- Wang, S. and R. A. Noe (2010). Knowledge Sharing: A Review and Directions for Future Research. *Human Resource Management Review 20*(2), 115 131.
- Ziv, A. (1993). Information Sharing in Oligopoly: The Truth-Telling Problem. The RAND Journal of Economics 24(3), 455–465.

## A Proofs for Section 3

**Proposition A.1.** The only Bayesian Nash Equilibrium of the game  $\Gamma_E$  consists of a strategy profile  $\{a_k^*()\}_{i=1}^2$  with

$$a_k^*(s_k) = s_k$$

Each player's ex-ante expected equilibrium payoff of the game  $\Gamma_E$  is

$$\pi^E = r(1 - \alpha)$$

*Proof.* Note that agent k's best response to each strategy of agent (3 - k) requires selecting the most likely correct action of nature given k's signal. Since

$$\mathbb{P}\left[s^* = s_k | s_k\right] = r,$$

while

$$\mathbb{P}\left[s^* = 1/2|s_k\right] = 1 - r$$

and then since r > 1/2, it must be that in equilibrium each agent selects an action that coincides with her signal and thus indeed  $a_k^*(s_k) = s_k$ .

Given the equilibrium strategies, each agent guesses the the correct action with probability *r* and thus the ex-ante equilibrium payoff of each player is  $\pi^E = r(1 - \alpha)$ .

**Proposition A.2.** Consider a benevolent third party that observes both signals, cares equally about the agents and solves for the first-best

$$\pi^{E,FB} = \max_{\{a_k(s)\}_{i=1}^2} \mathbb{E}_s \left[ \sum_k u_k (a_k, a_{3-k}, s) \right]$$

This problem is solved by  $a_k(s) = s^*$  and accordingly  $\pi^{E,FB} = 1 - \alpha$ .

*Proof.* Note that for each choice of  $\{a_k(s)\}_{i=1}^2$ 

$$\mathbb{E}_{s}\left[\sum_{k}u_{k}\left(a_{k},a_{3-k},s\right)\right] = \sum_{k}(p_{k}-\alpha p_{3-k}) = (1-\alpha)\sum_{k}p_{k},$$
(4)

where  $p_k$  is the unconditional probability of player k guessing the correct action  $s^*$  given  $a_k(s)$ . When  $r \in (1/2, 1)$  and  $\alpha \in (0, 1)$ , the expression in (4) is maximized by  $p_k = 1$  for each k. The only way to achieve  $p_k = 1$  is by setting  $a_k(s) = s^*$ . The corresponding expected payoff of each agent is  $\pi^{E,FB} = 1 - \alpha$ .

**Proposition A.3.** Consider the game  $\Gamma_E$  extended with a finite set  $W_k$  of messages that agent k can send to agent (3 - k) upon observing her private signal. Let  $\Gamma_D$  denote the extended game. In every weak Perfect Bayesian Equilibrium of  $\Gamma_D$  for each message  $w_{3-k} \in W_{3-k}$  player k chooses the action according to her private signal only:  $a_k(s_k, w_{3-k}) = s_k$ .

*Proof.* Note first that according to Proposition A.1, each agent chooses the same action as her signal in the absence of any additional information. Expected equilibrium payoffs of the players in the game with communication take the form

$$\pi_k^D = p_k - \alpha p_{3-k}$$
  
$$\pi_{3-k}^D = p_{3-k} - \alpha p_k,$$

where  $p_k$  is the probability of player k guessing the correct action. Since each player can guarantee herself a correct guess with probability r based on the private information only, it must be that  $p_k, p_{3-k} \ge r$ .

Now suppose that there exists a message  $w_{3-k}$  sent by type  $s_{3-k}$  of agent (3 - k) in equilibrium, such that some type  $s_k$  of agent k chooses action 1/2.

There are two cases: either (i) the type  $1 - s_k$  continues to choose action  $1 - s_k$  upon observing message  $w_{3-k}$  or (ii) the type  $1 - s_k$  chooses action 1/2 upon observing message  $w_{3-k}$ . In both of these cases player (3 - k) can deceive agent k to choose the correct action with probability less than r.

- (i) If type 1 sk continues to choose action 1 sk upon observing message w<sub>3-k</sub>, then by sending message w<sub>3-k</sub> irrespective of her own signal, agent (3 k) induces type sk of agent k to do action 1/2 both in case of coinciding or non-coinciding signals, while type 1 sk continues to act on private information only. Thus in case of (3 k) always sending message w<sub>3-k</sub>, the probability of k guessing the correct action is p̂k = 1/2 < r ≤ pk. Since for a fixed strategy of k the probability of (3 k) guessing the correct action p<sub>3-k</sub> remains constant, agent (3 k) has a profitable deviation.
- (ii) If type  $1 s_k$  chooses action 1/2 upon observing message  $w_{3-k}$ , then by sending message  $w_{3-k}$  irrespective of her signal, agent (3 k) induces all types of agent k to do action 1/2 and the probability of k guessing the correct action in case of such a deviation is  $\hat{p}_k = 1 r < r \le p_k$ . Again, since for a fixed strategy of k the probability of (3 k) guessing the correct action  $p_{3-k}$  remains constant, agent (3 k) has a profitable deviation.

Thus if there exists a message  $w_{3-k}$  that induces some type of agent k to choose action 1/2, agent (3 - k) necessarily has a profitable deviation. Therefore, it must be that in every every weak Perfect Bayesian Equilibrium of the game  $\Gamma_D$  with direct communication, each agent acts on her private information only.

## **B** Optimal mediation for the illustrative example

This appendix section finds the optimal mediation protocol for the illustrative example of Section 3. Recall that the baseline case of the example consists of the game  $\Gamma_E$ . In the game, each agent  $k \in \{1, 2\}$  obtains a binary signals  $s_k \in S = \{0, 1\}$  with the following joint distribution  $\mathbb{P}$  over  $S^2 = S \times S$  parametrized by  $r \in (1/2, 1)$ :

$$\frac{\pi}{s_2 = 0 \quad s_2 = 1}$$

$$s_1 = 0 \quad \frac{r}{2} \quad \frac{1 - r}{2}$$

$$s_1 = 1 \quad \frac{1 - r}{2} \quad \frac{r}{2}$$

Together, these signals determine the correct action

$$s^* = \frac{1}{2} \sum_k s_k.$$

Both agents would like to guess  $s^*$  by choosing an action in the set  $\mathcal{A}_k = \{0, 1/2, 1\}$ . The agents have a conflict of interest and prefer the opponent not to be able to guess the correct action. The payoffs representing such preferences are given by

$$u_k(a, s) = \mathbb{1} \{a_k = s^*\} - \alpha \times \mathbb{1} \{a_{3-k} = s^*\}$$

where  $\alpha \in (0, 1)$ .

Consider now a mediated game. That is, a mediation protocol is introduced that receives reports from the agents and sends messages back to the agents. For each possible report profile received from the agents the protocol specifies a distribution on the messages sent back to the agents.

Due to the revelation principle (see Myerson (1982, 1986) and Forges (1986)), attention can be restricted to *direct revelation* mediation protocols that take the type-reports from the agents and send back action recommendations to each player. A direct revelation mediation protocol M is defined as a function from the product type space into the joint distributions over the action recommendations  $M : S \rightarrow \Delta(A \times A)$ . That

is, each agent is asked to submit her signal received and conditional on the pair of reports is advised on an action, possibly in a random way. The mediation protocol should be such that the agents find it optimal to report their true types and follow the recommended action conditional on the other player reporting truthfully and following recommendations<sup>11</sup>.

Let  $\mathcal{M}$  be the set of all *incentive-compatible* mediation protocols. Let  $\pi_k^M$  be the ex-ante expected equilibrium payoff of agent k in game  $\Gamma_M$  with two agents, signals jointly distributed according to  $\mathbb{P}$  in (B), the correct action  $s^*$  determined as in (B), payoff-relevant actions be  $A_k = \{0, 1/2, 1\}$  and the mediation protocol  $M \in \mathcal{M}$ . Consider now the problem of optimal mediation protocol design faced by competitors

$$\max_{M \in \mathcal{M}} \left[ \sum_{k} \pi_{k}^{M} \right]$$
(5)

The following sequence of lemmas first simplifies the problem in (5) and leads to Theorem B.1 which presents the optimal mediation protocol.

Lemma B.1 below shows that, when solving (5), without loss of generality one can consider only mediation protocols that generate independent recommendations conditional on the pair of reports.

LEMMA B.1. For every  $M \in \mathcal{M}$  there exists  $M' \in \mathcal{M}$  such that

- (i)  $M' : S \to \Delta(A) \times \Delta(A)$
- (ii)  $M(\hat{s})$  and  $M'(\hat{s})$  have the same marginal distributions for every  $\hat{s}$
- (iii)  $\pi_k^{M'} = \pi_k^M$  for every k

*Proof.* Consider an incentive-compatible mediation protocol  $M \in \mathcal{M}$ . Define M' to be a mapping from S to  $\Delta(A) \times \Delta(A)$  such that for each  $\hat{s} \in S$ , for each vector of recommendations *a* 

$$M'(\hat{s})(a) = \prod_{k} \left( \sum_{\hat{a}_{3-k}} M(\hat{s})(a_k, \hat{a}_{3-k}) \right)$$

That is, M' is defined to be the product of the marginal distributions of M for each pair of reports s. By construction, (i)  $M' \in \Delta(A) \times \Delta(A)$  and (ii) M and M' have the same marginal distributions.

Now, since  $M \in M$ , it must also be that  $M' \in M$ . To see this observe first that for type  $s_k$  of agent k that submitted report  $\hat{s}_k$  and received a recommendation  $\hat{a}_k$ , agent k's posterior over the signal of agent (3 - k)[and thus also k's preferred action] is pinned down by the marginal of recommendation distributions. Thus for each report  $\hat{s}_k$  of type  $s_k$  of agent k under M', the resulting distribution of k's actions is the same as under M and consequently the probability of type  $s_k$  agent k making a correct guess is the same under M and M'. Similarly, for each report  $\hat{s}_k$  of agent k under M', the resulting distribution of (3 - k)'s actions is the same as under M. Since the distribution of both agents' actions pin down the expected payoff of agent k, type  $s_k$  of agent k has the same expected payoff for each report  $\hat{s}_k$  under M and M'. By assumption there were no profitable deviations from truth-telling under M, thus so is the case under M'. It is therefore proved that  $M' \in M$  and (iii)  $\pi_k^{M'} = \pi_k^M$ .

Exploiting Lemma B.1, one can restrict attention to mediation protocols with independent action recommendations when solving for the optimal communication protocol. Thus from now on  $\mathcal{M}$  is redefined to be the set of IC mediation protocols with independent action recommendations.

Next, Lemma B.2 shows that none of the IC mediation protocols recommend action  $1 - \hat{s}_k$  to agent k who reported type  $\hat{s}_k$ . That is, an action that is known to be incorrect by a particular agent is never recommended to this agent.

<sup>&</sup>lt;sup>11</sup>The solution concept corresponds to *information design with elicitation* in terminology of Bergemann and Morris (2018).

LEMMA B.2. For every  $M \in \mathcal{M}$  it must be that  $M(\hat{s}) \in \times_k \Delta(\{\hat{s}_k, 1/2\})$ .

*Proof.* Suppose that for some  $M \in \mathcal{M}$  action  $1 - \hat{s}_k$  is recommended to agent with a truthful report  $\hat{s}_k$ . Note that following such a recommendation yields a probability 0 of agent k guessing  $s^*$ . Obviously, agent k can deviate from following the recommendation, act on her private information only and guarantee herself at least a probability r of guessing  $s^*$ .

Every  $M \in \mathcal{M}$  is now completely summarized by a vector in  $[0,1]^8$  with typical values  $m_{i,j}^k \in [0,1]$ presented in the following table:

	М	
	$\hat{s}_2 = 0$	$\hat{s}_2 = 1$
$\hat{s}_1 = 0$	$m_{00}^1, m_{00}^2$	$m_{01}^1, m_{10}^2$
$\hat{s}_1 = 1$	$m_{10}^1, m_{01}^2$	$m_{11}^1, m_{11}^2$

In the table  $m_{i,j}^k$  is the probability of recommending the correct action to agent k, when agent k reported i and agent (3 - k) reported  $j^{12}$ . The vector m consisting of  $m_{i,j}^k \in [0, 1]$  is *feasible* if the corresponding M belongs to  $\mathcal{M}$ .

Lemma B.3 below simplifies the problem in (5) even further, stating that there is no loss in optimizing the weighted sum of payoffs with respect to just two variables: the probability of recommending a correct action in case of (i) coinciding and (ii) non-coinciding reports.

To prove Lemma B.3 the following claim is first established.

<u>Claim:</u> A mediation protocol *M* defined by a vector  $\left\{ \left\{ m_{i,j}^k \right\}_{i,j=0,0}^{1,1} \right\}_{k=1}^2$  belongs to the set of IC mediation protocols  $\mathcal{M}$  if and only if the following two sets of condition hold. First, for each k and i

$$m_{ii}^k r \ge (1 - m_{i,1-i}^k)(1 - r) \tag{6}$$

$$m_{i,1-i}^k(1-r) \ge (1-m_{i,i}^k)r$$
 (7)

Second, for each k and i

$$\begin{split} &(m_{ii}^{k} - \alpha m_{ii}^{3-k})r + (m_{i,1-i}^{k} - \alpha m_{1-i,i}^{3-k})(1-r) \geqslant \\ &\begin{cases} r(1 - \alpha(1 - m_{i,1-i}^{3-k})) + (1-r)(0 - \alpha(1 - m_{1-i,1-i}^{3-k})) & \text{if } (1 - m_{1-i,i}^{k})r > m_{1-i,1-i}^{k}(1-r) \\ r(m_{1-i,i}^{k} - \alpha(1 - m_{i,1-i}^{3-k})) + (1-r)(m_{1-i,1-i}^{k} - \alpha(1 - m_{1-i,1-i}^{3-k})) & \text{if } (1 - m_{1-i,i}^{k})r \leqslant m_{1-i,1-i}^{k}(1-r) \end{split}$$

Conditions (6)-(7) ensure that each agent finds it profitable to follow the mediation protocol's recommendation. Conditions (8) ensure that there are no profitable deviations from reporting truthfully to the mediation protocol.

Proof: The two sets of IC conditions (6)-(7) and (8) are established separately.

*"Following recommendation" conditions* Under an IC direct revelation mediation protocol, agent k finds it profitable to follow the mediation protocol's recommendation. In the problem at hand, one needs to establish conditions under which this is the case for every signal and recommendation obtained by agent k.

<sup>&</sup>lt;sup>12</sup>The correct action recommendation to agent k in this case is (i+j)/2 in case of truthful reporting, while the incorrect action recommendation to agent k is ((i + (1 - j))/2).

• Suppose agent k has signal i and obtained a recommendation R to do i. Such a recommendation can occur when the reported pair of types is either (i, i) or (i, 1 - i). The posterior probability that the state is (i, i) (and thus that the correct action is indeed i) is equal to

$$\mathbb{P}\left[s = (i,i)|R = i, s_k = i\right] = \mathbb{P}\left[s_{3-k} = i|R = i, s_k = i\right]$$
$$= \frac{\mathbb{P}\left[s_k = i, s_{3-k} = i, R = i\right]}{\mathbb{P}\left[s_k = i, R = i\right]}$$
$$= \frac{m_{ii}^k \cdot \frac{r}{2}}{m_{ii}^k \cdot \frac{r}{2} + (1 - m_{i,1-i}^k) \cdot \frac{1-r}{2}}$$

Agent k following recommendation means that s = (i, i) is more likely than s = (i, 1-i), which yields the first IC constraint

$$m_{ii}^k r \ge (1 - m_{i,1-i}^k)(1 - r)$$

and condition (6) is established.

• Suppose now agent k has signal i and obtained a recommendation R to do 1/2. Such a recommendation can occur when the reported pair of types is either (i, i) or (i, 1 - i). The posterior probability that the state is (i, 1 - i) (and thus the correct action is indeed 1/2) is equal to

$$\mathbb{P}\left[s = (i, 1 - i)|R = 1/2, s_k = i\right] = \mathbb{P}\left[s_{3-k} = 1 - i|R = 1/2, s_k = i\right]$$
$$= \frac{\mathbb{P}\left[s_k = i, s_{3-k} = 1 - i, R = 1/2\right]}{\mathbb{P}\left[s_k = i, R = 1/2\right]}$$
$$= \frac{m_{i,1-i}^k \cdot \frac{1-r}{2}}{m_{i,1-i}^k \cdot \frac{1-r}{2} + (1 - m_{ii}^1) \cdot \frac{r}{2}}$$

Agent 1 following recommendation means that (i, 1 - i) is more likely than (i, i), which yields the second IC constraint

$$m_{i,1-i}^k(1-r) \ge (1-m_{ii}^k)r$$

and condition (7) is established.

The "following recommendation" conditions are thus established.

"Truthful reporting" conditions Under an IC direct revelation mediation protocol, agent k finds it profitable to report truthfully to the mediation protocol and follow the recommendation rather than misreporting and doing some other action upon receiving a recommendation. In the problem at hand, one needs to establish conditions under which this is the case for every signal obtained by agent k.

• Suppose agent k received a signal i. If she reports [T]ruthfully and follows the recommendation by the mediation protocol (and so does agent (3 - k)), the expected payoff is

$$\pi_{k}^{T} = \mathbb{E}\left[\pi_{k}^{T}|s_{k}=i\right] = \mathbb{E}\left[\pi_{k}^{T}|s_{k}=i, s_{3-k}=i\right] \mathbb{P}\left[s_{3-k}=i|s_{k}=i\right] \\ + \mathbb{E}\left[\pi_{k}^{T}|s_{k}=i, s_{3-k}=1-i\right] \mathbb{P}\left[s_{3-k}=1-i|s_{k}=i\right] \\ = (m_{ii}^{k} - \alpha m_{ii}^{3-k})r + (m_{i,1-i}^{k} - \alpha m_{1-i,i}^{3-k})(1-r)$$

• If agent k misreports and sends 1 - i to the mediation protocol instead of i, it is possible to hear two recommendations in response: 1 - i or 1/2. The optimal actions in each of these cases are established below:

- What is the optimal action if 1 - i is recommended back by the mediation protocol? The conditional probability of agent (3 - k) having a signal *i* is equal to

$$\mathbb{P}_{U}\left[s_{3-k}=i|R=1-i, s_{k}=i\right] = \frac{\mathbb{P}_{U}\left[s_{k}=i, s_{3-k}=i, R=1-i\right]}{\mathbb{P}_{U}\left[R=1-i, s_{k}=i\right]}$$
$$= \frac{(1-m_{1-i,i}^{k}) \cdot \frac{r}{2}}{(1-m_{1-i,i}^{k}) \cdot \frac{r}{2} + m_{1-i,1-i}^{k} \cdot \frac{1-r}{2}},$$

where  $\mathbb{P}_U[]$  stands for the updated probabilities given an [U]ntruthful report. Similarly, the conditional probability of agent (3 - k) having a signal 1 - i is equal to

$$\mathbb{P}_{U}\left[s_{3-k} = 1 - i | R = 1 - i, s_{k} = i\right] = \frac{m_{1-i,1-i}^{k} \cdot \frac{1-r}{2}}{(1 - m_{1-i,i}^{k}) \cdot \frac{r}{2} + m_{1-i,1-i}^{k} \cdot \frac{1-r}{2}}$$

Consequently, if  $(1 - m_{1-i,i}^k)r > m_{1-i,1-i}^k(1 - r)$ , agent k will do action i upon receiving signal 1 - i and if  $(1 - m_{1-i,i}^k)r \le m_{1-i,1-i}^k(1 - r)$ , agent k will do action 1/2 upon receiving signal 1 - i.

- What is the optimal action if 1/2 is recommended back by the mediation protocol? The conditional probability of agent (3 - k) having a signal *i* is equal to

$$\mathbb{P}_{U}\left[s_{3-k}=i|R=1/2, s_{k}=i\right] = \frac{\mathbb{P}_{U}\left[s_{3-k}=i, s_{k}=i, R=1/2\right]}{\mathbb{P}_{U}\left[R=1/2, s_{k}=i\right]}$$
$$= \frac{m_{1-i,i}^{k} \cdot \frac{r}{2}}{m_{1-i,i}^{k} \cdot \frac{r}{2} + (1-m_{1-i,1-i}^{1}) \cdot \frac{1-r}{2}},$$

while the conditional probability of agent (3 - k) having a signal 1 - i is equal to

$$\mathbb{P}_{U}\left[s_{3-k}=1-i|R=1/2, s_{k}=i\right] = \frac{(1-m_{1-i,1-i}^{k})\cdot\frac{1-r}{2}}{m_{1-i,i}^{k}\cdot\frac{r}{2}+(1-m_{1-i,1-i}^{k})\cdot\frac{1-r}{2}},$$

Note that the IC constraint (7) and Assumption 1 imply that  $m_{1-i,i}^k r > (1 - m_{1-i,1-i}^k)(1 - r)$  and thus agent *k* prefers to do action *i* upon hearing a recommendation of 1/2 from the mediation protocol.

• Now [after trivially calculating the probabilities of agent (3 - k) making the correct guess], the expected payoff of agent k in case of misreporting and sending 1 - i instead of i can be computed. If  $(1 - m_{1-i,i}^k)r > m_{1-i,1-i}^k(1 - r)$ , agent k does i in any case and gets an expected utility of

$$\pi_k^{U,>} = r(1 - \alpha(1 - m_{i,1-i}^{3-k})) + (1 - r)(0 - \alpha(1 - m_{1-i,1-i}^{3-k}))$$

If  $(1 - m_{1-i,i}^k)r \le m_{1-i,1-i}^k(1 - r)$ , agent k does 1/2 in case of hearing a recommendation of 1 - i and i in case of recommendation of 1/2 and gets an expected utility of

$$\begin{split} \pi_k^{U,\leqslant} &= r(m_{1-i,i}^l \cdot 1 + (1 - m_{1-i,i}^k) \cdot 0 - \alpha(1 - m_{i,1-i}^k)) \\ &+ (1 - r)(m_{1-i,1-i}^k \cdot 1 + (1 - m_{1-i,1-i}^l) \cdot 0 - \alpha(1 - m_{1-i,1-i}^{3-k})) \\ &= r(m_{1-i,i}^k - \alpha(1 - m_{i,1-i}^{3-k})) + (1 - r)(m_{1-i,1-i}^k - \alpha(1 - m_{1-i,1-i}^{3-k})) \end{split}$$

⊲

Thus the IC constraint for agent k reporting truthfully upon observing signal i is

$$\begin{split} &(m_{ii}^k - \alpha m_{ii}^{3-k})r + (m_{i,1-i}^k - \alpha m_{1-i,i}^{3-k})(1-r) \geqslant \\ & \left\{ \begin{array}{ll} r(1 - \alpha(1 - m_{i,1-i}^{3-k})) + (1-r)(0 - \alpha(1 - m_{1-i,1-i}^{3-k})) & \text{ if } (1 - m_{1-i,i}^k)r > m_{1-i,1-i}^k(1-r) \\ r(m_{1-i,i}^k - \alpha(1 - m_{i,1-i}^{3-k})) + (1-r)(m_{1-i,1-i}^k - \alpha(1 - m_{1-i,1-i}^{3-k})) & \text{ if } (1 - m_{1-i,i}^k)r \leqslant m_{1-i,1-i}^k(1-r) \\ \end{split} \right. \end{split}$$

and "truthful reporting" conditions (8) are established. The proof of the claim is now completed.

LEMMA B.3. There exists a solution to the optimal mediation protocol design problem in eq. (5) such that  $\forall k, i \ m_{i,i}^k = p \in [0, 1] \text{ and } m_{i,1-i}^k = q \in [0, 1].$ 

Proof: Note first that the optimal mediation protocol design problem

$$\max_{M \in \mathcal{M}} \left[ \sum_{i} \pi_{i}^{M} \right]$$
(9)

can be written more explicitly as

.

$$\max_{\substack{m_{i,1-i}^{k} \\ m_{i,1-i}^{k}}} \left[ \frac{r}{2} \sum_{i} \sum_{k} \left( m_{ii}^{k} - \alpha m_{ii}^{3-k} \right) + \frac{1-r}{2} \sum_{i} \sum_{k} \left( m_{i,1-i}^{k} - \alpha m_{i,1-i}^{3-k} \right) \right]$$
(10)  
s.t. (6)-(8)

To prove the lemma, one needs to notice that (i) the constraint set defined by IC conditions is convex; (ii) the objective and the constraint set are symmetric with respect to players and states.

Convexity of the constraint set (6)-(8) To establish (i) one can show first that each IC condition defines a convex set. This is obvious for conditions (6)-(7) as these a linear in parameters. Now consider the typical truthful-reporting IC constraint

$$\begin{split} (m_{ii}^{k} - \alpha m_{ii}^{3-k})r + (m_{i,1-i}^{k} - \alpha m_{1-i,i}^{3-k})(1-r) \geqslant \\ \begin{cases} r(1 - \alpha (1 - m_{i,1-i}^{3-k})) + (1-r)(0 - \alpha (1 - m_{1-i,1-i}^{3-k})) & \text{if } (1 - m_{1-i,i}^{k})r > m_{1-i,1-i}^{k}(1-r) \\ r(m_{1-i,i}^{k} - \alpha (1 - m_{i,1-i}^{3-k})) + (1-r)(m_{1-i,1-i}^{k} - \alpha (1 - m_{1-i,1-i}^{3-k})) & \text{if } (1 - m_{1-i,i}^{k})r \le m_{1-i,1-i}^{k}(1-r) \\ \end{cases}$$
(11)

and consider two vectors t and s that satisfy those constraints.

If both t and s are such that  $\left[ (1 - t_{1-i,i}^k)r > l_{1-i,1-i}^k (1-r) \text{ and } (1 - s_{1-i,i}^k)r > s_{1-i,1-i}^k (1-r) \right]$  or  $\left[ (1 - t_{1-i,i}^k)r > s_{1-i,1-i}^k (1-r) \right]$  or  $\left[ (1 - t_{1-i,i}^k)r > s_{1-i,1-i}^k (1-r) \right]$  $\leq l_{1-i,1-i}^k(1-r)$  and  $(1-s_{1-i,i}^k)r \leq s_{1-i,1-i}^k(1-r)$ , then the convex combination  $u = \beta t + (1-\beta)s, \beta \in [0,1]$ also satisfies the constraint (11) with  $(1 - u_{1-i,i}^k)r > u_{1-i,1-i}^k(1-r)$  or  $(1 - u_{1-i,i}^k)r \le u_{1-i,1-i}^k(1-r)$  respectively by linearity of both RHS and LHS of the inequality in (11). Now suppose that

$$\begin{cases} (1 - t_{1-i,i}^{k})r > l_{1-i,1-i}^{k}(1 - r) \\ (t_{ii}^{k} - \alpha t_{ii}^{3-k})r + (t_{i,1-i}^{k} - \alpha t_{1-i,i}^{3-k})(1 - r) \ge r(1 - \alpha(1 - t_{i,1-i}^{3-k})) + (1 - r)(0 - \alpha(1 - t_{1-i,1-i}^{3-k})) \\ (1 - s_{1-i,i}^{k})r \le s_{1-i,1-i}^{k}(1 - r) \\ (s_{ii}^{k} - \alpha s_{ii}^{3-k})r + (s_{i,1-i}^{k} - \alpha s_{1-i,i}^{3-k})(1 - r) \ge r(s_{1-i,i}^{k} - \alpha(1 - s_{i,1-i}^{3-k})) + (1 - r)(s_{1-i,1-i}^{k} - \alpha(1 - s_{1-i,1-i}^{3-k})) \end{cases}$$

and let  $u = \beta t + (1 - \beta)s$  with  $\beta \in [0, 1]$ . First, either  $(1 - u_{1-i,i}^k)r > u_{1-i,1-i}^k(1-r)$  or  $(1 - u_{1-i,i}^k)r \le u_{1-i,1-i}^k(1-r)$ . Suppose  $(1 - u_{1-i,i}^k)r > u_{1-i,1-i}^k(1-r)$ .  $u_{1-i,1-i}^{k}(1-r)$ , then to show convexity of the set defined by constraint (11), one needs to establish that

$$(s_{ii}^{k} - \alpha s_{ii}^{3-k})r + (s_{i,1-i}^{k} - \alpha s_{1-i,i}^{3-k})(1-r) \ge r(1 - \alpha(1 - s_{i,1-i}^{3-k})) + (1-r)(0 - \alpha(1 - s_{1-i,1-i}^{3-k}))$$
(12)

Indeed, if this is the case, then also

$$(u_{ii}^k - \alpha u_{ii}^{3-k})r + (u_{i,1-i}^k - \alpha u_{1-i,i}^{3-k})(1-r) \ge r(1 - \alpha(1 - u_{i,1-i}^{3-k})) + (1 - r)(0 - \alpha(1 - u_{1-i,1-i}^{3-k}))$$

since u is a convex combination of t and s.

To show (12) note that

$$\begin{split} (s_{ii}^k - \alpha s_{ii}^{3-k})r + (s_{i,1-i}^k - \alpha s_{1-i,i}^{3-k})(1-r) &\geq r(s_{1-i,i}^k - \alpha (1-s_{i,1-i}^{3-k})) + (1-r)(s_{1-i,1-i}^k - \alpha (1-s_{1-i,1-i}^{3-k})) \\ &\geq r(1-\alpha (1-s_{i,1-i}^{3-k})) + (1-r)(0-\alpha (1-s_{1-i,1-i}^{3-k})), \end{split}$$

where the last inequality follows from  $(1-s_{1-i,i}^k)r \leq s_{1-i,1-i}^k(1-r)$ . The case of  $(1-u_{1-i,i}^k)r \leq u_{1-i,1-i}^k(1-r)$  is similar.

Since the intersection of convex sets is convex, the constraint set is convex itself.

Symmetry of the objective (10) and the constraint set (6)-(8) Note that if the maximization problem is solved by some vector  $m = (m^1, m^2)$  (where  $m^i$  denotes the subvector of probabilities related to agent *i*), then vector  $m' = (m^2, m^1)$  leads to the same value of the objective function and constraints are satisfied at *m'* by symmetry. Thus the same value of the objective is achieved at the average of *m*, *m'* and one can restrict attention to maximizing with 4-element vector  $m_{00}, m_{01}, m_{10}, m_{11}$ . Again, swapping  $m_{00}$  with  $m_{11}$ and  $m_{01}$  with  $m_{10}$  leads to the same value of the modified objective function and constraints being satisfied by symmetry. Thus one can restrict attention to maximizing with respect to 2-element vector (p, q) with  $p = m_{ii}$  and  $q = m_{i,1-i}$  and the proof of the lemma is now completed.

Utilizing the results from the preceding lemmas, Theorem B.1 provides an explicit solution to the problem of designing the optimal mediation protocol.

**Theorem B.1.** The optimal mediation protocol design problem in eq. (5) is solved by  $M^* \in \mathcal{M}$  such that  $\forall k, i \ (m^*)_{i,i}^k = p^* \text{ and } (m^*)_{i,1-i}^k = q^* \text{ with }$ 

$$\begin{pmatrix} p^* = 1, & q^* = \frac{2r - 1}{2r - 1 + \alpha} \end{pmatrix}, & if \ \alpha \le 1 - r \\ (p^* = 1, & q^* = 0), & if \ \alpha > 1 - r \end{cases}$$

Proof. Due to Lemma B.3 the optimal mediation protocol design problem is reduced to

$$\max_{p,q \in [0,1]^2} [pr + q(1 - r)]$$
subject to
$$pr \ge (1 - q)(1 - r) \tag{13}$$

$$q(1 - r) \ge (1 - p)r \tag{14}$$

$$(p - \alpha p)r + (q - \alpha q)(1 - r) \ge \begin{cases} r(1 - \alpha(1 - q)) + (1 - r)(0 - \alpha(1 - p)) & \text{if } (1 - q)r > p(1 - r) \\ r(q - \alpha(1 - q)) + (1 - r)(p - \alpha(1 - p)) & \text{if } (1 - q)r \le p(1 - r) \end{cases}$$

$$(15)$$

For a moment, ignore the first two constraints (13)-(14). It will be verified later that the solution of the relaxed maximization problem still satisfies these two constraints.

• Note that the value of the objective grows with p along the line (1-q)r = p(1-r) under Assumption 1. Indeed, substituting

$$q = 1 - p \frac{1 - r}{r}$$

into the objective yields coefficient equal to 2 - 1/r on the variable p and thus the objective grows in p. Having observed this, it is easy to see that the value of the objective in the region with (1-q)r > p(1-r) is not higher than at the point point on its boundary with the highest value of p, which is  $\left(1, \frac{2r-1}{r}\right)$ .



Figure 1: Typical constraint sets.

• Now also note that the objective grows with p along the line

$$(p-\alpha p)r+(q-\alpha q)(1-r)=r(q-\alpha(1-q))+(1-r)(p-\alpha(1-p))$$

Indeed, substituting

$$q = p \frac{1 - 2r + \alpha}{1 - 2r - \alpha} - \frac{\alpha}{1 - 2r - \alpha}$$

into the objective yields coefficient  $(\alpha + 1)(2r - 1)/(\alpha + 2r - 1) > 0$  on p. Thus if the point  $\left(1, \frac{2r-1}{2r+\alpha-1}\right)$  satisfies  $(1-q)r \le p(1-r)$ , it has the highest value of the objective in the region with  $(1-q)r \le p(1-r)$ .

- Note that α ≤ 1 − r simultaneously guarantees that (1, <sup>2r-1</sup>/<sub>2r+α-1</sub>) satisfies (1 − q)r ≤ p(1 − r) and has a higher value of the objective than (1, <sup>2r-1</sup>/<sub>r</sub>). Moreover, the point (1, <sup>2r-1</sup>/<sub>2r+α-1</sub>) satisfies the two omitted constraints (13)-(14) of the maximization problem. Thus the point that maximizes the objective for α ≤ 1 − r is (1, <sup>2r-1</sup>/<sub>2r+α-1</sub>).
- If in turn α > 1-r, then there are no points that satisfy the constraint with (1-q)r < p(1-r) : the set defined by (p αp)r + (q αq)(1-r) ≥ r(q α(1-q)) + (1-r)(p α(1-p)) has no intersection with the set defined by (1 q)r < p(1 r), which is easy to verify by comparing the values of the linear constraints at the boundary points of the constraint set p = 0 and p = 1. Moreover, the only point satisfying the constraint with (1-q)r ≥ p(1-r) is (1,0), which remains the only candidate for the optimal mediation protocol when α > 1 r. This point obviously satisfies the omitted constraints (13)-(14).
- For a graphical treatment, two typical constraint sets are shown in Figure 1.

The search for the optimal mediation protocol is thus completed

$$\begin{pmatrix} p^* = 1, & q^* = \frac{2r - 1}{2r - 1 + \alpha} \end{pmatrix}, & \text{if } \alpha \le 1 - r \\ (p^* = 1, & q^* = 0), & \text{if } \alpha > 1 - r \\ \end{cases}$$

and Theorem B.1 is proved.

**Remark B.1.**  $M^*$  is the unique solution of the optimal mediation protocol design problem in the set M of protocols with independent action recommendations.

*Proof.* Table 1 below explicitly presents  $M^*$  together with its gradient and the gradients of the 8 constraints that bind at the maximum found in Theorem B.1. The binding constraints are (i)  $\left\{\left\{m_{ii}^k \leq 1\right\}_{i=0}^{1}\right\}_{k=1}^{2}$ ; (ii) the "truthful reporting" constraints in (8).

$m_{ij}^k \mid M^* \mid$	∇ of obj.	i I				∇ of	constraints		
I I I I		l I	$m_{ii}^k$	≤ 1			"truthful	reporting"	
$m_{00}^1 + p^* +$	$r(1-\alpha)$	1	0	0	0	- <i>r</i>	$\alpha r$	1 – <i>r</i>	$\alpha(1-r)$
$m_{01}^{1} + q^{*} + ($	$(1-r)(1-\alpha)$	0	0	0	0	-(1-r)	$\alpha r$	r	$\alpha(1-r)$
$m_{10}^1 + q^* + 0$	$(1-r)(1-\alpha)$	0	0	0	0	r	$\alpha(1-r)$	-(1-r)	$\alpha r$
$m_{11}^{1^{\circ}} p^{*}$	$r(1-\alpha)$	0	1	0	0	1 - r	$\alpha(1-r)$	-r	$\alpha r$
$m_{00}^2 + p^* +$	$r(1-\alpha)$	0	0	1	0	$\alpha r$	-r	$\alpha(1-r)$	1 – <i>r</i>
$m_{01}^2 + q^* + ($	$(1-r)(1-\alpha)$	0	0	0	0	$\alpha r$	-(1-r)	$\alpha(1-r)$	r
$m_{10}^2 + q^* + ($	$(1-r)(1-\alpha)$	0	0	0	0	$\alpha(1-r)$	r	$\alpha r$	-(1-r)
$m_{11}^2 + p^* +$	$r(1-\alpha)$	0	0	0	1	$\alpha(1-r)$	1 – <i>r</i>	$\alpha r$	-r

<b>Table 1.</b> Value of objective, gradients of objective and constraints
--

Now note that the gradient of the objective is a linear combination of the gradients of the constraints with weights (-(2-1)(-2-1))

$$\lambda = \begin{pmatrix} -\frac{(2r-1)(\alpha^{2}-1)}{2r+\alpha-1} \\ -\frac{(2r-1)(\alpha^{2}-1)}{2r+\alpha-1} \\ -\frac{(2r-1)(\alpha^{2}-1)}{2r+\alpha-1} \\ -\frac{(2r-1)(\alpha^{2}-1)}{2r+\alpha-1} \\ \frac{(r-1)(\alpha-1)}{2r+\alpha-1} \\ \frac{(r-1)(\alpha-1)}{2r+\alpha-1} \\ \frac{(r-1)(\alpha-1)}{2r+\alpha-1} \\ \frac{(r-1)(\alpha-1)}{2r+\alpha-1} \end{pmatrix},$$
(16)

that are strictly positive when  $r \in (1/2, 1)$  and  $\alpha \in (0, 1)$ . Thus  $M^*$  is locally the unique solution to the maximization problem that defines the optimal mediation protocol. Due to the convexity of the constraint set and the linearity of the objective,  $M^*$  is also the unique global solution in the set  $\mathcal{M}$  of mediation protocols with independent action recommendations.

**Remark B.2.** It is easy to see that the mediation protocol of Theorem B.1 can be implemented by a generalized almost-truthful interim-biased mediation protocol with

$$\varepsilon_k(\hat{s}_k, \hat{s}_{3-k}) = \begin{cases} 0, & \text{if } \hat{s}_k = \hat{s}_{3-k} \\ \frac{\alpha}{2r - 1 + \alpha}, & \text{if } \hat{s}_k \neq \hat{s}_{3-k} \end{cases}$$

## C The case of payoffs non-separable in agents' actions

This appendix demonstrates that almost-truthful interim-biased mediation protocols can also facilitate communication in cases when agents' payoffs are not fully separable in actions, provided that the the payoff changes due to action interaction are relatively small compared to payoff changes induced by the own-action component  $v_k$ . Consider the following modification to the specification of agents' payoffs in eq. (2):

$$u_k(s,a) = v_k(a_k,s) - \alpha \times (c_k(a_{3-k},s) + z_k(a_k,a_{3-k},s)),$$
(17)

where the function  $z_k$  captures the effect of action interaction on agents' payoffs. Assume also that  $\alpha \in (0, 1]$  and suppose that Assumption 4.1 (which guarantees the uniqueness of the state-specific correct action) holds throughout this part of the appendix. Also assume that the interaction of actions leads to a smaller variation in payoffs relative to the  $v_k$ -component of the utility:

Assumption C.1. For every agent  $k \in \{1, 2\}$ , state  $s \in S$ , action pair  $a_k, a'_k \in \mathcal{A}_k$  and (3 - k)'s action  $a_{3-k} \in \mathcal{A}_{3-k}$  if  $v_k(a_k, s) - v_k(a'_k, s) > 0$  then

$$|z_k(a_k, a_{3-k}, s)) - z_k(a'_k, a_{3-k}, s))| < v_k(a_k, s) - v_k(a'_k, s)$$
(18)

This assumption essentially guarantees that agent k's optimal action in state s is not affected by agent (3-k)'s action.

Under this assumption, the main results of the paper continue to hold even when agents' actions interact in determining the payoffs as in (17). The intuition is similar to the main result. First, the mediator is able to control agent's beliefs by only rarely introducing distortions. With rare distortions, the mediator is able to induce the agents to take actions  $a^*(s_k, m)$  that are correct given their private signals and the mediator's message (these actions do not depend on the counterpart's action due to Assumption C.1 as demonstrated below). Given that the mediator is able to shift actions in this way, agents'  $v_k$ -component of the utility is maximized by truth-telling. This motive for truth-telling dominates any secondary payoff effects through the counterpart's action provided that the misalignment of interests  $\alpha$  is small enough. The formal proof is sketched below with some details omitted given the similarity to the main part of the paper.

LEMMA C.1. Consider agent k and state  $s \in S$ . For all  $a_{3-k} \in \mathcal{A}_{3-k}$  and  $\alpha \in (0, 1]$  the correct action  $a_k^*(s)$  maximizes  $v_k(a_k, s) - \alpha \times z_k(a_k, a_{3-k}, s)$  with respect to  $a_k$ .

*Proof:* To see this notice first that given the definition of  $a_k^*(s)$  and the uniqueness assumption (part (i) of Assumption 4.1), for every  $a'_k \in \mathcal{A}_k : a'_k \neq a^*_k(s)$ 

$$v_k(a_k^*(s), s) > v_k(a_k', s)$$
 (19)

Now, since  $v_k(a_k^*(s), s) - v_k(a_k', s) > 0$ , then by Assumption C.1 for every  $a_{3-k}$  and  $\alpha \in (0, 1]$ 

$$v_k(a_k^*(s), s) - v_k(a_k', s) > \alpha \times |z_k(a_k^*(s), a_{3-k}, s) - z_k(a_k', a_{3-k}, s)|$$
  
$$\geq \alpha \times (z_k(a_k^*(s), a_{3-k}, s) - z_k(a_k', a_{3-k}, s))$$

Thus for every  $a'_k \in \mathcal{A}_k : a'_k \neq a^*_k(s)$  and every  $a_{3-k}$ 

$$v_k(a_k^*(s), s) - \alpha \times z_k(a_k^*(s), a_{3-k}, s)) > v_k(a_k', s) - \alpha \times z_k(a_k', a_{3-k}, s)),$$

which implies that indeed for every  $a_{3-k}$  and  $\alpha \in (0, 1]$  the correct action  $a_k^*(s)$  maximizes  $v_k(a_k, s) - \alpha \times z_k(a_k, a_{3-k}, s))$  with respect to  $a_k$ .

Next, similarly to Lemma 4.1, notice that an action that is correct for a given signal of the counterpart will be chosen for high enough belief on this signal for any action of the counterpart.

LEMMA C.2. Let  $\tilde{\pi}_k$  be agent k's belief over  $S_{3-k}$ . There exists a  $\bar{\delta}_k^z < 1$  such that for all  $s_{3-k}$ ,  $\alpha \in (0,1]$  and any  $a_{3-k} \in \mathcal{A}_{3-k}$ , if  $\tilde{\pi}_k(s_{3-k}) \ge \bar{\delta}_k^z$ , then  $\arg \max_{a_k} \mathbb{E}_{\tilde{\pi}_k} \left[ v_k(a_k,s) - \alpha z_k(a_k,a_{3-k},s) \right] = a_k^*(s_k,s_{3-k})$ .

*Proof.* Notice that for  $\tilde{\pi}_k(s_{3-k}) = 1$  and  $a^* = a_k^*(s_k, s_{3-k})$  (which is unique by part (i) of Assumption 4.1)

$$\mathbb{E}_{\tilde{\pi}_k}\left[v_k(a^*, s) - \alpha z_k(a^*, a_{3-k}, s)\right] = \sum_{t_{3-k}} \tilde{\pi}_k(t_{3-k}) \left(v_k(a^*, s) - \alpha z_k(a^*, a_{3-k}, s)\right)$$

$$= v_{k}(a^{*}, s) - \alpha z_{k}(a^{*}, a_{3-k}, s)$$

$$> v_{k}(a', s) - \alpha z_{k}(a', a_{3-k}, s)$$

$$= \sum_{t_{3-k}} \tilde{\pi}_{k}(t_{3-k}) (v_{k}(a', s) - \alpha z_{k}(a', a_{3-k}, s))$$

$$= \mathbb{E}_{\tilde{\pi}_{k}} \left[ v_{k}(a', s) - \alpha z_{k}(a', a_{3-k}, s) \right]$$

for every  $a' \neq a_k^*(s_k, s_{3-k})$ , where the inequality for every  $a_{3-k}$  is implied by Lemma C.1. Thus by continuity of  $\sum_{t_{3-k}} \tilde{\pi}_k(t_{3-k}) v_k(a_k, s) - \alpha z_k(a_k, a_{3-k}, s)$  with respect to  $\tilde{\pi}_k(t_{3-k})$ , there exists a  $\bar{\delta}_k^z(s_{3-k}, a_{3-k}) < 1$  such that the same strict inequality holds for all  $\tilde{\pi}_k$  such that  $\tilde{\pi}_k(s_{3-k}) \ge \bar{\delta}_k^z(s_{3-k}, a_{3-k})$ . The proof of the lemma is completed by defining  $\bar{\delta}_k^z = \max_{s_{3-k}, a_{3-k}} \{ \bar{\delta}_k^z(s_{3-k}, a_{3-k}) \}$ .

Now, to see that the equilibrium with truth-telling exists even with the interplay of agents actions, notice that the remaining steps of the proof in the main part of the paper can be reproduced almost without changes. First, Lemma 4.2 showing that the mediator can create a belief weight arbitrarily close to 1 by using rare distortions does not depend on the interplay of the actions and goes through in exactly the same form. Next, the definition of almost-truthful mediation requires only a small change reflecting that the belief weights that are enough for each agent to behave as if the counterpart's signal is perfectly known now also depend on possible variation in the counterpart's action:

DEFINITION C.1. Let an *almost-truthful interim-biased mediation protocol* [for the case of payoff functions allowing for an interaction between agents' actions]  $m^{az}$  be a collection of random variables  $\{m_k^{az}\}_{k=1,2}$ 

with 
$$m_k^{az}(\hat{s}_k, \hat{s}_{3-k}) = m_k^b(\hat{s}_k, \hat{s}_{3-k})$$
 for some  $\varepsilon_k \in (0, \bar{\varepsilon}_k(\bar{\delta}_k^z)]$ 

The only difference between this definition and Definition 4.5 is that  $\bar{\delta}_k$  is replaced with  $\bar{\delta}_k^z$ . Then, Lemma 4.3 can be reformulated similarly:

LEMMA C.3. For each agent  $k, \alpha \in (0, 1]$ , signal  $s_k$ , signal report  $\hat{s}_k, (3 - k)$ 's action and realization m of mediator's message  $m_k^{az}(\hat{s}_k, s_{3-k})$ , agent k's optimal action coincides with the correct action in state  $(s_k, m)$ :  $\arg \max_{a_k} \mathbb{E}_k \left[ v_k(a_k, s) - \alpha z_k(a_k, a_{3-k}, s) | s_k, m_k^a(\hat{s}_k, s_{3-k}) = m \right] = a_k^*(s_k, m).$ 

*Proof.* By Definition C.1 agent k's posterior belief over agent (3 - k)'s signal places a higher than  $\bar{\delta}_k^z$  weight on *m*. By Lemma C.2,  $\arg \max_{a_k} \mathbb{E}_k \left[ v_k(a_k, s) - \alpha z_k(a_k, a_{3-k}, s) | s_k, m_k^a(\hat{s}_k, s_{3-k}) = m \right] = a_k^*(s_k, m).$ 

Next, since the agent's optimal behavior upon observing the mediator's message is known for any counterpart action (and coincides with the correct action in the corresponding state), Lemma 4.4 holds without any modifications, reproduced here for reference:

LEMMA C.4. Suppose that Assumption 4.4 holds. For each agent k, signals  $s_k \neq \hat{s}_k$  and mediation protocol  $m^{az}$ ,  $V_{m_{\mu}^{az}}(s_k, s_k) > V_{m_{\mu}^{az}}(s_k, \hat{s}_k)$ .

That is, the  $v_k$ -component of the utility (excluding the action interaction part and the  $c_k$ -component) is still strictly maximized by truth-telling.

Finally, the proof of the Theorem 4.1 requires only a small modification to show the existence of low enough misalignment of interest that ensures truth-telling:

**Theorem C.1.** Suppose that Assumption 4.4 holds. There exists an  $\bar{\alpha}$  such that for all  $\alpha \in (0, \bar{\alpha})$ , for each agent k, signals  $s_k \neq \hat{s}_k$  and almost-truthful mediation protocol  $m_k^{az}$ ,  $U_{m_k^{az}}(s_k, s_k) > U_{m_k^{az}}(s_k, \hat{s}_k)$ .

*Proof.* Define  $\Delta Z_{m_k^{az}}(s_k, \hat{s}_k)$ ,  $\Delta U_{m_k^{az}}(s_k, \hat{s}_k)$  and  $\Delta C_{m_k^{az}}(s_k, \hat{s}_k)$  analogously to  $\Delta V_{m_k^{az}}(s_k, \hat{s}_k)$ . Notice that  $\Delta Z_{m_k^{az}}(s_k, \hat{s}_k)$  and  $\Delta U_{m_k^{az}}(s_k, \hat{s}_k)$  are well-defined expectations, since the counterpart's action conditional on their private signal is uniquely pinned down by the mediator's message to the counterpart (i.e. there's a unique equilibrium in the post-communication game which only depends on the signals and mediator's messages). Also,

$$\Delta U_{m_{\nu}^{az}}(s_k, \hat{s}_k) = \Delta V_{m_{\nu}^{az}}(s_k, \hat{s}_k) - \alpha \times (\Delta C_{m_{\nu}^{az}}(s_k, \hat{s}_k) + \Delta Z_{m_{\nu}^{az}}(s_k, \hat{s}_k))$$

Let

$$\bar{\Delta}V_{m_k^{az}} = \min_{s_k, \hat{s}_k} \left[ \Delta V_{m_k^{az}}(s_k, \hat{s}_k) \right]$$

and

$$\bar{\Delta}H_{m_k^{az}} = \max\left\{0, \max_{s_k, \hat{s}_k} \left[\Delta C_{m_k^{az}}(s_k, \hat{s}_k) + \Delta Z_{m_k^{az}}(s_k, \hat{s}_k)\right]\right\}$$

Notice that  $\bar{\Delta}V_{m_k^{az}} > 0$  by Lemma C.4 and  $\bar{\Delta}H_{m_k^{az}} \ge 0$  by construction. Define

$$\bar{\alpha} = \begin{cases} 1 & \text{if } \forall k \,\bar{\Delta} H_{m_k^{az}} = 0\\ \min\left\{1, \min_k \left[\frac{\bar{\Delta} V_{m_k^{az}}}{\bar{\Delta} H_{m_k^{az}}}\right]\right\} & \text{if } \exists k \,\bar{\Delta} H_{m_k^{az}} \neq 0 \end{cases}$$

where the last line ensures that  $\bar{\alpha} \in (0, 1]$ . It remains to notice that for  $\alpha \in (0, \bar{\alpha})$ ,  $\Delta U_{m_k^{az}}(s_k, \hat{s}_k) > 0$  for every *k* and  $s_k \neq \hat{s}_k$ , which completes the proof.