

# Multi-Agent Persuasion: Leveraging Strategic Uncertainty\*

Tetsuya Hoshino<sup>†</sup>

## Abstract

A principal wishes to persuade multiple agents to take a particular action profile. Each agent cares about both a payoff-relevant state and other agents’ actions. The principal discloses information about the state to control the agents’ behavior by using their strategic uncertainty. We show that for any non-degenerate prior, the principal can persuade the agents to take an action profile as a unique rationalizable outcome if that action profile satisfies a generalization of risk dominance. Moreover, this result remains true even if each of the agents is allowed to strategically choose whether to receive information from the principal or not.

Keywords: persuasion, information design, information robustness.

JEL codes: C72, D82, D83.

## 1 Introduction

Bayesian persuasion offers insight into how to persuade a single agent (Kamenica and Gentzkow, 2011). The key idea is that a principal can induce any distribution of the agent’s posterior beliefs that satisfies the martingale property. Recent studies have analyzed how to persuade multiple agents (e.g., Bergemann and Morris, 2016a,b, 2019; Inostroza and Pavan, 2020; Li et al., 2020; Morris et al., 2020; Mathevet et al., 2020; Taneva, 2019).

Multi-agent persuasion has two features that are absent from single-agent persuasion, both of which arise from strategic interaction between agents. The first is that the principal controls not only the agents’ first-order beliefs but also their higher-order beliefs. The second is that there may be multiple rationalizable strategies for the agents, given the information disclosed by the principal. The multiplicity of rationalizable strategies makes it difficult for the principal to predict the agents’ actions and to evaluate the value of the disclosed information.

We take a “worst-case” approach: We assume that if there are multiple rationalizable strategies given disclosed information, the principal anticipates that the agents will take the “worst” rationalizable strategy profile, which minimizes the principal’s payoff. This approach is motivated as follows: The principal may be unable to coordinate the agents’ behavior on the principal’s most

---

\*First version: March, 2017; this version: September, 2021. An earlier version of this paper was circulated under the title “Using Strategic Uncertainty to Persuade Multiple Agents.” I am grateful to the co-editor, Masaki Aoyagi, and three anonymous referees for valuable suggestions. I thank Nageeb Ali, Yu Awaya, Kalyan Chatterjee, Andrei Gomberg, Soomin Jung, Vijay Krishna, Rohit Lamba, Fei Li, Stephen Morris, and Romans Pancs for their helpful comments. All errors are my own.

<sup>†</sup>ITAM. E-mail: [tetsuya.hoshino@itam.mx](mailto:tetsuya.hoshino@itam.mx)

preferred equilibrium, while this uncertainty about the agents’ behavior could motivate the principal to be cautious and to choose an information disclosure policy under the worst-case scenario.<sup>1</sup> This approach is becoming standard (e.g., [Inostroza and Pavan, 2020](#); [Li et al., 2020](#); [Morris et al., 2020](#)); however, since the usual revelation principle argument does not apply, we may rely on the structure of an underlying basic game to solve for an optimal information structure.

**Main Results** How much can the principal manipulate agents in multi-agent persuasion? In an  $N$ -agent game, given any  $\mathbf{p} = (p_1, p_2, \dots, p_N) \in (0, 1]^N$ , action profile  $a^0 = (a_1^0, a_2^0, \dots, a_N^0)$  is a  $\mathbf{p}$ -dominant equilibrium if each agent  $i$  strictly prefers action  $a_i^0$  when the other agents  $-i$  take actions  $a_{-i}^0$  with probability at least  $p_i$  (e.g., [Morris et al., 1995](#)). In this study, we assume that the agents play the following class of games ([Assumption 1](#)): (i) Action profile  $a^0$  is a  $\mathbf{p}$ -dominant equilibrium with  $\sum_i p_i \leq 1$  in all states. (ii) Action  $a_i^0$  is strictly dominant for each agent  $i$  in some state. Action  $a_i^0$  can be interpreted as a “safe” option that agent  $i$  may want to take when he is uncertain how the other agents  $-i$  behave.

This class of games has applications that are interesting and economically relevant. An example is a game of political revolution in which each agent chooses whether to attack a regime or not ([Section 2](#)). The regime is toppled only if more agents attack than its strength, and the strength is modeled as a state of nature. An agent will be punished if he attacks the regime but it survives. In this game, not attacking is the “safe” option for an agent and is optimal when he is uncertain what the other agents will do or when the regime is strong.<sup>2</sup> Other applications include bank runs and (joint) investment games, which we will detail later.

We show that for any non-degenerate prior, the principal can persuade the agents to take action profile  $a^0$  as a unique rationalizable outcome of a (suitably designed) Bayesian game ([Theorem 1](#)). The optimal persuasion perturbs both the agents’ first- and higher-order beliefs. This is reminiscent of information robustness (e.g., [Rubinstein, 1989](#); [Kajii and Morris, 1997](#)). Our study observes the tight connection between multi-agent persuasion (with adversarial selection) and information robustness. We do not only translate information robustness into multi-agent persuasion, but also combine the ideas from the two fields. The principal obfuscates the agents’ first-order beliefs about a state in the persuasion fashion, and perturbs their higher-order beliefs in the robustness fashion. The beliefs of different orders are defined consistently under the  $\mathbf{p}$ -dominance condition.

Are the agents willing to receive information from the principal? In the context of information robustness, an information structure is exogenous; therefore, we need not consider agents’ incentive to learn information. In the context of persuasion, however, an information structure is designed by a principal, and is thus endogenous. It is then unclear why the agents are willing to learn the information. To be more specific, consider a situation in which if they receive no information from

---

<sup>1</sup>If the principal could choose which rationalizable action profile the agents would play, she would be able to disclose full information about a state and then recommend to each agent the most preferable strategy. This is uninteresting and unrealistic.

<sup>2</sup>In this example, not attacking is strictly dominant at a strong state, but attacking is not at any state. This assumption is often made in political regime change (e.g., [Shadmehr and Bernhardt, 2011](#)).

the principal, they can take an action profile that all of them prefer to action profile  $a^0$ .<sup>3</sup> That is, the principal’s information makes them worse-off. We show that even if each agent is allowed to strategically choose whether to receive information or not, the principal can still induce action profile  $a^0$  as a unique rationalizable outcome (Theorem 2).

Additionally, we examine some assumptions for our results. First, we study what makes Theorem 1 different from existing results (e.g., concavification). The key is the cardinality of a signal space. Existing studies often assume finite signal spaces. In multi-agent persuasion, an infinite signal space is important for a principal to control agents’ higher-order beliefs even if the underlying state space is finite. Specifically, we show that if a signal space is finite, the principal cannot always persuade the agents to take action profile  $a^0$  (Proposition 1 and Theorem 3); moreover, this result holds true for general games that may not satisfy Assumption 1. Second, we show that the condition  $\sum_{i \in I} p_i \leq 1$  of Assumption 1 is a tight condition for Theorem 1. Using a simple example, we show that if  $\sum_i p_i > 1$ , a principal cannot always induce agents into action profile  $a^0$ , regardless of an information structure with a finite or infinite signal space (Proposition 2).

**Related Literature** This study is related to two strands of the literature: information robustness and multi-agent persuasion. We outlined the relationship to information robustness above, and presently discuss the relationship to multi-agent persuasion. [Mathevet et al. \(2020\)](#) characterize a principal’s optimal payoff, assuming that a signal space is arbitrarily large but finite. This assumption implies that, for any finite  $k$ , the principal can control the agents’ beliefs up to the  $k$ -th order. This does not restrict the principal’s ability to manipulate the agents if their behavior depends only on their finite order beliefs.<sup>4</sup> In contrast, we are interested in the case in which agents’ behavior may depend on their infinite order beliefs. Then, a principal can persuade, with probability 1, agents to take a particular action profile  $a^0$  if a signal space is infinite (Theorem 1); however, this result does not hold if the signal space is finite (Theorem 3). These contrasting results highlight the implications of the finiteness assumption of the signal space.

Our study complements the literature that has paid particular attention to persuasion in binary-action coordination games. [Inostroza and Pavan \(2020\)](#) consider persuasion in a continuum-agent global game, and show that an optimal information structure satisfies the perfect coordination property that all agents take the same action. [Li et al. \(2020\)](#) study persuasion in a continuum-agent coordination game, and characterize an optimal information structure. [Morris et al. \(2020\)](#) consider persuasion in a finite-agent supermodular game. They characterize adversarial-equilibrium implementable outcomes and provide sufficient conditions for the perfect coordination property. In contrast, we neither focus on coordination or supermodular games nor assume binary actions. It is true that the  $\mathbf{p}$ -dominance condition has a similar flavor to that of the coordination game structure, and that the class of games that satisfy the  $\mathbf{p}$ -dominance condition includes the binary-

<sup>3</sup>We provide such an example in Section 2.

<sup>4</sup>An example in [Mathevet et al. \(2020, Section 5\)](#) is a two-state, two-agent, two-action coordination game, in which each action is strictly dominant in one state. In this case, an agent who is certain about a state does not care about the other agent’s action or belief.

action games of regime change, as discussed above.<sup>5</sup> However, this class of games is not restricted to these games. Our study, as well as [Morris et al. \(2020\)](#), highlights the tight connection between the multi-agent persuasion and information robustness in [Kajii and Morris \(1997\)](#).<sup>6</sup> The existing persuasion studies focus on the endogeneity due to a principal’s design of information. In addition, we investigate another potential source of endogeneity, by asking whether agents are willing to receive manipulative information from the principal. For this potential concern, we show that the agent-related endogeneity does not restrict the principal’s influence.

**Layout** The remainder of this paper is organized as follows. Section 2 provides a simple example to illustrate our results. Section 3 builds a model, and Section 4 presents the main results. Section 5 discusses the tightness of key assumptions, followed by Section 6, which concludes.

## 2 Example

Using a simple example, we illustrate our results and discuss the tight connection between multi-agent persuasion (with adversarial selection) and information robustness.

**Set-up** There are two agents, denoted  $i \in I = \{1, 2\}$ . Each agent  $i$  chooses action  $a_i \in A_i = \{0, 1\}$ , where action 1 is interpreted as “attack” a regime and action 0 as “not attack.” There are two states, denoted  $\theta \in \Theta = \{0, 1\}$ , where state 1 is interpreted as a “weak” regime and state 0 as a “strong” one. There is a common prior  $\mu \in \Delta(\Theta)$ ; let  $\mu^\theta$  be the prior probability of a state  $\theta$ . Agent  $i$ ’s payoff  $u_i : A_1 \times A_2 \times \Theta \rightarrow \mathbb{R}$  is represented by the following tables:

$\theta = 0$	$a_2 = 0$	$a_2 = 1$	$\theta = 1$	$a_2 = 0$	$a_2 = 1$
$a_1 = 0$	0, 0	0, -2	$a_1 = 0$	0, 0	0, -2
$a_1 = 1$	-2, 0	-2, -2	$a_1 = 1$	-2, 0	1, 1

Table 1: the agents’ payoffs

In the state-1 complete-information game, there are two pure-strategy equilibria,  $a^1 = (1, 1)$  and  $a^0 = (0, 0)$ , where the former is Pareto dominant. In the state-0 complete-information game, there is a unique equilibrium  $a^0$ , as action 0 is strictly dominant.

The regime (a principal) aims to forestall a coordinated attack  $a^1$ . Its (state-independent) payoff is modeled by the function  $v : A_1 \times A_2 \rightarrow \mathbb{R}$  such that  $v(a_1, a_2) = 1 - a_1 a_2$ . It receives payoff 0 if the agents launch the coordinated attack  $a^1$  and payoff 1 otherwise. It designs an information structure  $(S, \pi)$ . That is, it chooses signal spaces  $S_1, S_2$  for the agents, and a distribution  $\pi(\cdot | \theta) \in \Delta(S)$  for each state  $\theta$ , where we write  $S = S_1 \times S_2$ .

<sup>5</sup>Another assumption of ours is that action  $a_i^0$ , which the principal wants agent  $i$  to take, is strictly dominant at some state. This assumption is maintained in the said studies.

<sup>6</sup>See also [Bergemann and Morris \(2019\)](#).

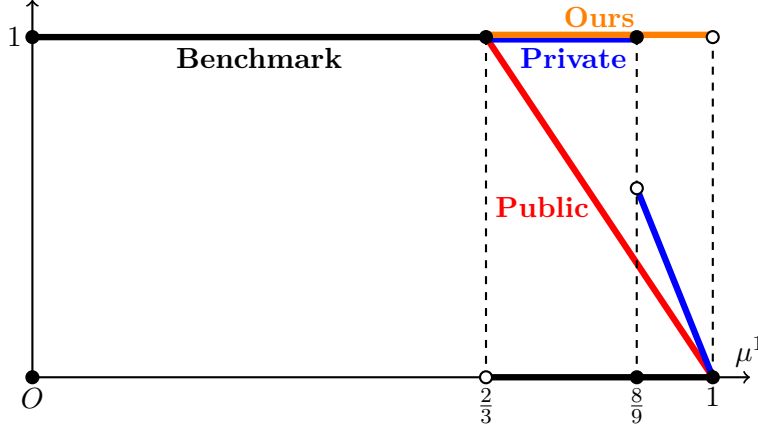


Figure 1: the regime’s maximum worst-case payoffs

Assume that the regime maximizes the worst-case payoff: If the agents have multiple rationalizable strategies given disclosed information, the regime anticipates that they will take the “worst” rationalizable strategy profile, which minimizes the regime’s payoff. We study the regime’s worst-case payoffs under different information structures. Figure 1 illustrates its maximum worst-case payoffs as a function of the prior  $\mu^1$  in cases that are discussed below.

**Benchmark** Suppose that the regime sends no informative signals. If we assume that the agents do not play weakly dominated strategies then the set of (pure-strategy) equilibria is  $\{a^0\}$  if  $\mu^1 \leq \frac{2}{3}$  and  $\{a^1, a^0\}$  if  $\mu^1 > \frac{2}{3}$ . Hence, the worst-case payoff is 1 if  $\mu^1 \leq \frac{2}{3}$  and 0 if  $\mu^1 > \frac{2}{3}$ .

**Binary Signals** Suppose that the regime sends binary signals. The signal space for each agent has two elements:  $|S_1| = |S_2| = 2$ . We assume that the agents do not play weakly dominated strategies. There are two cases to consider. First, we consider the case of *public* signals: The two agents receive the same signal realization. This case is equivalent to single-agent persuasion and thus the regime’s expected payoff is characterized by concavification.<sup>7</sup> Second, we consider the case of *private* signals: The two agents may receive different signal realizations. In this example, it is straightforward to find an optimal information structure, where the regime sends different signals with non-zero probabilities (Appendix B).

**Arbitrary Signals** Suppose that the regime sends arbitrary signals. Then, it can induce action profile  $a^0$ , given any prior  $\mu^1 \neq 1$  (Theorem 1). That is, its maximum payoff is independent of a prior  $\mu$  (except for the degenerate prior putting probability 1 on state 1). Unlike in the above cases, we need not assume that the agents do not play weakly dominated strategies.

Here is an optimal information structure. Since the regime achieves payoff 1 with no signal for any prior  $\mu^1 \leq [0, \frac{2}{3}]$ , we focus on a prior  $\mu^1 \in (\frac{2}{3}, 1)$ . Let  $S_1 = S_2 = \{0, 1, \dots\}$  denote signal spaces. To define a signal distribution, let  $P_\theta(s_1, s_2)$  denote the probability of a signal profile  $(s_1, s_2)$  being

<sup>7</sup>The concavification needs the assumption that the agents do not play weakly dominated strategies.

realized conditional on state  $\theta$ . We define the following distributions: At state  $\theta = 0$ ,  $P_0(0, 0) = 1$ . At state  $\theta = 1$ ,  $P_1(1, 0) = P_1(0, 1) = \nu$  and  $P_1(k, k - 1) = P_1(k - 1, k) = \rho^{k-2}\nu$  for each  $k \geq 2$ , where  $\nu = \frac{\mu^0}{\mu^1} \in (0, \frac{1}{2})$  and  $\rho = \frac{1-4\nu}{1-2\nu} < 1$ . The probabilities for any other signal profiles are zero.

We will now show that this information structure induces each agent to take action 0. First, suppose that agent 1 receives signal  $s_1 = 0$ . He assigns to state  $\theta = 0$  probability  $\mathbb{P}(\theta = 0 \mid s_1 = 0) = \frac{\mu^0}{\mu^0 + \mu^1 \nu} = \frac{1}{2}$ , which makes action 0 strictly dominant. Second, suppose that agent 1 receives signal  $s_1 = 1$ . He knows state  $\theta = 1$  but does not know what agent 2 knows. He assigns probability  $\mathbb{P}(s_2 = 0 \mid s_1 = 1) = \frac{P_1(1, 0)}{P_1(1, 0) + P_1(1, 2)} = \frac{1}{2}$  to agent 2 receiving signal  $s_2 = 0$  (and taking action 0). Thus, agent 1 takes action 0. Third, suppose that agent 1 receives signal  $s_1 = k \geq 2$ . He knows state  $\theta = 1$  but does not know what agent 2 knows. He assigns probability  $\mathbb{P}(s_2 = k - 1 \mid s_1 = k) = \frac{P_1(k, k-1)}{P_1(k, k-1) + P_1(k, k+1)} = \frac{\rho^{k-2}\nu}{\rho^{k-2}\nu + \rho^{k-1}\nu} > \frac{1}{2}$  to agent 2 receiving signal  $s_2 = k - 1$  (and taking action 0). Thus, agent 1 takes action 0. Consequently, agent 1 takes action 0, regardless of his signal. The same holds true for agent 2.

This logic may be reminiscent of information robustness (e.g., [Rubinstein, 1989](#); [Kajii and Morris, 1997](#)). In this study, we will discuss the tight connection between multi-agent persuasion (with adversarial selection) and information robustness, applying this kind of logic to multi-agent persuasion. Note that multi-agent persuasion entails high-order uncertainty. In our example, since the agents care about each other's behavior as well as the state, the regime may be able to leverage strategic uncertainty among the agents. It obfuscates the agents' (first-order) beliefs about a state in the persuasion fashion and perturbs their higher-order beliefs in the robustness fashion. This can be implemented if the regime sends private information to each agent. For example, online communication has made such private communication possible ([Arieli and Babichenko, 2019](#)).

There are two more questions that we will investigate. The first question is: Are the agents willing to receive the regime's signals, even though the signals make the agents worse-off? As seen above, if they receive the signals then they are induced to take action profile  $a^0$ , which yields payoff 0 for each agent; meanwhile, if they received no signals then they could take action profile  $a^1$  at any  $\mu^1 > \frac{2}{3}$ , which yields payoff  $\mu^1 - 2\mu^0 > 0$  for each agent. However, even if each agent can strategically choose whether to receive a signal or not, the principal can still induce action profile  $a^0$  as a unique rationalizable outcome (Theorem 2). The second question is: What makes our information structure different from binary information structures? In the binary (public or private) information structures, the regime's payoff is close to 0 if  $\mu^1$  is high; however, in our information structure, it is equal to 1 if  $\mu^1$  is non-degenerate (Figure 1). The key difference turns out to be the signal space cardinality. Binary signal spaces, of course, are finite; however, our signal space is infinite. For any finite (not necessarily binary) signal space, the regime's payoff must be close to 0 if  $\mu^1$  is high (Theorem 3).

### 3 Model

**Payoff Environment** There is a finite set of agents  $I = \{1, 2, \dots, N\}$  (for any  $N \geq 2$ ), where  $i$  denotes a generic agent and  $-i$  denotes all agents but  $i$ . There is a countable set of payoff-relevant states  $\Theta = \{0, 1, \dots\}$ , where  $\theta$  denotes a generic state. A **basic game**  $G$  consists of (i) for each  $i \in I$ , a finite set of actions  $A_i$  and a payoff function  $u_i : A \times \Theta \rightarrow \mathbb{R}$ , where  $A = \prod_{i \in I} A_i$  is the set of action profiles, and (ii) a common prior  $\mu \in \Delta(\Theta)$ , where  $\mu^\theta$  is the probability of a state  $\theta$ . That is,  $G = ((A_i, u_i)_{i \in I}, \Theta, \mu)$ . As usual, let  $A_{-i} = \prod_{j \neq i} A_j$ , where  $a_{-i}$  denotes a generic element.

Here is our key equilibrium concept (e.g., Morris et al., 1995; Kajii and Morris, 1997). Roughly speaking, action profile  $a^* = (a_1^*, a_2^*, \dots, a_N^*)$  is a strict **p**-dominant equilibrium if each agent  $i$  wants to take action  $a_i^*$  when agents  $-i$  take actions  $a_{-i}^*$  with probability at least  $p_i$ .

**Definition 1.** Given any  $\mathbf{p} = (p_1, p_2, \dots, p_N) \in (0, 1]^N$ , action profile  $a^* = (a_1^*, a_2^*, \dots, a_N^*)$  is a *strict p-dominant equilibrium* at a state  $\theta$  if for each  $i \in I$ , each  $a_i \in A_i \setminus \{a_i^*\}$ , and each  $\lambda \in \Delta(A_{-i})$  with probability  $\lambda(a_{-i}^*) \geq p_i$ , it holds that

$$\sum_{a_{-i} \in A_{-i}} \lambda(a_{-i}) u_i(a_i^*, a_{-i}, \theta) > \sum_{a_{-i} \in A_{-i}} \lambda(a_{-i}) u_i(a_i, a_{-i}, \theta).$$

We will mainly study the class of games such that (i) action profile  $a^0 = (a_1^0, a_2^0, \dots, a_N^0)$  is a strict **p**-dominant equilibrium with  $\sum_i p_i \leq 1$  in all states, and (ii) action  $a_i^0$  is strictly dominant for each agent  $i$  in a state, say,  $\theta = 0$ .

**Assumption 1.** *There exists an action profile, denoted  $a^0 \in A$ , such that it is a strict p-dominant equilibrium with  $\sum_{i \in I} p_i \leq 1$  in each state  $\theta \in \Theta$  and that action  $a_i^0$  is strictly dominant for each agent  $i$  in state  $\theta = 0$ .*

We interpret action  $a_i^0$  of Assumption 1 as a “safe” option for agent  $i$  that he wants to take whenever he is uncertain how agents  $-i$  behave—that is, whenever he assigns probability at least  $p_i$  to agents  $-i$  taking actions  $a_{-i}^0$ . The strict **p**-dominance with  $\sum_i p_i \leq 1$  is a many-player many-action generalization of the risk-dominance. Indeed, the strict  $(\frac{1}{2}, \frac{1}{2})$ -dominant equilibrium coincides with the risk-dominant equilibrium in a two-player two-action symmetric game.

**Remark 1.** The class of games satisfying Assumption 1 has many applications. Here are several examples. The first application is a game of political revolution, which we take as the example in Section 2. In this example, action 0 (not attack) is a safe option for each agent. More precisely, for any  $p > \frac{1}{3}$ , action profile  $a^0 = (0, 0)$  is a strict  $(p, p)$ -dominant equilibrium in both states, while action 0 is strictly dominant for each agent in state 0.

The second application is bank runs. By way of illustration, suppose that there are two depositors, and each has 1 unit of money in a bank and decides whether to withdraw her money (run) or wait for the maturity (wait). The state of nature corresponds to the bank’s solvency. The bank with high solvency keeps 2 units of money; thus, each depositor receives 1 if she runs and receives  $1 + r$  if she waits, where  $r > 0$  is the interest payment. The bank with low solvency has  $2d$  units of

money for  $d \in (1 - r, 1)$ ; thus, if one runs and the other waits, the running one receives 1, whereas the waiting one receives  $2d - 1$ , while if both run, each receives only  $d$  (i.e., a bank run). Even in the low solvency case, if both wait, each receives  $1 + r$ . This game satisfies Assumption 1 with action profile  $a^0$  of both waiting, because (i) for any  $p > \frac{1-d}{1-d+r}$  (where  $\frac{1-d}{1-d+r} < \frac{1}{2}$ ), action profile  $a^0$  is the strict  $(p, p)$ -dominant equilibrium for each solvency, and (ii) waiting is strictly dominant if the solvency is high. We note that this kind of information design can be interpreted as how to design a stress test for the bank (Inostroza and Pavan, 2020).

The third application is an investment game. There are two firms, and each decides whether to invest in a new technology (e.g., software) or not. The investment cost  $c > 0$  is known, while the software quality  $q$  is not. Given a quality  $q$ , the firm that invests in the software obtains utility  $q - c$  from “intra-firm” use, and gains an additional (known) spillover  $x$  from “inter-firm” use if both firms invest. Normalize the utility from not investing to zero. This game satisfies Assumption 1 with action profile  $a^0$  of both investing, if the best quality  $\bar{q}$  is such that  $\bar{q} - c > 0$  and the worst quality  $\underline{q}$  is such that  $x > 2(c - \underline{q}) > 0$ . This is because (i) for any  $p > \frac{c-q}{x}$  (where  $\frac{c-q}{x} < \frac{1}{2}$ ), action profile  $a^0$  of both investing is the strict  $(p, p)$ -dominant equilibrium for any quality  $q$  and (ii) investing is strictly dominant if the quality is  $\bar{q}$ .

Lastly, we note that these applications are well suited to multi-agent persuasion. A principal (e.g., a bank or a software company) may send private information to each agent. Such private communication is possible, for example, via online communication (Arieli and Babichenko, 2019). In the above applications, information about the bank’s solvency, the software quality, etc. is sent through private communication.  $\square$

**Information Environment** A principal wishes to persuade the agents to take action profile  $a^0$ , regardless of a state  $\theta$ . That is, her preference is state-independent. She designs an information structure  $(S, \pi)$ . That is, she chooses a signal space  $S_i$  for each  $i \in I$  and a distribution  $\pi(\cdot | \theta) \in \Delta(S)$  for each  $\theta \in \Theta$ , where we write  $S = \prod_{i \in I} S_i$ .

For expositional purposes, we will maintain the usual commitment assumption, which requires that the principal be able to commit to the information structure  $(S, \pi)$ . However, this assumption turns out to be *unnecessary* for our results.

**Worst-Case Analysis** The game proceeds as follows. First, nature draws a state  $\theta$ . Second, given an information structure  $(S, \pi)$ , which generates a signal profile  $s = (s_1, s_2, \dots, s_N)$  according to the distribution  $\pi(\cdot | \theta)$ , each agent  $i$  observes the signal realization  $s_i$  and then chooses action  $a_i$ . Assume that the principal maximizes the worst-case payoff. That is, if the agents have multiple rationalizable strategies given disclosed information, the principal anticipates that they will take the “worst” rationalizable strategy profile, which minimizes the principal’s payoff.

## 4 Main Results

## 4.1 Optimal Persuasion

We will now describe an optimal information structure whereby the principal persuades the agents to take action profile  $a^0$ . For expositional purposes, we describe it as if the principal were sending information sequentially. We will rewrite the sequential protocol into the usual information structure  $(S, \pi)$  later. The advantage of using the sequential protocol is that the connection to information robustness becomes more apparent. To minimize the risk of confusion, we say that the principal sends “messages” when we treat the principal as if she sent information sequentially, while she sends “signals” when she sends information only once according to an information structure  $(S, \pi)$ . In Section 4.1, we assume that the agents *must* receive any signals or messages sent from the principal; that is, we ignore their incentives about whether they receive the signals or messages. We will revisit the incentives in Section 4.2.

**Optimal Information Structure** The principal manipulates the agents’ (higher-order) beliefs by sending messages that are correlated to state  $\theta$ . The agents update their beliefs based on their messages and then take actions. Given parameters  $\varphi, \phi_1, \dots, \phi_N$ , with sum  $\phi_I = \sum_{j \in I} \phi_j \leq 1$ , we consider the following information structure in the sequential protocol:

**Case of State  $\theta = 0$ :** The principal sends messages 0 to all agents with probability 1 and then sends no more messages (to any agent).

**Case of State  $\theta \neq 0$ :** The principal sends messages according to the following rule:

**Round 0:** The principal sends messages 1 to all agents with probability  $1 - N\varphi$ , or she chooses agent  $i$  with equal probability  $\varphi$  for each of them and sends message 1 to (selected) agent  $i$  and messages 0 to agents  $-i$ . If she has sent messages 1 to all agents, she proceeds to Round 1; otherwise, she does not send any more message (to any agent).

**Round  $k \geq 1$ :** The principal sends messages  $k + 1$  to all agents with probability  $1 - \phi_I$ , or she chooses agent  $i$  with probability  $\phi_i$  for each of them and sends message  $k + 1$  to (selected) agent  $i$  and no messages to agents  $-i$ . If she has sent messages  $k + 1$  to all agents, she proceeds to Round  $k + 1$ ; otherwise, she does not send any more message (to any agent).

Figure 2 illustrates the protocol in the case of state  $\theta \neq 0$ . The protocol starts from the hollow node on the left. The numbers in parentheses denote message profiles. For example,  $(2, 3, 2, \dots, 2, 2)$  is a message profile that agent 2 receives message 3 and the other agents receive messages 2. The numbers by arrows denote transition probabilities. For example,  $\phi_2$  in Round 2 is the probability that the principal sends message 3 only to agent 2 conditional on having sent messages 2 to all agents. The principal continues to send messages until she reaches some message profile (with no arrow from it). For example, the probability of message profile  $(2, 3, 2, \dots, 2, 2)$  conditional on state  $\theta \neq 0$  is  $(1 - N\varphi)(1 - \phi_I)\phi_2$ .

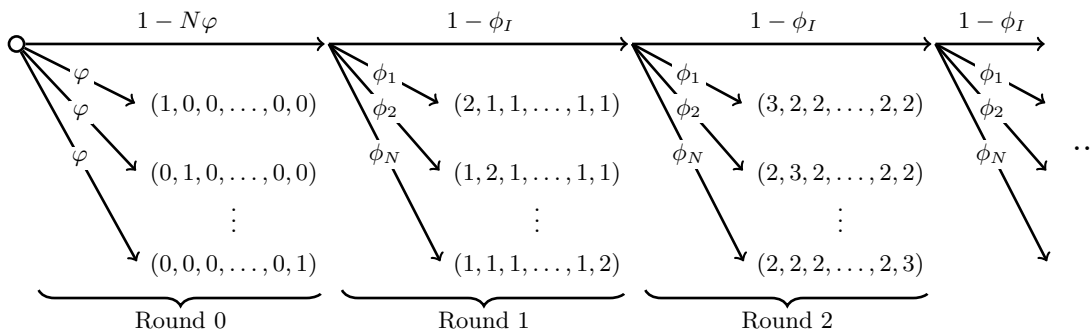


Figure 2: an optimal information structure at state  $\theta \neq 0$

**Optimal Information Structure Reformulated** We reformulate the sequential protocol into the form  $(S, \pi)$ , in which she sends information only once. Let  $S_i = \{0, 1, 2, \dots\}$  be agent  $i$ 's signal space for each  $i \in I$ , and let  $S = \prod_{i \in I} S_i$  be the set of signal profiles. It suffices to specify the distribution  $\pi(\cdot | \theta) \in \Delta(S)$  for each  $\theta \in \Theta$ . We say that for each  $k \in S_i$ , agent  $i$  receives signal  $s_i = k$  if he receives message  $k$  but not  $k + 1$  in the sequential protocol.

The optimal information structure, which is described above, identifies a signal distribution  $\pi(\cdot | \theta) \in \Delta(S)$  for each  $\theta \in \Theta$ :  $\pi(0, 0, \dots, 0 | \theta = 0) = 1$ ,  $\pi(0, \dots, 0, 1, 0, \dots, 0 | \theta \neq 0) = \varphi$ , and  $\pi(k, \dots, k, k + 1, k, \dots, k | \theta \neq 0) = (1 - N\varphi)(1 - \phi_I)^{k-1}\phi_i$ , where agent  $i$  is the agent who receives signal  $k + 1$ , for each  $k \geq 1$ . The probabilities for any other signal profiles are zero. Hence, we can translate the sequential protocol into the information structure  $(S, \pi)$  and vice versa.

Note that each agent's signal space is infinite. This is in contrast to most studies of persuasion, in which each agent's signal space is finite (e.g., [Mathevet et al., 2020](#)). In Section 5, we will see that the cardinality of signal spaces plays an important role in multi-agent persuasion.

Our first result is that the principal can persuade the agents to take action profile  $a^0$  with probability 1.

**Theorem 1.** *Let a basic game  $G$  satisfy Assumption 1. For any prior  $\mu^0 \neq 0$ , there exist parameters  $\varphi, \phi_1, \dots, \phi_N$  such that the corresponding information structure implements, with probability 1, action profile  $a^0$  as a unique rationalizable strategy profile.*

Theorem 1 is useful when the principal is concerned about the worst case, in which she assumes that the agents will take an adversarial rationalizable strategy profile of the Bayesian game  $\langle G, (S, \pi) \rangle$ , which consists of the basic game  $G$  and the principal's choice of information structure  $(S, \pi)$ . This result guarantees that the principal can uniquely achieve the agents' behavior that the principal desires, even under the cautious assumption about their behavior.

We will now sketch the proof of Theorem 1, relegating the details to Appendix A.

**Step 0:** Suppose that agent  $i$  receives signal  $s_i = 0$ . If  $\varphi$  is small enough, agent  $i$  will assign a high probability to state  $\theta = 0$ , thereby taking action  $a_i^0$ , which is the strictly dominant action at state  $\theta = 0$ .

**Step 1:** Suppose that agent  $i$  receives signal  $s_i = 1$ . Then, he knows that state  $\theta \neq 0$  is realized.<sup>8</sup>

He also knows that agent  $j$ 's signal is either 0, 1, or 2 for each  $j \neq i$ .<sup>9</sup> He does not know whether agents  $-i$  know that state  $\theta \neq 0$  is realized. If the principal takes  $\phi_I = \sum_{j \in I} \phi_j$  small enough then agent  $i$  assigns a high probability to agents  $-i$  receiving signals  $s_{-i} = 0$  (and taking action  $a_{-i}^0$ ). Hence, agent  $i$  chooses action  $a_i^0$ , which is, by Assumption 1, the best response to actions  $a_{-i}^0$  at state  $\theta \neq 0$ .

**Step 2:** Suppose that agent  $i$  receives signal  $s_i = 2$ . Then, he knows that state  $\theta \neq 0$  is realized and that agent  $j$ 's signal is either 1, 2, or 3 for each  $j \neq i$ .<sup>10</sup> Agent  $i$  knows that all agents know that state  $\theta \neq 0$  is realized, but does not know whether agents  $-i$  know that all agents know that state  $\theta \neq 0$  is realized. For carefully selected  $\phi_i$ , the principal can make agent  $i$  assign probability at least  $p_i$  to agents  $-i$  receiving signals  $s_{-i} = 1$  (and taking action  $a_{-i}^0$ ). Hence, agent  $i$  chooses action  $a_i^0$ , which is, by Assumption 1, the best response to actions  $a_{-i}^0$  at state  $\theta \neq 0$ .

**Step  $k \geq 3$ :** Suppose that agent  $i$  receives signal  $s_i = k$ . By the same logic as in Step 2, agent  $i$  will choose action  $a_i^0$ .

The optimal information structure combines the ideas of persuasion and (higher-order) belief perturbation. Step 0 manipulates the distribution of agents' (first-order) beliefs about state  $\theta$ . When agent  $i$  assigns a high probability to state  $\theta = 0$ , he will choose action  $a_i^0$ , which is strictly dominant. Hence, this step involves no strategic interaction among the agents and is analogous to (single-agent) persuasion. Since the principal does not have to tailor different (first-order) beliefs to different agents, we can take a parameter  $\varphi$  common to all the agents. Step  $k \geq 2$  perturbs the agents'  $(k + 1)$ th-order beliefs. In this step, agent  $i$  knows that state  $\theta \neq 0$  is realized and may not have a strictly dominant strategy. Since he cares about agents  $-i$ 's actions, the principal can persuade agent  $i$  into action  $a_i^0$  if she induces agent  $i$  to believe that agents  $-i$  will take actions  $a_{-i}^0$  with probability at least  $p_i$ . Since the principal has  $N$  objects to control (one for each agent), she has  $N$  parameters  $\phi_1, \dots, \phi_N$ . Once she chooses them, since all steps  $k \geq 2$  have the same structure, she can use the same parameters. This is why the choice of parameters  $\phi_1, \dots, \phi_N$  is independent of Step  $k \geq 2$ . Finally, Step 1 "connects" these two ideas. The principal adjusts the sum  $\phi_I$ , which is required to be small enough in Step 1, for all the steps to be compatible.

The commitment assumption (that the principal commits to her information structure) is not necessary for Theorem 1. The reason is that since every signal realization leads to the same action profile  $a^0$ , the principal has no incentive to misreport any signal realization.

---

<sup>8</sup>In this study, we say that an agent *knows* an event if he assigns probability 1 to the event and the event is true.

<sup>9</sup>If there are two agents, agent  $i$  receiving signal  $s_i = 1$  knows that agent  $-i$ 's signal is either 0 or 2.

<sup>10</sup>If there are two agents, agent  $i$  receiving signal  $s_i = 2$  knows that agent  $-i$ 's signal is either 1 or 3.

## 4.2 Agents' Strategic Decision about Whether to Receive Information

A principal can persuade agents into action profile  $a^0$  under the (implicit) assumption that all agents *must* receive the signals from the principal (Theorem 1). What if they can choose whether to receive the signals or not? In the example presented in Section 2, they are strictly better-off when they do not receive the signals than when they do. That is, the principal's signals make them worse-off. A natural question is: Are they willing to receive such signals? We show that even if each agent can strategically choose whether or not to receive his signal, the principal can still send the signals that each agent  $i$  chooses to receive, thereby achieving action profile  $a^0$ .<sup>11</sup>

**Modified Model** We modify the model so that each agent  $i$  can choose whether or not to receive signal  $s_i$ . Specifically, the modified model proceeds as follows. First, nature draws state  $\theta$ , while the principal designs an information structure  $(S, \pi)$ . Second, each agent  $i$  decides whether or not to receive signal  $s_i$ . Formally, he makes decision  $n_i \in \{\infty, \emptyset\}$ , where we interpret  $\infty$  as the agent receiving signal  $s_i$  and  $\emptyset$  as the agent not receiving it. Assume that agent  $i$ 's decision  $n_i$  is unobservable to the others.<sup>12</sup> Third, the principal sends signals  $s = (s_i)_i$  according to the information structure  $(S, \pi)$ . Since some agents may not receive signals, we distinguish the signal sent to agent  $i$  and the signal received by agent  $i$ . Let  $s_i^R$  denote the signal received by agent  $i$ , and let  $s_i$  denote for the signal sent to agent  $i$  (as before). For each  $s_i \in S_i$ , if  $n_i = \infty$  then agent  $i$  receives signal  $s_i^R = s_i$ , while if  $n_i = \emptyset$  then he receives no signal, denoted  $s_i^R = \emptyset$ . Lastly, each agent  $i$  takes action  $a_i$  based on his received signal  $s_i^R$ . Formally, a strategy for agent  $i$  is a pair  $(n_i, \alpha_i)$ , which comprises his signal-receiving decision  $n_i \in \{\infty, \emptyset\}$  and a function  $\alpha_i : S_i \cup \{\emptyset\} \rightarrow \Delta(A_i)$  that assigns to a received signal  $s_i^R$  his (random) action  $\alpha_i(s_i^R)$ . The original model corresponds to each agent  $i$  making decision  $n_i = \infty$ .

**Theorem 2.** *Let a basic game  $G$  satisfy Assumption 1. In the modified model, the unique rationalizable outcome is action profile  $a^0$  under the information structure  $(S, \pi)$  of Theorem 1.*

We discuss an intuition for this result. Consider the example in Section 2. There are two agents, each of whom chooses action  $a_i \in \{0, 1\}$ ; they will play action profile  $a^0 = (0, 0)$  if they receive the signals from the principal, but could play action profile  $a^1 = (1, 1)$  if they did not (at the prior  $\mu^1 > \frac{2}{3}$ ). We will see that each agent  $i$  chooses to receive signal  $s_i$ , regardless of whether agent  $j \neq i$  receives signal  $s_j$ . Suppose that agent  $j$  does not receive signal  $s_j$ . If agent  $i$  does not receive signal  $s_i$ , he acquires no information about state  $\theta$  and plays some  $\alpha_i(\emptyset)$ . However, if agent  $i$  receives signal  $s_i$ , he can decide which action to take conditional on a received signal  $s_i^R$ . Consider agent  $i$ 's strategy  $(\infty, \alpha'_i)$  such that (i) If  $s_i^R = 0$  (which makes him pretty sure of state

<sup>11</sup>Once an agent receives a signal, he automatically updates his belief based on the (realized) signal. Hence, the only way to avoid updating a belief is to not receive the signal.

<sup>12</sup>The unobservability assumption is natural in private persuasion. The setting that agents  $-i$  learn what agent  $i$  will know does not suit private persuasion. If their choice  $(n_1, n_2, \dots, n_N)$  were common knowledge then there could exist a trivial equilibrium such that the agents could play some action profile  $a \neq a^0$ . For example, in the game in Section 2, if we have a prior  $\mu^1 \geq \frac{2}{3}$ , then there is an equilibrium such that each agent  $i$  chooses  $n_i = \emptyset$  and plays  $a_i = 1$  if  $n_1 = n_2 = \emptyset$  and  $a_i = 0$  otherwise.

$\theta = 0$ ), he takes action  $a_i = 0$  and (ii) if  $s_i^R \geq 1$  (which makes him sure of state  $\theta = 1$ ), he plays  $\alpha'_i(s_i^R) = \alpha_i(\emptyset)$ . We see that strategy  $(\emptyset, \alpha_i)$  is strictly dominated by strategy  $(\infty, \alpha'_i)$ , if  $\alpha_i(\emptyset)$  does not assign probability 1 to action  $a_i = 0$ .<sup>13</sup> Hence, agent  $i$  receives signal  $s_i$ . Similarly, we see that even when agent  $j$  receives signal  $s_j$ , agent  $i$  receives signal  $s_i$ . In either case, since each agent  $i$  receives signal  $s_i$ , the situation is identical to that of Theorem 1; therefore, the agents play action profile  $a^0 = (0, 0)$ . Lastly, we note that the commitment assumption is not necessary.

**Modified Model with the Sequential Protocol** Theorem 2 assumes that the principal sends signals  $s = (s_i)_i$  once according to the information structure  $(S, \pi)$ . Another way of sending information is to send messages sequentially, as considered in Section 4.1. In this setting, each agent may “filter” which messages to receive. For example, an agent may receive messages up to 2 but not more. Since this setting allows the agents to select information more flexibly, it may seem that the agents can coordinate to receive some messages and play some action profile  $a \neq a^0$ .

Formally, we allow each agent  $i$  to choose  $n_i$  from the set  $\{1, 2, \dots, \infty\} \cup \{\emptyset\}$  (instead of the set  $\{\infty, \emptyset\}$ ). Here  $n_i \in \{1, 2, \dots, \infty\}$  represents the largest message he is willing to receive. Note that if  $n_i = 1$  then agent  $i$  receives either message 0 or 1. For any  $n_i \in \{1, 2, \dots, \infty\}$ , if the principal sends message  $k_i$  (but not  $k_i + 1$ ) to agent  $i$ , then agent  $i$  receives message  $\min\{k_i, n_i\}$  and takes action  $\alpha_i(\min\{k_i, n_i\})$ . All the other settings are maintained. Even in this setting, the unique rationalizable outcome is action profile  $a^0$ .<sup>14</sup>

**Theorem 2’.** Let a basic game  $G$  satisfy Assumption 1. In the modified model with the sequential protocol, the unique rationalizable outcome is action profile  $a^0$  when the principal sends messages sequentially according to the information structure  $(S, \pi)$  of Theorem 1.

## 5 Discussion

We discuss two key assumptions for our results.

### 5.1 Cardinalities of Signal Spaces

Consider the example of Section 2. Recall that the principal’s payoff from our information structure is significantly different from those from other information structures (Figure 1). What causes this difference? An information structure  $(S, \pi)$  induces a distribution  $\tau$  over agents’ posterior beliefs. The key is the cardinality of its support,  $|\text{supp}(\tau)|$ . The cardinality is infinite in our information

<sup>13</sup>If  $\alpha_i(\emptyset)$  assigns probability 1 to action  $a_i = 0$  then we can see that the agents play action profile  $a = (0, 0)$ .

<sup>14</sup>Binmore and Samuelson (2001) consider a version of Rubinstein’s (1989) electronic mail game in which communication is strategic. Each player  $i$  (out of two) is allowed to choose  $m_i \in \{0, 1, \dots, \infty\}$ , which is the largest message player  $i$  sends. In some equilibrium, the players send only a finite number of messages, and coordinate on a relevant action profile (if they receive the messages). In contrast, we consider communication between a principal and agents. In our model, each agent decides the largest message he receives. Our result says that in every equilibrium, the agents cannot coordinate on a relevant action profile.

structure, while it is finite in the other information structures. Many existing studies on multi-agent persuasion assume that the support is finite by restricting a signal space  $S$  to a finite set (e.g., Bergemann and Morris, 2016a; Mathevet et al., 2020; Taneva, 2019).

We examine the implication of the finite-support assumption. We demonstrate that the finiteness assumption restricts a principal’s ability to manipulate agents. To this end, it suffices to consider an information structure  $(S, \pi)$  with the signal space  $S$  finite. This is because, for any finite-support distribution  $\tau$  over agents’ posteriors, if it is induced by some information structure, there is another information structure with its signal space finite.<sup>15</sup> In the example of Section 2, we obtain the following proposition:

**Proposition 1.** *Suppose that a signal space  $S$  is finite. As a prior  $\mu^1$  tends to 1, the limit of the principal’s (worst-case) payoff is 0, for any distributions  $\{\pi(\cdot | \theta)\}_{\theta \in \Theta}$ .*

We prove Proposition 1 as a special case of a more general result (Theorem 3) below. Note that Proposition 1 fixes the cardinality of a signal space  $S$  and allows the principal to change a distribution  $\pi$  as the prior varies. If we allow the principal to change the cardinality of a signal space  $S$  (not only a distribution  $\pi$ ) according to the prior  $\mu^1$ , then the (worst-case) payoff may be close to 1 when  $\mu^1$  is close to 1.

**Generalization** The implication of this finiteness assumption (Proposition 1) is generalizable. For this analysis, we weaken Assumption 1.

**Assumption 2.** *Let  $G = ((A_i, u_i)_i, \Theta, \mu)$  be a basic game with a finite state space  $\Theta$  such that for each  $\theta \in \Theta$ , the state- $\theta$  complete-information game  $(A_i, u_i(\cdot, \theta))_i$  has a strict Nash equilibrium, denoted  $a^\theta$ .*

**Theorem 3.** *Let a basic game  $G$  satisfy Assumption 2, and fix any state  $\theta^* \in \Theta$ . Suppose that a signal space  $S$  is finite. For any  $\epsilon > 0$ , there exists some  $\delta > 0$  such that if we have a prior  $\mu(\theta^*) > 1 - \delta$ , then regardless of the choice of distributions  $\{\pi(\cdot | \theta)\}_{\theta \in \Theta}$ , the Bayesian game  $\langle G, (S, \pi) \rangle$  has a Bayesian Nash equilibrium in which  $a^{\theta^*}$  is played with probability at least  $1 - \epsilon$ .*

**Proof of Proposition 1.** The example of Section 2 satisfies Assumption 2. It has a strict Nash equilibrium  $a^0 = (0, 0)$  in the state-0 complete information game and two strict Nash equilibria  $a^1 = (1, 1)$  and  $a^0$  in the state-1 complete information game. From Theorem 3, it follows that as  $\mu^1 \rightarrow 1$ , regardless of the choice of distributions  $\{\pi(\cdot | \theta)\}_{\theta \in \Theta}$ , the corresponding Bayesian game has a Bayesian Nash equilibrium such that the probability that  $a^1$  is played converges to 1, so that the limit of the principal’s (worst-case) payoffs is  $v(1, 1) = 0$ . ■

**Remark 2.** We illustrate the role of the cardinality of a signal space  $S$  in Theorems 1 and 3. As illustrated above, the key to Theorem 1 is that when agent  $i$  receives signal  $s_i = k \geq 1$ , he assigns

<sup>15</sup>To see this, fix any finite-support distribution  $\tau$  that is induced by some information structure  $(S, \pi)$ , where the signal space  $S$  may be infinite. For any two posteriors  $t, t' \in \text{supp}(\tau)$ , if  $t \neq t'$  then there are two distinct realizations  $s, s' \in S$  of signal profiles. Hence,  $|\text{supp}(\tau)| \leq |S|$ . If  $|\text{supp}(\tau)| < \infty$  but  $|S| = \infty$ , then we can ignore all realizations  $s \in S$  that generate no posteriors (i.e., have probability zero under the distribution  $\pi$ ).

a high probability to agents  $-i$  receiving signals  $s_{-i} = k - 1$ . This makes agent  $i$  uncertain about agents  $-i$ 's behavior, driving agent  $i$  into the “safe” action  $a_i^0$ . To create such uncertainty, an infinite number of different signals are needed. This is why we assume that the signal space  $S$  is infinite in Theorem 1. In contrast, if a signal space is finite as in Theorem 3, there is the “last” signal  $\bar{K}$  such that if agent  $i$  receives signal  $s_i = \bar{K}$ , he assigns a high probability to agents  $-i$  receiving the same signals  $s_{-i} = \bar{K}$ . This implies that the agents would approximate (non-trivial) common knowledge about their signals and thus the state, which may lead to action profile  $a \neq a^0$ . Although this intuition is based on the special information structure of Theorem 1, we show, in the proof of Theorem 3, that the same argument works for all possible information structures.  $\square$

## 5.2 Assumption 1

Using the following variant of the example in Section 2, we show that the condition  $\sum_{i \in I} p_i \leq 1$  of Assumption 1 is a tight requirement for Theorem 1. The variant obtains merely by replacing the payoff table 1 with the payoff table 2 below. In this variant, each agent receives payoff  $g$  from action profile  $a^1 = (1, 1)$  at state  $\theta = 1$ . All the other settings are maintained.

$\theta = 0$	$a_2 = 0$	$a_2 = 1$	$\theta = 1$	$a_2 = 0$	$a_2 = 1$
$a_1 = 0$	0, 0	0, -2	$a_1 = 0$	0, 0	0, -2
$a_1 = 1$	-2, 0	-2, -2	$a_1 = 1$	-2, 0	$g, g$

Table 2: the agents’ payoffs for the modified example

Suppose that  $g > 2$ . Then, Assumption 1 fails to hold in this variant. At state  $\theta = 1$ , for any small  $\epsilon > 0$ , action profile  $a^0 = (0, 0)$  is a strict  $(p_1^0 + \epsilon, p_2^0 + \epsilon)$ -dominant equilibrium, where  $p_1^0 = p_2^0 = \frac{g}{2+g}$  (thus,  $\sum_i (p_i^0 + \epsilon) > 1$ ); in contrast, action profile  $a^1 = (1, 1)$  is a strict  $(p_1^1 + \epsilon, p_2^1 + \epsilon)$ -dominant equilibrium, where  $p_1^1 = p_2^1 = \frac{2}{2+g}$  (thus,  $\sum_i (p_i^1 + \epsilon) < 1$ ). That is,  $a^1$  is “safe” but  $a^0$  is not. In the following proposition, we show that, regardless of the choice of an information structure  $(S, \pi)$ , the principal (regime) cannot always prevent agents from playing action profile  $a^1$ .

**Proposition 2.** *If  $g > 2$ , then for any (finite or infinite) signal space  $S$  and any distributions  $\{\pi(\cdot | \theta)\}_{\theta \in \Theta}$ , the limit of the principal’s (worst-case) payoff is 0 as the prior  $\mu^1$  tends to 1.*

Using Proposition 2, we illustrate an intuitive reason why the condition  $\sum_i p_i \leq 1$  of Assumption 1 is needed for Theorem 1. The key idea of this theorem is that the principal perturbs the agents’ (higher-order) beliefs so that each agent is uncertain what the other agents know and thus how they behave. Under such beliefs, they want to play the “safe” action profile, which is the strict  $\mathbf{p}$ -dominant equilibrium with  $\sum_i p_i \leq 1$ . Recall that when  $g > 2$  (under which Assumption 1 fails to hold), in the state-1 complete-information game,  $a^1 = (1, 1)$  is “safe” but  $a^0 = (0, 0)$  is “unsafe”; therefore, whenever the agents can play  $a^0$  (unsafe) in an equilibrium, they can also play  $a^1$  (safe) in another equilibrium. Consequently, the principal’s (worst-case) prediction is  $a^1$ , which yields payoff 0 for her. In the proof of Proposition 2, we show that for any prior  $\mu^1$  close to 1, the agents will be commonly sure about state  $\theta = 1$  and thus can play  $a^1$ .

Proposition 2 implies that the condition  $\sum_i p_i \leq 1$  is tight: If action profile  $a$  is *not* the strict  $\mathbf{p}$ -dominant equilibrium for any  $\mathbf{p}$  with  $\sum_i p_i \leq 1$ , then the principal cannot always persuade the agents to uniquely play  $a$ . As discussed above, in the state-1 complete-information game, whenever the agents can play the unsafe actions in an equilibrium, they can also play the safe ones in another equilibrium. To see the tightness, we consider how threshold 1 of  $\sum_i p_i \leq 1$  is related to the safety of actions. Since  $(p_i^0 - \frac{1}{2})(p_i^1 - \frac{1}{2}) < 0$  for any  $g \neq 2$  (where  $p_i^0 = \frac{g}{2+g}$  and  $p_i^1 = \frac{2}{2+g}$ ), exactly one of  $p_i^0$  or  $p_i^1$  is less than  $\frac{1}{2}$ . This implies that if action  $a_i$  is safe (i.e.,  $p_i^{a_i} < p_i^{1-a_i}$ ) then  $p_i^{a_i} < \frac{1}{2}$ . That is, the threshold  $\frac{1}{2}$  determines which action is safe. Hence, if action profile  $a = (a_1, a_2)$  is not the strict  $\mathbf{p}$ -dominant equilibrium for any  $\mathbf{p}$  such that  $\sum_i p_i^{a_i} \leq 1$ , at least one agent  $i$  has action  $a_i$  such that  $p_i^{a_i} > \frac{1}{2}$  and it is not safe for him; thus, action profile  $a$  could not be played uniquely.

## 6 Conclusion

We have studied a persuasion model in which a principal endogenously designs multiple agents' information structure. The principal wants the agents to take action profile  $a^0$  that consists of “safe” options for the agents (Assumption 1). The principal can persuade the agents to take  $a^0$  as a unique rationalizable outcome, regardless of a (non-degenerate) prior (Theorem 1). This result is built on our observation of the tight connection between multi-agent persuasion and information robustness of equilibrium. Even if the agents can strategically choose whether to receive information from the principal or not, they are still induced into  $a^0$  (Theorem 2). Additionally, we have examined the underlying assumptions for these results. In particular, the cardinality of a signal space can affect the principal's ability to manipulate the agents. It is essential to our results that the signal space is infinite; if it were finite, our results would not hold (Theorem 3). This result highlights the implications of finite signal spaces, which are often assumed in the literature.

Future research can be considered along several dimensions. First, since many applications of interest are intrinsically dynamic, with agents learning from an unknown state from past actions (e.g., Basu et al., 2020), it would be interesting to extend the analysis. Second, it would also be interesting to allow agents to communicate. In this extension, a principal must design information structures, while predicting how the agents communicate their information.

## A Omitted Proofs

### A.1 Theorem 1

The following notations are useful:  $s = k$  means that each agent  $i$  receives signal  $s_i = k$ , while  $s \geq k$  means that each agent  $i$  receives signal  $s_i \geq k$ . Similarly,  $s_{-i} = k$  means that each agent  $j \neq i$  receives signal  $s_j = k$ . Moreover, we write  $\pi^\theta(\cdot) = \pi(\cdot \mid \theta)$  for each  $\theta$ . Note that  $\pi^0(s = 0) = 1$  and  $\pi^\theta = \pi^1$  for each  $\theta \neq 0$ .

**Step 0:** Suppose that agent  $i$  receives  $s_i = 0$ . Then, his belief is:

$$\mathbb{P}(\theta = 0 \mid s_i = 0) = \frac{\mu^0}{\mu^0 + (1 - \mu^0)\pi^1(s_i = 0)}, \quad (1)$$

where  $\pi^1(s_i = 0) = (N - 1)\varphi$ . This belief is arbitrarily close to 1 if we take  $\varphi > 0$  small enough. Let  $\varphi < \frac{1}{N}$ , in particular. Hence,  $a_i^0$  is strictly dominant for agent  $i$ .

**Step 1:** Suppose that agent  $i$  receives  $s_i = 1$ . Then, his belief is:

$$\mathbb{P}(s_{-i} = 0 \mid s_i = 1) = \frac{\varphi}{\varphi + (1 - N\varphi)(\phi_I - \phi_i)}. \quad (2)$$

Since  $\varphi < \frac{1}{N}$ , this belief is arbitrarily close to 1 if we take  $\phi_I = \sum_{j \in I} \phi_j$  small enough. That is, he assigns probability at least  $p_i$  to agents  $-i$  taking  $a_{-i}^0$ , and thus chooses  $a_i^0$ .

**Step  $k \geq 2$ :** Suppose that agent  $i$  receives  $s_i = k$ . Then, his belief is:

$$\mathbb{P}(s_{-i} = k - 1 \mid s_i = k) = \frac{\phi_i}{\phi_i + (1 - \phi_I)(\phi_I - \phi_i)}. \quad (3)$$

Let  $p_I = \sum_{j \in I} p_j$ . Then,  $p_I \leq 1$  by Assumption 1. If we take  $\phi_i = (p_i/p_I)\phi_I$  for each  $i \in I$  then this belief is at least  $p_i$ .<sup>16</sup> That is, he assigns probability at least  $p_i$  to agents  $-i$  taking  $a_{-i}^0$ , and thus chooses  $a_i^0$ .

Lastly, we show that  $\mathbb{P}(s_i < \infty) = 1$ . Since  $\mathbb{P}(s \geq k) = (1 - \phi_I)^{k-1}\mathbb{P}(s \geq 1)$  for each  $k \geq 1$ , there exist  $j \in I$  and  $K < \infty$  with probability 1 such that  $s_j = K$ . Then,  $s_{-i} = K$  for all except for one, say  $i \neq j$ . Thus,  $s_i = K + 1$ .

## A.2 Theorem 2

We use the same notations as in the proof of Theorem 1. In addition, for each  $\alpha_i \in \Delta(A_i)$  and each  $a_i \in A_i$ , we use the following notations:

- $\alpha_i(s_i^R) = a_i$  means that  $\alpha_i(s_i^R)$  assigns probability 1 to action  $a_i$ .
- $\alpha_i(s_i^R) \neq a_i$  means that  $\alpha_i(s_i^R)$  does not assign probability 1 to action  $a_i$ .

We also write  $(n_i, \alpha_i^{n_i})$  for agent  $i$ 's strategy when we want to be explicit about his decision  $n_i$ .

We prove the theorem by showing that each agent  $i$ 's strategy  $(\emptyset, \alpha_i^\emptyset)$  is strictly dominated whenever  $\alpha_i^\emptyset(\emptyset) \neq a_i^0$ . Given any strategies  $(n_{-i}, \alpha_{-i})$  for agents  $-i$ , if agent  $i$  plays  $(\emptyset, \alpha_i^\emptyset)$  then his payoff is

$$\mu^0 \mathbb{E} \left[ u_i(\alpha_i^\emptyset(s_i^R), \alpha_{-i}, \theta) \mid \theta = 0 \right] + (1 - \mu^0) \mathbb{E} \left[ u_i(\alpha_i^\emptyset(s_i^R), \alpha_{-i}, \theta) \mid \theta \neq 0 \right], \quad (4)$$

In this proof, we will omit the input  $s_{-i}^R$  of the function  $\alpha_{-i}$ , as it causes no confusion. Recall that  $\mu^0$  is the prior probability of  $\theta = 0$ . If agent  $i$  plays another strategy  $(\infty, \alpha_i^\infty)$  such that  $\alpha_i^\infty(0) = a_i^0$

<sup>16</sup>If  $p_I > 1$ , the belief would be smaller than  $p_i$  for any small  $\phi_I > 0$ .

and  $\alpha_i^\infty(s_i) = \alpha_i^\emptyset(\emptyset)$  for each  $s_i \geq 1$ , then his payoff is

$$\mu^0 \mathbb{E} \left[ u_i(\alpha_i^\infty(s_i^R), \alpha_{-i}, \theta) \mid \theta = 0 \right] + (1 - \mu^0) \mathbb{E} \left[ u_i(\alpha_i^\infty(s_i^R), \alpha_{-i}, \theta) \mid \theta \neq 0 \right]. \quad (5)$$

It suffices to prove that (5) – (4) > 0 regardless of agents  $-i$ 's strategies  $(n_{-i}, \alpha_{-i})$ , where

$$(5) - (4) = \mu^0 \overbrace{\left( \mathbb{E} \left[ u_i(\alpha_i^\infty(s_i^R), \alpha_{-i}, \theta) \mid \theta = 0 \right] - \mathbb{E} \left[ u_i(\alpha_i^\emptyset(s_i^R), \alpha_{-i}, \theta) \mid \theta = 0 \right] \right)}^{\equiv (6')} + (1 - \mu^0) \underbrace{\left( \mathbb{E} \left[ u_i(\alpha_i^\infty(s_i^R), \alpha_{-i}, \theta) \mid \theta \neq 0 \right] - \mathbb{E} \left[ u_i(\alpha_i^\emptyset(s_i^R), \alpha_{-i}, \theta) \mid \theta \neq 0 \right] \right)}_{\equiv (6'')}. \quad (6)$$

At state  $\theta = 0$ , the principal sends  $s_i = 0$  by the construction of  $(S, \pi)$ . If  $n_i = \infty$  then since  $s_i^R = 0$ , we have  $\alpha_i^\infty(s_i^R) = a_i^0$ . If  $n_i = \emptyset$  then since  $s_i^R = \emptyset$ , we have  $\alpha_i^\emptyset(s_i^R) = \alpha_i^\emptyset(\emptyset)$ . Hence

$$(6') = \mathbb{E} \left[ u_i(a_i^0, \alpha_{-i}, \theta) \mid \theta = 0 \right] - \mathbb{E} \left[ u_i(\alpha_i^\emptyset(\emptyset), \alpha_{-i}, \theta) \mid \theta = 0 \right].$$

Since  $\alpha_i^\emptyset(\emptyset) \neq a_i^0$ , there exist some  $\eta \in [0, 1)$  and some  $\beta_i \in \Delta(A_i)$  with support  $\text{supp}(\beta_i) \subseteq A_i \setminus \{a_i^0\}$  such that  $\alpha_i^\emptyset(\emptyset)$  plays  $a_i^0$  with probability  $\eta$  and  $\beta_i$  with probability  $1 - \eta$ . Then,

$$(6') = (1 - \eta) \left( \mathbb{E} \left[ u_i(a_i^0, \alpha_{-i}, \theta) \mid \theta = 0 \right] - \mathbb{E} \left[ u_i(\beta_i, \alpha_{-i}, \theta) \mid \theta = 0 \right] \right). \quad (7)$$

At state  $\theta \neq 0$ , the principal sends  $s_i = 0$  with probability  $(N - 1)\varphi$  and  $s_i \geq 1$  with probability  $1 - (N - 1)\varphi$ . If  $n_i = \infty$  then since  $s_i^R = s_i$ , it must be that  $\alpha_i^\infty(s_i^R)$  plays  $a_i^0$  with probability  $(N - 1)\varphi$  and  $\alpha_i^\emptyset(\emptyset)$  with probability  $1 - (N - 1)\varphi$ . Hence,

$$\mathbb{E} \left[ u_i(\alpha_i^\infty(s_i^R), \alpha_{-i}, \theta) \mid \theta \neq 0 \right] = (N - 1)\varphi \mathbb{E} \left[ u_i(a_i^0, \alpha_{-i}, \theta) \mid \theta \neq 0, s_i = 0 \right] + (1 - (N - 1)\varphi) \mathbb{E} \left[ u_i(\alpha_i^\emptyset(\emptyset), \alpha_{-i}, \theta) \mid \theta \neq 0, s_i \geq 1 \right].$$

If  $n_i = \emptyset$  then since  $s_i^R = \emptyset$  regardless of signal  $s_i$ , it follows that

$$\mathbb{E} \left[ u_i(\alpha_i^\emptyset(s_i^R), \alpha_{-i}, \theta) \mid \theta \neq 0 \right] = (N - 1)\varphi \mathbb{E} \left[ u_i(\alpha_i^\emptyset(\emptyset), \alpha_{-i}, \theta) \mid \theta \neq 0, s_i = 0 \right] + (1 - (N - 1)\varphi) \mathbb{E} \left[ u_i(\alpha_i^\emptyset(\emptyset), \alpha_{-i}, \theta) \mid \theta \neq 0, s_i \geq 1 \right].$$

Hence,

$$(6'') = (N - 1)\varphi \left( \mathbb{E} \left[ u_i(a_i^0, \alpha_{-i}, \theta) \mid \theta \neq 0, s_i = 0 \right] - \mathbb{E} \left[ u_i(\alpha_i^\emptyset(\emptyset), \alpha_{-i}, \theta) \mid \theta \neq 0, s_i = 0 \right] \right).$$

Since  $\alpha_i^\emptyset(\emptyset)$  plays  $a_i^0$  with probability  $\eta$  and  $\beta_i$  with probability  $1 - \eta$ , it follows that

$$(6'') = (1 - \eta)(N - 1)\varphi \left( \mathbb{E} \left[ u_i(a_i^0, \alpha_{-i}, \theta) \mid \theta \neq 0, s_i = 0 \right] - \mathbb{E} \left[ u_i(\beta_i, \alpha_{-i}, \theta) \mid \theta \neq 0, s_i = 0 \right] \right). \quad (8)$$

Note that (6) =  $\mu^0 \times (7) + (1 - \mu^0) \times (8)$ . Since  $\eta < 1$ , we have (6) > 0 if and only if (9) > 0, where

$$\begin{aligned} & \mu^0 \left( \overbrace{\mathbb{E}\left[u_i(a_i^0, \alpha_{-i}, \theta) \mid \theta = 0\right] - \mathbb{E}\left[u_i(\beta_i, \alpha_{-i}, \theta) \mid \theta = 0\right]}^{\equiv (9')} \right) + (1 - \mu^0)(N - 1)\varphi \\ & \times \left( \overbrace{\mathbb{E}\left[u_i(a_i^0, \alpha_{-i}, \theta) \mid \theta \neq 0, s_i = 0\right] - \mathbb{E}\left[u_i(\beta_i, \alpha_{-i}, \theta) \mid \theta \neq 0, s_i = 0\right]}^{\equiv (9'')} \right). \end{aligned} \quad (9)$$

Since  $a_i^0$  is strictly dominant at state  $\theta = 0$  and  $\beta_i$  assigns zero probability to  $a_i^0$ , it follows that (9') is bounded away from zero. That is, there exists some  $\bar{u} > 0$  such that (9') >  $\bar{u}$  for each  $\alpha_{-i}$ . Since the action spaces are finite, (9'') is finite for each  $\beta_i$  and each  $\alpha_{-i}$ . If  $\varphi > 0$  is small enough then (9) > 0. Hence, agent  $i$ 's strategy  $(\emptyset, \alpha_i^\emptyset)$  is strictly dominated whenever  $\alpha_i^\emptyset(\emptyset) \neq a_i^0$ .

Since each agent  $i$ 's every rationalizable strategy  $(n_i, \alpha_i)$  is such that either (i)  $n_i = \emptyset$  and  $\alpha_i^\emptyset(\emptyset) = a_i^0$  or (ii)  $n_i = \infty$ , it suffices to show that his uniquely rationalizable action is action  $a_i^0$  in case (ii), where he plays strategy  $(\infty, \alpha_i^\infty)$ . First, if  $s_i^R = 0$ , his belief  $\mathbb{P}(\theta = 0 \mid s_i^R = 0)$  is equal to the right-hand side of equation (1) with  $\pi^1(s_i = 0) = (N - 1)\varphi$ , and is close to 1, which leads him to take action  $a_i^0$ . Hence,  $\alpha_i^\infty(0) = a_i^0$  for each  $i$ . Second, if  $s_i^R = 1$ , his belief  $\mathbb{P}(s_{-i} = 0 \mid s_i^R = 1)$  is equal to the right-hand side of equation (2), and is close to 1. Each agent  $j \neq i$  takes action  $a_j^0$  when  $s_j = 0$  is sent, regardless of whether agent  $j$  takes (i)  $n_j = \emptyset$  and  $\alpha_j^\emptyset(\emptyset) = a_j^0$  or (ii)  $n_j = \infty$  and  $\alpha_j^\infty(0) = a_j^0$ . When receiving  $s_i^R = 1$ , agent  $i$  assigns probability close to 1 to agents  $-i$  taking actions  $a_{-i}^0$ , which leads agent  $i$  to take action  $a_i^0$ . Hence,  $\alpha_i^\infty(1) = a_i^0$  for each  $i$ . Third, if  $s_i^R = k \geq 2$ , his belief  $\mathbb{P}(s_{-i} = k - 1 \mid s_i^R = k)$  is equal to the right-hand side of equation (3), and is at least  $p_i$ . Each agent  $j \neq i$  takes action  $a_j^0$  when  $s_j = k - 1$  is sent, regardless of whether agent  $j$  takes (i)  $n_j = \emptyset$  and  $\alpha_j^\emptyset(k - 1) = a_j^0$  or (ii)  $n_j = \infty$  and  $\alpha_j^\infty(k - 1) = a_j^0$ . Hence, when receiving  $s_i^R = k$ , agent  $i$  assigns probability at least  $p_i$  to agents  $-i$  taking actions  $a_{-i}^0$ , which leads agent  $i$  to take action  $a_i^0$ . By induction,  $\alpha_i^\infty(k) = a_i^0$  for each  $i$  and each  $k$ .

### A.3 Theorem 2'

For brevity, we use the notation  $s_i$  and  $s_i^R$ :  $s_i = k$  means that the principal sends message  $k$  but not  $k + 1$  to agent  $i$ , while  $s_i^R = k$  means that agent  $i$  receives message  $k$  but not  $k + 1$ . In particular, if the principal sends message  $k$  but not  $k + 1$  to agent  $i$  and he makes decision  $n_i \in \{1, 2, \dots, \infty\}$  then he receives  $s_i^R = \min\{n_i, k\}$ . We use the same notation as in the proof of Theorem 2.

**Step 1** For any small enough  $\varphi > 0$ , we show that each agent  $i$ 's strategy  $(\emptyset, \alpha_i^\emptyset)$  is strictly dominated whenever  $\alpha_i^\emptyset(\emptyset) \neq a_i^0$ . Fix any strategies  $(n_{-i}, \alpha_{-i})$  for agents  $-i$ . Since agent  $i$  receives no message, denoted  $s_i^R = \emptyset$ , he plays  $\alpha_i^\emptyset(\emptyset)$ , and his payoff from strategy  $(\emptyset, \alpha_i^\emptyset)$  is

$$\sum_{l \in S_i \setminus \{0\}} \mathbb{P}(s_i = l) \mathbb{E}\left[u_i(\alpha_i^\emptyset(\emptyset), \alpha_{-i}, \theta) \mid s_i = l\right] + \mathbb{P}(s_i = 0) \mathbb{E}\left[u_i(\alpha_i^\emptyset(\emptyset), \alpha_{-i}, \theta) \mid s_i = 0\right]. \quad (10)$$

In this proof, we will omit the input  $s_{-i}^R$  of the function  $\alpha_{-i}$ , as it causes no confusion. If agent  $i$  plays strategy  $(\infty, \alpha_i^\infty)$  such that  $\alpha_i^\infty(0) = a_i^0$  and  $\alpha_i^\infty(l) = \alpha_i^\emptyset(\emptyset)$  for each  $l \in S_i \setminus \{0\}$ , then his payoff is

$$\sum_{l \in S_i \setminus \{0\}} \mathbb{P}(s_i = l) \mathbb{E} \left[ u_i(\alpha_i^\infty(l), \alpha_{-i}, \theta) \mid s_i = l \right] + \mathbb{P}(s_i = 0) \mathbb{E} \left[ u_i(\alpha_i^\infty(0), \alpha_{-i}, \theta) \mid s_i = 0 \right]. \quad (11)$$

Since  $\alpha_i^\infty(l) = \alpha_i^\emptyset(\emptyset)$  for each  $l \in S_i \setminus \{0\}$ , the first terms of payoffs (10) and (11) are equal to each other. Since  $\mathbb{P}(s_i = 0) > 0$  and  $\alpha_i^\infty(0) = a_i^0$ , it follows that (11) > (10) if and only if

$$\mathbb{E} \left[ u_i(a_i^0, \alpha_{-i}, \theta) \mid s_i = 0 \right] > \mathbb{E} \left[ u_i(\alpha_i^\emptyset(\emptyset), \alpha_{-i}, \theta) \mid s_i = 0 \right]. \quad (12)$$

By construction,  $\mathbb{P}(\theta = 0 \mid s_i = 0)$  is equal to the right-hand side of equation (1) with  $\pi^1(s_i = 0) = (N-1)\varphi$ , and is close to 1 if  $\varphi$  is small enough. This implies that  $a_i^0$  is strictly dominant when  $s_i = 0$ . Hence, (12) holds, and strategy  $(\emptyset, \alpha_i^\emptyset)$  is strictly dominated whenever  $\alpha_i^\emptyset(\emptyset) \neq a_i^0$ . Moreover, as suggested by this argument, it holds that for each  $i \in I$  and his every rationalizable strategy  $(n_i, \alpha_i^{n_i})$ , if  $n_i \geq 1$  then  $\alpha_j(0) = a_j^0$ .<sup>17</sup>

In summary, if  $(n_j, \alpha_j)$  is agent  $j$ 's rationalizable strategy, then it satisfies either:

- (1a)  $n_j = \emptyset$  and  $\alpha_j(\emptyset) = a_j^0$ ; or
- (1b)  $n_j \geq 1$  and  $\alpha_j(0) = a_j^0$ .

**Step 2** For any small enough  $\phi_I > 0$ , we show that each agent  $i$ 's strategy  $(1, \alpha_i^1)$  is (iteratively) strictly dominated whenever  $\alpha_i^1(1) \neq a_i^0$ . Fix any rationalizable strategies  $(n_{-i}, \alpha_{-i})$  for agents  $-i$ , which satisfy either (1a) or (1b). Then, agent  $i$ 's payoff is

$$\sum_{l \in S_i \setminus \{1\}} \mathbb{P}(s_i = l) \mathbb{E} \left[ u_i(\alpha_i^1(\min\{1, l\}), \alpha_{-i}, \theta) \mid s_i = l \right] + \mathbb{P}(s_i = 1) \mathbb{E} \left[ u_i(\alpha_i^1(1), \alpha_{-i}, \theta) \mid s_i = 1 \right], \quad (13)$$

where, since  $n_i = 1$ , agent  $i$  receives  $s_i^R = \min\{1, l\}$  when  $s_i = l$  is sent. Let  $(\infty, \alpha_i^\infty)$  now denote agent  $i$ 's strategy such that  $\alpha_i^\infty(1) = a_i^0$  and  $\alpha_i^\infty(l) = \alpha_i^1(\min\{1, l\})$  for each  $l \in S_i \setminus \{1\}$ , then his payoff is

$$\sum_{l \in S_i \setminus \{1\}} \mathbb{P}(s_i = l) \mathbb{E} \left[ u_i(\alpha_i^\infty(l), \alpha_{-i}, \theta) \mid s_i = l \right] + \mathbb{P}(s_i = 1) \mathbb{E} \left[ u_i(\alpha_i^\infty(1), \alpha_{-i}, \theta) \mid s_i = 1 \right]. \quad (14)$$

Since  $\alpha_i^\infty(l) = \alpha_i^1(\min\{1, l\})$  for each  $l \in S_i \setminus \{1\}$ , the first terms of payoffs (13) and (14) are equal to each other. Since  $\mathbb{P}(s_i = 1) > 0$  and  $\alpha_i^\infty(1) = a_i^0$ , it follows that (14) > (13) if and only if

$$\mathbb{E} \left[ u_i(a_i^0, \alpha_{-i}, \theta) \mid s_i = 1 \right] > \mathbb{E} \left[ u_i(\alpha_i^1(1), \alpha_{-i}, \theta) \mid s_i = 1 \right]. \quad (15)$$

---

<sup>17</sup>When  $s_i^R = 0$ , agent  $i$  assigns high probability (1) to state  $\theta = 0$ , thereby taking action  $\alpha_i^{n_i}(0) = a_i^0$ .

Every agent  $j$ 's rationalizable strategy  $(n_j, \alpha_j)$  satisfies either (1a) or (1b). When  $s_j = 0$  is sent, agent  $j$  takes action  $a_j^0$  regardless of whether (1a) holds (so that  $s_j^R = \emptyset$  and  $\alpha_j(\emptyset) = a_j^0$ ) or (1b) holds (so that  $s_j^R = 0$  and  $\alpha_j(0) = a_j^0$ ). Agent  $i$ 's belief  $\mathbb{P}(\alpha_{-i}(s_{-i}^R) = a_{-i}^0 \mid s_i = 1) \geq \mathbb{P}(s_{-i} = 0 \mid s_i = k)$  is at least the right-hand side of equation (2), which is close to 1 if  $\phi_I > 0$  is small enough. Thus,  $\mathbb{P}(\alpha_{-i}(s_{-i}^R) = a_{-i}^0 \mid s_i = 1)$  is close to 1, which makes action  $a_i^0$  uniquely optimal. Hence, (15) holds, and strategy  $(1, \alpha_i^1)$  is (iteratively) strictly dominated whenever  $\alpha_i^1(1) \neq a_i^0$ . Moreover, as implied by this argument, for each agent  $i$ 's every rationalizable strategy  $(n_i, \alpha_i^{n_i})$ , if  $n_i \geq 2$  then  $\alpha_i^{n_i}(l) = a_i^0$  for each  $l = 0, 1$ .

In summary, if  $(n_j, \alpha_j)$  is agent  $j$ 's rationalizable strategy, then it satisfies either:

- (2a)  $n_j = \emptyset$  and  $\alpha_j(\emptyset) = a_j^0$ ; or
- (2b)  $n_j \geq 1$  and  $\alpha_j(l) = a_j^0$  for each  $l = 0, 1$ .

**Step 3** By induction, we show that for each  $i \in I$  and for each  $k \geq 2$ , if  $\alpha_i^k(k) \neq a_i^0$  then  $(k, \alpha_i^k)$  is (iteratively) strictly dominated. Given any  $k \geq 2$ , fix any rationalizable strategies  $(n_{-i}, \alpha_{-i})$  for agents  $-i$ , where for each  $j \neq i$ , agent  $j$ 's rationalizable strategy  $(n_j, \alpha_j)$  satisfies either condition:

- (3a)  $n_j = \emptyset$  and  $\alpha_j(\emptyset) = a_j^0$ ; or
- (3b)  $n_j \geq 1$  and  $\alpha_j(\min\{k, l\}) = a_j^0$  for each  $l = 0, 1, \dots, k-1$ .

Note that if  $k = 2$ , (3a) and (3b) are identical to (2a) and (2b). Agent  $i$ 's payoff from strategy  $(k, \alpha_i^k)$  is

$$\sum_{l \in S_i \setminus \{k\}} \mathbb{P}(s_i = l) \mathbb{E}[u_i(\alpha_i^k(\min\{k, l\}), \alpha_{-i}, \theta) \mid s_i = l] + \mathbb{P}(s_i = k) \mathbb{E}[u_i(\alpha_i^k(k), \alpha_{-i}, \theta) \mid s_i = k], \quad (16)$$

where, since  $n_i = k$ , agent  $i$  receives  $s_i^R = \min\{k, l\}$  when  $s_i = l$  is sent. Let  $(\infty, \alpha_i^\infty)$  now denote agent  $i$ 's strategy such that  $\alpha_i^\infty(k) = a_i^0$  and  $\alpha_i^\infty(l) = \alpha_i^1(\min\{k, l\})$  for each  $l \in S_i \setminus \{k\}$ , then his payoff is

$$\sum_{l \in S_i \setminus \{k\}} \mathbb{P}(s_i = l) \mathbb{E}[u_i(\alpha_i^\infty(l), \alpha_{-i}, \theta) \mid s_i = l] + \mathbb{P}(s_i = k) \mathbb{E}[u_i(\alpha_i^\infty(k), \alpha_{-i}, \theta) \mid s_i = k]. \quad (17)$$

Since  $\alpha_i^\infty(l) = \alpha_i^k(\min\{k, l\})$  for each  $l \in S_i \setminus \{k\}$ , the first terms of payoffs (16) and (17) are equal to each other. Since  $\mathbb{P}(s_i = k) > 0$  and  $\alpha_i^\infty(k) = a_i^0$ , it follows that (17)  $>$  (16) if and only if

$$\mathbb{E}[u_i(a_i^0, \alpha_{-i}(s_{-i}^R), \theta) \mid s_i = k] > \mathbb{E}[u_i(\alpha_i^k(k), \alpha_{-i}(s_{-i}^R), \theta) \mid s_i = k]. \quad (18)$$

Every agent  $j$ 's rationalizable strategy  $(n_j, \alpha_j)$  satisfies either (3a) or (3b). When  $s_j = k-1$  is sent, agent  $j$  takes action  $a_j^0$  regardless of whether (3a) holds (so that  $s_j^R = \emptyset$  and  $\alpha_j(\emptyset) = a_j^0$ ) or (3b) holds (so that  $s_j^R = \min\{k-1, n_j\}$  and  $\alpha_j(\min\{k-1, n_j\}) = a_j^0$ ). Agent  $i$ 's belief  $\mathbb{P}(\alpha_{-i}(s_{-i}^R) = a_{-i}^0 \mid s_i = k) \geq \mathbb{P}(s_{-i} = k-1 \mid s_i = k)$  is at least the right-hand side of equation (3), which is at least  $p_i$ . Thus,  $\mathbb{P}(\alpha_{-i}(s_{-i}^R) = a_{-i}^0 \mid s_i = k) \geq p_i$ , which makes a ction  $a_i^0$  uniquely optimal. Hence, (18) holds, and strategy  $(k, \alpha_i^k)$  is (iteratively) strictly dominated whenever  $\alpha_i^k(k) \neq a_i^0$ .

Finally, we note that each agent  $i$ 's rationalizable strategy  $(n_i, \alpha_i)$  satisfies one of the following conditions: (i)  $n_i = \emptyset$  and  $\alpha_i(\emptyset) = a_i^0$ ; (ii)  $n_i = 1, 2, \dots$  and  $\alpha_i(l) = a_i^0$  for each  $l = 0, 1, \dots, n_i$ ; or (iii)  $n_i = \infty$  and  $\alpha_i(l) = a_i^0$  for each  $l = 0, 1, \dots$ . In any case, agent  $i$  takes action  $a_i^0$  regardless of his received message.

#### A.4 Theorem 3

**Preliminaries** Let  $\Omega = \Theta \times S$  denote a (finite) state space, with a generic element  $\omega = (\theta, s_i, s_{-i})$ . Each agent  $i$  learns only his signal realization  $s_i$  at a state  $\omega$ . For any event  $E \subset \Omega$ , let  $B_i^q(E) = \{\omega : \mathbb{P}(E \mid s_i) \geq q\}$  be the event that agent  $i$  assigns to the event  $E$  probability at least  $q$ , and let  $B^q(E) = \bigcap_{i \in I} B_i^q(E)$ . Let  $(B^q)^1(E) = B^q(E)$  and  $(B^q)^k(E) = B^q((B^q)^{k-1}(E))$  for each  $k \geq 2$ . Define an event  $C^q(E) = \bigcap_{k=1}^{\infty} (B^q)^k(E)$ , and we say that the event  $E$  is common  $q$ -belief when a state  $\omega \in C^q(E)$  is realized. To ease notation, let  $C^q(\theta) = C^q(\{\theta\} \times S)$ . Analogously, let  $B_i^q(\theta) = B_i^q(\{\theta\} \times S)$  and  $B^q(\theta) = B^q(\{\theta\} \times S)$ .

In this proof, we will use the following result due to [Monderer and Samet \(1989\)](#):

**Lemma A.0.** *In the setting of Theorem 3, the following two hold:*

- (i) *Let  $\mathcal{G}$  be the Bayesian game defined by the basic game  $G$  with an information structure  $(S, \pi)$ . There exists some  $q \in (0, 1)$  such that the game  $\mathcal{G}$  has a Bayesian Nash equilibrium that plays action profile  $a^{\theta^*}$  on the event  $C^q(\theta^*)$ .<sup>18</sup>*
- (ii) *For each  $q$ ,  $\omega \in C^q(\theta^*)$  if and only if there exists an event  $F \ni \omega$  such that  $F \subset B^q(F)$  and  $F \subset B^q(\theta^*)$ .*

**Proof of Theorem 3** By Lemma A.0, it suffices to construct an event  $F$  such that  $F \subset B^q(F)$  and  $F \subset B^q(\theta^*)$  and that  $\mathbb{P}(F) \rightarrow 1$  as  $\mu(\theta^*) \rightarrow 1$ .

We define the event  $F$ . For each  $i \in I$  and each  $k \in \mathbb{N}$ , let  $S_i^0 = S_i$  and  $S_i^k = S_i^{k-1} \setminus D_i^k$ , where

$$D_i^k = \left\{ s_i \in S_i^{k-1} : \mathbb{P}(\{\theta^*\} \times S^{k-1} \mid s_i) \leq q \right\}.$$

As usual, let  $S^{k-1} = \prod_{i \in I} S_i^{k-1}$  for each  $k \in \mathbb{N}$ . It is immediate that  $S^k \subset S^{k-1}$  for each  $k \in \mathbb{N}$ . Since  $\Theta \times S$  is finite, there exists some  $K \geq 0$  such that  $S^{K+1} = S^K$ . Then, define  $F \equiv \{\theta^*\} \times S^K$ .

We show the desired properties mentioned above. First, we show that  $F \subset B^q(F)$ . Since  $S^{K+1} = S^K$ , it must be that  $D^{K+1} = \emptyset$ , which implies that for each  $s = (s_i)_i \in S^K$ ,  $\mathbb{P}(F \mid s_i) > q$ , where we recall  $F = \{\theta^*\} \times S^K$ . That is,  $F \subset B^q(F)$ . Second, we note that  $F \subset B^q(\theta^*)$ . By construction,  $S_i^K \subset S_i^1$  and  $F \subset B_i^1(\theta^*)$  for each  $i \in I$ .<sup>19</sup> Hence,  $F \subset B^q(\theta^*)$ . Third, we show that  $\mathbb{P}(F) \rightarrow 1$  as  $\mu(\theta^*) \rightarrow 1$ . Now we use the following lemma, which we will prove later.

**Lemma A.1.** *For each  $i \in I$ ,  $\mathbb{P}(\bigcup_{k=1}^K D_i^k) \rightarrow 0$  as  $\mu(\theta^*) \rightarrow 1$ .*

<sup>18</sup>Recall that action profile  $a^{\theta^*}$  is a strict Nash equilibrium of the state- $\theta^*$  complete-information game  $(A_i, u_i(\cdot, \theta^*))_i$ .

<sup>19</sup>If  $K = 0$  then by definition  $S_i^1 = S_i^0$ ; if  $K \geq 1$  then  $S_i^K \subset S_i^{K-1} \subset \dots \subset S_i^0$ .

Note that

$$\mathbb{P}(F) \geq 1 - \sum_{i \in I} \mathbb{P}\left(\{\theta^*\} \times \bigcup_{k=1}^K D_i^k \times S_{-i}\right) \geq 1 - \sum_{i \in I} \mathbb{P}\left(\bigcup_{k=1}^K D_i^k\right).$$

In view of Lemma A.1, this inequality implies that  $\mathbb{P}(F) \rightarrow 1$  as  $\mu(\theta^*) \rightarrow 1$ . Hence, we obtain the desired result.  $\blacksquare$

**Proof of Lemma A.1.** It suffices to prove that  $\mathbb{P}(D_i^k) \rightarrow 0$  as  $\mu(\theta^*) \rightarrow 1$ , for each  $k = 1, 2, \dots, K$ . The proof is by induction.

**Case  $k = 1$ :** For each  $i \in I$ , as  $\mu(\theta^*) \rightarrow 1$ , we have  $\mathbb{P}(D_i^1 \times S_{-i} \mid \theta^*) \rightarrow 0$  and thus

$$\mathbb{P}\left(D_i^1\right) = \mu(\theta^*)\mathbb{P}\left(D_i^1 \times S_{-i} \mid \theta^*\right) + \sum_{\theta \neq \theta^*} \mu(\theta)\mathbb{P}\left(D_i^1 \times S_{-i} \mid \theta\right) \rightarrow 0.$$

**Case  $k \geq 2$ :** Suppose, for induction, that  $\mathbb{P}(D_{-i}^{k'}) \rightarrow 0$  as  $\mu(\theta^*) \rightarrow 1$ , for each  $i \in I$  and each  $k' \leq k - 1$ . If agent  $i$  observes  $s_i \in D_i^k$ , then his belief is  $\mathbb{P}(\{\theta^*\} \times \{s_i\} \times S_{-i}^{k-1} \mid s_i) \leq q$ , where  $S_{-i}^{k-1} = \prod_{j \neq i} S_j^{k-1}$ . By Bayes' rule,  $\frac{\mathbb{P}(\{\theta^*\} \times \{s_i\} \times S_{-i}^{k-1})}{\mathbb{P}(\Theta \times \{s_i\} \times S_{-i})} \leq q$ , which is equivalent to

$$\mathbb{P}\left(\Theta \times \{s_i\} \times S_{-i}\right) \leq \frac{1}{1-q} \left\{ \mathbb{P}\left(\Theta \times \{s_i\} \times S_{-i}\right) - \mathbb{P}\left(\{\theta^*\} \times \{s_i\} \times S_{-i}^{k-1}\right) \right\}. \quad (19)$$

Since  $S_{-i}^{k-1} = S_{-i} \setminus \bigcup_{k'=1}^{k-1} D_{-i}^{k'}$ , it follows that

$$\begin{aligned} & \left(\Theta \times \{s_i\} \times S_{-i}\right) \setminus \left(\{\theta^*\} \times \{s_i\} \times S_{-i}^{k-1}\right) \\ & \subset \left(\Theta \times \{s_i\} \times \bigcup_{k'=1}^{k-1} D_{-i}^{k'}\right) \cup \left(\left(\Theta \setminus \{\theta^*\}\right) \times \{s_i\} \times S_{-i}\right). \end{aligned}$$

Hence,

$$(19) \leq \frac{1}{1-q} \left\{ \mathbb{P}\left(\Theta \times \{s_i\} \times \bigcup_{k'=1}^{k-1} D_{-i}^{k'}\right) + \mathbb{P}\left(\left(\Theta \setminus \{\theta^*\}\right) \times \{s_i\} \times S_{-i}\right) \right\}. \quad (20)$$

By the induction assumption, (20)  $\rightarrow 0$  as  $\mu(\theta^*) \rightarrow 1$ . Hence,  $\mathbb{P}(\Theta \times \{s_i\} \times S_{-i}) \rightarrow 0$ . Since  $D_i^k$  is finite, it must be that  $\mathbb{P}(D_i^k) \rightarrow 0$  as  $\mu(\theta^*) \rightarrow 1$ .  $\blacksquare$

## A.5 Proposition 2

We use the same notation as in Appendix A.4 (cf. the paragraph of Preliminaries). Note that action profile (1, 1) is a strict Nash equilibrium in the state-1 complete-information game. Analogously to Lemma A.0, any Bayesian game  $\mathcal{G}$ , defined by the basic game  $G$  of this variant and an information structure  $(S, \pi)$ , has a Bayesian Nash equilibrium that plays action profile (1, 1) on the event  $C^p(\theta = 1)$ , where for any small enough  $\epsilon > 0$ ,  $p = \frac{2}{2+\epsilon} + \epsilon < \frac{1}{2}$ . Hence, we are interested in the probability of  $C^p(\theta = 1)$  given  $(S, \pi)$ . From the critical path result (Kajii and Morris, 1997), it

follows that

$$\mathbb{P}(C^p(\theta = 1)) \geq 1 - (1 - \mu^1) \frac{1 - p}{1 - 2p}.$$

This inequality is true regardless of the agents' higher-order belief specifications and thus regardless of  $(S, \pi)$ . Note that  $\mathbb{P}(C^p(\theta = 1)) \rightarrow 1$  as  $\mu^1 \rightarrow 1$ . Therefore, the limit of the principal's (worst-case) payoff is  $v(1, 1) = 0$  as  $\mu^1 \rightarrow 1$ .

## B Details of Example in Section 2

In this appendix, we examine the case of binary signals. Fix the signal space  $S_i = \{0, 1\}$  for each agent  $i = 1, 2$ . If we have  $\mu^1 \leq \frac{2}{3}$ , the regime can forestall the coordinated attack and achieve the optimal payoff 1 by sending no information. Hence, it suffices to consider  $\mu^1 > \frac{2}{3}$ .

First, we consider the case of public signals. Since the agents observe the same signal realization, it must be that  $\pi(s_1 \neq s_2 | \theta) = 0$ . The optimal distribution is as follows: If  $\theta = 0$ , the regime sends  $(s_1, s_2) = (0, 0)$  with probability 1; if  $\theta = 1$ , it sends  $(s_1, s_2)$  according to Table 3.

$\theta = 1$	$s_2 = 0$	$s_2 = 1$
$s_1 = 0$	$2\nu$	$0$
$s_1 = 1$	$0$	$1 - 2\nu$

Table 3: public information structure at state  $\theta = 1$  ( $\nu = \mu^0/\mu^1 \in (0, \frac{1}{2})$ )

If agent  $i$  receives  $s_i = 0$ , he assigns to  $\theta = 1$  probability  $\frac{\mu^1 \times 2\nu}{\mu^0 \times 1 + \mu^1 \times 2\nu} = \frac{2}{3}$  and thus only action  $a_i = 0$  is rationalizable. If agent  $i$  receives  $s_i = 1$ , he assigns to  $\theta = 0$  with probability 0 and thus both actions  $a_i = 0, 1$  are rationalizable. Hence, the regime's payoff under the adversarial selection is  $\mu^0 + \mu^1 \times 2\nu = 3(1 - \mu^1)$ .

Second, we consider the case of private signals. Since the agents may observe different signal realizations, it may be that  $\pi(s_1 \neq s_2 | \theta) \geq 0$ . The optimal distribution is as follows: If  $\theta = 0$ , the regime sends  $(s_1, s_2) = (0, 0)$  with probability 1; if  $\theta = 1$ , it sends  $(s_1, s_2)$  according to Table 4.

$\theta = 1$	$s_2 = 0$	$s_2 = 1$
$s_1 = 0$	$0$	$2\nu$
$s_1 = 1$	$2\nu$	$1 - 4\nu$

Table 4: private information structure at state  $\theta = 1$  ( $\nu = \mu^0/\mu^1 \in (0, \frac{1}{2})$ )

If agent  $i$  receives  $s_i = 0$ , he assigns to  $\theta = 1$  probability  $\frac{\mu^1 \times 2\nu}{\mu^0 \times 1 + \mu^1 \times 2\nu} = \frac{2}{3}$  and thus action  $a_i = 1$  is weakly dominated. If agent  $i$  receives  $s_i = 1$ , he assigns to  $\theta = 1$  probability 1 and assigns to agent  $-i$  receiving  $s_{-i} = 1$  probability  $\frac{\mu^1 \times (1 - 4\nu)}{\mu^1 \times (1 - 4\nu) + \mu^1 \times 2\nu} = \frac{1 - 4\nu}{1 - 2\nu}$ . He may take action  $a_i = 1$  only if  $\frac{1 - 4\nu}{1 - 2\nu} \geq \frac{2}{3}$ , or equivalently  $\mu^1 \geq \frac{8}{9}$ . The regime's payoff under the adversarial selection is 1 if

$\frac{2}{3} < \mu^1 \leq \frac{8}{9}$  and  $1 - \mu^1(1 - 4\nu) = 5(1 - \mu^1)$  if  $\mu^1 > \frac{8}{9}$ . If  $\mu^1 > \frac{8}{9}$ , the probability that the agents take different actions is non-zero.

## References

- I. Arieli and Y. Babichenko. Private bayesian persuasion. *Journal of Economic Theory*, 182:185–217, 2019.
- P. Basu, K. Chatterjee, T. Hoshino, and O. Tamuz. Repeated coordination with private learning. *Journal of Economic Theory*, 190:105106, 2020.
- D. Bergemann and S. Morris. Bayes correlated equilibrium and the comparison of information structures in games. *Theoretical Economics*, 11(2):487–522, 2016a.
- D. Bergemann and S. Morris. Information design and Bayesian persuasion information design, Bayesian persuasion, and Bayes correlated equilibrium. *American Economic Review*, 106(5):586–591, 2016b.
- D. Bergemann and S. Morris. Information design: A unified perspective. *Journal of Economic Literature*, 57(1):1–57, 2019.
- K. Binmore and L. Samuelson. Coordinated action in the electronic mail game. *Games and Economic Behavior*, 35(1-2):6–30, 2001.
- N. Inostroza and A. Pavan. Persuasion in global games with application to stress testing. 2020.
- A. Kajii and S. Morris. The robustness of equilibria to incomplete information. *Econometrica*, 65(6):1283–1309, 1997.
- E. Kamenica and M. Gentzkow. Bayesian persuasion. *American Economic Review*, 101(6):2590–2615, 2011.
- F. Li, Y. Song, and M. Zhao. Global manipulation by local obfuscation. 2020.
- L. Mathevet, J. Perego, and I. Taneva. On information design in games. *Journal of Political Economy*, 128(4):1370–1404, 2020.
- D. Monderer and D. Samet. Approximating common knowledge with common beliefs. *Games and Economic Behavior*, 1(2):170–190, 1989.
- S. Morris, R. Rob, and H. S. Shin. p-dominance and belief potential. *Econometrica*, 63(1):145–157, 1995.
- S. Morris, D. Oyama, and S. Takahashi. Implementation via information design in binary-action supermodular games. 2020.
- A. Rubinstein. The electronic mail game: Strategic behavior under “almost common knowledge”. *American Economic Review*, 79(3):385–391, 1989.
- M. Shadmehr and D. Bernhardt. Collective action with uncertain payoffs: coordination, public signals, and punishment dilemmas. *American Political Science Review*, pages 829–851, 2011.
- I. Taneva. Information design. *American Economic Journal: Microeconomics*, 11(4):151–85, 2019.