

Almost-Truthful Mediation Enables Information Exchange between Agents with Opposing Interests

Dmitry Sedov*

*Northwestern University, Evanston, USA

May 12, 2021

Abstract

Information necessary for decision-making is often distributed among agents with opposing 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 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 messages back, almost always transmitting the received reports without any distortions. The rare distorted messages, however, are deliberately designed to prevent deviations from truth-telling. Specifically, the mediator's distortions encourage each agent to take the action that is interim-optimal given her private information as reflected by her signal report. Thus, in case of distorted messages, a truthful agent acts on private information only, while the action of a deviating agent is shifted away from the truly interim-optimal one. 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. Through its distinctive way of eliciting the truth, almost-truthful mediation highlights the role of information design in facilitating communication.

JEL classification: C72; D82; D83.

Keywords: Communication; Competition; Information; Mediation; Cheap talk; Mechanism design.

E-mail address: dsedov@u.northwestern.edu

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 Francisco Poggi and Quitzé Valenzuela-Stookey. The author is also grateful for the reviews received on an earlier version of this paper.

This paper previously circulated under the title "Information Exchange between Competitors". The author gratefully acknowledges financial support received from the Russian Science Foundation [project №15-18-30081] when that version of the paper was in progress.

1 Introduction

Information relevant to decision-making is often distributed among agents with opposing interests. For illustration, consider the following examples. Pharmaceutical companies possess only partial data on drugs under development, but prefer to be the sole owners of the final formula. In knowledge-intensive organizations complementary information is fragmented and held by employees who are frequently motivated by relative evaluation. Intelligence agencies gather incomplete material about the suspects during joint 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 opposing interests. The model features a one-stage game without monetary transfers, in which agents possess private signals, take actions and receive payoffs that depend on the combination of signals and actions. Agents lack commitment power and may have opposing interests regarding each other's actions. Moreover, the actions are neither strategic complements nor strategic substitutes, which eliminates coordination motives for truthful communication. The usual incentives for information transmission are thus ruled out, direct information exchange cannot be sustained in equilibrium, and information design is isolated as the sole means of enabling communication.

This paper shows that a special class of almost-truthful mediation protocols can facilitate communication between the agents. Protocols in this class are designed to receive private signal reports from the agents and to send messages back so that the response message to a given agent almost always contains the actual signal report of the other agent. With a small probability, however, the response message to a given agent is distorted. The distorted message depends on that agent's own signal report only and does not contain any additional information. The two types of messages (accurate and distorted) take value in the same set, and thus the agents can not distinguish them with certainty. In fact, by the almost-truthful design of the protocol, each agent is almost certain the response message contains the other agent's actual signal report. The undistorted messages lead to complete information revelation. The distorted messages are designed to encourage each agent to take an action that is interim-optimal given her private information as reflected by her signal report. When the mediator uses a distorted message, a non-deviating agent essentially acts on her private information only. A deviating agent, however, shifts her action away from the action that is interim-optimal with respect to her true private signal. As a result, the distorted message hurts the agent more if she deviates rather than not. This ensures that revealing private information truthfully is optimal from the perspective of agent's own actions. While a deviating agent may still benefit from the change in the other party's action caused by the deviation, incentives for truthful communication dominate if the misalignment of interests between the agents is small enough. Thus, an equilibrium with truth-telling exists. Importantly, 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. 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. 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. These three details (a design feature and two assumptions) are crucial for the mediated communication scheme developed in this

paper.

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, proves the non-existence of direct communication, characterizes the class of almost-truthful mediation protocols and demonstrates that such protocols enable information transmission. Section 5 concludes by summarizing the results and discussing applications to the settings of sharing clinical trials data, communication in organizations and exchange of information between intelligence agencies.

2 Literature review

Two branches of research on transmission of unverifiable information are relevant in the context of the present paper.

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). 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 decision-maker. That is, informed agents prefer to tell the truth, since the corresponding action is preferred. 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 interests¹. The only incentive for communication is the higher benefit of information received through the mediation protocol in case of truth-telling (i.e. lower downside of the distorted messages).

Second, communication between partially informed agents, who also take actions, is studied in a number of papers. 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 oligopolists², who may share unverifiable private

¹In 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.

²Further 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 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 [Goltzman 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 optimal 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 is equivalent to public persuasion, which is not the case in the present paper: the mediation protocol *does* need to condition her recommendation on agents’ reports in order to sustain informative communication.

To sum up, the present paper contributes to the literature on communication between agents with misaligned interests by considering a setting that differs from the ones discussed above in terms of assumptions, results and intuition.

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 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 Γ . Each of the two agents $i \in \{1, 2\}$ obtains a binary signal B_i . Let $S = \{0, 1\}^2$ with the following joint distribution c over $S^2 = S \times S$ parametrized by

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.

$A \geq 1/2 - 1^\circ$:

	C	
	$s_2 = 0$	$s_2 = 1$
$s_1 = 0$	$\frac{A}{2}$	$\frac{1-A}{2}$
$s_1 = 1$	$\frac{1-A}{2}$	$\frac{A}{2}$

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

$$U_1(s_1, s_2) = 1 - 2|s_1 - B| \quad U_2(s_1, s_2) = 2|s_1 - B| \tag{1}$$

where $U \in [0, 1]$ parametrizes the degree of interest misalignment between the agents: the higher U 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 $O_i = B$. 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 $A(1-U)$ to $1-U$. 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 mediation

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

Mediation setup The almost-truthful mediation protocol receives signal reports from the agents and sends messages back to them. It almost always sends agent 1's signal report to agent 2 without any distortion. However, with a positive probability γ the mediation protocol returns to agent 1 a message that coincides with agent 2's most likely signal given 1's own report. Formally, let $\{p_1, p_2\}$ be a collection of binary distributions over mediator's messages to agent 1 when the mediator receives reports $\{s_1, s_2\}$ from the agents:

$$p_1(s_1, s_2) = \begin{cases} p_2 & \text{with probability } 1 - \gamma \\ p_1 & \text{with probability } \gamma \end{cases}$$

where

$$\gamma \geq 0 - \frac{1-A}{A}$$

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 1 reports truthfully, agent 2 chooses the average between the private signal and the mediator's message as her action irrespective of whether 1 herself reports truthfully. To see this, note that the posterior probability on the event $B_3 = c$ where c is the observed mediator's message is greater than $1/2$. Indeed, if the mediator's message does not coincide with 1's report, $c < B_1$ - then agent 2 for sure knows that $c = B_3$ and $P_{B_3} (B_3 = c | B_1 < B_3) = 1$. If the mediator's message does coincide with 1's report, $c = B_1$ - then

$$P_{B_3} (B_3 = c | B_1 < B_3) = \frac{P (c = B_3 | B_1 < B_3) = c < B_1}{P (c = B_3 | B_1 < B_3) + P (c = B_1 | B_1 < B_3)}$$

$$= \frac{1}{1 + \frac{A}{Y}}$$

$$= \frac{1}{2} \frac{1 + Y}{1 + Y + A}$$

since

$$Y > \frac{1}{A}$$

As a result, for low values of Y the agent optimally chooses action $O_2 = \frac{1}{2} (B_1 + B_3)$, the average of the private signal and mediator's message.

Second, since 1'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 1 chooses action B . If the mediator distorts the information agent 1 chooses the interim-optimal action B_1 in case she doesn't deviate and chooses action $1/2$ in case she does deviate. That is, if 1 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 $1 - AY$. If 1 deviates, she makes a mistake when the mediator distorts and her interim-optimal action was actually ex-post optimal, which happens with probability AY . Thus the expected partial payoff from 1's own actions in case of a truthful report equals $U(B) = 1 - AY$ and equals $U(1/2) = 1 - AY$ in case of a deviation.

It remains to notice that the expected partial payoff in case of reporting truthfully is higher: $U(B) - U(1/2) = AY > 0$. Provided that U is low enough, the positive partial difference $U(B) - U(1/2)$ dominates the total payoff difference, and reporting truthfully is optimal for each agent. The existence of a truthful equilibrium is thus shown.

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. 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 $A > 1/2$ - which imposes positive correlation between the agents' signals. If $A = 1/2$ (signals are uncorrelated), then the mediator can not shift the action away from the interim-optimal action³. 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 $A = 1/2$ is problematic in two ways: (i) non-uniqueness of interim-optimal action, (ii)

³Which is not unique: both the private signal B_1 and $1/2$ are optimal actions for every agent.

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.

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 mediation enables communication, while direct communication is impossible. Sufficient conditions for the almost-truthful 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.

Baseline Formally, consider a 2-agents setup with a finite state space $S = S_1 \times S_2$. Each agent i learns the realization of $B_i \in S_i$ (the signal), but does not learn B_{-i} . Let $C^{i,0}$ be the common prior over S . Also, let $C^{i,1} | B_i^0$ and $E_i | B_i^0$ be agent i 's posterior and expectation operator upon learning B_i , respectively. Each agent i chooses an action O_i from a finite action space A_i . It is assumed that the agents' preferences over action profiles are represented by a state-dependent utility function

$$U_i(B_i, O_i, O_{-i}) = E_i(O_i | B_i^0) - B_i^0 + U_{-i}(O_i, O_{-i} | B_i^0)$$

$E_i(O_i | B_i^0)$ captures agent i 's value from taking action O_i in a given state $B_i \in S_i$. Without loss of generality assume $E_i(O_i | B_i^0) > 0$ for every i , O_i and B_i . The preferences of agent i with respect to O_{-i} 's actions in a given state are captured by the cost function $U_{-i}(O_i, O_{-i} | B_i^0)$. The cost function U_{-i} captures the idea of agents' opposing interests: as in the illustrative example above, a change in agent 1 's action that is beneficial for agent 2 is allowed to increase the cost function U_{-i} (and thus to be harmful) for agent 1 . The cost component of the utility function is further parametrized by the parameter $U_i \in [0, 1]$ that captures the degree of the misalignment of interests⁴. It

⁴This paper is primarily motivated by the situations where agents have opposing 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), but does not depend on such an alignment.

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 B in S each agent i can maximize her payoff by taking the same action for all actions of j . That is, due to the assumption of separability of agents' preferences with respect to each other's actions, agent i can choose the action that maximizes the E_i -component of her utility. Definition 4.1 below introduces formal notation for such state-specific correct actions:

DEFINITION 4.1. For agent i in state $B \in S$ let $O_i^1(B)$ be the set of state-specific correct actions. That is, $O_i^1(B) = \arg \max_{O_i} E_i(O_i, -B)$.

Similarly, for every privately observed state $B_i \in S_i$ the interim-correct actions are defined:

DEFINITION 4.2. For each agent i and signal $B_i \in S_i$ let $O_i^1(B_i)$ be the set of interim-correct actions. That is, $O_i^1(B_i) = \arg \max_{O_i} E_i(O_i, -B_j | B_i)$.

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 i and signal $B_i \in S_i$ let $\mathcal{F}_i^1(B_i)$ be the set of interim-correct counterpart signals. That is, $\mathcal{F}_i^1(B_i) = \{B_j \in S_j \mid O_i^1(B_i) = O_i^1(B_i, -B_j)\}$.

Specifically, if agent j 's signal was revealed to belong to the set $\mathcal{F}_i^1(B_i)$ then agent i endowed with signal B_i would have no incentive to take an action other than the interim-correct one. At this point, $\mathcal{F}_i^1(B_i)$ 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{c}_i be agent i 's belief over S_j . There exists a \bar{X}_i such that for all $B_j \in S_j$ if $\tilde{c}_i(B_j) > \bar{X}_i$ then $\arg \max_{O_i} E_{\tilde{c}_i}(E_i(O_i, -B^0) = O_i^1(B_i, -B_j)$.

Proof. Notice that for $\tilde{c}_i = 1$ and every $O_i \in O_i^1(B_i, -B_j)$

$$\begin{aligned} E_{\tilde{c}_i}(E_i(O_i, -B^0)) &= \int_{S_j} \tilde{c}_i(B_j) E_i(O_i, -B^0) \\ &= E_i(O_i, -B^0) \\ &= \int_{S_j} E_i(O_i, -B^0) \\ &= \int_{S_j} \tilde{c}_i(B_j) E_i(O_i, -B^0) \\ &= E_{\tilde{c}_i}(E_i(O_i, -B^0)) \end{aligned}$$

for every $O_i \in O_i^1(B_i, -B_j)$ by Definition 4.1. Thus by continuity of $\int_{S_j} \tilde{c}_i(B_j) E_i(O_i, -B^0)$ with respect to $\tilde{c}_i(B_j)$ there exists a $\bar{X}_i(B_j)$ such that the same strict inequality holds for all \tilde{c}_i such that $\tilde{c}_i(B_j) > \bar{X}_i(B_j)$. The proof of the lemma is completed by defining $\bar{X}_i = \max_{B_j \in S_j} \bar{X}_i(B_j)$.

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) $O_i^1 B^0$ is a singleton for every i and $B \in S$. (ii) For every agent i and signal B_i - if $B_3^0 < B_3$ - then $O_i^1 B_i - B_3^0 < O_i^1 B_i - B_3$.

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) $O_i^1 B_i^0$ is a singleton for every i and $B_i \in S_i$. (ii) For every agent i and signal B_i there exists B_3 such that $O_i^1 B_i^0 = O_i^1 B_i - B_3$.

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 C_i are close enough to degenerate ones: $C_i^1 j B_i^0 > 0$ for the true signal B_3 only⁵. 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 can not 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 ($A_i = O_i^1 B^0 \cup B \in C$ for every i and if the E_i -component is monotone with respect to own-signal (for every agent i - signals $B_i - B_3$ and $B_3^0 - E_i^1 O_i^1 B_i - B_3^0 - B_3 > E_i^1 O_i^1 B_i - B_3^0 - B_3$). Note that Assumption 4.1 and Assumption 4.2 jointly ensure that the set of interim-correct counterpart signals $\mathcal{F}_i^1 B_i^0$ is a singleton for all agents i and private states B_i .

Finally, the assumption regarding the information structure is made.

ASSUMPTION 4.3. (i) $|S_1| = |S_2|$. (ii) $C^1 B^0 > 0$ for all $B \in S$.

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 signals⁶. 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.

⁵To see this formally, one can utilize Lemma 4.1.

⁶The setup under part (i) of Assumption 4.3 can be interpreted as follows. Agents agree on the possible states of the world: $S_1 = S_2 = C$, and receive noisy signals about the true state with the signal space coinciding with the states of the world set.

4.3 Impossibility of direct communication

Can the two agents help each other directly by simultaneously sending messages that are at least partially informative? Intuitively, this should not be the case. Conditional on receiving a message from the counterpart and the corresponding posterior over the counterpart’s signal (which governs the action choice), each agent has no incentives to send truthful information. This section establishes the impossibility of direct communication.

Proposition 4.1. *There exist cost functions $2: 1 - 0$ such that for any direct communication extension \mathbb{V} : of the baseline setup only interim-optimal actions are taken.*

Proof. Consider the case of “extreme” conflict: $2: 1 0_3 \text{ :- } B^0 = E_3 \text{ :- } 1 0_3 \text{ :- } B^0$. Consider an extended game in which agents can simultaneously send each other a message from a predetermined set \mathbb{V} : upon observing the private signal. Suppose that there is an equilibrium in which agent $13 \text{ :- } 0^0$ ’s expected value of the own-choice component $E_3 \text{ :-}$ is strictly greater than under no communication. Then there exists a type of agent :- and a message $F \text{ :-}$ sent in equilibrium such that some type of agent $13 \text{ :- } 0^0$ is expected to shift her action away from the interim-optimal one conditional on receiving such a message. Then consider the following deviation: send $F \text{ :-}$ irrespective of own type and follow the same action plan (possibly contingent on received messages) as in equilibrium. Such a deviation is profitable, as it ensures that the expectation of $E_3 \text{ :-}$ is strictly lower than under no communication, and thus the expectation of $U_2 \text{ :-} = UE_3 \text{ :-}$ is higher. The expectation of $E_3 \text{ :-}$ is decreased under such a deviation since at least one type of $13 \text{ :- } 0^0$ shifts the action away from the interim-optimal irrespective of :- ’s true information.

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 :- to agent :- coincides with the set $S_3 \text{ :-}$. This captures the notion that the mediator may share some of $13 \text{ :- } 0^0$ ’s private information with agent :- .

Interim-biased mediation protocols introduced in Definition 4.4 below are randomized for each pair of signal reports. Such protocols transmit agent $13 \text{ :- } 0^0$ ’s report to agent :- with probability $1 - Y \text{ :-}$ and send a message with :- ’s interim-correct counterpart signal with the complementary probability $Y \text{ :-}$.

DEFINITION 4.4. Let an *interim-biased mediation protocol* $\langle 1 \text{ :-}$ be a collection of random variables $\langle 1 \text{ :-} \text{ :-} 1-2$ with

$$\langle 1 \text{ :-} B_3 \text{ :-} B_3 \text{ :-} 0^0 = \begin{cases} B_3 \text{ :-} & \text{with probability } 1 - Y \text{ :-} \\ \mathcal{F} \text{ :-} B_3 \text{ :-} 0^0 & \text{with probability } Y \text{ :-} \end{cases} \tag{2}$$

for some $Y \text{ :-} 2 \gg 0-1 \frac{1}{4}$.

The following lemma establishes that there exist positive probabilities $Y \text{ :-}$ such that, conditional on agent $13 \text{ :- } 0^0$ reporting truthfully, agent :- believes that the mediator’s message coincides with $13 \text{ :- } 0^0$ ’s report irrespective of :- ’s own report. Formally,

LEMMA 4.2. *For every agent :- and $1 \bullet 2 \check{Y} X \check{Y} 1-$ there exists an $Y \text{ :-} 1 X^0 \text{ :-} 0$ such that if an interim-biased mediation protocol satisfies $Y \text{ :-} 6 Y \text{ :-} 1 X^0-$ then $P \text{ :-} B_3 \text{ :-} = \langle j B_3 \text{ :-} \langle 1 \text{ :-} B_3 \text{ :-} 0^0 = \langle \text{ :-} \rangle X$ for all $B \text{ :-} B \text{ :-} 2 S \text{ :-}$ and $\langle 2 S_3 \text{ :-}$.*

Proof. Consider the beliefs of agent i upon reporting signal B_i and receiving message $<$ from the mediator.

1. If $< \in \mathcal{F}_i^1(B_i, \cdot)$ then the posterior probability that agent i 's signal report is equal to $<$ can be computed as follows.

$$P(B_i = < | B_i, <) = \frac{P(< | B_i, <) \cdot P(B_i = <)}{P(< | B_i, <) \cdot P(B_i = <) + \sum_{B_i \neq <} P(< | B_i, B_i) \cdot P(B_i = B_i)}$$

$$= \frac{1 \cdot Y_i}{P(< | B_i, <) + \sum_{B_i \neq <} P(< | B_i, B_i)}$$

Notice that given Definition 4.4, $P(< | B_i, <) = 1 \cdot Y_i$ and thus $P(B_i = < | B_i, <) = 1 > X$. Consequently, for the lemma to be true only the case of $< \in \mathcal{F}_i^1(B_i, \cdot)$ needs to be considered.

2. If $< \in \mathcal{F}_i^0(B_i, \cdot)$ then

$$P(B_i = < | B_i, <) = \frac{1}{1 + \frac{1}{Y_i \cdot c^1 \ell_3} \cdot \sum_{B_i \neq <} P(B_i = B_i | B_i, <)}$$

where

$$c^1 \ell_3 = \frac{C^1(B_i = <)}{C^1(B_i \neq <)}$$

Note that $c^1 \ell_3$ is a positive finite number for every $<$ and B_i under Assumption 4.3. Thus if Y_i satisfies

$$Y_i \geq \frac{1}{1 + \frac{1}{c^1 \ell_3} \cdot \sum_{B_i \neq <} P(B_i = B_i | B_i, <)}$$

$$(3)$$

then $P(B_i = < | B_i, <) > X$ for every B_i .

The proof is completed by observing that $Y_i \geq 1 - \epsilon$.

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 $1 - Y_i$ of truthful transmission of each agent's report is close to 1. Formally,

DEFINITION 4.5. Let an *almost-truthful interim-biased mediation protocol* $\langle \cdot \rangle^0$ be a collection of random variables $\langle \cdot \rangle_{i=1,2}^0$ with $\langle \cdot \rangle_{i=1,2}^0 = \langle \cdot \rangle_{i=1,2}^1$ for some $Y_i \geq 1 - \epsilon$.

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

LEMMA 4.3. For each agent i - signal B_i - signal report B_i and realization $<$ of mediator's message $\langle \cdot \rangle_{i=1,2}^0$ - agent i 's optimal action coincides with the correct action in state $(B_i, <)$: $\arg \max_{O_i} E_i(O_i | B_i, <) = O_i(B_i, <)$.

Proof. By Definition 4.5 agent i 's posterior belief over agent j 's signal places a higher than X weight on $<$. By Lemma 4.2, $\arg \max_{O_i} E_i(O_i | B_i, <) = O_i(B_i, <)$.

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. The additional required assumption ensures enough variation in intermediate beliefs across different types of agents.

First, the additional assumption is stated.

ASSUMPTION 4.4. For each agent i and pair of signals $B_i^0 < B_i - \mathcal{F}_i^1 B_i^0 < \mathcal{F}_i^1 B_i^0$.

This assumption is interpreted as follows: the interim-correct actions are rationalized by different counterpart's signals. Notice that Assumption 4.4 implies $|S_{1j}| = |S_{2j}|$ - 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.

Next, the E_i -component of the utility is shown to be strictly maximized by truthful reporting. Let $+_{<^0} B_i - B_i^0$ be the expectation of agent i 's E_i -component of the utility conditional on reporting B_i to an almost-truthful mediation protocol $<^0$ when the true signal is B_i and agent i reports truthfully. Notice that by Lemma 4.3,

$$+_{<^0} B_i - B_i^0 = \sum_{\mathcal{C}_i} C_i^1 \mathcal{C}_i : j B_i^0 \quad 1 \quad Y_i^0 \quad E_i^1 O_i^1 B_i - \mathcal{C}_i^0 - B_i - \mathcal{C}_i^0 :^{00}$$

$$\sum_{\mathcal{C}_i} C_i^1 \mathcal{C}_i : j B_i^0 \quad Y_i \quad E_i^1 O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} - B_i - \mathcal{C}_i^0 :^{00}$$

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

LEMMA 4.4. Suppose that Assumption 4.4 holds. For each agent i - signals $B_i < B_i$ and mediation protocol $<^0$ - $+_{<^0} B_i - B_i^0 \geq +_{<^0} B_i - B_i^0$.

Proof. Define $+_{<^0} B_i - B_i^0 = +_{<^0} B_i - B_i^0 +_{<^0} B_i - B_i^0$. Notice that

$$+_{<^0} B_i - B_i^0 =$$

$$= \sum_{\mathcal{C}_i} C_i^1 \mathcal{C}_i : j B_i^0 \quad E_i^1 O_i^1 B_i - \mathcal{C}_i^0 - B_i - \mathcal{C}_i^0 :^{00} \gg Y_i \quad Y_i^0$$

$$\sum_{\mathcal{C}_i} C_i^1 \mathcal{C}_i : j B_i^0 \quad Y_i \quad E_i^1 O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} - B_i - \mathcal{C}_i^0 :^{00} \quad E_i^1 O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} - B_i - \mathcal{C}_i^0 :^{00}$$

$$= Y_i \sum_{\mathcal{C}_i} C_i^1 \mathcal{C}_i : j B_i^0 \quad E_i^1 O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} - B_i - \mathcal{C}_i^0 :^{00} \quad E_i^1 O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} - B_i - \mathcal{C}_i^0 :^{00}$$

Rearranging,

$$+_{<^0} B_i - B_i^0 = Y_i \sum_{\mathcal{C}_i} C_i^1 \mathcal{C}_i : j B_i^0 \quad E_i^1 O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} - B_i - \mathcal{C}_i^0 :^{00}$$

$$Y_i \sum_{\mathcal{C}_i} C_i^1 \mathcal{C}_i : j B_i^0 \quad E_i^1 O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} - B_i - \mathcal{C}_i^0 :^{00}$$

First notice that $O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} = O_i^1 B_i^0$ 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 $\mathcal{F}_i^1 B_i^0 < \mathcal{F}_i^1 B_i^0$. Also, by Assumption 4.1 the optimal action varies for different signals of the counterpart and thus $O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} < O_i^1 B_i - \mathcal{F}_i^1 B_i^{00}$. Since by Assumption 4.2 there

is a unique interim-correct action, $O_i^1 B_i - \mathcal{F}_i^1 B_i^{00} < O_i^1 B_i^0 = O_i^1 B_i - \mathcal{F}_i^1 B_i^{00}$ Combining, for $O_i^0 = O_i^1 B_i - \mathcal{F}_i^1 B_i^{00}$,

$$+_{<^0} O_i^1 B_i - \mathcal{B}_i^0 = \frac{Y_i}{i^0} \left\{ \frac{E_{C_i} E_i^1 O_i^1 B_i^0 - \mathcal{B}_i^0}{E_{C_i} E_i^1 O_i^0 - \mathcal{B}_i^0} \right\} \quad i^0$$

where the first inequality $Y_i > 0$ is ensured by Definition 4.5 of almost-truthful mediation protocols, and the second inequality $E_{C_i} E_i^1 O_i^1 B_i^0 - \mathcal{B}_i^0 > E_{C_i} E_i^1 O_i^0 - \mathcal{B}_i^0 > 0$ 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 simultaneously choose actions. Let $*_{<^0} O_i^1 B_i - \mathcal{B}_i^0$ be the expectation of agent i ’s utility conditional on reporting \hat{B}_i to an almost-truthful mediation protocol $<^0$ when the true signal is B_i . Let $<^0 O_i^1 B_i - \mathcal{B}_i^0$ be the expectation of agent i ’s 2.-component of the utility conditional on reporting \hat{B}_i to an almost-truthful mediation protocol $<^0$ when the true signal is B_i .

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

Theorem 4.1. *Suppose that Assumption 4.4 holds. There exists an U such that for all $U \geq U^0$ -for each agent i -signal $B_i < \hat{B}_i$ and almost-truthful mediation protocol $<^0$ - $*_{<^0} O_i^1 B_i - \mathcal{B}_i^0 > <^0 O_i^1 B_i - \mathcal{B}_i^0$.*

Proof. Define $*_{<^0} O_i^1 B_i - \mathcal{B}_i^0$ and $<^0 O_i^1 B_i - \mathcal{B}_i^0$ analogously to $+_{<^0} O_i^1 B_i - \mathcal{B}_i^0$. Notice that

$$*_{<^0} O_i^1 B_i - \mathcal{B}_i^0 = +_{<^0} O_i^1 B_i - \mathcal{B}_i^0 \quad U > <^0 O_i^1 B_i - \mathcal{B}_i^0$$

Let

$$-_{<^0} = \min_{B_i - \hat{B}_i} \quad h \quad i \quad +_{<^0} O_i^1 B_i - \mathcal{B}_i^0$$

and

$$-_{<^0} = \max_{B_i - \hat{B}_i} \quad h \quad i \quad 0 - \max_{<^0} O_i^1 B_i - \mathcal{B}_i^0$$

Notice that $-_{<^0} > 0$ by Lemma 4.4 and $-_{<^0} > 0$ by construction. Define

$$U = \begin{cases} \infty & \text{if } \exists: -_{<^0} = 0 \\ \min_{<^0} \frac{+_{<^0}}{-_{<^0}} & \text{if } \exists: -_{<^0} < 0 \end{cases}$$

It remains to notice that for $U \geq U^0$ - $*_{<^0} O_i^1 B_i - \mathcal{B}_i^0 > 0$ for every i and $B_i < \hat{B}_i$ - 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 interim-optimal action. 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.

On information structure and almost-truthful mediation As stated in Section 4.2, the assumptions imposed on the setting in order to show that almost-truthful mediation enables information exchange require that the action space is of intermediate coarseness. A similar observation can be made regarding the prior distribution C . Intuitively, C should exhibit an intermediate degree of dependence across its two dimensions, S_1 and S_2 - in order for a given almost-truthful mediation protocol to work.

To see this, first fix an U and an almost-truthful mediation protocol characterized by Y that allows for information exchange. Consider a small increase in $C: \mathcal{F}^1 B; \circ | B; \circ$ for all $B; \circ$. Notice that Y ; defined in (3) necessarily goes down. There is thus no guarantee that the Y -protocol ensures a high enough posterior belief weight on its message among the agents. Consequently, there is no guarantee that the Y -protocol still allows information exchange under the new prior C . On the other hand, consider the case when almost-truthful mediation is impossible due to the failure of Assumption 4.4. Increasing the conditional beliefs $C: \mathcal{F}^1 B; \circ$ on the most likely counterpart's signal given $B; \circ$ can restore Assumption 4.4 and thus lead to the possibility of communication through a given almost-truthful mediation protocol.

Increasing the degree of dependence between the two dimensions of C can thus make a given protocol less trustworthy, but can also make the interim beliefs more sensitive to signals. Informally, this makes intermediate degree of dependence across C 's two dimensions most suitable for almost-truthful mediation.

4.6 Almost-truthful mediation can be optimized with a linear program

The almost-truthful mediation can be generalized, allowing the probability Y ; 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

$$\langle \mathcal{F}^1 B; \circ - \mathcal{B}_3; \circ \rangle = \begin{cases} \mathcal{B}_3; \circ & \text{with probability } 1 - Y; \mathcal{F}^1 B; \circ - \mathcal{B}_3; \circ \\ \mathcal{F}^1 B; \circ & \text{with probability } Y; \mathcal{F}^1 B; \circ - \mathcal{B}_3; \circ \end{cases}$$

such that $Y; \mathcal{F}^1 B; \circ - \mathcal{B}_3; \circ \geq 0 - Y; \mathcal{F}^1 \bar{X}; \circ$ for all $B; \circ - \mathcal{B}_3; \circ$.

Proposition 4.2. *The Pareto-optimal generalized almost-truthful interim-biased mediation protocol can be found in polynomial time with a linear programming problem.*

Proof. Similarly to previous cases, let $\langle \mathcal{F}^1 B; \circ - \mathcal{B}; \circ \rangle$ be the expectation of agent i 's utility conditional on reporting $\mathcal{B}; \circ$ to a generalized almost-truthful mediation protocol $\langle \mathcal{F}^1 \rangle$ when the true signal is $B; \circ$ and agent i reports truthfully. Notice that

$$\begin{aligned} \langle \mathcal{F}^1 B; \circ - \mathcal{B}; \circ \rangle &= C; \mathcal{F}^1 \mathcal{L}_3; \circ | B; \circ \cdot 1 - Y; \mathcal{F}^1 B; \circ - \mathcal{L}_3; \circ \cdot E; \mathcal{F}^1 O; \mathcal{F}^1 B; \circ - \mathcal{L}_3; \circ - \mathcal{B}; \circ - \mathcal{L}_3; \circ \\ &\quad + C; \mathcal{F}^1 \mathcal{L}_3; \circ | B; \circ \cdot Y; \mathcal{F}^1 B; \circ - \mathcal{L}_3; \circ \cdot E; \mathcal{F}^1 O; \mathcal{F}^1 B; \circ - \mathcal{F}^1 \mathcal{B}; \circ - \mathcal{B}; \circ - \mathcal{L}_3; \circ \\ &\quad + C; \mathcal{F}^1 \mathcal{L}_3; \circ | B; \circ \cdot 1 - Y_3; \mathcal{F}^1 \mathcal{L}_3; \circ - \mathcal{B}; \circ \cdot 2; \mathcal{F}^1 O; \mathcal{F}^1 B; \circ - \mathcal{L}_3; \circ - \mathcal{B}; \circ - \mathcal{L}_3; \circ \\ &\quad + C; \mathcal{F}^1 \mathcal{L}_3; \circ | B; \circ \cdot Y_3; \mathcal{F}^1 \mathcal{L}_3; \circ - \mathcal{B}; \circ \cdot 2; \mathcal{F}^1 O; \mathcal{F}^1 \mathcal{F}_3; \mathcal{F}^1 \mathcal{L}_3; \circ - \mathcal{L}_3; \circ - \mathcal{B}; \circ - \mathcal{L}_3; \circ \end{aligned}$$

is linear in the Y -parameters. Thus the Pareto-optimization problem with weights λ_i :

$$\begin{aligned} \max_{\{Y_i\}_{i=1}^n} \quad & \sum_{i=1}^n \lambda_i C_i^1 B_i^0 - \sum_{i=1}^n \lambda_i C_i^0 B_i^1 \\ \text{subject to} \quad & \sum_{i=1}^n \lambda_i C_i^1 B_i^0 - \sum_{i=1}^n \lambda_i C_i^0 B_i^1 \leq 0 \end{aligned}$$

is a linear program and can be solved by the standard methods (and in polynomial time), see [Vanderbei \(2014\)](#) and references therein.

Remark 4.1. *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 interim-biased 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.*

5 Conclusion

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

5.1 Applications

The mediation protocol introduced in Definition 4.5 suggests potentially useful insights for arranging information exchange in settings where competition between different parties is impossible (or too costly) to avoid by changing an organizational structure altogether.

Consider the case of pharmaceutical companies mentioned in the introduction. While in the recent years there has been some movement both by regulators and private entities in the direction of sharing clinical trials data⁷, it is still recognized that there are significant challenges to such a practice, in particular, there exist risks of commercial misuse of the shared data and the proper policies to mitigate such risks are still to be developed. For example, a report by the [Institute of Medicine \(2015, p. 142\)](#)⁸ acknowledges that (i) “sponsors of clinical trials have serious concerns about competitors copying data packages that lack strong regulatory data protection” and (ii) “what types of policies might be implemented to protect data from “unfair commercial use” is a difficult question”, listing registration-only or online-only access to shared information and data request reviews “by an independent third party” as potential solutions. The intuition developed in this paper suggests that a central authority gathering the data from competing companies and distributing the results with a small bias to each company’s private report may provide appropriate motivation to share data in addition to the policies recommended by the [Institute of Medicine \(2015\)](#).

Many organizations try to encourage knowledge sharing among their employees⁹. Appropriately

⁷See, for example, the Wall Street Journal article by [Burton \(2016\)](#) on the expansion of requirements for publication of clinical studies’ results and the New York Times publication by [Thomas \(2012\)](#) on GlaxoSmithKline’s effort to open up its drug research, which resulted in the clinical trials data request system being set up at www.clinicalstudydatarequest.com.

⁸The Institute of Medicine was renamed to the National Academy of Medicine in 2015.

⁹See, for example, the Harvard Business Review discussion of knowledge sharing by [Myers \(2015\)](#) and the Wall Street Journal article on knowledge hiding by [Deal \(2018\)](#).

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 indicates that organizational leaders can also act as key communication intermediaries for their subordinates, who possess some complementary knowledge and engage in competition, and hints at a strategy of information dissemination. In particular, transmitting full information in most cases, and sometimes encouraging the employees to act on private information may work as the proper mediation tactic. While some information is transferred, the rare encouragement to act on private information only can ensure that the employees who engage in knowledge hiding are put to a disadvantage. Such implicit penalty can thus prevent improper communication, importantly, even if the organizational leader is unable to establish the quality of the reports directly. Finally, 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 \(2002\)](#). 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. The concern about the difficulty of providing proper incentives for sharing information across competing agencies is still relevant, see the discussions by [Garicano and Posner \(2005, p. 161\)](#) and the current head of the U.S. State Department’s Bureau of Counterterrorism, [Sales \(2010, p. 281\)](#). Modest implications of the results in this paper point to (i) the disincentives to share complete intelligence information in the presence of inter-agency competition; and (ii) Information Sharing Environment’s potential use as an imperfect information transmission filter, establishing proper incentives for knowledge sharing. When considering these applications, one should keep in mind that the mediation protocol class constructed in this paper relies on the assumption of the separability of agents’ utility functions in each other’s actions. Thus the insights apply best to settings where agents can not mitigate each other’s actions directly but are affected by such actions (companies operating in separate markets, employees working on different projects, intelligence agencies investigating parallel cases). At the same time, one should also keep in mind that the construction presented in this paper is primarily theoretical and demonstrates that appropriately designing an information environment can facilitate communication.

5.2 Concluding remarks

This paper studies communication between partially informed agents with opposing 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 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) suggests applications of the results to the settings of clinical trials data exchange, organizational knowledge hiding and sharing of intelligence information.

Three generalizations of this paper’s results can be interesting. First, relaxing the separability of preferences in agents’ actions may be desirable. Provided that the component of utility related to the interplay of agents’ actions is sufficiently small compared to the $E \cdot \mathbb{1}_{O \cdot} - B^p$ action-state component of the utility, a similar argument to Theorem 4.1 seems feasible. Next, 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 possibility result to multi-player settings. Finally, 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 can’t be a non-random function of the agent’s report); (ii) the deviating agent is shifted away from the interim-optimal action *more*¹⁰ than the truthful agent when distorted messages are sent. While formalizing these arguments is interesting, it is outside of the scope of the present paper, which provides a possibility result and highlights intuition.

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.
- Burton, T. M. (2016, September). New U.S. Rule to Expand Requirements for Publication of Clinical Trials. *The Wall Street Journal*. Accessed: 2018-09-26. Archived at <https://web.archive.org/web/20171113163510/https://www.wsj.com/articles/new-u-s-rule-to-expand-requirements-for-publication-of-clinical-trials-1474040304>.
- 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/background/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-stop-employees-from-deliberately-hiding-information/>.
- Farrell, J. and R. Gibbons (1989). Cheap Talk with Two Audiences. *The American Economic Review* 79(5), 1214 – 1223.

¹⁰In the continuous action space both agents’ optimal actions conditional on the available information will reflect the possibility that the mediator’s message was distorted.

- 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.
- Institute of Medicine (2015). *Sharing Clinical Trial Data: Maximizing Benefits, Minimizing Risk*. Washington, D.C.: The National Academies Press.
- 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.
- Myers, C. G. (2015, November). [Is Your Company Encouraging Employees to Share What They Know?](https://web.archive.org/web/20181015154019/https://hbr.org/2015/11/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-your-company-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.
- 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.
- Thomas, K. (2012, October). [Glaxo Opens Door to Data on Research](https://web.archive.org/web/20180928103418/https://www.nytimes.com/2012/10/11/business/glaxo-opens-door-to-data-on-its-research.html). *The New York Times*. Accessed: 2018-09-27. Archived at <https://web.archive.org/web/20180928103418/https://www.nytimes.com/2012/10/11/business/glaxo-opens-door-to-data-on-its-research.html>.
- 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.

Vanderbei, R. (2014). *Linear programming: Foundations and Extensions.* New York: Springer.

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 consists of a strategy profile $\theta_{\beta=1}^{10}$ with*

$$O_{\beta=1}^1 B_{\beta=1}^0 = B_{\beta=1}^0$$

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

$$C = A^1 1 U^0$$

Proof. Note that agent 1's best response to each strategy of agent 2 $\theta_{\beta=1}^0$ requires selecting the most likely correct action of nature given 1's signal. Since

$$P(B|B) = P(A|B) = A$$

while

$$P(B|A) = 1 - P(A|B) = 1 - A$$

and then since $A > 1/2$ it must be that in equilibrium each agent selects an action that coincides with her signal and thus indeed $O_{\beta=1}^1 B_{\beta=1}^0 = B_{\beta=1}^0$.

Given the equilibrium strategies, each agent guesses the the correct action with probability A and thus the ex-ante equilibrium payoff of each player is $C = A^1 1 U^0$.

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

$$C^* = \max_{\theta_{\beta=1}^2} E_B \left[D_{\beta=1}^1 O_{\beta=1}^0 - O_{\beta=1}^0 \right] = B^0$$

This problem is solved by $O_{\beta=1}^1 B_{\beta=1}^0 = B_{\beta=1}^0$ and accordingly $C^ = 1 U^0$.*

Proof. Note that for each choice of $f: \mathcal{O} \rightarrow \mathcal{B}^{\circ}$

$$E_B \left[\sum_{\theta=1}^{\#} \tilde{O}_\theta \cdot D_\theta \cdot \tilde{B}_\theta \right] = \sum_{\theta=1}^{\#} \tilde{O}_\theta \cdot U_\theta \cdot \tilde{B}_\theta \quad (4)$$

where \tilde{B}_θ is the unconditional probability of player i guessing the correct action B given \tilde{O}_θ . When $A \geq \frac{1}{2}$ and $U \geq 0$ -the expression in (4) is maximized by $\tilde{B}_\theta = 1$ for each θ . The only way to achieve $\tilde{B}_\theta = 1$ is by setting $\tilde{O}_\theta = B$. The corresponding expected payoff of each agent is $C_i = 1 - U$.

Proposition A.3. Consider the game Γ extended with a finite set \mathcal{W} of messages that agent i can send to agent j upon observing her private signal. Let $\tilde{\Gamma}$ denote the extended game. In every weak Perfect Bayesian Equilibrium of $\tilde{\Gamma}$ for each message $F_j \in \mathcal{W}$ player i chooses the action according to her private signal only: $\tilde{O}_i(F_j) = B_i$.

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

$$C_i = \sum_{\theta=1}^{\#} \tilde{O}_\theta \cdot U_\theta$$

$$C_j = \sum_{\theta=1}^{\#} \tilde{O}_\theta \cdot U_\theta$$

where \tilde{O}_θ is the probability of player i guessing the correct action. Since each player can guarantee herself a correct guess with probability A based on the private information only, it must be that $\tilde{O}_\theta \geq A$.

Now suppose that there exists a message F_j sent by type B_j of agent j in equilibrium, such that some type B_i of agent i chooses action $1 \neq B_i$.

There are two cases: either (i) the type $1 \neq B_i$ continues to choose action $1 \neq B_i$ upon observing message F_j or (ii) the type $1 \neq B_i$ chooses action $1 \neq B_i$ upon observing message F_j . In both of these cases player j can deceive agent i to choose the correct action with probability less than A .

- (i) If type $1 \neq B_i$ continues to choose action $1 \neq B_i$ upon observing message F_j - then by sending message F_j irrespective of her own signal, agent j induces type B_i of agent i to do action $1 \neq B_i$ both in case of coinciding or non-coinciding signals, while type $1 \neq B_i$ continues to act on private information only. Thus in case of j always sending message F_j , the probability of i guessing the correct action is $\tilde{O}_i = 1 - A < A$. Since for a fixed strategy of i the probability of j guessing the correct action \tilde{O}_j remains constant, agent j has a profitable deviation.
- (ii) If type $1 \neq B_i$ chooses action $1 \neq B_i$ upon observing message F_j - then by sending message F_j irrespective of her signal, agent j induces all types of agent i to do action $1 \neq B_i$ and the probability of i guessing the correct action in case of such a deviation is $\tilde{O}_i = 1 - A < A$. Again, since for a fixed strategy of i the probability of j guessing the correct action \tilde{O}_j remains constant, agent j has a profitable deviation.

Thus if there exists a message F_j that induces some type of agent i to choose action $1 \neq B_i$ - agent j necessarily has a profitable deviation. Therefore, it must be that in every every weak Perfect Bayesian Equilibrium of the game $\tilde{\Gamma}$ 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 Γ . In the game, each agent $i \in \{1, 2\}$ obtains

a binary signals $B: 2^S = \{0, 1\}^2$ with the following joint distribution P over $S^2 = S \times S$ parametrized by $A \in [0, 1]^2$:

	c	
	$s_2 = 0$	$s_2 = 1$
$s_1 = 0$	$\frac{A}{2}$	$\frac{1-A}{2}$
$s_1 = 1$	$\frac{1-A}{2}$	$\frac{A}{2}$

Together, these signals determine the correct action

$$B = \frac{1}{2} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \cdot B.$$

Both agents would like to guess B by choosing an action in the set $A = \{0, 1\}^2$. 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

$$D: U_1(a, B) = \mathbb{1}\{a = B\} \quad U_2(a, B) = \mathbb{1}\{a \neq B\} \tag{5}$$

where $U \in [0, 1]^2$.

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 μ is defined as a function from the product type space into the joint distributions over the action recommendations $\mu: S \times S \rightarrow \Delta(A)$. 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 recommendations¹¹.

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

$$\max_{\mu \in \mathcal{M}} \sum_{i=1}^2 c_i^* \tag{6}$$

The following sequence of lemmas first simplifies the problem in (6) and leads to Theorem B.1 which presents the optimal mediation protocol.

Lemma B.1 below shows that, when solving (6), 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 $\mu \in \mathcal{M}$ there exists $\mu^0 \in \mathcal{M}$ such that

(i) $\mu^0: S \times S \rightarrow \Delta(A \times A)$

¹¹The solution concept corresponds to *information design with elicitation* in terminology of Bergemann and Morris (2018).

(ii) $\mu^1_{\mathcal{B}} \text{ and } \mu^0_{\mathcal{B}}$ have the same marginal distributions for every \mathcal{B}

(iii) $c^0_i = c^1_i$ for every i

Proof. Consider an incentive-compatible mediation protocol $\mu \in \mathcal{M}$. Define μ^0 to be a mapping from \mathcal{S} to $\mathcal{A}^1 \times \mathcal{A}^0$ such that for each $\mathcal{B} \in \mathcal{S}$ - for each vector of recommendations O

$$\mu^0_{\mathcal{B}}(O) = \prod_{i \in \mathcal{I}} \mu^i_{\mathcal{B}}(O_i)$$

That is, μ^0 is defined to be the product of the marginal distributions of μ for each pair of reports \mathcal{B} . By construction, (i) $\mu^0 \in \mathcal{M}$ and (ii) μ and μ^0 have the same marginal distributions.

Now, since $\mu \in \mathcal{M}$ - it must also be that $\mu^0 \in \mathcal{M}$. To see this observe first that for type \mathcal{B}_i of agent i that submitted report \mathcal{B}_i and received a recommendation O_i - agent i 's posterior over the signal of agent $j \neq i$ [and thus also i 's preferred action] is pinned down by the marginal of recommendation distributions. Thus for each report \mathcal{B}_i of type \mathcal{B}_i of agent i under μ^0 - the resulting distribution of i 's actions is the same as under μ and consequently the probability of type \mathcal{B}_i agent i making a correct guess is the same under μ and μ^0 . Similarly, for each report \mathcal{B}_j of agent j under μ^0 - the resulting distribution of j 's actions is the same as under μ . Since the distribution of both agents' actions pin down the expected payoff of agent i - type \mathcal{B}_i of agent i has the same expected payoff for each report \mathcal{B}_i under μ and μ^0 . By assumption there were no profitable deviations from truth-telling under μ - thus so is the case under μ^0 . It is therefore proved that $\mu^0 \in \mathcal{M}$ and (iii) $c^0_i = c^1_i$.

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 $1 - \mathcal{B}_i$ to agent i who reported type \mathcal{B}_i . That is, an action that is known to be incorrect by a particular agent is never recommended to this agent.

LEMMA B.2. For every $\mu \in \mathcal{M}$ it must be that $\mu^1_{\mathcal{B}_i}(1 - \mathcal{B}_i) = 0$.

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

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

		μ	
		$\hat{\mathbf{s}}_2 = \mathbf{0}$	$\hat{\mathbf{s}}_2 = \mathbf{1}$
$\hat{\mathbf{s}}_1 = \mathbf{0}$	$\langle \mu_{00} \rangle - \langle \mu_{01} \rangle$	$\langle \mu_{00} \rangle - \langle \mu_{01} \rangle$	$\langle \mu_{01} \rangle - \langle \mu_{10} \rangle$
$\hat{\mathbf{s}}_1 = \mathbf{1}$	$\langle \mu_{10} \rangle - \langle \mu_{11} \rangle$	$\langle \mu_{10} \rangle - \langle \mu_{11} \rangle$	$\langle \mu_{11} \rangle - \langle \mu_{11} \rangle$

In the table $\langle \cdot \rangle_{\theta, \theta}$ is the probability of recommending the correct action to agent i when agent i reported θ and agent j reported θ ¹². The vector $\langle \cdot \rangle$ consisting of $\langle \cdot \rangle_{\theta, \theta}$ is *feasible* if the corresponding μ belongs to M .

Lemma B.3 below simplifies the problem in (6) 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.

Claim: A mediation protocol μ defined by a vector $\langle \cdot \rangle_{\theta, \theta}$ belongs to the set of IC mediation protocols M if and only if the following two sets of condition hold. First, for each i and θ

$$\langle \cdot \rangle_{\theta, \theta} A_i > 1 \quad \langle \cdot \rangle_{\theta-1, \theta}^{\theta} 1 \quad A_i^0 \tag{7}$$

$$\langle \cdot \rangle_{\theta-1, \theta} 1 \quad A_i^0 > 1 \quad \langle \cdot \rangle_{\theta, \theta}^{\theta} A_i \tag{8}$$

Second, for each i and θ

$$\begin{aligned} & \langle \cdot \rangle_{\theta, \theta} U_i < \langle \cdot \rangle_{\theta, \theta}^{\theta} A_i \quad \langle \cdot \rangle_{\theta-1, \theta} U_i < \langle \cdot \rangle_{\theta-1, \theta}^{\theta} 1 \quad A_i^0 > \\ & A_i^1 1 \quad U_i^1 < \langle \cdot \rangle_{\theta-1, \theta}^{\theta} 1 \quad A_i^0 1 \quad U_i^1 < \langle \cdot \rangle_{\theta-1, \theta}^{\theta} 1 \quad A_i^0 \quad \text{if } 1 \quad \langle \cdot \rangle_{\theta, \theta}^{\theta} A_i < \langle \cdot \rangle_{\theta-1, \theta} 1 \quad A_i^0 \\ & A_i^1 < \langle \cdot \rangle_{\theta, \theta} U_i^1 < \langle \cdot \rangle_{\theta-1, \theta}^{\theta} 1 \quad A_i^0 < \langle \cdot \rangle_{\theta-1, \theta} U_i^1 < \langle \cdot \rangle_{\theta-1, \theta}^{\theta} 1 \quad A_i^0 \quad \text{if } 1 \quad \langle \cdot \rangle_{\theta, \theta}^{\theta} A_i < \langle \cdot \rangle_{\theta-1, \theta} 1 \quad A_i^0 \end{aligned} \tag{9}$$

Conditions (7)-(8) ensure that each agent finds it profitable to follow the mediation protocol’s recommendation. Conditions (9) ensure that there are no profitable deviations from reporting truthfully to the mediation protocol.

Proof: The two sets of IC conditions (7)-(8) and (9) are established separately.

“Following recommendation” conditions Under an IC direct revelation mediation protocol, agent i 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 i .

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

$$\begin{aligned} P(B = \theta-\theta' | \theta' = \theta) &= P(B_3 = \theta | \theta' = \theta) = P(B_3 = \theta | \theta' = \theta) \\ &= \frac{P(B_3 = \theta - B_3 = \theta - \theta' = \theta)}{P(B_3 = \theta - \theta' = \theta)} \\ &= \frac{\langle \cdot \rangle_{\theta, \theta} \frac{A_i}{2}}{\langle \cdot \rangle_{\theta, \theta} \frac{A_i}{2} + 1 \cdot \langle \cdot \rangle_{\theta-1, \theta}^{\theta} \frac{1-A_i}{2}} \end{aligned}$$

Agent i following recommendation means that $B = \theta-\theta'$ is more likely than $B = \theta-1-\theta'$, which yields the first IC constraint

$$\langle \cdot \rangle_{\theta, \theta} A_i > 1 \quad \langle \cdot \rangle_{\theta-1, \theta}^{\theta} 1 \quad A_i^0$$

and condition (7) is established.

¹²The correct action recommendation to agent i in this case is $\theta-\theta'$ in case of truthful reporting, while the incorrect action recommendation to agent i is $\theta-1-\theta'$.

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

$$\begin{aligned}
 P(B_3 = \theta-1 \mid \theta, \theta') &= \frac{P(\theta-1, \theta')}{P(\theta-1, \theta') + P(\theta, \theta')} = \frac{1}{1 + \frac{P(\theta, \theta')}{P(\theta-1, \theta')}} \\
 &= \frac{1}{1 + \frac{P(\theta, \theta-1)}{P(\theta-1, \theta-1)}} = \frac{1}{1 + \frac{1}{2} \frac{A}{\theta-1}} \\
 &= \frac{2(\theta-1)}{2(\theta-1) + A}
 \end{aligned}$$

Agent 1 following recommendation means that $(\theta-1, \theta')$ is more likely than (θ, θ') , which yields the second IC constraint

$$\frac{2(\theta-1)}{2(\theta-1) + A} > \frac{1}{2}$$

and condition (8) is established.

The “following recommendation” conditions are thus established.

“Truthful reporting” conditions Under an IC direct revelation mediation protocol, agent 1: 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 1:.

- Suppose agent 1: received a signal θ . If she reports [T]ruthfully and follows the recommendation by the mediation protocol (and so does agent 3:), the expected payoff is

$$\begin{aligned}
 E(c_1 \mid \theta) &= E(c_1 \mid B_3 = \theta) = E(c_1 \mid B_3 = \theta-1) = \theta \\
 &= \frac{1}{2} \theta + \frac{1}{2} (\theta-1) = \frac{2\theta-1}{2}
 \end{aligned}$$

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

- What is the optimal action if $\theta-1$ is recommended back by the mediation protocol? The conditional probability of agent 3: having a signal θ is equal to

$$\begin{aligned}
 P^*(B_3 = \theta-1 \mid \theta-1) &= \frac{P(\theta-1, \theta-1)}{P(\theta-1, \theta-1) + P(\theta, \theta-1)} \\
 &= \frac{1}{1 + \frac{1}{2} \frac{A}{\theta-1}}
 \end{aligned}$$

where P^* stands for the updated probabilities given an [U]ntruthful report.

Similarly, the conditional probability of agent 3: having a signal $\theta-1$ is equal to

$$P^*(B_3 = \theta-1 \mid \theta-1) = \frac{\frac{1}{2} \frac{A}{\theta-1}}{1 + \frac{1}{2} \frac{A}{\theta-1}}$$

Consequently, if $\frac{1}{2} \frac{A}{\theta-1} > \frac{1}{2}$ - agent 1: will do action θ upon receiving signal $\theta-1$ and if $\frac{1}{2} \frac{A}{\theta-1} < \frac{1}{2}$ - agent 1: will do action 1•2 upon receiving signal $\theta-1$.

- What is the optimal action if 1•2 is recommended back by the mediation protocol? The conditional probability of agent 1•3 having a signal θ is equal to

$$P^*_{B_3 : = \theta} = \frac{P^*_{B_3 : = \theta} \cdot P_{1 \cdot 2 : = \theta}}{P^*_{1 \cdot 2 : = \theta}} = \frac{P^*_{B_3 : = \theta} \cdot \frac{A}{2}}{P^*_{1 \cdot 2 : = \theta} \cdot \frac{A}{2}}$$

while the conditional probability of agent 1•3 having a signal 1 - θ is equal to

$$P^*_{B_3 : = 1 - \theta} = \frac{P^*_{B_3 : = 1 - \theta} \cdot P_{1 \cdot 2 : = 1 - \theta}}{P^*_{1 \cdot 2 : = 1 - \theta}} = \frac{P^*_{B_3 : = 1 - \theta} \cdot \frac{A}{2}}{P^*_{1 \cdot 2 : = 1 - \theta} \cdot \frac{A}{2}}$$

Note that the IC constraint (8) and Assumption 1 imply that $P^*_{B_3 : = \theta} > P^*_{B_3 : = 1 - \theta}$ and thus agent 1•3 prefers to do action θ upon hearing a recommendation of 1•2 from the mediation protocol.

- Now [after trivially calculating the probabilities of agent 1•3 making the correct guess], the expected payoff of agent 1•3 in case of misreporting and sending 1 - θ instead of θ can be computed. If agent 1•3 does θ in any case and gets an expected utility of

$$C^*_{1 \cdot 3 : \theta} = A^1 U^1_{1 \cdot 3 : \theta} + A^0 U^0_{1 \cdot 3 : \theta}$$

If agent 1•3 does 1/2 in case of hearing a recommendation of 1 - θ and θ in case of recommendation of 1•2 and gets an expected utility of

$$C^*_{1 \cdot 3 : 1/2} = A^1 \left(\frac{1}{2} U^1_{1 \cdot 3 : \theta} + \frac{1}{2} U^1_{1 \cdot 3 : 1 - \theta} \right) + A^0 \left(\frac{1}{2} U^0_{1 \cdot 3 : \theta} + \frac{1}{2} U^0_{1 \cdot 3 : 1 - \theta} \right)$$

Thus the IC constraint for agent 1•3 reporting truthfully upon observing signal θ is

$$C^*_{1 \cdot 3 : \theta} > C^*_{1 \cdot 3 : 1/2} \iff A^1 U^1_{1 \cdot 3 : \theta} + A^0 U^0_{1 \cdot 3 : \theta} > A^1 \left(\frac{1}{2} U^1_{1 \cdot 3 : \theta} + \frac{1}{2} U^1_{1 \cdot 3 : 1 - \theta} \right) + A^0 \left(\frac{1}{2} U^0_{1 \cdot 3 : \theta} + \frac{1}{2} U^0_{1 \cdot 3 : 1 - \theta} \right)$$

and “truthful reporting” conditions (9) are established.

The proof of the claim is now completed. z

LEMMA B.3. *There exists a solution to the optimal mediation protocol design problem in eq. (6) such that $P^*_{1 \cdot 2 : = \theta} > P^*_{1 \cdot 2 : = 1 - \theta}$ and $P^*_{B_3 : = \theta} > P^*_{B_3 : = 1 - \theta}$.*

Proof: Note first that the optimal mediation protocol design problem

$$\max_{\theta} C^*_{1 \cdot 3 : \theta} \tag{10}$$

can be written more explicitly as

$$\max_{\theta} \frac{A}{2} \left(U^1_{1 \cdot 3 : \theta} + U^0_{1 \cdot 3 : \theta} \right) > \frac{A}{2} \left(U^1_{1 \cdot 3 : 1 - \theta} + U^0_{1 \cdot 3 : 1 - \theta} \right) \tag{11}$$

s. t. (7)-(9)

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

$$\begin{aligned}
 & \langle \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \rangle \\
 & A^1 \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \circ 1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \quad \text{if } \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0 \\
 & A^1 \langle \cdot \rangle_{\delta-1 \delta}^1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \circ 1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \quad \text{if } \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0
 \end{aligned} \tag{12}$$

and consider two vectors ℓ and B that satisfy those constraints.

If both ℓ and B are such that $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ and $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ or $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ and $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ - then the convex combination $D = V\ell + V'B$ with $V \geq 0, V' \geq 0, V + V' = 1$ also satisfies the constraint (12) with $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ or $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ respectively by linearity of both RHS and LHS of the inequality in (12).

Now suppose that

$$\begin{aligned}
 & \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0 \\
 & \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \rangle A^1 \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \circ 1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \\
 & \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0 \\
 & \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \rangle A^1 \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \circ 1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0
 \end{aligned}$$

and let $D = V\ell + V'B$ with $V \geq 0, V' \geq 0, V + V' = 1$.

First, either $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ or $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$. Suppose $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ - then to show convexity of the set defined by constraint (12), one needs to establish that

$$\langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \rangle A^1 \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \circ 1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \tag{13}$$

Indeed, if this is the case, then also

$$\langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \rangle A^1 \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \circ 1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0$$

since D is a convex combination of ℓ and B .

To show (13) note that

$$\begin{aligned}
 & \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \rangle A^1 \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \circ 1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \\
 & \qquad \qquad \qquad > A^1 \langle \cdot \rangle_{\delta\delta}^1 \quad U \langle \cdot \rangle_{\delta\delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0 \circ 1 \quad U \langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ 1 \quad A^0
 \end{aligned}$$

where the last inequality follows from $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$. The case of $\langle \cdot \rangle_{\delta-1 \delta}^1 \quad \langle \cdot \rangle_{\delta-1 \delta}^3 \rangle \circ A \succ \langle \cdot \rangle_{\delta-1 \delta}^1 \quad A^0$ is similar.

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

Symmetry of the objective (11) and the constraint set (7)-(9) Note that if the maximization problem is solved by some vector $\langle \cdot \rangle = \langle \cdot \rangle^1 - \langle \cdot \rangle^0$ (where $\langle \cdot \rangle^\delta$ denotes the subvector of probabilities related to agent δ), then vector $\langle \cdot \rangle^0 = \langle \cdot \rangle^2 - \langle \cdot \rangle^1$ leads to the same value of the objective function and constraints are satisfied at $\langle \cdot \rangle^0$ by symmetry. Thus the same value of the objective is achieved at the average of $\langle \cdot \rangle - \langle \cdot \rangle^0$ and one can restrict attention to maximizing with 4-element vector $\langle \cdot \rangle_{00} - \langle \cdot \rangle_{01} - \langle \cdot \rangle_{10} - \langle \cdot \rangle_{11}$. Again, swapping $\langle \cdot \rangle_{00}$ with $\langle \cdot \rangle_{11}$

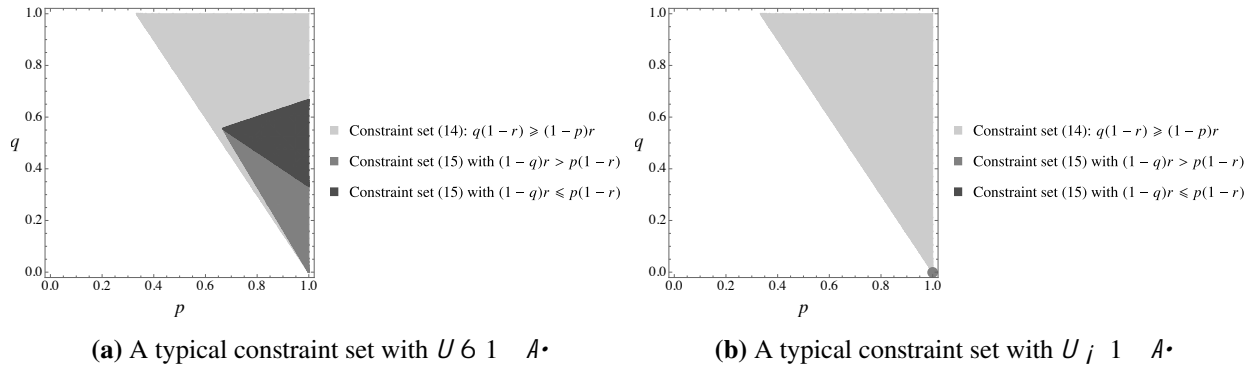


Figure 1: Typical constraint sets.

and α_1 with α_0 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 (α, β) with $\alpha = \alpha_0$ and $\beta = \alpha_1$ and the proof of the lemma is now completed. \square

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. (6) is solved by $(\alpha, \beta) \in M$ such that $\alpha = \alpha_0$ and $\beta = \alpha_1$ with*

$$\alpha = 1 - \beta = \frac{2A - 1}{2A - 1 - U} \quad \text{if } U < 1 - A$$

$$\alpha = 1 - \beta = 0 \quad \text{if } U > 1 - A$$

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

$$\max_{\alpha, \beta \in [0, 1]^2} \alpha A + \beta (1 - A)$$

subject to

$$\alpha A + \beta (1 - A) > \alpha A + \beta (1 - A) \tag{14}$$

$$\alpha A + \beta (1 - A) > \alpha A + \beta (1 - A) \tag{15}$$

$$\alpha A + \beta (1 - A) > \alpha A + \beta (1 - A) \tag{16}$$

For a moment, ignore the first two constraints (14)-(15). 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 α along the line $\alpha A + \beta (1 - A) = \alpha A + \beta (1 - A)$ under Assumption 1. Indeed, substituting

$$\beta = 1 - \alpha \frac{1 - A}{A}$$

into the objective yields coefficient equal to $2 - \frac{1 - A}{A}$ on the variable α and thus the objective grows in α . Having observed this, it is easy to see that the value of the objective in the region with $\alpha A + \beta (1 - A) = \alpha A + \beta (1 - A)$ is not higher than at the point point on its boundary with the highest value of α , which is $1 - \frac{2A - 1}{A}$.

- Now also note that the objective grows with α along the line

$$\alpha A + \beta (1 - A) = \alpha A + \beta (1 - A)$$

Indeed, substituting

$$\alpha = \frac{1 - 2A - U}{1 - 2A - U} \frac{U}{1 - 2A - U}$$

into the objective yields coefficient $1 - \frac{2A - 1}{2A - U}$ on α . Thus if the point $(1 - \frac{2A - 1}{2A - U}, \frac{U}{1 - 2A - U})$ satisfies $1 - \alpha \geq \beta$, it has the highest value of the objective in the region with $1 - \alpha \geq \beta$.

- Note that $U \leq 1 - A$ simultaneously guarantees that $(1 - \frac{2A - 1}{2A - U}, \frac{U}{1 - 2A - U})$ satisfies $1 - \alpha \geq \beta$ and has a higher value of the objective than $(1 - \frac{2A - 1}{A}, 0)$. Moreover, the point $(1 - \frac{2A - 1}{2A - U}, \frac{U}{1 - 2A - U})$ satisfies the two omitted constraints (14)-(15) of the maximization problem. Thus the point that maximizes the objective for $U \leq 1 - A$ is $(1 - \frac{2A - 1}{2A - U}, \frac{U}{1 - 2A - U})$.
- If in turn $U > 1 - A$ then there are no points that satisfy the constraint with $1 - \alpha \geq \beta$: the set defined by $1 - \alpha \geq \beta$ has no intersection with the set defined by $1 - \alpha \geq \beta$ - which is easy to verify by comparing the values of the linear constraints at the boundary points of the constraint set $\alpha = 0$ and $\alpha = 1$. Moreover, the only point satisfying the constraint with $1 - \alpha \geq \beta$ is $(1 - 0, 0)$ - which remains the only candidate for the optimal mediation protocol when $U > 1 - A$. This point obviously satisfies the omitted constraints (14)-(15).
- For a graphical treatment, two typical constraint sets are shown in Appendix B.

The search for the optimal mediation protocol is thus completed

$$\alpha = 1 - \frac{2A - 1}{2A - U} \quad \text{if } U \leq 1 - A$$

$$\alpha = 0 \quad \text{if } U > 1 - A$$

and Theorem B.1 is proved.

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

Proof. Table 1 below explicitly presents α^* 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) $\alpha \leq 1 - \beta$; (ii) the “truthful reporting” constraints in (9).

α	"	r of obj.	r of constraints								
			$\alpha \leq 1 - \beta$				"truthful reporting"				
$\alpha = 1$?	$A^1 1 - U^0$	1	0	0	0	0	A	UA	$1 - A$	$U^1 1 - A^0$
$\alpha = 1 - \frac{2A - 1}{2A - U}$	@	$1 - A^0 1 - U^0$	0	0	0	0	0	$1 - A^0$	UA	A	$U^1 1 - A^0$
$\alpha = 1 - \frac{2A - 1}{2A - U}$	@	$1 - A^0 1 - U^0$	0	0	0	0	0	A	$U^1 1 - A^0$	$1 - A^0$	UA
$\alpha = 1 - \frac{2A - 1}{2A - U}$?	$A^1 1 - U^0$	0	1	0	0	0	$1 - A$	$U^1 1 - A^0$	A	UA
$\alpha = 1 - \frac{2A - 1}{2A - U}$?	$A^1 1 - U^0$	0	0	1	0	0	UA	A	$U^1 1 - A^0$	$1 - A$
$\alpha = 1 - \frac{2A - 1}{2A - U}$	@	$1 - A^0 1 - U^0$	0	0	0	0	0	UA	$1 - A^0$	$U^1 1 - A^0$	A
$\alpha = 1 - \frac{2A - 1}{2A - U}$	@	$1 - A^0 1 - U^0$	0	0	0	0	0	$U^1 1 - A^0$	A	UA	$1 - A^0$
$\alpha = 1 - \frac{2A - 1}{2A - U}$?	$A^1 1 - U^0$	0	0	0	1	0	$U^1 1 - A^0$	$1 - A$	UA	A

Table 1: Value of objective, gradients of objective and constraints.

