

# Persuasion via Weak Institutions

---

Elliot Lipnowski

*Columbia University*

Doron Ravid

*University of Chicago*

Denis Shishkin

*University of California San Diego*

A sender commissions a study to persuade a receiver but influences the report with some probability. We show that increasing this probability can benefit the receiver and can lead to a discontinuous drop in the sender's payoffs. To derive our results, we geometrically characterize the sender's highest equilibrium payoff, which is based on the concavification of a capped value function.

## I. Introduction

Many institutions routinely collect and disseminate information. Although the collected information is instrumental to its consumers, often the main goal of dissemination is to persuade. Persuading one's audience, however,

Lipnowski and Ravid acknowledge support from the National Science Foundation (grant SES-1730168). We thank Roland Bénabou, Ben Brooks, Joyee Deb, Eddie Dekel, Wouter Dessein, Jon Eguia, Emir Kamenica, Navin Kartik, Stephen Morris, Pietro Ortoleva, Wolfgang Pesendorfer, Carlo Prato, Marzena Rostek, Evan Sadler, Zichang Wang, Richard Van Weelden, Leat Yariv, and various audiences for useful suggestions. We also thank

Electronically published August 18, 2022

*Journal of Political Economy*, volume 130, number 10, October 2022.

© 2022 The University of Chicago. This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License (CC BY-NC 4.0), which permits non-commercial reuse of the work with attribution. For commercial use, contact [journalpermissions@press.uchicago.edu](mailto:journalpermissions@press.uchicago.edu). Published by The University of Chicago Press.

<https://doi.org/10.1086/720462>

requires the audience to believe what one says. In other words, the institution must be credible, meaning it must be capable of delivering both good and bad news. Yet if the institution is not independent from its superiors, delivering unfavorable news might be especially difficult. This paper studies how an institution's credibility influences its persuasiveness and the quality of information it provides.

For concreteness, consider a head of state who wants to sway a firm to invest as much as possible in her country's economy. The firm can make a large investment (2), a small investment (1), or no investment (0). Whereas the country's leader wants to maximize the firm's expected investment, the firm's net benefit from investing depends on the state of the economy, which can be either *good* or *bad*. When the economy is *good*, the firm makes a profit of 1 from a large investment and 3/4 from a small investment. Investing in a *bad* economy results in losses, yielding the firm a payoff of  $-1$  and  $-1/4$  from a large and small investment, respectively. Not investing always generates a payoff of 0 to the firm, regardless of the state. Therefore, the firm will make a large (no) investment whenever it assigns a probability of at least 3/4 to the economy being *good* (*bad*). For intermediate beliefs, the firm makes a small investment. The firm and the leader share a prior belief of  $\mathbb{P}(\textit{good}) = 1/2$  (fig. 1).

To persuade the firm to invest, the leader commissions a report by the country's central bank. By specifying the report's parameters—its data, methods, assumptions, focus, and so on—the leader controls what information the report is supposed to convey. Formally, the commissioned report is a signal structure,  $\xi(\cdot|\textit{good})$  and  $\xi(\cdot|\textit{bad})$ , specifying a distribution over messages that the firm observes conditional on the state if the report is conducted as announced. To execute the report as planned, however, the bank must withstand the leader's behind-the-scenes pressures; that is, the firm observes a message drawn from  $\xi$  only if the bank is independent, which occurs with probability  $\chi$ . With complementary probability, the bank is influenced, meaning it releases a message of the leader's choice. Once the message is realized, the firm observes it and chooses how much to invest without knowing whether the report is influenced.

When the central bank is fully credible,  $\chi = 1$ , it is committed to the official report. As such, the leader can communicate any information she chooses, and so this example falls within the framework of Kamenica and Gentzkow (2011). Using their results, one can deduce that the policy maker optimally chooses a symmetric binary signal,

$$\begin{aligned}\xi_1^*(g|\textit{good}) &= 3/4, \xi_1^*(g|\textit{bad}) = 1/4, \\ \xi_1^*(b|\textit{good}) &= 1/4, \xi_1^*(b|\textit{bad}) = 3/4.\end{aligned}$$

---

Chris Baker, Sulagna Dasgupta, Takuma Habu, and Elena Istomina for excellent research assistance. This paper was edited by Emir Kamenica.

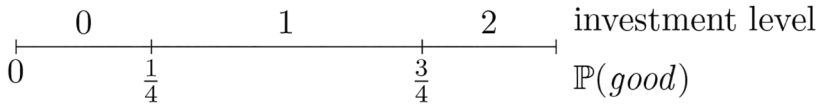


FIG. 1.—Firm’s best response in central bank example.

Under this signal structure, the firm is willing to invest 2 following a *g* signal, and 1 following a *b* signal. Ex ante, the two signals occur with equal probability, leading the firm to invest 3/2 on average.

If the central bank were weaker, its messages would be less persuasive because the firm would no longer take them at face value. To illustrate, suppose that  $\chi = 2/3$  and that the leader commissioned the same report as under full credibility. In this case, the firm could not possibly make a large investment after seeing *g*; otherwise, the leader would always send *g* when influencing the report, which would make a small investment strictly better for the firm. Thus, when  $\chi = 2/3$ , the leader’s full-commitment report is not sufficiently persuasive to increase the firm’s involvement in the local economy beyond its no-information investment of 1.

The leader can, however, overcome the firm’s skepticism by asking the bank to release more information. In fact, when  $\chi = 2/3$ , commissioning a fully revealing report that sends *g* if and only if the economy is good is optimal for the leader. In the resulting equilibrium, the leader always sends *g* when influencing the report, whereas the firm makes a large investment when seeing *g* and invests nothing otherwise. The reason the firm invests 2 upon seeing *g* is that the bank’s official report is so informative that a *g* message results in the firm believing the economy is good with probability 3/4 despite the leader’s possible interference. Because the firm sees the *g* message with probability 2/3, it invests 4/3 on average in the leader’s economy.

Since a weaker central bank results in the leader commissioning a more informative report, the firm may benefit from a reduction in the bank’s credibility. To illustrate, observe that when  $\chi = 1$ , the firm is no better off with the leader’s report than it was without it: in either case, the firm expects a profit of 1/4. By contrast, when  $\chi = 2/3$ , the firm strictly benefits from the leader’s communications, making an expected profit of 1/2 from investing 2 after seeing *g* and not investing otherwise. On average, the firm’s profit equals 1/3. Thus, the leader responds to the central bank’s decreased credibility by commissioning a report whose informativeness more than compensates the firm for the central bank’s increased susceptibility.

To understand examples such as the one above, we study a general model of strategic communication between a receiver (he) and a sender (she) who cares about only the receiver’s action. The receiver’s preferences over

his actions depend on an unknown state,  $\theta$ . To learn about  $\theta$ , the receiver relies on information provided by an institution under the sender's control. The game begins with the sender publicly announcing an official reporting protocol, which is an informative signal about the state. With probability  $\chi$ , the sender's institution is independent, delivering the receiver a message drawn according to the originally announced protocol. With complementary probability, the report is influenced: the sender learns the state and chooses what message to send to the receiver. Seeing the message (but not its origin), the receiver takes an action. Thus,  $\chi$  represents the credibility of the sender's institution, that is, the institution's ability to resist interference by its superiors.

At the extremes, our framework specializes to two prominent models of information transmission. When  $\chi = 1$ , the sender can never influence the report, so our setting reduces to one in which the sender publicly commits to her communication protocol at the beginning of the game. In other words, under full credibility, our model is equivalent to Bayesian persuasion (Kamenica and Gentzkow 2011). When  $\chi = 0$ , the receiver knows the sender is choosing the report's message *ex post*. Because messages are costless, they are just cheap talk (Crawford and Sobel 1982; Green and Stokey 2007), meaning that our no-credibility case corresponds to a cheap-talk game with state-independent preferences (Chakraborty and Harbaugh 2010; Lipnowski and Ravid 2020).

The corner cases of our model lend themselves to geometric analysis. Let the sender's *value function* be the highest value the sender can obtain from the receiver responding optimally at a given posterior belief. Kamenica and Gentzkow (2011) show that concavifying this function gives the sender's largest equilibrium payoff in the Bayesian persuasion model. More recently, Lipnowski and Ravid (2020) observe that as long as the sender cares about only the receiver's actions, quasiconcavifying the sender's value function delivers her highest equilibrium payoff under cheap talk.

Our theorem 1 uses the aforementioned geometric approach to characterize the sender's maximal equilibrium value in the intermediate credibility case,  $\chi \in (0, 1)$ . To do so, the theorem partitions the sender's equilibrium messages into two sets: messages the sender willingly sends when influencing the report (e.g.,  $g$  in the above example) and messages communicated only by the official report. One might guess that concavification and quasiconcavification characterize the sender's payoffs from official and influenced reporting, respectively. However, we show that whereas quasiconcavification characterizes the sender's payoffs from influenced reporting, one cannot find the sender's utility from official reporting using concavification alone. The reason is that the sender's payoff from a message cannot surpass the utility she obtains under compromised reporting: if it did, the sender would have a profitable

deviation. To account for this incentive constraint, one must cap the sender's value function at her utility from influenced reporting before concavifying it.

Using theorem 1, we explore how the use of weaker institutions affects persuasion. Proposition 1 identifies situations in which the receiver does better with a less credible sender. In particular, the proposition shows that such productive mistrust can occur when the sender wants to reveal intermediate information under full credibility. In such circumstances, a less credible sender may choose to commission a report that releases more news that is bad for her, so that the receiver believes messages that are good for the sender. We see this case in the central bank example above: when  $\chi = 1$ , the bank never fully reveals any state, whereas under  $\chi = 2/3$ , the report must occasionally reveal that the economy is bad in order to ensure that the firm invests 2 when seeing  $g$ .

Our next result, proposition 2, shows that small decreases in credibility can lead to large drops in the sender's value. More precisely, we show that such a collapse occurs at some full-support prior and some credibility level if and only if the sender can benefit from persuasion. Such a collapse is present in the above example: whenever  $\chi < 2/3$ , the leader cannot induce the firm to invest 2 even when she chooses to commission a fully revealing report. Thus, the best she can do when  $\chi < 2/3$  is to get an investment of 1 for sure by communicating no information—a drop of  $1/3$  from the  $4/3$  average investment the leader obtains when  $\chi$  is exactly  $2/3$ .

One may wonder if such collapses may occur at full credibility. Our proposition 3 shows that such a discontinuity can occur but only in knife-edge cases. Thus, although the sender's value often drops at some prior and some  $\chi$  because of small decreases in credibility, it rarely does so at  $\chi = 1$ .

*Related literature.*—This paper contributes to the literature on strategic information transmission. To place our work, consider two extreme benchmarks: full credibility and no credibility. Our full-credibility case is the model used in the Bayesian persuasion literature (Aumann and Maschler 1995; Kamenica and Gentzkow 2011; Kamenica 2019), which studies sender-receiver games in which a sender commits to an information transmission strategy. The no-credibility specialization of our model reduces to cheap talk (Crawford and Sobel 1982; Green and Stokey 2007). In particular, we build on Lipnowski and Ravid (2020), who use the belief-based approach to study cheap talk under state-independent sender preferences.

Two recent papers (Fréchette, Lizzeri, and Perego 2022; Min 2021) study closely related models. Fréchette, Lizzeri, and Perego (2022) test experimentally the connection between the informativeness of the sender's communication and her credibility in the binary state, binary action version of our model. Min (2021) looks at a generalization of our model in which the sender's preferences can be state dependent. He shows that

the sender weakly benefits from a higher commitment probability. Applying Blume, Board, and Kawamura's (2007) results on noisy communication, Min (2021) also shows that allowing the sender to commit with positive rather than zero probability strictly helps both players in Crawford and Sobel's (1982) uniform quadratic example.

Other thematically related work studies games of information transmission while varying the (exogenous or endogenous) limits to communication. Some such work focuses on games of direct communication, showing how some manner of commitment power can be sustained (for either a sender or a receiver) via lying costs (e.g., Kartik 2009; Guo and Shmaya 2021; Nguyen and Tan 2021), repeated interactions (e.g., Mathevet, Pearce, and Stacchetti 2022; Best and Quigley 2022), verifiable information (e.g., Glazer and Rubinstein 2006; Sher 2011; Hart, Kremer, and Perry 2017; Ben-Porath, Dekel, and Lipman 2019), informational control (e.g., Ivanov 2010; Luo and Rozenas 2018), or mediation (e.g., Goltzman et al. 2009; Salamanca 2021). Other work considers models in which a sender chooses an experiment *ex ante*, asking how persuasion can be shaped by exogenous experiment constraints (e.g., Ichihashi 2019; Perez-Richet and Skreta 2022) or by signaling motives (e.g., Perez-Richet 2014; Hedlund 2017; Alonso and Câmara 2018).

More broadly, weak institutions often serve as a justification for examining mechanism design under limited commitment (e.g., Bester and Strausz 2001; Skreta 2006). We complement this literature by relaxing a principal's commitment power in the control of information rather than incentives.

## II. A Weak Institution

We analyze a game with two players: a sender (she) and a receiver (he). Whereas both players' payoffs depend on the receiver's action,  $a \in A$ , the receiver's payoff also depends on an unknown state,  $\theta \in \Theta$ . Thus, the sender and the receiver have objectives  $u_S : A \rightarrow \mathbb{R}$  and  $u_R : A \times \Theta \rightarrow \mathbb{R}$ , respectively, and each aims to maximize expected payoffs.

The game begins with the sender commissioning a report,  $\xi : \Theta \rightarrow \Delta M$ , to be delivered by a research institution. The state then realizes, and the receiver sees a message  $m \in M$  (without observing  $\theta$ ). Given any  $\theta$ , the sender is credible with probability  $\chi$ , meaning  $m$  is drawn according to the official reporting protocol,  $\xi(\cdot|\theta)$ . With probability  $1 - \chi$ , the sender is not credible, in which case the sender decides which message to send after privately observing  $\theta$ . Only the sender learns her credibility type, and she learns it only after announcing the official reporting protocol.<sup>1</sup>

<sup>1</sup> In the appendix, we show that our payoff results are unchanged if the sender learns her credibility type before choosing the official report.

We impose some technical restrictions on our model.<sup>2</sup> Both  $A$  and  $\Theta$  are finite spaces with at least two elements. The state,  $\theta$ , follows some prior distribution  $\mu_0 \in \Delta\Theta$ , which is known to both players. Finally, we assume that  $M$  is rich enough to ensure that the sender faces no exogenous constraints on communication.<sup>3</sup>

We now define an equilibrium, which consists of four objects: the sender’s official reporting protocol,  $\xi : \Theta \rightarrow \Delta M$ , executed whenever the sender is credible; the strategy that the sender employs when not committed, that is, the sender’s influencing strategy,  $\sigma : \Theta \rightarrow \Delta M$ ; the receiver’s strategy,  $\alpha : M \rightarrow \Delta A$ ; and the receiver’s belief map,  $\pi : M \rightarrow \Delta\Theta$ , assigning a posterior belief to each message. A  $\chi$ -equilibrium is an official reporting policy announced by the sender,  $\xi$ , together with a perfect Bayesian equilibrium of the subgame following the sender’s announcement. Formally, a  $\chi$ -equilibrium is a tuple  $(\xi, \sigma, \alpha, \pi)$  of maps such that it is consistent with Bayesian updating, and both the receiver and the sender behave optimally; that is,

1. *Bayesian updating*: the belief map  $\pi : M \rightarrow \Delta\Theta$  satisfies Bayes’s rule given prior  $\mu_0$  and the message policy

$$\chi\xi + (1 - \chi)\sigma : \Theta \rightarrow \Delta M.$$

2. *Receiver optimality*: every  $m \in M$  has  $\alpha(m)$  supported on

$$\operatorname{argmax}_{a \in A} \sum_{\theta \in \Theta} u_R(a, \theta) \pi(\theta | m).$$

3. *Sender optimality*: every  $\theta \in \Theta$  has  $\sigma(\theta)$  supported on

$$\operatorname{argmax}_{m \in M} \sum_{a \in A} u_S(a) \alpha(a | m).$$

We view the sender as a principal capable of steering the receiver toward her favorite  $\chi$ -equilibria. In Lipnowski, Ravid, and Shishkin (2022), we define the notion of perfect Bayesian  $\chi$ -equilibrium in which we explicitly model the sender’s incentives at the experiment choice stage. By appropriately completing off-path play, that paper shows that the sender’s highest  $\chi$ -equilibrium payoff coincides with her highest perfect Bayesian  $\chi$ -equilibrium payoff.

<sup>2</sup> We view every topological space as a measurable space with its Borel field. For any measurable space  $Y$ , we denote by  $\Delta Y$  the set of all probability measures over  $Y$ . For any measurable spaces  $X, Y$ , a map  $X \rightarrow Y$  is a measurable function  $X \rightarrow Y$ .

<sup>3</sup> For example, we could take  $M = [0, 1]$  (see appendix). Moreover, corollary 1 in the appendix implies that the sender’s optimal equilibrium payoff would remain unchanged if  $M$  were instead finite with  $|M| \geq \min\{|A|, 2|\Theta| - 1\}$ .

### III. Persuasion with Partial Credibility

In this section, we characterize the sender's maximal  $\chi$ -equilibrium payoff. Our analysis applies the belief-based approach (Kamenica 2019; Forges 2020). Within an equilibrium, each message  $m$  that the sender communicates to the receiver induces a posterior belief  $\mu = \pi(m) \in \Delta\Theta$  and an expected sender utility from the receiver's (potentially mixed) action  $s = \sum_{a \in A} u_s(a)\alpha(a|m) \in \mathbb{R}$ . By replacing each message with its associated  $\mu$  and  $s$ , one can transform the equilibrium distribution of messages into its induced joint distribution  $\mathbf{P}$  of the receiver's beliefs and the sender's continuation payoffs. We refer to  $(\mu, s) \in \Delta\Theta \times \mathbb{R}$  as an *outcome*, and to a distribution  $\mathbf{P} \in \Delta(\Delta\Theta \times \mathbb{R})$  as an *outcome distribution*, and we define a  $\chi$ -equilibrium *outcome distribution* to be an outcome distribution induced by a  $\chi$ -equilibrium.

#### A. The Extreme Cases

We now review existing results that cover the extreme cases of our model. These cases serve as building blocks for proving our main theorem, which covers the case in which  $\chi$  is intermediate.

##### 1. Full Credibility

When  $\chi = 1$ , the sender's official announcement is binding, and so our model reduces to the Bayesian persuasion model of Kamenica and Gentzkow (2011). We now review some of their results. With full credibility, the sender is hampered by only two constraints. The first constraint is that the receiver updates his beliefs using Bayes's rule, which is equivalent to the receiver's posterior belief averaging to his prior. That is,  $\mathbf{P}$  must satisfy

$$\int \mu \, d\mathbf{P}(\mu, s) = \mu_0. \quad (\text{Bayes})$$

The second constraint is that the receiver must be best responding: for any belief the receiver holds, he must take only actions he finds optimal. To formalize this requirement, define the sender's *value correspondence* to be the correspondence mapping each posterior belief to the set of payoffs the sender can attain from the receiver-optimal behavior,<sup>4</sup>

$$V : \Delta\Theta \rightrightarrows \mathbb{R}$$

$$\mu \mapsto \text{co } u_s \left( \underset{a \in A}{\text{argmax}} \sum_{\theta \in \Theta} u_R(a, \theta) \mu(\theta) \right).$$

<sup>4</sup> The reason for the convex hull in  $V$ 's definition is that the receiver may choose to mix in the event that he has multiple best responses to a given belief.

Then,  $\mathbf{P}$  is compatible with the receiver’s incentive constraint if and only if  $\mathbf{P}$  is supported on the graph of  $V$ ; that is, a message can induce an outcome  $(\mu, s)$  only if  $s \in V(\mu)$ . Letting  $\text{gr } V := \{(\mu, s) : s \in V(\mu)\}$  denote the graph of  $V$ , we can state this constraint formally as

$$\mathbf{P}(\text{gr } V) = 1. \tag{R-IC}$$

As noted by Kamenica and Gentzkow (2011), the conditions (R-IC) and (Bayes) are together necessary and sufficient for an outcome distribution  $\mathbf{P}$  to arise from some 1-equilibrium. Denote the subset of  $\Delta(\Delta\Theta \times \mathbb{R})$  that satisfy these conditions for a prior  $\mu_0$  and value correspondence  $V$  by

$$\text{BP}(\mu_0, V) = \{\mathbf{P} \in \Delta(\Delta\Theta \times \mathbb{R}) : \mathbf{P} \text{ satisfies (Bayes) and (R-IC)}\}.$$

One can characterize the sender’s highest 1-equilibrium payoff using her *value function*,

$$\begin{aligned} v : \Delta\Theta &\rightarrow \mathbb{R} \\ \mu &\mapsto \max V(\mu), \end{aligned}$$

which maps every belief to the utility the sender obtains if the receiver chooses optimally and breaks ties in the sender’s favor given multiple best responses. Specifically, one can show that the sender’s utility in her favorite 1-equilibrium equals  $\hat{v}(\mu_0)$ , where

$$\hat{v} := \text{cav}(v)$$

is the lowest concave function that is everywhere above  $v$  (e.g., Aumann and Maschler 1995; Kamenica and Gentzkow 2011). The function  $\hat{v}$  is known as  $v$ ’s concavification.

Figure 2 illustrates the above in the context of the central bank example from the introduction. Because the state is binary, we identify the receiver’s posterior belief  $\mu$  with the probability it assigns to the economy being good. The left panel in figure 2 plots the sender’s value correspondence, taking  $\mu$

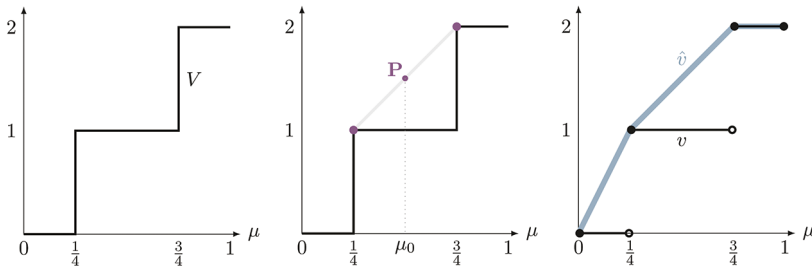


FIG. 2.—Value correspondence  $V$ , sender’s best 1-equilibrium outcome  $\mathbf{P}$ , and value function  $v$  with its concavification  $\hat{v}$  in central bank example.

as an input. For  $\mu < 1/4$ , the sender can only get a payoff of 0, whereas when  $\mu \in (1/4, 3/4)$ , she can only get 1, and when  $\mu > 3/4$ , she can only get 2. The sender can attain any payoff between 0 and 1 when  $\mu = 1/4$  and any payoff between 1 and 2 when  $\mu = 3/4$ . The middle panel depicts the sender's best 1-equilibrium outcome distribution  $\mathbf{P}$ , which assigns equal weight to the points  $(\mu, s) = (1/4, 1)$  and  $(\mu, s) = (3/4, 2)$ . As can be seen, both points lie on the graph of  $V$ , meaning that this distribution satisfies (R-IC). This distribution also satisfies (Bayes) because the average probability assigned to  $\theta = \textit{good}$  equals  $1/2$ , which is the probability assigned to that state by the prior. One can visually verify that this distribution is indeed sender optimal by examining the right panel, which shows the sender's value function along with its concave envelope,

$$v(\mu) = \begin{cases} 0 & \text{if } \mu \leq 1/4, \\ 1 & \text{if } \mu \in [1/4, 3/4], \\ 2 & \text{if } \mu \geq 3/4, \end{cases} \quad \hat{v}(\mu) = \begin{cases} 4\mu & \text{if } \mu \leq 1/4, \\ 1 + 2(\mu - 1/4) & \text{if } \mu \in [1/4, 3/4], \\ 2 & \text{if } \mu \geq 3/4. \end{cases} \quad (1)$$

As seen in the figure, the outcome distribution  $\mathbf{P}$  gives the sender an expected payoff of  $3/2$ , which is also the value of  $\hat{v}(\mu_0)$ , thereby confirming that  $\mathbf{P}$  is indeed sender optimal.

## 2. No Credibility

We now turn to the  $\chi = 0$  case, in which the receiver knows the sender is choosing  $m$  after observing the state. Being freely chosen, the sender's communication is cheap talk (Crawford and Sobel 1982; Green and Stokey 2007) and thus needs to satisfy the sender's incentive constraints. Our assumption that the sender's preferences are state independent simplifies these constraints considerably: the sender must be indifferent between all on-path messages. The reason is that if the sender's payoffs across two distinct messages differ, the sender will never (in any state) want to send the lower-payoff message. As such, the sender's payoff from all outcomes in the support of a 0-equilibrium outcome distribution must be the same. In other words, every 0-equilibrium outcome distribution  $\mathbf{P}$  must satisfy

$$\mathbf{P}\{\Delta\Theta \times \{s_i\}\} = 1 \text{ for some } s_i \in \mathbb{R}. \quad (\text{CP})$$

Combining (CP) with the restrictions imposed by Bayesian updating (Bayes) and the receiver incentives (R-IC), one obtains a full characterization of the attainable outcome distributions under no credibility (see Aumann and Hart 2003; Lipnowski and Ravid 2020). It follows that the sender's highest 0-equilibrium payoff is given by

$$\max_{\mathbf{P} \in \text{BP}(\mu_0, V)} \int s \, d\mathbf{P}(\mu, s) \text{ subject to (CP).} \tag{CT}$$

Lipnowski and Ravid (2020) show that this maximal payoff is equal to  $\bar{v}(\mu_0)$ , where

$$\bar{v} = \text{qcav}(v)$$

is  $v$ 's quasiconcavification, that is, the lowest quasiconcave function that is everywhere above  $v$ .

Figure 3 depicts  $v$ 's quasiconcavification and concavification, respectively, for some function  $v$ . These functions describe the sender's ability to benefit from communication by connecting points on the graph of the sender's value correspondence. With full credibility, the sender can connect such points using any affine segment. When  $\chi = 0$ , the sender's incentive constraints dictate that her payoff coordinate must remain constant; that is, the sender can use only flat segments.

Let us revisit the example from the introduction when  $\chi = 0$ . Observe that the optimal 1-equilibrium outcome distribution in this example does not satisfy (CP), because it generates two outcomes with different sender payoffs and so cannot be induced by a 0-equilibrium (see fig. 2, middle panel). We now argue that the sender cannot attain any value above 1 in any 0-equilibrium. One way of seeing this fact is to observe that the sender's value function in this example is quasiconcave and is therefore equal to its quasiconcavification. Alternatively, observe that (Bayes) requires every 0-equilibrium outcome distribution  $\mathbf{P}$  to induce at least one outcome with  $\mu \leq 1/2$ , whereas (R-IC) requires the sender's payoff from all such beliefs to be below 1. Because the sender's payoff must be constant over  $\mathbf{P}$ 's support by (CP), it follows that  $\mathbf{P}$  cannot induce a sender payoff strictly larger than 1.

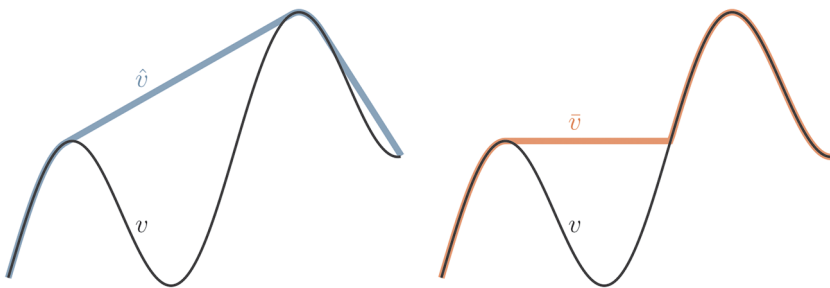


FIG. 3.—Value function  $v$  and its quasiconcavification  $\hat{v}$  and concavification  $\bar{v}$ .

*B. The Intermediate Credibility Case*

This section presents theorem 1, which geometrically characterizes the sender's optimal  $\chi$ -equilibrium value for our general model.

Suppose that credibility is not extreme ( $0 < \chi < 1$ ) so that both the official reporting protocol and the sender's influencing strategy are relevant, and let  $\mathbf{P}$  be a  $\chi$ -equilibrium outcome distribution. Notice that the receiver optimality and the Bayesian-updating conditions are as in the full- and no-credibility cases, and so  $\mathbf{P}$  must satisfy (Bayes) and (R-IC); that is,  $\mathbf{P} \in \text{BP}(\mu_0, V)$ . We now use these conditions to derive an upper bound on the sender's value from  $\mathbf{P}$ .

We begin by decomposing  $\mathbf{P}$  into two distributions. To do so, let

$$s_{\max} := \max\{s : (\mu, s) \in \text{supp}(\mathbf{P})\}$$

be the highest payoff in the support of  $\mathbf{P}$ , and let  $k \in [0, 1]$  denote the  $\mathbf{P}$ -probability of sender payoffs strictly below  $s_{\max}$ . In what follows, we focus on the case in which  $0 < k < 1$ .<sup>5</sup> Let  $\mathbf{G}$  be the distribution over outcomes induced by  $\mathbf{P}$  conditional on  $s = s_{\max}$ , and let  $\mathbf{B}$  be the outcome distribution conditional on  $s < s_{\max}$ . By construction,

$$\mathbf{P} = (1 - k)\mathbf{G} + k\mathbf{B}.$$

For an example, consider the optimal 1-equilibrium outcome distribution  $\mathbf{P}$  from the central bank example, which generates the outcomes  $(\mu, s) = (1/4, 1)$  and  $(\mu, s) = (3/4, 2)$  with equal probability. In this case,  $s_{\max} = 2$  and  $k = 1/2$ , whereas  $\mathbf{G}$  and  $\mathbf{B}$  are degenerate on  $(3/4, 2)$  and  $(1/4, 1)$ , respectively.

We now bound the sender's payoff from  $\mathbf{P}$  from above by applying the results of the extreme cases of our model to the above decomposition. We begin by bounding the value the sender obtains from  $\mathbf{G}$ . To do so, note that because  $\mathbf{P}$  satisfies (R-IC),  $\mathbf{G}$  is supported on the graph of  $V$ . It follows that  $\mathbf{G} \in \text{BP}(\gamma, V)$ , where  $\gamma = \int \mu d\mathbf{G}(\mu, s)$  is the receiver's expected posterior under  $\mathbf{G}$ . Moreover, observe that  $\mathbf{G}$  satisfies the constant sender payoff condition (CP): by construction,  $\mathbf{G}$  only induces outcomes that give the sender a payoff of  $s_{\max}$ . Hence, given the above characterization of feasible distributions for the no-credibility case,  $\mathbf{G}$  is compatible with a 0-equilibrium for the game with modified prior  $\gamma$ . Therefore, we can bound the sender's expected payoff from  $\mathbf{G}$  using the quasiconcavification of the sender's value function:

$$s_{\max} = \int s d\mathbf{G}(\mu, s) \leq \bar{v}(\gamma).$$

<sup>5</sup> It will be apparent that in the cases of  $k = 0$  and  $k = 1$ , the payoff upper bound we derive will remain an upper bound.

Next, we use concavification to bound from above the sender’s expected payoff from  $\mathbf{B}$ . Toward this goal, for every payoff  $\bar{s}$ , define the correspondence  $V_{\bar{s}} : \Delta\Theta \rightrightarrows \mathbb{R}$  that censors  $V(\mu)$  from above by  $\bar{s}$ :

$$V_{\bar{s}}(\mu) = \{\min\{s, \bar{s}\} : s \in V(\mu)\}.$$

Figure 4 illustrates  $V_{\bar{s}}$ . The graph of this correspondence is constructed by reducing to  $\bar{s}$  the payoff coordinate of every outcome  $(\mu, s)$  in  $V$ ’s graph whose  $s$  is above  $\bar{s}$ . Other outcomes in  $V$ ’s graph are kept unchanged.

To understand why  $V_{\bar{s}}$  is a useful correspondence, observe that  $\mathbf{B}$  is supported on the graph of  $V$  and that, by definition,  $\mathbf{B}$  never yields a sender payoff above  $s_{\max}$ . In other words, for any  $\bar{s}$  larger than  $s_{\max}$ ,  $\mathbf{B}$  only generates outcomes from the graph of  $V$  that are also in the graph of  $V_{\bar{s}}$ . Hence, whenever  $\bar{s} \geq s_{\max}$ , the outcome distribution  $\mathbf{B}$  is in the set  $\text{BP}(\beta, V_{\bar{s}})$ , where  $\beta = \int \mu d\mathbf{B}(\mu, s)$  is the receiver’s average posterior under  $\mathbf{B}$ . Therefore,  $\mathbf{B}$  must give the sender a utility below the maximal payoff that the sender can get from some distribution in this set. As we explained in section III.A, one can find this maximal payoff using concavification. Specifically, let

$$v_{\bar{s}} : \Delta\Theta \rightarrow \mathbb{R}$$

$$\mu \mapsto \max V_{\bar{s}}(\mu) = \min\{v(\mu), \bar{s}\}$$

be the function that assigns every belief  $\mu$  with the highest sender utility in  $V_{\bar{s}}(\mu)$ , and let  $\hat{v}_{\bar{s}}$  be the concavification of  $v_{\bar{s}}$ . Then,  $\hat{v}_{\bar{s}}(\beta)$  is the highest payoff the sender can obtain from any distribution in  $\text{BP}(\beta, V_{\bar{s}})$ . Because  $\bar{v}(\gamma) \geq s_{\max}$ , setting  $\bar{s} = \bar{v}(\gamma)$  delivers that  $\mathbf{B}$  gives the sender an expected payoff below  $\hat{v}_{\bar{v}(\gamma)}$ . To ease notational burden, we use

$$\hat{v}_{\Lambda\gamma} := \text{cav}(v_{\Lambda\bar{v}(\gamma)})$$

as shorthand for  $\hat{v}_{\Lambda\bar{v}(\gamma)}$ .

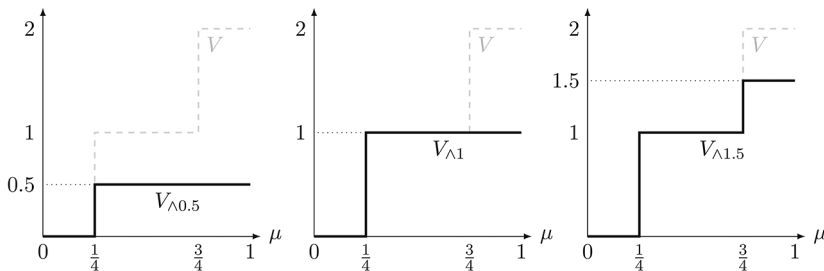


FIG. 4.—Construction of  $V_{\bar{s}}$  for  $\bar{s} = 0.5, 1, 1.5$  in central bank example.

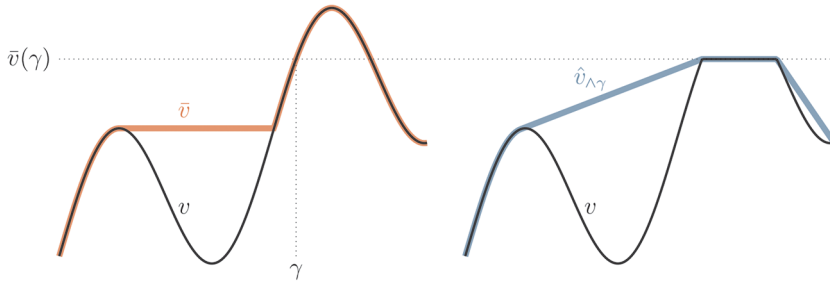


FIG. 5.—Construction of concavification of value function capped at some  $\gamma$ .

Figure 5 illustrates the construction of  $\hat{v}_{\Lambda\gamma}$ . The first step in the construction is to find  $\bar{v}(\gamma)$ , the value of the quasiconcavification of  $v$  at an arbitrary  $\gamma$ . Using this value, one then caps the sender's value function so that no belief results in a payoff higher than  $\bar{v}(\gamma)$ . The result is the function  $v_{\Lambda\gamma}(\cdot) = \min\{v(\cdot), \bar{v}(\gamma)\}$ , which is the same function one obtains by mapping every belief  $\mu$  to the maximal value in  $V_{\Lambda\bar{v}(\gamma)}$ . Concavifying this function delivers  $\hat{v}_{\Lambda\gamma}$ .

Collecting the above observations allows us to bound the sender's payoff from a fixed  $\chi$ -equilibrium outcome distribution  $\mathbf{P}$ ,

$$\begin{aligned} \int s d\mathbf{P}(\mu, s) &= k \int s d\mathbf{B}(\mu, s) + (1 - k) \int s d\mathbf{G}(\mu, s) \\ &\leq k\hat{v}_{\Lambda\gamma}(\beta) + (1 - k)\bar{v}(\gamma). \end{aligned}$$

Of course, the above bound holds only for  $\mathbf{P}$ , the  $\chi$ -equilibrium outcome distribution we started from. To attain an upper bound across all  $\chi$ -equilibria, we maximize the right-hand side of the above equation over all  $(\beta, \gamma, k)$  satisfying two restrictions necessary for a  $\chi$ -equilibrium outcome distribution. For the first restriction, recall that  $\mathbf{P}$  must satisfy the Bayesian updating constraint (Bayes), and so

$$\mu_0 = \int \mu d\mathbf{P}(\mu, s) = k \int \mu d\mathbf{B}(\mu, s) + (1 - k) \int \mu d\mathbf{G}(\mu, s).$$

Because  $\int \mu d\mathbf{B}(\mu, s) = \beta$  and  $\int \mu d\mathbf{G}(\mu, s) = \gamma$ , it follows that  $(\beta, \gamma, k)$  must satisfy the Bayesian splitting constraint

$$k\beta + (1 - k)\gamma = \mu_0. \quad (\text{BS})$$

For the second restriction, observe that an influencing sender only sends messages whose induced outcome results in a sender payoff of  $s_{\max}$ . Indeed, she never attains a higher payoff, since no on-path message leads to a payoff above  $s_{\max}$ , and she cannot find sending a message yielding a lower payoff optimal, because then she would prefer to deviate to a

message generating a payoff of  $s_{\max}$ . Hence, for each state  $\theta$ , the probability the state is  $\theta$  and the sender obtains a payoff of  $s_{\max}$  is at least the probability the state is  $\theta$  and reporting is influenced—that is,  $(1 - \chi)\mu_0(\theta)$ . Expressing this inequality directly in terms of  $\mathbf{P}$  and using the definitions of  $k$  and  $\mathbf{G}$  gives

$$(1 - \chi)\mu_0(\theta) \leq \int_{\{(\mu,s) : s=s_{\max}\}} \mu(\theta) d\mathbf{P}(\mu, s) = (1 - k) \int \mu(\theta) d\mathbf{G}(\mu, s).$$

Recalling that  $\int \mu d\mathbf{G}(\mu, s) = \gamma$  delivers that  $(\beta, \gamma, k)$  must satisfy the credibility constraint

$$(1 - k)\gamma(\theta) \geq (1 - \chi)\mu_0(\theta) \quad \forall \theta \in \Theta. \tag{\chi C}$$

Thus, we have obtained the following upper bound on the sender’s maximal  $\chi$ -equilibrium value:

$$v_\chi^*(\mu_0) := \max_{\beta, \gamma \in \Delta\Theta, k \in [0,1]} \{k\hat{v}_{\chi\gamma}(\beta) + (1 - k)\bar{v}(\gamma)\} \tag{*}$$

subject to (BS) and ( $\chi$ C).

Our main theorem shows that this bound is also tight when  $\chi$  is intermediate.

**THEOREM 1.** Some  $\chi$ -equilibrium exists in which the sender’s value is  $v_\chi^*(\mu_0)$ . Moreover, any such  $\chi$ -equilibrium is sender optimal.

Our proof uses a  $(\beta, \gamma, k)$  that solves the program (\*) to construct a  $\chi$ -equilibrium yielding the sender a value of  $v_\chi^*(\mu_0)$ . Intuitively, one pastes together a sender-optimal equilibrium of a cheap talk game with prior  $\gamma$  and a Bayesian persuasion solution with prior  $\beta$ . We give an informal description of this construction in appendix A and a formal proof in appendix B.

We now apply the theorem to the introduction’s central bank example. To solve the program for  $v_\chi^*(\mu_0)$ , first note that setting  $(\beta, \gamma, k) = (\mu_0, \mu_0, 0)$  is always feasible, and hence  $v_\chi^*(\mu_0) \geq \bar{v}(\mu_0) = 1$ . But what form must a solution  $(\beta, \gamma, k)$  take if  $v_\chi^*(\mu_0) > 1$ ? First, because the objective is bounded above by  $\bar{v}(\gamma)$ , it must be that  $\bar{v}(\gamma) > 1$ . Equivalently,  $\gamma \geq 3/4$ . Constraint (BS) then requires  $\beta \leq 1/2$  and further gives us an exact formula for  $k$  in terms of  $(\beta, \gamma)$ :

$$k = k_{\beta,\gamma} := \frac{\gamma - \mu_0}{\gamma - \beta}.$$

In what follows, we treat the program as an optimization over  $(\beta, \gamma)$ , taking for granted that  $k$  will be set to  $k_{\beta,\gamma}$ .

Observe that we can (still under the hypothesis that  $v_\chi^*(\mu_0) > \bar{v}(\mu_0)$ ) take  $\gamma = 3/4$ . Indeed, moving  $\gamma \in [3/4, 1]$  closer to the prior—hence,

lowering  $k$  to preserve (BS)—always preserves ( $\chi$ C).<sup>6</sup> Meanwhile, because  $\hat{v}_{\lambda\gamma}(\beta) \leq \bar{v}(\gamma)$  by definition, such a modification raises the program's objective if the modification does not alter the value of  $\bar{v}(\gamma)$ . Therefore, because  $\bar{v}$  is constant on  $[3/4, 1]$ , any solution  $(\beta, \gamma, k)$  such that  $\gamma \geq 3/4$  can be replaced with one that has  $\gamma = 3/4$ .

Thus, we have argued that the program (\*) always admits a solution of the form  $(\beta, 3/4, k_{\beta, 3/4})$  for  $\beta \in [0, 1/2]$ . Restricted to solutions of this form, the program (\*) reduces to a univariate constrained maximization program, which can be solved in three exhaustive cases. If  $\chi \geq 3/4$ , the triplet  $(1/4, 3/4, 1/2)$  is a feasible  $(\beta, \gamma, k)$  that delivers the sender her full commitment value of  $v_x^*(\mu_0) = 3/2$ , meaning that said triplet is optimal. If  $2/3 \leq \chi < 3/4$ , it is optimal to set  $\beta$  equal to

$$\beta_x^* := \frac{3\chi - 2}{4\chi - 2},$$

which is the highest  $\beta$  for which  $k_{\beta, 3/4}$  and  $\gamma = 3/4$  satisfy the constraint ( $\chi$ C). The sender's utility in this case is  $v_x^*(\mu_0) = 2\chi$ . Finally, if  $\chi < 2/3$ , no  $\beta \in [0, 1/2)$  can satisfy the constraints required to support  $\gamma = 3/4$ , and so we cannot improve upon feasible solution  $(\beta, \gamma, k) = (1/2, 1/2, 0)$ , which yields value  $v_x^*(\mu_0) = 1$ ; that is, the sender can do no better than a babbling equilibrium. To summarize, the sender's maximal equilibrium payoff is given by

$$v_x^*(\mu_0) = \begin{cases} 1 & \text{if } \chi < 2/3, \\ 2\chi & \text{if } \chi \in [2/3, 3/4], \\ 3/2 & \text{if } \chi \geq 3/4. \end{cases}$$

Figure 6 illustrates the calculation of this value for some  $\chi \in (2/3, 3/4)$ .

The way the sender obtains the above value—following the construction described after the proof of theorem 1—depends on  $\chi$ . When  $\chi < 2/3$ , it is best for the sender to leave the receiver uninformed. When  $\chi = 1$ , the sender is best commissioning the report described in the introduction,  $\xi_1^*$ . To obtain her full-credibility payoff when  $\chi \in [3/4, 1)$ , the sender commissions a report that induces the same information about  $\theta$  in equilibrium, but the official report is itself more informative than  $\xi_1$  to compensate for the fact that an influencing sender always sends the high message. When  $\chi \in (2/3, 3/4)$ , the sender commissions a report that sends three different messages. The low and medium messages, which induce posterior beliefs 0 and  $1/4$ , respectively, are only ever sent under

<sup>6</sup> In the presence of (BS), the constraint ( $\chi$ C) is equivalent to requiring  $h\beta(\theta) \leq \chi\mu_0(\theta)$  for every state  $\theta$ , a constraint that relaxes as  $k$  decreases.

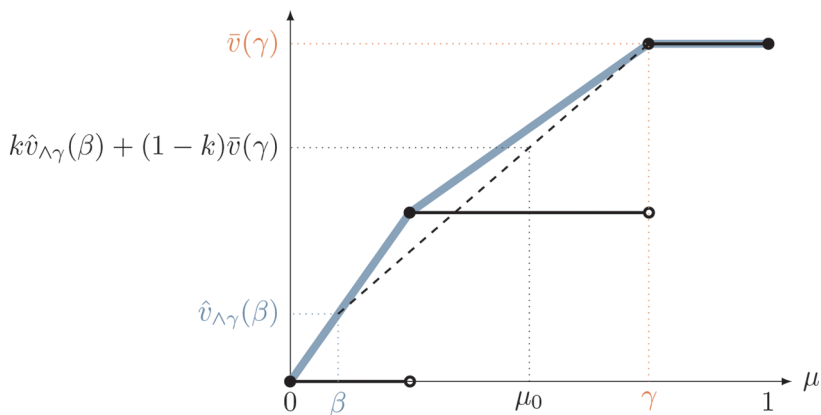


FIG. 6.—Calculating sender value for feasible  $\beta$  and  $\gamma$  in central bank example.

official reporting. The high message would induce a belief strictly higher than  $3/4$  if it were known to come from official reporting, but when taking into account that influenced reporting sends this message in either state, its induced receiver belief is exactly  $3/4$ . Finally, the case of  $\chi = 3/4$  is a limiting version of the latter case in which the medium message is never sent; in this case, the official report is fully informative.

#### IV. Varying Credibility

This section uses theorem 1 to conduct general comparative statics. First, we study how a decrease in the sender’s credibility affects the receiver’s value. In particular, we provide sufficient conditions for the receiver to benefit from a less credible sender. Second, we show that small reductions in the sender’s credibility can often lead to a large drop in the sender’s payoffs. Finally, we note that these drops rarely occur at full credibility. In other words, the full-credibility value is usually robust to small imperfections in the sender’s commitment power.

##### A. Productive Mistrust

We now study how a decrease in the sender’s credibility affects the receiver’s value and the informativeness of the sender’s equilibrium communication. In general, the less credible the sender, the smaller the set of equilibrium outcome distributions.<sup>7</sup> However, that the set of outcome distributions shrinks does not mean that less information is transmitted in the sender’s preferred equilibrium. Our introductory example is a

<sup>7</sup> Given credibility levels  $\chi' < \chi$  and a  $\chi'$ -equilibrium  $(\xi, \sigma, \alpha, \pi)$ , one can construct a  $\chi$ -equilibrium that generates the same outcome distribution, e.g.,  $((\chi'/\chi)\xi + [1 - (\chi'/\chi)]\sigma, \sigma, \alpha, \pi)$ .

case in point, showing that lowering the sender's credibility can result in a more informative equilibrium (à la Blackwell 1953). Moreover, in that example, the receiver uses this additional information, obtaining a strictly higher value when the sender's credibility is lower. In what follows, we refer to this phenomenon as productive mistrust and provide sufficient conditions for it to occur.

Our key sufficient condition involves the sender's optimal outcome distribution under full credibility. For a state  $\theta$ , let  $\delta_\theta \in \Delta\Theta$  be the degenerate belief that generates  $\theta$  with probability 1. Given prior  $\mu$ , an outcome distribution  $\mathbf{P} \in \text{BP}(\mu, V)$  is a *show-or-best* (SOB) outcome distribution if every supported receiver belief lies in

$$\{\delta_\theta\}_{\theta \in \Theta} \cup \underset{\mu' \in \Delta[\text{supp}(\mu)]}{\text{argmax}} v(\mu').$$

In words,  $\mathbf{P}$  is an SOB distribution if it either reveals the state to the receiver or brings the receiver to a posterior belief that attains the sender's best feasible value. Say the sender is a *two-faced SOB* if for every binary support prior  $\mu \in \Delta\Theta$ , every  $\mathbf{P} \in \text{BP}(\mu, V)$  is outperformed by an SOB distribution  $\mathbf{P}' \in \text{BP}(\mu, V)$ ; that is,  $\int s d\mathbf{P}(\mu', s) \leq \int s d\mathbf{P}'(\mu', s)$ . Figure 7 depicts an example in which the sender is a two-faced SOB. Note that productive mistrust cannot occur in this example: one can show that if the sender's favorite equilibrium outcome distribution changes as credibility

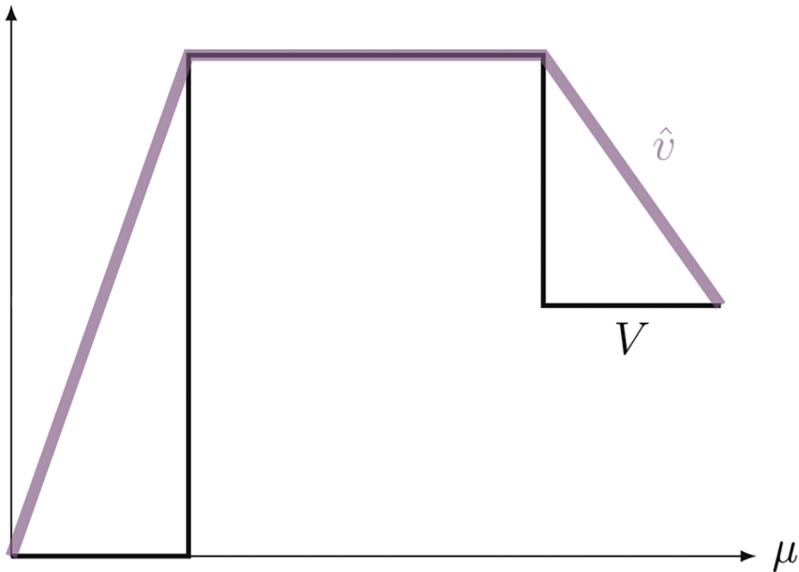


FIG. 7.—Sender is a two-faced SOB.

declines, no information must become sender optimal.<sup>8</sup> As such, the receiver need not benefit from a less credible sender.

Finally, say a model is *generic* if the receiver is (1) not indifferent between any two actions at any degenerate belief and (2) not indifferent between any three actions at any binary support belief.<sup>9</sup>

Proposition 1 below shows that in generic settings, the sender not being a two-faced SOB is sufficient for productive mistrust to occur for some full-support priors at some credibility levels. Intuitively, the sender being an SOB means that a highly credible sender has no bad information to hide: under full credibility, the sender's bad messages are maximally informative, subject to keeping the receiver's posterior fixed following the sender's good messages. The sender not being an SOB at some prior means her bad messages optimally hide some instrumental information. By reducing the sender's credibility just enough to make the full-credibility solution infeasible, one can push her to reveal some of that information to the receiver. In other words, the sender commits to potentially revealing more extreme bad information in order to preserve the credibility of her good messages. Proposition 1 below formalizes this intuition.

**PROPOSITION 1.** Consider a generic model in which the sender is not a two-faced SOB. Then, a full-support prior and credibility levels  $\chi' < \chi$  exist such that every sender-optimal  $\chi'$ -equilibrium is strictly better for the receiver than every sender-optimal  $\chi$ -equilibrium.<sup>10</sup>

The proposition builds on the binary state case, extending to the general case via a continuity argument. We now sketch the binary state argument. To follow the argument, consulting figure 8, which depicts the relevant objects for the central bank example, is useful. Because the model is generic,  $\bar{v}$  has a nondegenerate interval of maximizers (which correspond to beliefs in  $[3/4, 1]$  in fig. 8). Fixing a prior near this interval but toward the nearest kink, we then find the lowest  $\chi \in [0, 1]$  at which the sender still obtains her full-credibility value. In the central bank example, one can use any prior in  $(1/4, 3/4)$ . If we choose  $\mu_0 = 1/2$ , we

<sup>8</sup> For an explanation, observe that the claim is obvious for priors that allow the sender to attain her first-best under no information. For other priors, a feasible  $(\beta, \gamma, k)$  exists that improves on the sender's no-information payoff if and only if a feasible  $(\beta, \gamma, k)$  exists that gives the sender her full-credibility payoff.

<sup>9</sup> Given a fixed finite  $A$  and  $\Theta$ , genericity holds for (Lebesgue) almost every  $u_R \in \mathbb{R}^{A \times \Theta}$ . In particular, it holds if  $u_R(a, \theta) \neq u_R(a', \theta)$  for all distinct  $a, a' \in A$  and all  $\theta \in \Theta$ , and  $(u_R(a_1, \theta_1) - u_R(a_2, \theta_1)) / (u_R(a_1, \theta_2) - u_R(a_2, \theta_2)) \neq (u_R(a_2, \theta_1) - u_R(a_3, \theta_1)) / (u_R(a_2, \theta_2) - u_R(a_3, \theta_2))$  for all distinct  $a_1, a_2, a_3 \in A$  and all distinct  $\theta_1, \theta_2 \in \Theta$ .

<sup>10</sup> Two additional remarks are in order. First, when  $|\Theta| = 2$ , every sender-optimal  $\chi'$ -equilibrium is more Blackwell informative than every sender-optimal  $\chi$ -equilibrium.

Second, with more than two states, one can also find payoff environments in which every sender-optimal 0-equilibrium is strictly better for the receiver than every sender-optimal 1-equilibrium.



*B. Collapse of Trust*

Theorem 1 immediately implies that lowering the sender’s credibility can only decrease her value.<sup>11</sup> Below, we show that this decrease is discontinuous for many payoff specifications of our model. In other words, small decreases in the sender’s credibility can result in a large drop in the sender’s benefits from communication.

PROPOSITION 2. The following are equivalent:

- i. A collapse of trust never occurs:

$$\lim_{\chi' > \chi} v_{\chi'}^*(\mu_0) = v_{\chi}^*(\mu_0)$$

for every  $\chi \in [0, 1]$  and every full-support prior  $\mu_0$ .

- ii. Commitment is of no value:  $v_1^* = v_0^*$ .
- iii. No conflict occurs:  $v(\delta_{\theta}) = \max v(\Delta\Theta)$  for every  $\theta \in \Theta$ .

Let us sketch proposition 2’s proof. To this end, notice that two of the proposition’s three implications are immediate. First, whenever no conflict occurs, the sender can reveal the state in an incentive-compatible way while obtaining her first-best payoff (given the receiver’s incentives), meaning commitment is of no value; that is, point iii implies point ii. Second, because the sender’s highest equilibrium value increases with her credibility, commitment having no value means that the sender’s best equilibrium value is constant (and, a fortiori, continuous) in the credibility level; that is, point ii implies point i.

To show that point i implies point iii, we show that any failure of point iii implies the failure of point i. To do so, we fix a full-support prior  $\mu_0$  at which  $\bar{v}$  is minimized. Because conflict occurs,  $\bar{v}$  is nonconstant and thus takes values strictly greater than  $\bar{v}(\mu_0)$ . By theorem 1, one has that  $v_{\chi}^*(\mu_0) > \bar{v}(\mu_0)$  if and only if a feasible triplet  $(\beta, \gamma, k)$  with  $k < 1$  exists such that  $\bar{v}(\gamma) > \bar{v}(\mu_0)$ . Using upper semicontinuity of  $\bar{v}$ , we show that such a triplet is feasible for credibility  $\chi$  if and only if  $\chi$  is weakly greater than some strictly positive  $\chi^*$ . We thus have

$$v_{\chi^*}^*(\mu_0) \geq k\bar{v}(\mu_0) + (1 - k)\bar{v}(\gamma) > \bar{v}(\mu_0) = \max_{\chi \in [0, \chi^*]} v_{\chi}^*(\mu_0),$$

where the first inequality follows from  $\mu_0$  minimizing  $\bar{v}$ ; that is, a collapse of trust occurs.

<sup>11</sup> In app. sec. B.1.4, we show that credibility increases have a continuous payoff effect: a sufficiently small increase in the sender’s credibility never results in a large gain in the sender’s benefits from communication. Thus, the sender’s value is an upper-semicontinuous function of  $\chi$ . Proposition 2 implies that lower semicontinuity is frequently violated.

C. *Robustness of the Commitment Case*

Given the large and growing literature on optimal persuasion with commitment, one may wonder whether the commitment solution is robust to small decreases in the sender's credibility. Proposition 3 shows the answer is almost always.

PROPOSITION 3. The following are equivalent:

- i. The full-commitment value is robust:  $\lim_{\chi \rightarrow 1} v_\chi^*(\mu_0) = v_1^*(\mu_0)$  for every full-support  $\mu_0$ .
- ii. The sender receives the benefit of the doubt: every  $\theta \in \Theta$  is in the support of some member of  $\operatorname{argmax}_{\mu \in \Delta\Theta} v(\mu)$ .

Thus, the proposition shows that the sender's full-credibility value is robust if and only if the sender can persuade the receiver to take her favorite action without ruling out any states. A sufficient condition for the latter is that the receiver is willing to take the sender's preferred undominated action at some full-support belief, a property that holds generically.<sup>12</sup> Hence, although small decreases in credibility often lead to a collapse in the sender's value, these collapses rarely occur at  $\chi = 1$ .

The argument behind proposition 3 establishes a four-way equivalence between

- a. the sender getting the benefit of the doubt,
- b.  $\bar{v}$  being maximized by a full-support prior  $\gamma$ ,
- c. a full-support  $\gamma$  existing such that  $\hat{v}_{\gamma}$  and  $\hat{v}$  agree over all full-support priors, and
- d. robustness to limited credibility.

To see that point a implies point b, notice that whenever the sender receives the benefit of the doubt, one can find a full-support prior in the convex hull of the beliefs in which the receiver is willing to give the sender her first-best action. Splitting this prior across those beliefs gives an outcome distribution in  $\operatorname{BP}(\mu_0, V)$  that delivers the sender her highest feasible payoff for every supported outcome, meaning the sender can attain this payoff using cheap talk. For the converse direction, one can use the fact that  $\max \bar{v}(\Delta\Theta) = \max v(\Delta\Theta)$ . Specifically, this fact implies  $\bar{v}$  is maximized at a full-support prior  $\gamma$  if and only if one can split  $\gamma$  in a way

<sup>12</sup> More precisely, proposition 3 implies that the sender's full-credibility value is robust whenever a sender-best action among those not strictly dominated for the receiver is a best reply for some full-support belief. It follows from lemma 1 in Lipnowski, Ravid, and Shishkin (2022) that this property holds for Lebesgue-almost every preference specification.

that attains  $v$ 's maximal value at all posteriors, because  $\bar{v}$  gives the sender's highest cheap-talk payoff for every prior. The sender receiving the benefit of the doubt then follows from  $\gamma$  having full support.

For the equivalence of points b and c, note that  $\hat{v}$  and  $\hat{v}_{\lambda\gamma}$  are both continuous because  $A$  and  $\Theta$  are finite. Therefore, the two functions agree over all full-support priors if and only if they are equal, which is equivalent to the cap on  $v_{\lambda\bar{v}(\gamma)}$  being nonbinding; that is,  $\gamma$  maximizes  $\bar{v}$ .

To see why point c is equivalent to point d, fix some full-support  $\mu_0$  and consider two questions about theorem 1's program. First, which beliefs can serve as  $\gamma$  for  $\chi < 1$  large enough? Second, how do the optimal  $(\beta, k)$  for a given  $\gamma$  change as  $\chi$  goes to 1? For the first question, the answer is that  $\gamma$  is feasible for some  $\chi < 1$  if and only if  $\gamma$  has full support.<sup>13</sup> For the second question, one can show that it is always optimal to choose  $(\beta, k)$  so as to make  $(\chi C)$  bind while still satisfying (BS).<sup>14</sup> Direct computation reveals that as  $\chi$  goes to 1, every such  $(\beta, k)$  must converge to  $(\mu_0, 1)$ . Combined, one obtains that as  $\chi$  increases, the sender's optimal value converges to  $\max_{\gamma \in \text{int}(\Delta\Theta)} \hat{v}_{\lambda\gamma}(\mu_0)$ . Thus, the sender's value is robust to limited credibility if and only if some full-support  $\gamma$  exists for which  $\hat{v}_{\lambda\gamma} = \hat{v}$  for all full-support priors; that is, point c is equivalent to point d. The proposition follows.

**V. Conclusion**

This paper studies a model of persuasion through a weak institution whose messages are compromised. Our model has certain features that are worth further discussion.

Throughout the paper, we assumed that the sender's credibility is independent of the state of the world. However, in many scenarios, it is natural for the sender's credibility to be correlated with the state. For example, an autocrat may be more likely to influence the media in a rich economy with abundant resources than in a country where resources are scarce (e.g., Egorov, Guriev, and Sonin 2009). One can capture such correlation by supposing that when the state is  $\theta$ , the message is drawn

<sup>13</sup> It is easy to see that every full-support  $\gamma$  admits some  $\beta$  and  $k < 1$  that make (BS) hold. Moreover,  $(\chi C)$  is also satisfied at  $(\beta, \gamma, k)$  for all sufficiently high  $\chi$ , because  $(\chi C)$ 's right-hand side converges to zero as  $\chi \rightarrow 1$ . Conversely, observe that if  $\gamma(\theta) = 0$ ,  $(\chi C)$  is violated at  $\theta$  for all  $\chi < 1$ , because  $\mu_0$  has full support.

<sup>14</sup> To see why, for any feasible  $(\beta, \gamma, k)$ , a  $(\beta', \gamma, k')$  exists such that  $(\beta', \gamma, k')$  is feasible,  $(\chi C)$  binds, and  $k' \geq k$ . By (BS),  $\beta' = (k/k')\beta + (1 - k/k')\gamma$ . Because  $\hat{v}_{\lambda\gamma}$  is concave and  $\hat{v}_{\lambda\gamma}(\gamma) = \bar{v}(\gamma)$ ,

$$\begin{aligned} k'\hat{v}_{\lambda\gamma}(\beta') + (1 - k')\bar{v}(\gamma) &= k'\hat{v}_{\lambda\gamma}\left(\frac{k}{k'}\beta + \left(1 - \frac{k}{k'}\right)\gamma\right) + (1 - k')\bar{v}(\gamma) \\ &\geq k\hat{v}_{\lambda\gamma}(\beta) + (k' - k)\hat{v}_{\lambda\gamma}(\gamma) + (1 - k')\bar{v}(\gamma) = k\hat{v}_{\lambda\gamma}(\beta) + (1 - k)\bar{v}(\gamma). \end{aligned}$$

from the sender's official report with probability  $\chi(\theta)$ . Theorem 1 generalizes to this case with a minor modification. For a bounded and measurable  $f: \Theta \rightarrow \mathbb{R}$  and  $\mu \in \Delta\Theta$ , let  $f_\mu$  denote the measure on  $\Theta$  given by  $f_\mu(\hat{\Theta}) := \int_{\hat{\Theta}} f d\mu$ . Then, appendix B shows that some sender-favorite equilibrium exists, and the sender's value in this equilibrium is given by

$$v_\chi^*(\mu_0) = \max_{\beta, \gamma \in \Delta\Theta, k \in [0,1]} k\hat{v}_{\chi, \gamma}(\beta) + (1-k)\bar{v}(\gamma) \quad (2)$$

$$\text{subject to } k\beta + (1-k)\gamma = \mu_0,$$

$$(1-k)\gamma \geq (1-\chi)\mu_0. \quad (\chi C)$$

With the above characterization in hand, the propositions of section IV extend to the state-dependent credibility model in a straightforward manner; see the appendix for precise statements.

We also assumed that the sender announces her official report before knowing whether the announcement is credible. In practice, the sender may be privy to institutional features that affect her chances of influencing the report before she commissions it. To understand such situations, appendix C considers a modified model in which the sender learns her credibility type before announcing the official reporting protocol. We show that this modification has no impact on the sender's equilibrium payoffs, and so the sender's maximal equilibrium value remains unchanged.

Finally, we formulated our model as having a finite number of actions and states. However, many applications admit infinite states, infinite actions, or both (e.g., Gentzkow and Kamenica 2016; Kolotilin et al. 2017; Dworzak and Martini 2019). To accommodate such applications, the appendix considers a more general model in which both the action and the state space are compact metrizable. As we show there, our characterization of sender-optimal equilibrium payoffs generalizes to this case in a straightforward manner.

## References

- Alonso, Ricardo, and Odilon Câmara. 2018. "On the Value of Persuasion by Experts." *J. Econ. Theory* 174:103–23.
- Aumann, Robert J., and Sergiu Hart. 2003. "Long Cheap Talk." *Econometrica* 71 (6): 1619–60.
- Aumann, Robert J., and Michael Maschler. 1995. *Repeated Games with Incomplete Information*. Cambridge, MA: MIT Press.
- Ben-Porath, Elchanan, Eddie Dekel, and Barton L. Lipman. 2019. "Mechanisms with Evidence: Commitment and Robustness." *Econometrica* 87 (2): 529–66.
- Best, James, and Daniel Quigley. 2022. "Persuasion for the Long Run." Working paper.

- Bester, Helmut, and Roland Strausz. 2001. "Contracting with Imperfect Commitment and the Revelation Principle: The Single Agent Case." *Econometrica* 69 (4): 1077–98.
- Blackwell, David. 1953. "Equivalent Comparisons of Experiments." *Ann. Math. Statist.* 24 (2): 265–72.
- Blume, Andreas, Oliver J. Board, and Kohei Kawamura. 2007. "Noisy Talk." *Theoretical Econ.* 2 (4): 395–440.
- Chakraborty, Archishman, and Rick Harbaugh. 2010. "Persuasion by Cheap Talk." *A.E.R.* 100 (5): 2361–82.
- Crawford, Vincent P., and Joel Sobel. 1982. "Strategic Information Transmission." *Econometrica* 50 (6): 1431–51.
- Dworczak, Piotr, and Giorgio Martini. 2019. "The Simple Economics of Optimal Persuasion." *J.P.E.* 127 (5): 1993–2048.
- Egorov, Georgy, Sergei Guriev, and Konstantin Sonin. 2009. "Why Resource-Poor Dictators Allow Freer Media: A Theory and Evidence from Panel Data." *American Polit. Sci. Rev.* 103 (4): 645–68.
- Forges, Françoise. 2020. "Games with Incomplete Information: From Repetition to Cheap Talk and Persuasion." *Ann. Econ. and Statist.* (137): 3–30.
- Fréchette, Guillaume, Alessandro Lizzeri, and Jacopo Perego. 2022. "Rules and Commitment in Communication: An Experimental Analysis." *Econometrica*, forthcoming.
- Gentzkow, Matthew, and Emir Kamenica. 2016. "A Rothschild-Stiglitz Approach to Bayesian Persuasion." *A.E.R.* 106 (5): 597–601.
- Glazer, Jacob, and Ariel Rubinstein. 2006. "A Study in the Pragmatics of Persuasion: A Game Theoretical Approach." *Theoretical Econ.* 4 (1): 395–410.
- Goltsman, Maria, Johannes Hörner, Gregory Pavlov, and Francesco Squintani. 2009. "Mediation, Arbitration and Negotiation." *J. Econ. Theory* 144 (4): 1397–420.
- Green, Jerry R., and Nancy L. Stokey. 2007. "A Two-Person Game of Information Transmission." *J. Econ. Theory* 135 (1): 90–104.
- Guo, Yingni, and Eran Shmaya. 2021. "Costly Miscalibration." *Theoretical Econ.* 16 (2): 477–506.
- Hart, Sergiu, Ilan Kremer, and Motty Perry. 2017. "Evidence Games: Truth and Commitment." *A.E.R.* 107 (3): 690–713.
- Hedlund, Jonas. 2017. "Bayesian Persuasion by a Privately Informed Sender." *J. Econ. Theory* 167:229–68.
- Ichihashi, Shota. 2019. "Limiting Sender's Information in Bayesian Persuasion." *Games and Econ. Behavior* 117:276–88.
- Ivanov, Maxim. 2010. "Informational Control and Organizational Design." *J. Econ. Theory* 145 (2): 721–51.
- Kamenica, Emir. 2019. "Bayesian Persuasion and Information Design." *Annual Rev. Econ.* 11 (1): 249–72.
- Kamenica, Emir, and Matthew Gentzkow. 2011. "Bayesian Persuasion." *A.E.R.* 101 (6): 2590–615.
- Kartik, Navin. 2009. "Strategic Communication with Lying Costs." *Rev. Econ. Studies* 76 (4): 1359–95.
- Kolotilin, Anton, Tymofiy Mylovanov, Andriy Zapechelnyuk, and Ming Li. 2017. "Persuasion of a Privately Informed Receiver." *Econometrica* 85 (6): 1949–64.
- Lipnowski, Elliot, and Doron Ravid. 2020. "Cheap Talk with Transparent Motives." *Econometrica* 88 (4): 1631–60.
- Lipnowski, Elliot, Doron Ravid, and Denis Shishkin. 2022. "Perfect Bayesian Persuasion." Working paper.

- Luo, Zhaotian, and Arturas Rozenas. 2018. "Strategies of Election Rigging: Trade-Offs, Determinants, and Consequences." *Q. J. Polit. Sci.* 13 (1): 1–28.
- Mathevet, Laurent, David Pearce, and Ennio Stacchetti. 2022. "Reputation for a Degree of Honesty." Working paper.
- Min, Daehong. 2021. "Bayesian Persuasion under Partial Commitment." *Econ. Theory* 72:743–64.
- Nguyen, Anh, and Teck Yong Tan. 2021. "Bayesian Persuasion with Costly Messages." *J. Econ. Theory* 193:105212.
- Perez-Richet, Eduardo. 2014. "Interim Bayesian Persuasion: First Steps." *A.E.R.* 104 (5): 469–74.
- Perez-Richet, Eduardo, and Vasiliki Skreta. 2022. "Test Design under Falsification." *Econometrica* 90 (3): 1109–42.
- Salamanca, Andrés. 2021. "The Value of Mediated Communication." *J. Econ. Theory* 192:105191.
- Sher, Itai. 2011. "Credibility and Determinism in a Game of Persuasion." *Games and Econ. Behavior* 71 (2): 409–19.
- Skreta, Vasiliki. 2006. "Sequentially Optimal Mechanisms." *Rev. Econ. Studies* 73 (4): 1085–111.